

# Building a Large Proof Repair Dataset

Author One

Affiliation One

Author Two

Affiliation Two

Author Three

Affiliation Three

Talia Ringer

University of Illinois Urbana-Champaign

---

## Abstract

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

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

**Keywords and phrases** TODO: TODO keywords

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

**Supplementary Material** TODO URL to any code or data we can share

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

## 1 Introduction

Machine learning is coming for proofs. Recent years have seen a surge in interest in machine learning for proofs—one reflected by the many recent research venues [4, 5, 6], papers [21, 34, 27], tools [7, 12, 25], industrial research groups [36, 2], and funding opportunities [1, 3] centering or prominently featuring machine learning for proofs. The surge in interest blurs the line between proofs and data so that any proof development, once released, may itself become data to improve proof automation for future proof developments.

We in the proof engineering community have agency in how this surge of interest plays out. We can close our eyes and hope that we are sufficiently different from the many fields that machine learning has changed. Or we can embrace that change and adapt; make sure that our voices are heard. We can develop datasets and benchmark suites that steer the machine-learning community toward the tasks that matter most. We can build infrastructure that makes it easy to develop those datasets and benchmark suites, or to work on those tasks. And we can build evaluation methodologies that measure success on those tasks in ways that truly matter, so that state-of-the-art results on benchmarks will transfer smoothly to real-world improvements in proof automation.



© Author One, Author Two, Author Three, 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. 23:1–23:60  
Leibniz International Proceedings in Informatics



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

## 23:2 Building a Large Proof Repair Dataset

This paper takes a step in that direction. In particular, it presents a dataset and benchmark suite for an important proof automation task in Coq. The task of choice is *proof repair* [38]: the automatic repair of proofs in response to breaking changes in programs or specifications. Proof repair is a useful task for which data is scarce and challenging to collect [41]. This paper presents a new proof repair dataset and benchmark suite while also highlighting the challenges involved in building it, with an emphasis on how the proof engineering community can adapt to those challenges as machine learning becomes increasingly relevant.

Its contributions are the following:

1. a proof repair dataset and benchmark suite for Coq that is accessible to machine learning experts (Section 3),
2. tools for building and extracting information from Coq projects reusable for future dataset efforts (Section 4), and
3. a forward-facing discussion of the challenges we encountered in the course of this effort and how to overcome them (Section 5)

Our overarching goal is to build the infrastructure and proof assistant community support we need to steer the machine learning community toward the tasks that matter most before it is too late (Section 2).

## 2 Overview

A proof is a serialization of a truth, and assuming that our logical foundations are sound, truth is inviolable: what is true will always be true. This is great in isolation, but in practice, the software and hardware that form our proof’s hypotheses about the world change over time, and so do the truths we need to prove. All software, even that which is formally verified, must receive maintenance to survive sustained contact with reality.

The necessity and difficulty of proof maintenance has been borne out empirically. A recent user study of eight intermediate through expert proof engineers showed that maintenance happened constantly for participating proof engineers during proof development [41], and that even experts sometimes gave up in the face of change [38].

Consider, for example, the change in Figure 1, in which a user study participant updated a lemma statement in response to a change in a dependency. As noted in the user study paper, this was part of a larger change, with 10 other definitions or lemma statements changing in analogous ways. Furthermore, this change broke at least five proofs, four of which the user study participant—an expert proof engineer—admitted or aborted rather than repair.

The ubiquity of maintenance and the challenges of repair have been largely neglected in machine-learning tools for proofs. Machine-learning tools for proofs have instead historically fixated on development tasks like predicting tactics [45, 12, ?], synthesizing proofs from scratch [?, ?, ?], **TODO: TOM: rephrase this re slack messages** and automatically formalizing informal natural language statements [?, ?, ?].

If maintenance is so ubiquitous and repair is so challenging, what stands in the way of good machine-learning tools for proof repair? One of the central challenges is building a good dataset and benchmark suite, as datasets and benchmark suites can really drive results in machine learning for any domain. If the datasets and benchmark suites are not fit for maintenance tasks, the machine-learning community may neglect those tasks entirely, instead chasing state-of-the-art results on tasks for which existing benchmark suites suffice.

Existing datasets and benchmark suites are not sufficient for training and evaluating proof repair models. The main dataset of use in Coq is CoqGym [51], a collection of 124 proof developments in Coq. The CoqGym data for proof developments is static in that

```

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 user study paper [41].

it includes the projects frozen in time, rather than the histories or changes that would be needed to train a model for proof repair. In line with this, the CoqGym benchmarks measure success on a proof synthesis task in terms of the number of proofs synthesized in the suite for a given set of theorems. The data from the REPLICA [41] user study of Coq proof engineers, in contrast, includes very granular changes useful for repair, but is far too small in size to use to train dedicated proof repair models.

Our experiences interacting with machine-learning experts and building datasets ourselves suggest that the choice of datasets and benchmark suites for a domain is not driven solely by what is likely to be useful—it is also driven by barriers imposed by infrastructure, lack of domain expertise, or social factors. Here are three examples of these barriers:

1. In building a model for a brand new task with an industrial team, one of the authors—a Coq expert—chose Isabelle/HOL over Coq because of the utility of the Archive of Formal Proofs. They were unable to build a dataset for this task in Coq on a reasonable timeline because Coq lacked a **centralized archive** with **well-documented versioned histories**.
2. One of the authors witnessed several machine-learning experts across industry and academia erase meaningful information from data or make simplifying decisions about data and metrics for success that may corrupt results because of the difficulties in automating **parsing** and **building** Isabelle/HOL proof developments, and in **checking** proofs. These machine-learning experts largely preferred to treat proofs as plaintext data, and to interact with the proof assistant as little as possible.
3. The REPLICA dataset of atomic edits in Coq was fundamentally limited in size by the difficulties in **recruiting enough participants** in the user study [41].

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

## 23:4 Building a Large Proof Repair Dataset

We find this especially prudent given that the danger of chasing benchmarks that may not transfer to real life workflows has been realized quite dramatically in other domains—from incorrect patches to programs [35] to incorrect clinical interpretation of x-ray results [53]. Furthermore, these challenges can influence not just the tasks that the machine-learning community chooses to tackle, but even the very proof assistants the machine-learning community chooses to build tools in support of. It is in the community’s best interest to drive strong, practical results for useful tasks in a way that is equitable across proof assistants.

Our hopes are twofold. First, we hope that our dataset will be immediately useful for proof repair. Second, and perhaps more prudently, we hope our discussion of the challenges involved in building it will serve as a call for action to build better infrastructure. That way, we in the proof engineering community can ensure that machine-learning experts focus on the proof automation tasks that matter most to our community, that they measure success on those tasks in ways that are likely to transfer smoothly to real-world usefulness, and that they do so equitably across proof assistants.

### 3 A Proof Repair Dataset

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

#### 3.1 The Task: Proof Repair

In machine learning, a *task* refers to a high-level input/output specification of what is being learned. A dataset and benchmark suite typically organizes itself around a particular task while remaining agnostic to the details of the model implementation.

We define proof repair as a machine learning task as follows:

1. **Inputs:** old theorem, proof of old theorem, and new theorem
2. **Outputs:** proof of new theorem

Note that this particular proof repair task assumes that we already know how to repair the theorem statement and all of its dependencies. We could also consider a second proof repair task that allows the model to also repair the specification itself. This second proof repair task would be harder to evaluate, so we do not focus our benchmark suite on it at this time, though our dataset supports this. We discuss our plans for measuring success on the second task more in Section 3.3.

#### Inputs

**TODO: Show an Example.** We aim to provide sufficient context in the data to support a wide range of machine-learning approaches and repair model architectures. At a high level, the input to the machine-learning repair model is the entire state of a project where the

---

<sup>1</sup> We plan to release the dataset before the camera ready. The research team includes government contractors who are subject to US government approval processes, and cannot release datasets and code without a month’s notice. We have included a small sample as supplementary material in our submission, which we were able to approve a month in advance.

approach dictates how much of this state (and in what form) actually reaches the model. More precisely, we expect the input to comprise the statement of the theorem whose proof should be repaired, any contextual definitions on which it depends, the step-by-step goals and hypotheses for each sentence in the old proof, and known changes to the project up to and including changes in the theorem statement and its dependencies. For each of these components, we supply raw text representations (as one would observe in source code or in feedback messages within *CoqIDE*), abstract syntax trees (ASTs), and identifiers from which one could hypothetically compute data flow graphs. In the case of goals and hypotheses, serialized Coq kernel representations supply detailed internal proof states. In addition, environmental dependencies such as Coq compiler version are captured for errors induced by external application programming interface (API) changes.

## Outputs

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

## 3.2 The Data: Aligned Git Commits

The training and testing data comprise of aligned Git commits for a selection of realistic Coq projects. In total, it includes  $n$  commits for  $m$  Coq projects, with  $p$  repair examples in total. **TODO: Summary of other important stats goes here.** A full summary of all projects and commits in the dataset is in Table 6 in Appendix A.

## Choice of Coq Projects

Projects were originally selected by querying GitHub's API for projects that

1. contained Coq source code,
2. had a file called 'Makefile' in the project's root directory, and
3. had at least 100 commits to mine repair data from.

Eventually, we included CoqGym [51] and filtered to projects that were listed in opam repositories. Some projects were excluded for not containing any proofs or for having ulterior motives in their builds (e.g., projects that intended to test the performance of the Coq compiler *coqc* rather than build fascinating proofs).

**TODO: How do we know these are reasonable Coq projects to use for this task? Also note we are trying to broaden beyond what we currently have, but infra challenges, tease to challenges section, discuss later.**

**TODO: Enforceable note about right to delete**

## Repair Examples

Within each Coq project, the data comprises a number of repair examples—that is, changes to definitions or proofs. A repair example is constructed by comparing a definition or its proof before and after a change. Since sentences, identifiers, and files may be moved, renamed, added, or deleted between commits, they must first be *aligned* to ensure the right changes

## 23:6 Building a Large Proof Repair Dataset

196 are compared. At a high level, this means that Vernacular commands in one commit are  
197 assigned one-by-one to commands in another commit, where these assignments may cross file  
198 boundaries. Note that each command may not get a partner, indicating that it was either  
199 deleted or added. We describe this in more detail in Section 4.

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

### 208 Data Split

209 Machine-learning datasets and benchmark suites often include a data split between training  
210 data, validation data, and testing data. We do not commit to a single split ahead of time,  
211 but we consider two different ways of splitting data:

- 212 1. across projects, and
- 213 2. chronologically within projects.

214 These two splits test two different kinds of generalization beyond the training data, so we  
215 may consider them two different benchmarks for the same proof repair task. We find it  
216 reasonable to consider both of these kinds of generalization in measuring success of a proof  
217 repair model, so we plan to include defaults for both splits in the final release of the dataset  
218 and benchmark suite.

### 219 Across Projects

220 The first split—across projects—chooses distinct sets of Coq projects to use for training,  
221 validation, and test data. This is standard in existing benchmark suites for proof generation  
222 like CoqGym [51], and effectively measures generalization of the learned model to new  
223 projects not seen at training time.

### 224 Chronologically

225 The second split—chronologically within projects—uses the same set of projects for training  
226 and test data. However, these two sets in particular are split chronologically, so that training  
227 data includes earlier commits, and testing data includes later commits for the same projects.  
228 This effectively measures generalization of the learned model to new changes within a given  
229 project, when the model was trained on older data for that project.

## 230 3.3 The Metrics: Proof Checking

231 Changes in proof developments that break proofs can be fixed in two ways: by repairing the  
232 proofs themselves, or by repairing some other definition, like a program or specification [38].  
233 Our dataset includes both kinds of changes. For the sake of metrics for evaluation, we focus  
234 our benchmark suite on the former (repairing proofs), as the metric for success is immediately  
235 clear. We hope our benchmark suite will also be useful for the latter (repairing definitions),  
236 but we believe the problem of choosing a good metric for success for repairing definitions to  
237 be an open research problem.

## 238 Repairing Proofs

239 We focus our initial benchmarks on the problem of repairing a proof script assuming that the  
 240 statement of the repaired theorem to be proven is already known. In this case, checking the  
 241 correctness of the repaired proof amounts to using Coq’s kernel to proof check the type of  
 242 the repaired proof against the type that represents the desired repaired theorem statement.

243 The metric of proof checking is the same as that used for the standard proof generation  
 244 benchmark suite for Coq—CoqGym [51]. This metric is sound and complete (up to the  
 245 correctness of Coq’s trusted kernel, and designating nonterminating proof scripts as incorrect  
 246 proof scripts), as:

- 247 1. any proof that checks with the desired type is a proof of the theorem the type encodes  
 248 (**soundness**), and
- 249 2. all proofs that prove that theorem will check with the desired type (**completeness**).

250 In this sense, for this flavor of proof repair, we are able to take advantage of the fact  
 251 that proof checking is a perfect oracle when the theorem statement is known. The presence  
 252 of a perfect oracle has been hugely beneficial for existing machine learning work for proof  
 253 generation [21], and for early symbolic work on proof repair when specifications do not  
 254 change [42]. It continues to benefit machine learning for proof repair.

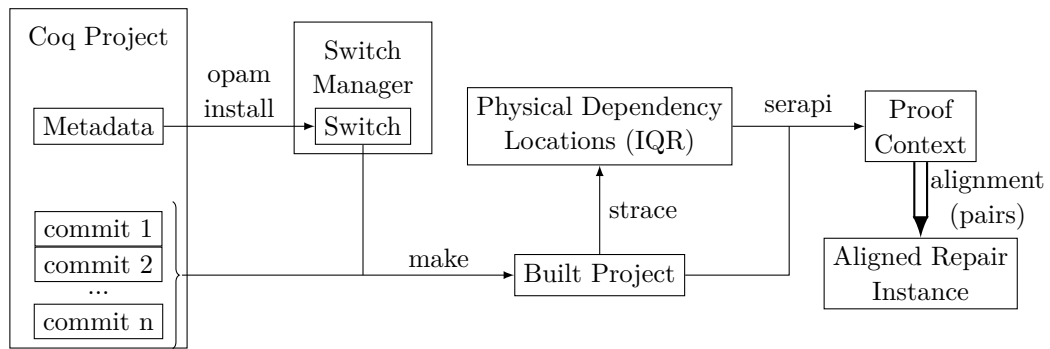
## 255 Repairing Definitions

256 The metric of proof checking is useful when the repaired specification and all of its dependen-  
 257 cies are known in advance. But what when they are not? Helping users fix the definitions that  
 258 the specification itself depends on—or the theorem statement itself—is also highly desirable.  
 259 In fact, a recent user study of eight intermediate through expert proof engineers in Coq  
 260 found that 75% of the time those proof engineers fixed a broken proof throughout the course  
 261 of the user study, they did so by fixing something else, like a program or specification [41].  
 262 So there is reason to suspect that supporting this use case may actually be *more* helpful to  
 263 proof engineers than supporting the original flavor of proof repair.

264 We believe the problem of finding a suitable metric for this use case to be an open research  
 265 problem. Existing metrics are insufficient:

- 266 ■ The metric of **proof checking** is no longer sufficient when the repaired specification  
 267 is not fully known in advance, as showing that the proof type checks is not meaningful  
 268 unless we know the type it ought to have.
- 269 ■ The conservative metric of **exact equality** with an expected repaired specification or  
 270 definition is sound, but far from complete, as there are many equivalent ways to state the  
 271 same theorems or write the same definitions.
- 272 ■ Loosenings of exact equality to use, for example, common notions of **definitional**  
 273 **equality** or **propositional equality** in Coq are slightly less conservative, but are still  
 274 far from complete.
- 275 ■ Proof repair across type equivalences [40] defines a more complete notion of preserving  
 276 the specification correctly across a change in datatype in terms of **univalent transport**,  
 277 but this notion is too narrow to express some changes in specifications that add or  
 278 remove behaviors or information. Furthermore, checking this metric automatically in a  
 279 benchmark would be difficult, as it is undecidable without a generated proof of correctness  
 280 of the change, and requiring proofs of correctness of the change would place a very large  
 281 burden on repair tools.





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

282 ■ In other domains, it is common to use **distance metrics**, but common distance metrics  
 283 like BLEU [31] are poor measures of success in program synthesis tasks **TODO: cite.**  
 284 **TODO: (AG) There is also Code-BLEU, but its computation is complicated by the**  
 285 **difficulty in obtaining ASTs and data flows** As formal proof developments are in some  
 286 sense programs, we expect existing distance metrics like BLEU score to be inadequate for  
 287 proofs as well. Distance metrics do have the distinct advantage of being continuous in  
 288 some sense—they make it possible to define what it means to be “close to correct.” But  
 289 being “close to correct” does not mean much when the notion of correctness used is itself  
 290 incorrect.

291 For now, we largely focus on the proof repair task where the repaired specification and all  
 292 of its dependencies are known in advance, as it is easy to measure success on this task. Still,  
 293 the data for the task of repairing definitions like programs and specifications is present in the  
 294 dataset. If you choose to train and evaluate a model for the task of repairing definitions, we  
 295 recommend using a conservative metric like exact or definitional equality, so as to avoid the  
 296 danger of chasing misleading benchmarks alluded to in Section 2. We hope the community  
 297 will come together to develop a suitable less conservative metric going forward.

## 298 4 Building the Proof Repair Dataset

300 Now that we have introduced the proof repair dataset, we take a step back and describe  
 301 the processes behind our data collection efforts. The process of generating repair data from  
 302 a project is summarized in Figure 2 and comprises the following steps:

- 303 ■ We start with a collection of aggregated open-source **Coq Projects**, which form the  
 304 foundation of our dataset (Section 4.1).
- 305 ■ We determine how to build each project using the **Switch Manager (SwiM)**, which  
 306 takes the dependencies of a Coq project and gives a copy of the cached *switch* (OCaml  
 307 Package Manager (opam) sandbox or virtual environment) that satisfies as many of them  
 308 as possible and in which the remainder can be installed (Section 4.2).
- 309 ■ Once we have a switch, we run the build command in the generated switch to produce a  
 310 **Built Project** (Section 4.3).
- 311 ■ We **strace** the build process to scrape the **Physical Dependency Locations (IQR**  
 312 **flags)** of each document in the project, which lets us know the mapping between Coq  
 313 imports and physical disk locations (Section 4.4).



314 ■ Using the IQR flags for each document in a built project along with the Coq serializer  
315 SerAPI [8], we extract a **Proof Context** corresponding to each intermediate proof state  
316 for the project by querying Coq’s state during the execution of proofs (Section 4.5).  
317 ■ Finally, we align changed proofs across commits and save those along with their interme-  
318 mediate proof contexts to arrive at an **Aligned Repair Instance**, which is a final product  
319 in our dataset (Section 4.6).

320 Building this infrastructure was a significant undertaking with many challenges en-  
321 countered along the way; we discuss these challenges in Section 5. Our hope is that the  
322 infrastructure we have built will make it easier to collect similar datasets in the future.

## 323 4.1 Coq Projects

324 The foundation of the dataset comprises open-source Coq projects found in publicly available  
325 repositories on web hosting services (primarily Github). The choice of projects is summarized  
326 in Table 6 of Appendix A. Mining the commits of these projects eventually yields examples  
327 of refactors or repairs. Each project is also accompanied by per-commit metadata containing  
328 project dependencies, source URL, and build commands that is in parts manually curated  
329 and programmatically inferred.

## 330 4.2 Switch Manager

331 To extract information about these projects like intermediate proof states, we need to be able  
332 to build them. This requirement is nontrivial because different projects and even different  
333 commits of the same project can include different versions of dependencies—including different  
334 versions of Coq or of the OCaml compiler.

335 To resolve dependencies and make it possible to build many different commits of the  
336 same project, we introduce a novel SwiM capability that works in tandem with opam—an  
337 OCaml package manager that acts as the primary distributor of Coq projects—to extract  
338 a Python object that models the build environment for a given commit of a given project.  
339 The Python object in particular models an opam switch, which is the way opam represents  
340 dependencies along with OCaml compiler versions.

341 This capability subverts the typical opam workflow that requires one to manually create  
342 a switch and to set environment variables for the current shell to activate that switch. This  
343 manual workflow would be intractable at the scale of hundreds of commits for each of  
344 hundreds of projects. With the SwiM, we can automate this functionality and extract a  
345 dataset at scale.

346 A number of benefits arise from the design of the SwiM. First, the SwiM enables  
347 sandboxing of project builds by providing temporary isolated switch clones that last for the  
348 duration of the commit’s extraction. Furthermore, the SwiM minimizes the time required to  
349 obtain a clone by maintaining a pool of switches with pre-installed packages and choosing  
350 the one upon request that is closest to satisfying a commit’s requirements.

351 As new commits are built, switches containing their dependencies are added to the  
352 managed pool of switches. Since switches range in size from hundreds of megabytes to  
353 a few gigabytes, a least-recently-used cache maintains the total disk consumption below  
354 an implicit limit by deleting stale switches that are not frequently used. A shared switch  
355 manager facilitates multiprocessing with a shared pool of switches by queuing requests from  
356 concurrent threads/processes. Implementation of this capability required reflection of opam’s  
357 dependency formula parsing and evaluation logic from OCaml to Python.

### 358 4.3 Built Project

359 Using the switch provided by the SwiM and a build command from the metadata, we may  
 360 be able to build the project. Confounding issues that may prevent building include undefined  
 361 opam variables within dependency formulas, for which we only provide default values of  
 362 “True” for each of the `build`, `post`, and `dev` builtin variables. In practice, we have so far seen  
 363 a build failure rate of about 68%. We attempt to build each commit with seven different  
 364 major versions of Coq ranging from 8.9 to 8.15 corresponding to versions of SerAPI that  
 365 support capabilities deemed necessary for data extraction. Since Coq releases are rarely  
 366 backwards compatible, many of the build failures can be explained by the fact that each  
 367 commit can only be expected to build without error for the single Coq version for which it  
 368 was written. Furthermore, since we pin one of the seven Coq versions in the switch supplied  
 369 by the SwiM, conflicting version requirements may yield an opam command that has no  
 370 solution. Consequently, some build errors are inevitable.

371 However, other errors are due to mistakes or missing information in the human-sourced  
 372 metadata. We plan to address this latter class of build errors over time by fixing problems  
 373 in the metadata through automated inference mechanisms. If the project build fails, we  
 374 hope in the future to be able to recover proofs from the documents that built before the  
 375 failure as well as subsequent independent proofs. We are also exploring ways to automatically  
 376 recover from simple build errors such as dependency mismatches between the switch and the  
 377 project’s requirements by using the date the commit was made as a version hint.

### 378 4.4 Physical Dependency Location (IQR)

379 In order to run any of the Coq or SerAPI tools (e.g., `coqc`, `coqtop`, `sertop`) on a given Coq  
 380 source file, many Coq source files require one or more flags to be passed to these commands  
 381 to specify the physical location of dependencies. These flags are described below:

- 382 ■ The `-I` flag allows a directory to be added to the OCaml loadpath
- 383 ■ The `-Q` flag adds a physical directory to the loadpath and binds it to a specified logical  
 384 path
- 385 ■ The `-R` flag does the same as the `-Q` flag, except that subdirectories are also made available  
 386 recursively

387 In publicly available Coq projects, these flags (informally referred to as “IQR” flags from  
 388 here on) are specified in one or more build or configuration files. There are several common  
 389 approaches for specifying IQR flags, but there is no single standardized approach, making it  
 390 difficult to automatically infer these flags from configuration and build files alone. While  
 391 projects will be able to build successfully without our inferring of these flags, we need to  
 392 infer them in order to use SerAPI tools in other stages of our framework.

393 Our solution to this problem builds off an approach developed in IBM’s PyCoq <sup>2</sup>. Following  
 394 PyCoq, we use `strace` to inspect the actual commands run during the build process for a  
 395 Coq project. Each build command is captured and any present IQR flags are extracted using  
 396 regular expressions. In some projects, build files may be nested, and IQR flags may specify  
 397 physical paths that are relative to the nested directories. We need to ensure that the inferred  
 398 IQR flags are relative to the project root directory, so before we store the inferred IQR flags,  
 399 their paths are resolved to the project root directory.

---

<sup>2</sup> <https://github.com/IBM/pycoq>

## 4.5 Proof Context

Once the project has been built, the individual Coq source files are parsed into sentences and then interactively executed with `sertop` to capture intermediate proof states. A parser (`sercomp`) is available through SerAPI, but it can only be used on Coq source files whose dependencies have already been compiled, which prohibits its use in planned recoveries from partial builds. Furthermore, `sercomp` introduces significant overhead through computation that will be redundant with `sertop`.

From this, we can collect thorough context from the document, such as whitespace-normalized text, ASTs, command types, and intermediate proof steps with goals and hypotheses. Each command is accompanied by inferred identifiers of the command itself (e.g., the name of an inductive type and its constructors) and a list of fully-qualified identifiers referenced within the command, which enables models to more easily incorporate relevant local context or apply graph-based approaches. Any errors that are encountered in the execution of a command are cached as well for repair models that can adapt their approach based on the error. Accompanying source code locations allow for accurate provenance of data and application of proposed repairs to appropriate destinations for testing.

## 4.6 Aligned Repair Instance

The final step in our data collection process is extracting proof repair examples from different versions of the projects. To extract these examples, we need a way to tell how different versions of definitions and proofs across different commits correspond to each other.

We approach this as an alignment problem. Because lines in one commit are not necessarily located at the same index as in the other commit, we first have to establish a robust mapping between lines that can handle the changes usually associated with a commit: lines can be changed, deleted, added, and rearranged. This is similar to the ‘diff’ utility, which essentially tries to explain what changes were made between two files, but it lacks capacities we deem necessary, such as matching lines that were rearranged and matching lines that were only changed slightly. In particular, we wish to establish a ‘diff’ in terms of Coq Vernacular commands treated as atomic entities that can capture whether an individual definition, lemma, inductive or other identifiable entity was moved, renamed, or otherwise altered during a refactor.

A traditional order-preserving alignment between two sequences **TODO**: e.g., **CITE** is not quite an appropriate approach to resolve this issue as it cannot correctly align two independent definitions whose order has been reversed during a refactor (perhaps due to an introduced dependency). Instead, we must consider permutations of commands within and across file boundaries. We approach the problem as a bipartite matching or *assignment* between the elements of two sets such that the overall similarity of matched elements is maximized. We can formally specify the desired assignment between two commits  $X$  and  $Y$  considered as respective sets of  $m$  and  $n$  commands across one or more files as the solution

## 23:12 Building a Large Proof Repair Dataset

438 to the following optimization problem:

$$\begin{aligned}
& \underset{f}{\text{minimize}} && \sum_{x \in X} \sum_{y \in Y} C(x, y) f(x, y) \\
& \text{subject to} && \sum_{y \in Y} f(x, y) \leq 1, \\
& && \sum_{x \in X} f(x, y) \leq 1, \\
& && \sum_{x \in X} \sum_{y \in Y} f(x, y) = \min\{m, n\}, \\
& && f(x, y) \in \{0, 1\}
\end{aligned} \tag{1}$$

440 Here,  $C(x, y)$  is a non-negative cost function that measures the *distance* between the com-  
441 mands  $x$  and  $y$ , and  $f(x, y)$  indicates the binary assignment of  $x$  to  $y$ . The constraints, in  
442 order, state that no element of  $X$  can be assigned to more than one command of  $Y$ , no  
443 element of  $Y$  can be assigned to more than one command of  $X$ , and as many commands as  
444 possible must be assigned (i.e., choosing to assign nothing in an attempt to minimize the  
445 cost is not allowed). This optimization is an instance of the well-known *assignment problem*,  
446 which one can solve exactly in polynomial time according to the Hungarian algorithm [30]  
447 (or one of its many alternatives) or approximately in linear time according to the Sinkhorn  
448 algorithm [17] (which also admits GPU acceleration).

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

$$451 \quad C(x, y) = \frac{2E(x, y)}{\|x\| + \|y\| + E(x, y)}, \tag{2}$$

452 where  $\|x\|$  and  $\|y\|$  give the character lengths of  $x$  and  $y$  considered as text (not including  
453 proof bodies). This normalization is an instance of the biotope or Steinhaus transform [18],  
454 which preserves the metric properties (such as the triangle inequality) of the edit distance.  
455 We further threshold the distance by a constant  $t$  such that  $C_t(x, t) = \min\{C(x, y), t\}$ , which  
456 also preserves metric properties [32]. After solving for  $f$ , commands  $x$  and  $y$  assigned to one  
457 another ( $f(x, y) = 1$ ) with a cost of  $t$  are considered to be unassigned (i.e., we determine  
458 that  $x$  was dropped between commits and  $y$  was added). We choose  $t = 0.4$ , which roughly  
459 corresponds to 50% of a command’s text being changed before it is considered to have been  
460 dropped. **TODO: (AG): This is a point for discussion as we could also introduce some**  
461 **alternative heuristics for checking whether two assigned commands are the “same”.**

462 Solving this assignment problem for two entire commits can be costly: solving for the  
463 assignment exactly is cubic complexity, and the edit distance must be calculated for all pairs  
464 of commands from both commits, which is necessarily quadratic complexity. Furthermore,  
465 the assignments produced may be somewhat spurious, especially in the event of multiple  
466 global optima. Though we cannot enforce preservation of order as an absolute constraint, we  
467 do wish to apply it as a pseudo-constraint to resolve ambiguity. We mitigate these issues by  
468 applying the assignment problem only to those commands known to have changed in some  
469 manner between the commits according to their intersection with a (Git) ‘diff’. The final  
470 resolution of the problem is thus somewhere in between alignment and assignment.

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

## 5 Challenges

Machine learning for proofs will likely continue to flourish as a field whether or not the core proof assistant community participates. But if we do participate, perhaps we can work to ensure that it grows in the right directions—towards the tasks and benchmarks that matter, equitably across proof assistants.

We have a long way to go before that is true. Collecting and building datasets and benchmark suites for many tasks is still extremely challenging, and it is challenging in a way that is not at all equitable across proof assistants.

Here, we discuss challenges we faced in building this dataset and benchmark suite that fall into three categories:

1. Archival & Versioning (Section 5.1),
2. Package Management & Distribution (Section 5.2), and
3. Parsing & Serialization (Section 5.3).

For each of these categories, we discuss our experiences dealing with those challenges, what we believe the Coq community can learn from other proof assistant communities, and how the proof assistant community at large could address them more sustainably going forward.

### 5.1 Archival & Versioning

One of the largest barriers to building this dataset was the lack of a centralized archive for Coq proof data, combined with the difficulties of dealing with different versions of different projects and their dependencies. These interacted in complex ways, the end result of which was a scattered source of data that was difficult to aggregate in a way that would be useful for downstream machine learning tools, especially for the proof repair task.

#### 5.1.1 Our Experiences

We faced four challenges related to archival and versioning:

1. the **lack of centralized archival** of Coq proof developments,
2. frequent changes in **proof assistant versioning** that broke the abstraction barrier,
3. significant **build system variation** across different proof developments, and
4. tracking information across **changing definitions & proofs** within proof developments.

#### Lack of Centralized Archival

Coq lacks a centralized archive that it is standard for proof engineers to use. There are a few archives that have existed at various points in time, but there is not a strong norm of using them. Consequentially, we had to aggregate proof developments from public sources like GitHub. While this worked nicely, it also meant that proof developments lacked useful associated metadata, or were formatted in different ways, making automatic aggregation, parsing, and building more challenging. This perhaps is at the root of many other challenges we faced and had to overcome.

#### Proof Assistant Versioning

One of the barriers to centralized archival in Coq is the large variation in proof assistant versioning, while the up-to-date archive that does exist requires proof developments to be on the latest version. There was also notably not a uniform way to denote which proof

assistant version a given proof development was on, nor was there a uniform way to interact with Coq across proof assistant versions. The tools we used for serialization and parsing, for example, were highly version-dependent (see Section 5.3). This meant that we had to infer proof assistant versions using imperfect proxies, and build data collection tools that inherently depended on those inferred versions.

## Build System Variation

Over the years, the recommended build system for Coq proof developments has been in flux. In 2019, for example, the Coq development team urged proof engineers to move their proof developments to Dune<sup>3</sup>. This did not fully succeed, and the documentation for the latest Coq version includes instructions for both Dune and the native Coq build system<sup>4</sup>. The native Coq build system itself has also changed over time in significant ways, losing compatibility with previous versions of the same build system. For example, the arguments to the `coq_makefile` command changed over time. Because of this fragmented build infrastructure, we had to employ extremely abstract methods to extract information whenever a build system was involved, such as using `strace` to grab flags passed to Coq’s compilation command `coqc` during builds (described in Section 4.4), while making almost no assumptions about the process invoking `coqc`.

## Changing Definitions & Proofs

Building a proof repair dataset meant mining not static proof developments, but rather changes in proof developments over time. Figuring out which definitions and proofs corresponded to one another across versions was a significant challenge. This challenge is one that the REPLICA user study also faced; rather than custom alignment, they relied on GitHub’s diff command, but needed a human to manually crawl the result to find repair examples [41]. Since our data was less granular than REPLICA’s, we were able to make the simplifying assumption that all repair examples occurred commits that occurred immediately after one another, without intermediate commits. We then mined repair examples automatically using the alignment algorithm discussed in Section 4.6.

### 5.1.2 Other Proof Assistants

Other proof assistants can learn from Isabelle/HOL’s rich archival culture. In Isabelle/HOL, it is standard to upload proof developments to the Archive of Formal Proofs (AFP). Isabelle/HOL’s AFP is highly centralized, and also makes it easy to associate proof developments with metadata that may be useful for machine learning. It is also neatly versioned, with all proof developments updated for every minor and major version of Isabelle/HOL, and with major version semantically grouped in different folders.

At the time of writing, the AFP includes 725 proof developments<sup>5</sup>. The AFP already forms the basis of a dataset for Isabelle/HOL for static data [23]. We suspect it will also make for a very strong basis for a proof repair dataset due to its neat versioning. We believe that other proof assistants should take influence from Isabelle/HOL in developing a large centralized archive that is standard to use and neatly versioned, and that makes it easy to associate proof developments with metadata useful for machine learning tasks.

<sup>3</sup> <https://coq.discourse.group/t/proposal-a-custom-build-tool-for-coq-projects/239/2>

<sup>4</sup> <https://coq.github.io/doc/master/refman/practical-tools/utilities.html>

<sup>5</sup> <https://www.isa-afp.org/statistics/>

### 5.1.3 Recommendations

There are still a number of ways that even the best archival and versioning infrastructures for proof assistants fall short. For example, while Isabelle/HOL's AFP makes a natural data source for machine learning tools for proofs, it does not include any processes for informed consent—and the datasets that build on it assume that all publicly available data is fair game. While this is standard in machine learning for code and proofs, it is not ideal. Consent should be considered at every step in the process, and archival is a natural place to consider that.

Another problem with archives is keeping them up-to-date. Is this sustainable as archives continue to grow, while the cost of repair is still so high? We are not sure how the AFP maintainers manage this burden, but we assume that not all communities can manage it, and lowering the cost of repair will be central to better archival across all proof assistants, and to getting proof developments on the latest version to archive to begin with.

Finally, while archival makes it possible to associate metadata with proof developments, there is little consideration for adding metadata that associates definitions and proofs *across versions* of a proof development. This kind of metadata would make building repair datasets and benchmark suites much easier, and make it simpler to address the problems of change as they interact with archival in the future.

Based on our observations, we make the following recommendations for the proof assistant community going forward:

1. Work with the Isabelle/HOL community to learn how to build **centralized archives** as successful as the AFP for other proof assistants.
2. Include an **informed consent form** in any centralized archive that allows proof engineers to opt in or out of their developments being used for machine learning tools.
3. Create **better and more consistent standards** for tools like build systems.
4. Help proof engineers **port legacy proof developments** to meet those standards.
5. Continue to work on tools for **proof repair and reuse** that ease the burden of porting said legacy developments.
6. Determine what **kinds of metadata** both within and across proof developments would ease the development of machine learning for proofs tools for the tasks that matter most.
7. Make it easy to **track that metadata** inside of any centralized archive.
8. Create standard ways of **associating that metadata** with proof developments even outside of centralized archives.

## 5.2 Package Management & Distribution

The beating heart of our collection effort is package management, which gives us a programmatic interface to building compatible environments for the data-set's constituent projects. In our case, we use opam: a popular OCaml package manager and the primary distributor for Coq packages due to the close relationship between the OCaml and Coq communities. **TODO: OPAM isn't supposed to be in caps, even though it's an acronym.** The opam package manager was nominally designed to service individual developers using a few switches, not use-cases like our own where we have to spin up dozens of switches efficiently. This required us to reimplement and expand upon some of opam's capabilities.

### 5.2.1 Our Experiences

We faced three challenges related to package management & distribution:



## 23:16 Building a Large Proof Repair Dataset

- 599 1. Packages in opam have very **expressive dependencies** which complicate efficient installs.
- 600 2. The opam package manager has no mechanism to cache builds, so we need to **copy**
- 601 **switches** to avoid rebuilding the same packages.
- 602 3. It is difficult to use opam with packages that have **inconsistent dependencies**, though
- 603 this is desirable to us.

### 604 Expressive Dependencies

605 The opam package manager provides a powerful and expressive syntax (package formulae)  
606 for packages to specify dependencies over other packages. Package formulae allow package  
607 developers to restrict the versions of dependencies that can be installed, to conjunct and  
608 disjunct formulae into more complicated expressions, and to refer to variables declared  
609 elsewhere in the environment. This benefits the library developer, who can precisely specify  
610 the environment that their code will run in, but for our purposes package formulae pose a  
611 challenge: packages can be very picky about their environments and force opam to rebuild  
612 existing libraries. Since we need to install many versions of many packages, we needed to  
613 look into efficient ways to create and select switch environments they could be installed in,  
614 which meant interpreting these expressions. As a result, we had to reimplement a majority  
615 of opam's package formula features, including parsing the custom grammar for package  
616 formulae and implementing package version comparison, in order to reason about what  
617 existing switches would be the install environment for a package where opam would do the  
618 least amount of work to install the new package.

619

### 620 Copying Switches

621 If multiple opam packages conflict with each other, they can be installed in separate switches.  
622 Since we were installing packages from various points in the past, packages frequently  
623 required different version of dependencies, so they frequently conflicted. This meant we had  
624 to spin up many switches, but spinning up a switch required us to rebuild packages from  
625 source each time. We reasoned it would be better to build a package once and deploy it in  
626 multiple switches. However, many executables built by opam contain their absolute path as  
627 a hardcoded variable, which means these executables stop working if the name or location of  
628 the switch changes. In other words, a compiled opam package only necessarily works in one  
629 switch.

630 Our workaround was to copy switch directories without informing opam, which allowed  
631 us to start a new switch using another one as a base, saving valuable build time. Whenever  
632 a copied switch needed to be used, we bind-mounted the copied switch over the original  
633 switches location, so a copied switch effectively shared the same path as its parent. We  
634 used the **bwrap** utility to allow unprivileged processes to bind-mount, which is already used  
635 internally within opam to sandbox package builds. Of course, handling these cloned switches  
636 required additional bookkeeping and infrastructure, complicating our efforts.

637

### 638 Inconsistent Dependencies

639 For our repair dataset, we may like to see how a package breaks when dependencies are  
640 updated, since this breakage is another repair instance. Naturally, opam does not want to let  
641 us install packages it knows are incompatible. It is possible to force opam to install packages  
642 with inconsistent dependencies anyways through the use of the **ignore-dependencies-on**

flag, but this puts the solver in an inconsistent state. Trying to use opam after inconsistent packages are installed typically results in opam trying to change the versions of packages as soon as it can, even without being prompted to do so. Attempting to pin both packages to inconsistent versions renders opam inconsolable, refusing to do anything with dependencies until the inconsistency is resolved. Finally, it is possible to continue using opam as long as you use the previously stated `ignore-dependencies-on` flag once again to disregard all dependencies on the wrong-versioned-package, but if you install another package that depends on the package whose dependencies are being ignored, it will simply install the newest version, not necessarily one that is compatible (ignoring all dependencies, not just the one exception created initially).

A cleaner way to tamper with dependencies is to edit the opam package file of the package of interest directly so that opam's dependency solver isn't involved in the process. opam even provides a mechanism to pin local, modified versions of packages and a command line tool (`opam pin edit`) to open a package's metadata in an editor, which are useful for package tinkering and allow inconsistent packages to be solved without causing the above issues. Unfortunately, this is only a complete solution for humans editing a few packages. An automatic tool, such as would be useful for our efforts, still has to parse the package metadata fields, parse the dependencies, modify them, and paste them back in the right spot.

### 5.2.2 Other Proof Assistants

Agda has its own library management system, which it uses in combination with hackage (to install Agda itself). Anecdotally, researchers we have spoken to cite installation difficulties as a barrier for picking up Agda. Lean similarly has its own package manager, but lacks advanced features for the kinds of problems we faced. Isabelle/HOL's AFP serves as a de facto group of packages to install. But Isabelle/HOL in general takes a very IDE-centric approach to builds and other tooling [39, 47] one of the authors has found that students learning Isabelle/HOL in a proof automation course struggle to understand how to build dependencies.

In summary, we are not aware of an elegant and effective solution to package management for other proof assistants. We are eager to hear from experts in other proof assistant communities about whether the outlook is any better in other proof assistants, and about what we can learn from any solutions that do exist. If the situation is in fact universally bleak, perhaps we ought to come together to innovate.

**TODO:** It's also possible to use nix for Coq, and this uses online caches so you don't need to rebuild. Why don't people use this? Not sure where to drop this but it's an important note I think. Also, Isabelle sessions are a thing, and can inherit from other sessions, solving some but not all of the caching issue. Also the `-R` option. See Tweet responses and consider.

### 5.2.3 Recommendations

There is a legitimate argument to be made that package management by default targets a very different use case from ours, and perhaps existing tools are sufficient for that use case. If that is the case, it may make more sense to build shared high-level libraries and tools on top of existing package managers to support bulk cases like our own, especially since these use cases are common when building machine learning datasets and tools for formal proof.

Nonetheless, there may be room for improvement of the package managers themselves. For example, the problem we encountered of copying switches was due to poor caching of build dependencies, which itself was due to some degree to hard-coding of paths. It may be

## 23:18 Building a Large Proof Repair Dataset

688 better if applications could consult something more malleable than a compile time string  
689 to determine if the applications consulted something more malleable than a compile time  
690 string to determine these paths, like an environment variable, which opam already configures  
691 several of over the course of normal switch operation. This would likely require effort on the  
692 part of each of the applications, but may help with the problems we encountered and more.

693 In all, we make the following recommendations for the proof assistant community going  
694 forward:

- 695 1. Consider building **shared libraries and tools** optimized for the problems of bulk builds.
- 696 2. Consider limiting the scope of possible **dependency complexity**, or else building better  
697 tools to resolve complex dependencies efficiently.
- 698 3. Work on better **build caching** across multiple or bulk builds, for example by avoiding  
699 the hard-coded paths seen in opam.
- 700 4. Provide more effective ways of **overriding conflict declarations** that are overly conser-  
701 vative.
- 702 5. Work together across proof assistant communities to determine **common lessons and**  
703 **infrastructure** to resolve these problems.
- 704 6. Consider opening conversations with **language developers and companies** outside of  
705 verification about their package management and distribution solutions, as this problem  
706 is pervasive across all software.

### 707 5.3 Parsing & Serialization

708 **TODO: Radiance folks and Tom: I need help with 5.3.0 and 5.3.1 before it's in a state such**  
709 **that I can write a good 5.3.2 and 5.3.3. I'd appreciate if you all can take a good solid pass at**  
710 **this.**

711 **TODO: Intro paragraph that is higher-level than what we had before.** No matter how  
712 sophisticated the build system, we cannot get detailed data about individual proofs without  
713 parsing the Coq files and serializing proof state to text. SerAPI [8] is the de facto standard  
714 for serializing interactive Coq data, providing a query protocol for exposing aspects of  
715 internal Coq program data including definitions in the global environment, syntax trees,  
716 goals, types, and more. We used the CoqGym [51] Python wrapper as a starting point for  
717 our implementation, taking care to decouple it from CoqGym's custom versions of Coq and  
718 SerAPI since we need to support multiple versions of each coinciding with chosen projects'  
719 Git histories. This need to support multiple versions of Coq exacerbated challenges arising  
720 from gaps in SerAPI's query protocol, requiring us to implement workarounds using the  
721 most public and arguably stable interface Coq possesses: its Vernacular query commands.

#### 722 5.3.1 Our Experiences

723 We faced five challenges related to parsing & serialization:

- 724 1. Executing a file one Coq sentence at a time requires accurately **parsing sentence**  
725 **boundaries**, but parsing requires execution: a catch-22.
- 726 2. **Calculating command dependencies** requires identification of novel definitions and  
727 full qualification of embedded references with library prefixes, neither of which are  
728 provided through SerAPI queries.
- 729 3. **Determining the scope of a conjecture** is complicated by the potential presence  
730 of nested proofs/definitions, arbitrary grammar extensions, and the lack of obligatory  
731 syntax to delineate proof boundaries.

- 732 4. Some proof developments consume an enormous amount of memory and processor time,  
 733 requiring **resource limits** to keep data collection tractable.
- 734 5. SerAPI is experimental software, which leads to breaking **changes between versions**.

735 **TODO:** Add challenge for error handling and recovery with discussion of regression  
 736 proving and asynchronous proof checking and why despite its availability as a feature of  
 737 Coq, it is not an option for us (would require invasive modification of build processes for  
 738 every commit, deemed impractical; also does not work for all projects). Also, not clear  
 739 how asynchronous proof checking would help with extraction of intermediate proof states  
 740 except that it implies there is a way to determine conjecture scopes without evaluating them  
 741 (actually, the way would be to use the Vernacular classification functions alluded to later).

## 742 Parsing Sentence Boundaries

743 A Coq statement or ‘sentence’ ends with a period (.), but Coq also uses the symbol for  
 744 import paths and module members so that one cannot identify sentences in a file merely by  
 745 splitting on periods. To further complicate matters, Coq boasts an extensible syntax that  
 746 can greatly improve the read and writability of Coq expressions, but allows a user to define  
 747 syntax that allows periods to show up in even more places.

748 An example of a problematic extension of ‘.’ is shown in 3. This snippet defines syntax to  
 749 create tuples using a period rather than a comma, which complicates sentence splitting to  
 750 the point where the latest version of *CoqIDE*, the official editor for Coq, cannot correctly  
 751 parse and run this code even though Coq itself can.

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

■ **Figure 3** This example was found (here).

752 We did not discover any public or officially supported mechanism to extract the sentences  
 753 of a Coq document. Though one could hypothetically parse the output of *coqdoc*’s HTML  
 754 or LaTeX output, correlating the sentences back to their original source code locations  
 755 (line, column, and character numbers) brings additional complexity. We opted to implement  
 756 our own heuristic regular-expression-based parser in Python for simplicity and maximal  
 757 portability between build environments, noting that it inherits weaknesses exhibited by  
 758 *CoqIDE* but reasoning them equally acceptable as the official editor. A custom string subclass  
 759 tracks locations of individual characters across splits, slices, and concatenations of the parsed  
 760 source code.

## 761 Calculating Command Dependencies

762 To identify the dependent definitions referenced in a statement, we must associate with each  
 763 applicable command the identifier(s) that it defines and references. No SerAPI query provides  
 764 the name(s) introduced by a given command, nor is there any unambiguous syntax one can  
 765 use to identify a name that will be visible in the global environment that also generalizes  
 766 across builtin Vernacular commands (and Coq versions) or unforeseen grammar extensions.  
 767 Furthermore, identifiers within ASTs yielded from SerAPI typically match their qualification  
 768 in the source code and thus are not guaranteed to unambiguously identify referenced or  
 769 shadowed definitions when examined out of context.

770 Instead, we rely upon parsing user-level feedback and queries for detecting definitions.  
 771 Though Coq typically emits feedback of the form “ $X$  is defined” when a new function  
 772 or inductive type is introduced, it is not guaranteed to do so for a proposition when its  
 773 proof is completed (whether it does or not depends upon the proposition type and Coq  
 774 version). Consequently, a Vernacular **Print All** command is inserted when no such feedback  
 775 is detected to print the names and types of all definitions introduced since the start of the  
 776 interactive session. Monitoring changes in the set of all locally defined names then allows  
 777 one to deduce which definition, if any, was introduced by a given command.

778 We qualify identifiers in an AST through SerAPI’s **Locate** query, which requires a separate  
 779 query for each identifier. These repetitious queries add a significant amount of unavoidable  
 780 overhead. Care must be taken to ensure locally bound variables within binders, patterns, or  
 781 other sub-expressions do not get incorrectly qualified as any top-level definition that they  
 782 may shadow. Our solution is ultimately limited by the lack of accurate scope rules that  
 783 would require access to or replication of Coq’s internal name resolution capabilities. Aside  
 784 from being unable to disambiguate locally bound variables sharing the same name, we also  
 785 note one restriction on resolving globally bound identifiers: if a new identifier shadows an  
 786 existing one, then the defining command cannot reference the shadowed identifier. Violation  
 787 of this assumption is possible (consider a recursive function **nat** that expects arguments of  
 788 type **nat**) but not expected to be a significant risk as it is unlikely in the first place and  
 789 would generally be considered poor practice. If the restriction is violated, then the shadowed  
 790 identifier will simply be shadowed starting with the violating statement.

791 Note that all of the above solutions require parsing identifiers to some degree, which poses  
 792 its own challenge due to Coq’s versatile support for Unicode. The character set comprising  
 793 Coq identifiers intersects each of the usual meta characters used to denote whitespace  
 794 (non-breaking space 0x00A0), word characters, and other symbols. The exact character  
 795 set is not publicly documented. Consequently, one cannot accurately construct a regular  
 796 expression for parsing Coq identifiers without duplicating the Unicode tables and implied  
 797 regular expressions within Coq’s source code, which is what we ultimately resorted to as we  
 798 could not eliminate false positive and false negative detection of identifiers with handcrafted  
 799 regular expressions.

## 800 Determining Conjecture Scope

801 Determination of conjecture scope decomposes into two subchallenges: attribution of proof  
 802 steps to the correct conjecture and detection of proof (conjecture) completion. Though one  
 803 can query the open goals at any given time, the conjecture name cannot be derived from  
 804 the information (and identical goals can appear in different proofs). Furthermore, steps of  
 805 distinct proofs may be intermingled or nested as shown in 4.

806 A Vernacular command—**Show Conjectures**—again provides the solution in the absence  
 807 of a SerAPI query. This command lists the names of currently stated but unproved conjectures  
 808 and by all observations is guaranteed to list the conjecture actively being proved first, though  
 809 no such guarantee is stated in its documentation. We rely upon this presumed order to  
 810 identify the name of the current conjecture, accumulating proof steps in stacks associated  
 811 with each open conjecture.

812 The accuracy of this approach depends upon the assumption that no conjecture fails to  
 813 enter proof mode after its initial sentence is executed. The only known exceptions to this  
 814 rule comprise Programs, which do not enter proof mode until their first Obligation’s proof is  
 815 begun. Behavior due to violations of this rule is untested as all known plugins are compliant.

816 Special handling is required to associate each Obligation with the correct Program since

```

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

```

■ **Figure 4** A simple example showing that proofs may be interleaved and that multiple proofs (obligations) may be associated with one term.

817 **Show Conjectures** reveals a unique name for each Obligation. These Obligation names  
 818 are reliably constructed from the Program's name as `<program_name>_obligation_<id>`,  
 819 where `id` is an integer. We can thus easily accumulate proof stacks for each Obligation  
 820 and maintain their association with one another. However, the special handling means any  
 821 grammar extension that defines its own Obligation or Program equivalents (i.e., multi-block  
 822 proofs) cannot be serialized to the same level of accuracy. If any extension does so, then  
 823 each Obligation is expected to be serialized as an unrelated theorem.

824 We rely upon detection of definitions to determine when and if a conjecture was proved,  
 825 assuming that no conjecture emits an identifier before it is defined (i.e., before it is proved).  
 826 More precisely, neither verbose user feedback nor **Print All** command should indicate the  
 827 conjecture is defined before its proof(s) are complete. The only allowed exceptions to this  
 828 rule comprise subproofs (generally delimited by braces and bullets) whose identifiers are  
 829 similarly structured to those of Obligations and are filtered with a regular expression. Due  
 830 to nesting, one cannot assume that a detected definition in the midst of a proof corresponds  
 831 to the conjecture. Nor can one assume that the name of the conjecture once defined will  
 832 actually match its name as returned by **Show Conjectures** during the proof since Vernacular  
 833 commands such as **Save <name>** exist, which introduce the conjecture into the environment  
 834 under the name given.

835 We ultimately detect the completion of a proof by requiring two conditions: a change in  
 836 the currently detected conjecture and the detection of a new definition. This rule necessarily  
 837 invokes an additional assumption: a change in the current conjecture implies that either  
 838 a new proof has begun or the current proof has ended (but not both). Mutual exclusivity  
 839 follows from the mutually exclusive classification of Vernacular commands into proof starters

## 23:22 Building a Large Proof Repair Dataset

or ends within Coq’s source code<sup>6</sup>. Since we assume that a conjecture cannot emit an identifier before it is completed, we deduce that the emission of an identifier upon the change of the current conjecture implies the completion of the prior conjecture.

Finally, if the conjecture is aborted, then it will never be detected as a definition at all even though its proof has ended. Detecting an aborted proof requires a different approach for which no alternative other than explicitly checking the type of command was forthcoming. We thus assume that no grammar extension defines its own `Abort` or `Abort All` equivalents, i.e., no other type of command concludes a proof without emitting an identifier.

### Limiting Resources

Efficiently serializing the files of each project and commit in the dataset requires robust automation and multiprocessing infrastructure. **TODO: Wesley: Mention project that consumed 100s of GBs of RAM and several hours: coqrel, RelDefinitions.v, 532be10aa3d14a5a775c99926b813729011f54cf, 45 GB after 5 minutes, over 200 GB some time after (crashed entire extraction process)** **TODO: Incorporate link to open issue about Coq memory usage: <https://github.com/coq/coq/issues/12487>** **TODO: Speak to solution to limiting time and memory: `subprocess.run` with `timeout` arg and `preexec_fn` arg with resource `RLIMIT_AS`** **TODO: Speak to solution to CPU limitation: semaphore to limit number of files being serialized** **TODO: Speak to SwiM resource limitations as well including LRU cache; or not, seems to be covered in Section 4.2**

### Serialization and Version Changes

SerAPI was in theory supposed to help with some of the proof assistant versioning problems mentioned in Section 5.1. In practice, though, SerAPI itself depends on the version of Coq, and we found we had to break the SerAPI abstraction barrier often as the Coq version changed. Some useful features that exist in newer versions of SerAPI are not backported to previous versions. Furthermore, while SerAPI provides a convenient interface to expose certain Coq internals, those internals are not necessarily stable. For example, SerAPI had “can’t-fix” bugs involving nested proofs because the serialization errors occur in the Coq codebase itself.<sup>7</sup> There seems to be an implicit assumption that one is always using the latest version of Coq, which made it challenging to collect data for older versions.

**TODO: integrate any remaining parsing points:**

- Different representations, sometimes you want both and need to query for more info like paths (definitely came up in Passport, too) **TODO: This may be indirectly addressed in the Calculating Command Dependencies section**

## 5.3.2 Other Proof Assistants

**TODO: What other proof assistants do. Here citing PISA is a good call? maybe SerAPI is nice but library on top of it is needed? at least Coq has a parser though, Isabelle really doesn’t have a canonical parser and that’s a huge issue ... Parsing is *even harder* in Isabelle, though things like PISA have helped with this a lot, especially for non-experts**

<sup>6</sup> See the `vernac_classification` type.

<sup>7</sup> <https://github.com/ejgallego/coq-serapi/issues/117>



### 5.3.3 Recommendations

TODO: What the community can do to make things easier going forward. In some sense SerAPI also kind of assumes you're a PL expert, but for ML we may want a very different interface, since folks writing ML datasets and tools are often not PL experts. Or maybe we want to make ML more accessible to PL experts, too. TODO: Exposure of Vernac classification in SerAPI (VtStartProof, VtProofStep, etc.) would help substantially. In particular, we believe that 100% accurate serialization requires access to Coq's Vernacular classification functions that determine the class (e.g., proof starter, proof ender, side-effect, etc.) of a given Vernacular command, whether the command be built-in or a grammar extension. Structure the language to be machine-readable human/compiler-out-of-the-loop or at least provide a public API for parsing it into its major components (for example, C provides braces, Python requires indentation, etc.). To this end, serapi is "still a research, experimental project, and it is expected to evolve considerably"<sup>8</sup>, and there are already ongoing plans for further developments that would expose new features that tackle some existing challenges. For instance, there is recently a plan to rebase serapi on a different project<sup>9</sup>, which exposes features like document overviews that appear to list all the definitions in the file and the ability to fold proofs, implying that it has the capability to list theorems and gather the associated lines— one of our current challenges.

## 6 Related Work

We describe related work in datasets and benchmark suites for programs and proofs, in proof repair, and in machine learning for proofs.

### Datasets & Benchmark Suites

The REPLICA [41] user study collected incremental edit data from eight proof engineers over the course of a month. The data collected for REPLICA was more incremental than the data we collected via Git commits. On the other hand, due to the difficulty in finding volunteers for the user study, the amount of data collected was too small for many modern data-hungry machine-learning tools. In contrast, our dataset is less incremental but much larger. The REPLICA data may make a useful supplement to our current dataset in the future, especially when more incremental data is useful.

There have been a number of datasets and benchmark suites released recently targeting *autoformalization*: the automatic translation of natural language specifications and proofs to formal specifications and proofs in a proof assistant. Recent autoformalization datasets consisting of aligned natural and formal language examples include ProofNet [9] for Lean, and the Isabelle Parallel Corpus [14] for Isabelle/HOL. The MiniF2F [54] benchmark suites includes a number of math Olympiad problems formalized across different proof assistants, and is used to measure success on both autoformalization and proof synthesis tasks.

There are a few datasets, benchmark suites, and proving environments specifically for proof synthesis tasks. Examples include CoqGym [51] for Coq and HOList [11] for HOL Light. The distinguishing feature of our dataset is that it includes repair examples rather than static data. We hope to expand on CoqGym in the near future using static data from the latest versions of projects in our repair dataset.

<sup>8</sup> <https://github.com/ejgallego/coq-serapi/issues/252>

<sup>9</sup> <https://github.com/ejgallego/coq-serapi/issues/252#issuecomment-1365510329>

## 23:24 Building a Large Proof Repair Dataset

918 There is a large amount of work on machine learning for code in recent years, summarized  
919 for example in a recent survey paper on neurosymbolic programming [15]. We expect there  
920 is much to be learned from the lessons, infrastructure, and metrics designed for machine  
921 learning for code, especially given the similarities between code and proofs. For example,  
922 perhaps distance metrics like CodeBLEU [37] will be reasonable proxies for proximity to  
923 correctness in machine learning tasks for formal proof.

924 In the field of software engineering, accessible datasets facilitate new research. For  
925 example, Defects4 [24] is a collection of bugs and patches in Java. It has been frequently  
926 used as a benchmark for automated program repair [19, 46, 28], even recently, despite  
927 nearly a decade having passed since its release. We hope that producing an accessible proof  
928 repair dataset will spur new research in proof repair. We also hope that, by focusing on  
929 good benchmarks and metrics for success early on, we can avoid some of the methodology  
930 challenges faced in program repair [35].

### 931 Proof Repair

932 The machine-learning task that our dataset and benchmark suite focuses on is proof repair,  
933 which is summarized in the namesake thesis [38]. There is not yet published work we are  
934 aware of for machine learning for proof repair, though we are aware of ongoing machine  
935 learning for proof repair projects by other teams in parallel in proof assistants other than  
936 Coq. We plan to train and evaluate at least two distinct proof repair models in Coq using  
937 this dataset, and we hope that this dataset makes it easy for others to do the same.

938 Proof repair is closely related to work in proof reuse [20, 43, 13], proof refactoring [49, 48,  
939 44], and proof transformation [33]. These and other related topics in proof engineering have  
940 a long history, described in detail in the proof engineering survey paper QED at Large [39],  
941 as well as in the proof repair namesake thesis [38].

942 Proof repair can be viewed as program repair [29, 22] for proofs. There is a large amount  
943 of recent work on learning to repair programs, both symbolically (for example, by way of  
944 anti-unification in Getafix [10]) and neurally (for example, by way of unsupervised learning  
945 in Break-It-Fix-It [52]). This recent work may provide useful insights when building machine  
946 learning datasets, benchmark suites, and models for proof repair, though care must be taken  
947 to consider the differences between typical programs and formal proof developments [38].

### 948 Machine Learning for Proofs

949 Advances in machine learning have had a transformative effect on many fields, and theorem  
950 provers are not excluded. Examples of recent work on machine learning for synthesizing  
951 formal proofs include GPT-f [34] and HTPS [26] for Metamath and Lean; Proverbot9001 [45],  
952 ASTactic [51], Tactician [12], and DIVA [21] for Coq; and DeepHOL [11] for HOL Light.  
953 Also of note is recent work on autoformalization in Isabelle/HOL [50], Lean [9], and Coq [16].  
954 More machine learning work for proofs can be found in QED at Large [39]. Our main goal  
955 is to expand the scope of tasks covered in machine learning for proofs, reaching important  
956 tasks not previously explored.

## 957 7 Conclusions & Future Work

958 **TODO: Conclusion will go here.**

959 Moving forward, our immediate plan is to build machine learning models for proof repair  
960 in Coq. We would also like to develop better metrics for measuring success at repairing

961 definitions. Finally, we hope to work with the rest of the proof assistant community to  
 962 address the many challenges we have highlighted, so that we may steer machine learning for  
 963 proofs in the right direction.

## 964 — References —

- 965 1 Europroofnet. URL: <https://europroofnet.github.io/>.
- 966 2 Openai. URL: <https://openai.com/>.
- 967 3 Proof engineering, adaptation, repair, and learning for software (pearls). URL: <https://sam.gov/opp/da84366306554cc981f37f703a78c698/view>.
- 968 4 AI for Theorem Proving, 2016-2022. URL: <http://aitp-conference.org/>.
- 969 5 2nd MATH-AI Workshop at NeurIPS'22, 2021-2022. URL: <https://mathai2022.github.io/>.
- 970 6 Beyond Bayes: Paths Towards Universal Reasoning Systems, 2022. URL: <https://beyond-bayes.github.io/>.
- 971 7 Arpan Agrawal, Emily First, Zhanna Kaufman, Tom Reichel, Shizhuo Zhang, Timothy Zhou,  
 974 Alex Sanchez-Stern, and Talia Ringer. Proofster. URL: <https://www.alexsanchezstern.com/papers/proofster.pdf>.
- 975 8 Emilio Jesús Gallego Arias. SerAPI: Machine-Friendly, Data-Centric Serialization for  
 976 COQ. Technical Report hal-01384408, HAL, 2016. URL: <http://dml.mathdoc.fr/item/hal-01384408/>.
- 977 9 Zhangir Azerbayev, Bartosz Piotrowski, and Jeremy Avigad. Proofnet: A benchmark for  
 978 autoformalizing and formally proving undergraduate-level mathematics problems. In *Second*  
 979 *MATH-AI Workshop*, 2022. URL: <https://mathai2022.github.io/>.
- 980 10 Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. Getafix: Learning to fix  
 981 bugs automatically. *Proc. ACM Program. Lang.*, 3(OOPSLA), oct 2019. doi:10.1145/3360585.
- 982 11 Kshitij Bansal, Sarah M Loos, Markus N Rabe, Christian Szegedy, and Stewart Wilcox. Holist:  
 983 An environment for machine learning of higher order logic theorem proving. In *ICML*, 2019.
- 984 12 Lasse Blaauwbroek, Josef Urban, and Herman Geuvers. The tactician: A seamless, interactive  
 985 tactic learner and prover for coq. In *Intelligent Computer Mathematics: 13th International*  
 986 *Conference, CICM 2020, Bertinoro, Italy, July 26–31, 2020, Proceedings*, page 271–277, Berlin,  
 987 Heidelberg, 2020. Springer-Verlag. doi:10.1007/978-3-030-53518-6\_17.
- 988 13 Olivier Boite. Proof reuse with extended inductive types. In *Theorem Proving in Higher Order*  
 989 *Logics: 17th International Conference, TPHOLs 2004, Park City, Utah, USA, September*  
 990 *14–17, 2004. Proceedings*, pages 50–65, Berlin, Heidelberg, 2004. Springer. doi:10.1007/  
 991 978-3-540-30142-4\_4.
- 992 14 Anthony Bordg, Yiannos A Stathopoulos, and Lawrence C Paulson. A parallel corpus of  
 993 natural language and isabelle artefacts. In *7th Conference on Artificial Intelligence and*  
 994 *Theorem Proving (AITP)*, 2022. URL: [http://aitp-conference.org/2022/abstract/AITP\\_2022\\_paper\\_8.pdf](http://aitp-conference.org/2022/abstract/AITP_2022_paper_8.pdf).
- 995 15 Swarat Chaudhuri, Kevin Ellis, Oleksandr Polozov, Rishabh Singh, Armando Solar-Lezama,  
 996 Yisong Yue, et al. Neurosymbolic programming. *Foundations and Trends® in Programming*  
 997 *Languages*, 7(3):158–243, 2021.
- 998 16 Garrett Cunningham, Razvan C. Bunescu, and David Juedes. Towards autoformalization  
 999 of mathematics and code correctness: Experiments with elementary proofs, 2023. URL:  
 1000 <https://arxiv.org/abs/2301.02195>, doi:10.48550/ARXIV.2301.02195.
- 1001 17 Marco Cuturi. Sinkhorn Distances: Lightspeed Computation of Optimal Trans-  
 1002 port. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q.  
 1003 Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages  
 1004 2292–2300. Curran Associates, Inc., 2013. URL: <http://papers.nips.cc/paper/4927-sinkhorn-distances-lightspeed-computation-of-optimal-transport.pdf>.

- 1009 **18** Michel Marie Deza and Elena Deza. *Encyclopedia of Distances*. Springer Berlin Heidelberg,  
1010 Berlin, 3 edition, October 2014. URL: [https://link.springer.com/book/10.1007/](https://link.springer.com/book/10.1007/978-3-642-30958-8)  
1011 [978-3-642-30958-8](https://link.springer.com/book/10.1007/978-3-642-30958-8).
- 1012 **19** Thomas Durieux, Matias Martinez, Martin Monperrus, Romain Sommerard, and Jifeng  
1013 Xuan. Automatic repair of real bugs: An experience report on the defects4j dataset. *CoRR*,  
1014 [abs/1505.07002](https://arxiv.org/abs/1505.07002), 2015. URL: <http://arxiv.org/abs/1505.07002>, [arXiv:1505.07002](https://arxiv.org/abs/1505.07002).
- 1015 **20** Amy Felty and Douglas Howe. Generalization and reuse of tactic proofs. In *Logic Programming*  
1016 *and Automated Reasoning: 5th International Conference, LPAR '94*, pages 1–15, Berlin,  
1017 Heidelberg, 1994. Springer. doi:10.1007/3-540-58216-9\_25.
- 1018 **21** Emily First and Yuriy Brun. Diversity-driven automated formal verification. In *Proceedings*  
1019 *of the 44th International Conference on Software Engineering (ICSE)(22–27)*. Pittsburgh, PA,  
1020 USA. <https://doi.org/10.1145/3510003.3510138>, 2022.
- 1021 **22** Luca Gazzola, Daniela Micucci, and Leonardo Mariani. Automatic software repair: A survey. In  
1022 *Proceedings of the 40th International Conference on Software Engineering, ICSE '18*, page 1219,  
1023 New York, NY, USA, 2018. Association for Computing Machinery. doi:10.1145/3180155.  
1024 3182526.
- 1025 **23** Albert Jiang, Wenda Li, Jesse Han, and Wu Yuhuai. Lisa: Language models of isabelle proofs.  
1026 In *6th Conference on Artificial Intelligence and Theorem Proving (AITP)*, 2021.
- 1027 **24** René Just, Darioush Jalali, and Michael D. Ernst. Defects4J: A Database of existing faults to  
1028 enable controlled testing studies for Java programs. In *ISSTA 2014, Proceedings of the 2014*  
1029 *International Symposium on Software Testing and Analysis*, pages 437–440, San Jose, CA,  
1030 USA, July 2014. Tool demo.
- 1031 **25** Daniel Kühlwein, Jasmin Christian Blanchette, Cezary Kaliszyk, and Josef Urban. Mash:  
1032 Machine learning for sledgehammer. In Sandrine Blazy, Christine Paulin-Mohring, and David  
1033 Pichardie, editors, *Interactive Theorem Proving*, pages 35–50, Berlin, Heidelberg, 2013. Springer  
1034 Berlin Heidelberg.
- 1035 **26** Guillaume Lample, Marie-Anne Lachaux, Thibaut Lavril, Xavier Martinet, Amaury Hayat,  
1036 Gabriel Ebner, Aurélien Rodriguez, and Timothée Lacroix. Hypertree proof search for neural  
1037 theorem proving. *arXiv preprint arXiv:2205.11491*, 2022.
- 1038 **27** Guillaume Lample, Marie-Anne Lachaux, Thibaut Lavril, Xavier Martinet, Amaury Hayat,  
1039 Gabriel Ebner, Aurélien Rodriguez, and Timothée Lacroix. Hypertree proof search for neural  
1040 theorem proving, 2022. URL: <https://arxiv.org/abs/2205.11491>, doi:10.48550/ARXIV.  
1041 2205.11491.
- 1042 **28** Thibaud Lutellier, Hung Viet Pham, Lawrence Pang, Yitong Li, Moshi Wei, and Lin Tan.  
1043 Coconut: Combining context-aware neural translation models using ensemble for program  
1044 repair. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software*  
1045 *Testing and Analysis*, ISSTA 2020, page 101–114, New York, NY, USA, 2020. Association for  
1046 Computing Machinery. doi:10.1145/3395363.3397369.
- 1047 **29** Martin Monperrus. Automatic software repair: A bibliography. *ACM Comput. Surv.*, 51(1),  
1048 jan 2018. doi:10.1145/3105906.
- 1049 **30** James Munkres. Algorithms for the Assignment and Transportation Problems. *Journal of*  
1050 *the Society for Industrial and Applied Mathematics*, 5(1):32–38, 1957. Publisher: Society for  
1051 Industrial and Applied Mathematics. URL: <https://www.jstor.org/stable/2098689>.
- 1052 **31** Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A method for automatic  
1053 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for*  
1054 *Computational Linguistics*, ACL '02, page 311–318, USA, 2002. Association for Computational  
1055 Linguistics. doi:10.3115/1073083.1073135.
- 1056 **32** Ofir Pele and Michael Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th*  
1057 *International Conference on Computer Vision*, pages 460–467, 2009. doi:10.1109/ICCV.2009.  
1058 5459199.
- 1059 **33** Frank Pfenning. *Proof Transformations in Higher-Order Logic*. PhD thesis, Carnegie Mellon  
1060 University Pittsburgh, 1987.

- 1061 34 Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem  
1062 proving. *CoRR*, abs/2009.03393, 2020. URL: <https://arxiv.org/abs/2009.03393>, arXiv:  
1063 2009.03393.
- 1064 35 Zichao Qi, Fan Long, Sara Achour, and Martin Rinard. An analysis of patch plausibility and  
1065 correctness for generate-and-validate patch generation systems. In *Proceedings of the 2015*  
1066 *International Symposium on Software Testing and Analysis*, ISSTA 2015, page 24–36, New  
1067 York, NY, USA, 2015. Association for Computing Machinery. doi:10.1145/2771783.2771791.
- 1068 36 Markus N. Rabe and Christian Szegedy. Towards the automatic mathematician. In André  
1069 Platzter and Geoff Sutcliffe, editors, *Automated Deduction – CADE 28*, pages 25–37, Cham,  
1070 2021. Springer International Publishing.
- 1071 37 Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming  
1072 Zhou, Ambrosio Blanco, and Shuai Ma. Codebleu: a method for automatic evaluation of  
1073 code synthesis. *CoRR*, abs/2009.10297, 2020. URL: <https://arxiv.org/abs/2009.10297>,  
1074 arXiv:2009.10297.
- 1075 38 Talia Ringer. *Proof Repair*. PhD thesis, University of Washington, 2021.
- 1076 39 Talia Ringer, Karl Palmskog, Ilya Sergey, Milos Gligoric, and Zachary Tatlock. QED at large:  
1077 A survey of engineering of formally verified software. *CoRR*, abs/2003.06458, 2020. URL:  
1078 <https://arxiv.org/abs/2003.06458>, arXiv:2003.06458.
- 1079 40 Talia Ringer, RanDair Porter, Nathaniel Yazdani, John Leo, and Dan Grossman. Proof repair  
1080 across type equivalences. PLDI 2021, page 112–127, New York, NY, USA, 2021. Association  
1081 for Computing Machinery. doi:10.1145/3453483.3454033.
- 1082 41 Talia Ringer, Alex Sanchez-Stern, Dan Grossman, and Sorin Lerner. REPLica: REPL  
1083 instrumentation for Coq analysis. In *Proceedings of the 9th ACM SIGPLAN International*  
1084 *Conference on Certified Programs and Proofs*, CPP 2020, page 99–113, New York, NY, USA,  
1085 2020. Association for Computing Machinery. doi:10.1145/3372885.3373823.
- 1086 42 Talia Ringer, Nathaniel Yazdani, John Leo, and Dan Grossman. Adapting proof automation to  
1087 adapt proofs. In *Proceedings of the 7th ACM SIGPLAN International Conference on Certified*  
1088 *Programs and Proofs*, CPP 2018, page 115–129, New York, NY, USA, 2018. Association for  
1089 Computing Machinery. doi:10.1145/3167094.
- 1090 43 Talia Ringer, Nathaniel Yazdani, John Leo, and Dan Grossman. Ornaments for proof reuse in  
1091 coq. In *Interactive Theorem Proving*, 2019.
- 1092 44 Valentin Robert. *Front-end tooling for building and maintaining dependently-typed functional*  
1093 *programs*. PhD thesis, UC San Diego, 2018.
- 1094 45 Alex Sanchez-Stern, Yousef Alhessi, Lawrence Saul, and Sorin Lerner. Generating cor-  
1095 rectness proofs with neural networks, 2019. URL: <https://arxiv.org/abs/1907.07794>,  
1096 doi:10.48550/ARXIV.1907.07794.
- 1097 46 Ming Wen, Junjie Chen, Rongxin Wu, Dan Hao, and Shing-Chi Cheung. Context-aware patch  
1098 generation for better automated program repair. In *Proceedings of the 40th International*  
1099 *Conference on Software Engineering*, ICSE ’18, page 1–11, New York, NY, USA, 2018.  
1100 Association for Computing Machinery. doi:10.1145/3180155.3180233.
- 1101 47 Makarius Wenzel. Isabelle/jedit — a prover IDE within the PIDE framework. *CoRR*,  
1102 abs/1207.3441, 2012. URL: <http://arxiv.org/abs/1207.3441>, arXiv:1207.3441.
- 1103 48 Iain Johnston Whiteside. *Refactoring proofs*. PhD thesis, University of Edinburgh, November  
1104 2013. URL: <http://hdl.handle.net/1842/7970>.
- 1105 49 Karin Wibergh. Automatic refactoring for agda. Master’s thesis, Chalmers University of  
1106 Technology and University of Gothenburg, 2019.
- 1107 50 Yuhuai Wu, Albert Q Jiang, Wenda Li, Markus N Rabe, Charles Staats, Mateja Jamnik,  
1108 and Christian Szegedy. Autoformalization with large language models. *arXiv preprint*  
1109 *arXiv:2205.12615*, 2022.
- 1110 51 Kaiyu Yang and Jia Deng. Learning to prove theorems via interacting with proof assistants.  
1111 In *International Conference on Machine Learning (ICML)*, Long Beach, CA, USA, 2019. URL:  
1112 <http://proceedings.mlr.press/v97/yang19a/yang19a.pdf>.

## 23:28 Building a Large Proof Repair Dataset

- 1113 52 Michihiro Yasunaga and Percy Liang. Break-it-fix-it: Unsupervised learning for program  
1114 repair. In *International Conference on Machine Learning*, pages 11941–11952. PMLR, 2021.
- 1115 53 John R. Zech, Marcus A. Badgeley, Manway Liu, Anthony B. Costa, Joseph J. Titano, and  
1116 Eric Karl Oermann. Variable generalization performance of a deep learning model to detect  
1117 pneumonia in chest radiographs: A cross-sectional study. *PLOS Medicine*, 15(11):e1002683,  
1118 November 2018. doi:10.1371/journal.pmed.1002683.
- 1119 54 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. minif2f: a cross-system benchmark for  
1120 formal olympiad-level mathematics. In *International Conference on Learning Representations*,  
1121 2022. URL: <https://openreview.net/forum?id=9ZPegFuFTFv>.

## 1122 A Dataset Project Sources

1123 TODO: will reformat table (may need help), and may not include Github links if we'll release  
1124 the dataset itself, since they make things harder to format

1125 TODO: Need to denote what is currently in the dataset and what is not yet there because  
1126 of issues

1127 TODO: ITP submission guidelines say not to include an appendix, not sure what that  
1128 means, may need to move our appendix table online or something? May need to ask  
1129 conference chair.

■ **Table 1** The complete list of projects under consideration for either repair or pretraining datasets. Note that the proof and sentence counts are estimated based upon the heuristic parsing described in TODO: TODO.

Project	Proof Count	Sentence Count	Repository URL
Abel	429	4063	<a href="https://github.com/math-comp/Abel">https://github.com/math-comp/Abel</a>
additions	103	1425	<a href="https://github.com/coq-contribs/additions">https://github.com/coq-contribs/additions</a>
ails	225	4605	<a href="https://github.com/coq-contribs/ails">https://github.com/coq-contribs/ails</a>
algebra	559	6609	<a href="https://github.com/coq-contribs/algebra">https://github.com/coq-contribs/algebra</a>
AML-Formalization	1532	43978	<a href="https://github.com/harp-project/AML-Formalization">https://github.com/harp-project/AML-Formalization</a>
amm11262	28	1154	<a href="https://github.com/coq-contribs/amm11262">https://github.com/coq-contribs/amm11262</a>
analysis	4663	39651	<a href="https://github.com/math-comp/analysis">https://github.com/math-comp/analysis</a>
angles	93	2013	<a href="https://github.com/coq-contribs/angles">https://github.com/coq-contribs/angles</a>
area-method	779	10603	<a href="https://github.com/coq-contribs/area-method">https://github.com/coq-contribs/area-method</a>
argosy	304	3614	<a href="https://github.com/mit-pdos/argosy">https://github.com/mit-pdos/argosy</a>
asn1fpcq	194	2982	<a href="https://github.com/digamma-ai/asn1fpcq">https://github.com/digamma-ai/asn1fpcq</a>
atbr	819	11089	<a href="https://github.com/coq-community/atbr">https://github.com/coq-community/atbr</a>
automata	344	5833	<a href="https://github.com/coq-contribs/automata">https://github.com/coq-contribs/automata</a>
axiomatic-abp	399	4494	<a href="https://github.com/coq-contribs/axiomatic-abp">https://github.com/coq-contribs/axiomatic-abp</a>



banach_tarski	479	8608	<a href="https://github.com/roglo/banach_tarski">https://github.com/roglo/banach_tarski</a>
bbv	653	7343	<a href="https://github.com/mit-plv/bbv">https://github.com/mit-plv/bbv</a>
bdds	232	15510	<a href="https://github.com/coq-contribs/bdds">https://github.com/coq-contribs/bdds</a>
bedrock	4423	66779	<a href="https://github.com/mit-plv/bedrock">https://github.com/mit-plv/bedrock</a>
bedrock2	1233	32668	<a href="https://github.com/mit-plv/bedrock2">https://github.com/mit-plv/bedrock2</a>
bedrock-mirror-shard	1726	18771	<a href="https://github.com/gmalecha/bedrock-mirror-shard">https://github.com/gmalecha/bedrock-mirror-shard</a>
bellantonicook	498	7515	<a href="https://github.com/davidnowak/bellantonicook">https://github.com/davidnowak/bellantonicook</a>
bigenough	5	43	<a href="https://github.com/math-comp/bigenough">https://github.com/math-comp/bigenough</a>
bignums	638	9191	<a href="https://github.com/coq-community/bignums">https://github.com/coq-community/bignums</a>
bits	436	4895	<a href="https://github.com/coq-community/bits">https://github.com/coq-community/bits</a>
buchberger	753	9705	<a href="https://github.com/coq-community/buchberger">https://github.com/coq-community/buchberger</a>
cage	456	7425	<a href="https://github.com/gstew5/cage">https://github.com/gstew5/cage</a>
Categories	519	7757	<a href="https://github.com/amintimany/Categories">https://github.com/amintimany/Categories</a>
category-theory	1377	15827	<a href="https://github.com/jwiegley/category-theory">https://github.com/jwiegley/category-theory</a>
cecoa	1281	33151	<a href="https://github.com/davidnowak/cecoa">https://github.com/davidnowak/cecoa</a>
celsius	211	4470	<a href="https://github.com/clementbladeau/celsius">https://github.com/clementbladeau/celsius</a>
ceramist	427	8223	<a href="https://github.com/verse-lab/ceramist">https://github.com/verse-lab/ceramist</a>
cerise	1026	35029	<a href="https://github.com/logsem/cerise">https://github.com/logsem/cerise</a>
certicoq	4350	133553	<a href="https://github.com/CertiCoq/certicoq">https://github.com/CertiCoq/certicoq</a>
CertiGraph	3522	106982	<a href="https://github.com/CertiGraph/CertiGraph">https://github.com/CertiGraph/CertiGraph</a>
cgraphs	812	16559	<a href="https://github.com/julesjacobs/cgraphs">https://github.com/julesjacobs/cgraphs</a>
ChargeCore	299	3750	<a href="https://github.com/jesper-bengtson/ChargeCore">https://github.com/jesper-bengtson/ChargeCore</a>
checker	4	55	<a href="https://github.com/coq-contribs/checker">https://github.com/coq-contribs/checker</a>
cheerios	83	1078	<a href="https://github.com/uwplse/cheerios">https://github.com/uwplse/cheerios</a>
chinese	137	1769	<a href="https://github.com/coq-contribs/chinese">https://github.com/coq-contribs/chinese</a>
circuits	220	2723	<a href="https://github.com/coq-contribs/circuits">https://github.com/coq-contribs/circuits</a>
ClassicalReal	1303	48201	<a href="https://github.com/QinxiangCao/ClassicalReal">https://github.com/QinxiangCao/ClassicalReal</a>
cls-coq	611	20968	<a href="https://github.com/combinators/cls-coq">https://github.com/combinators/cls-coq</a>



## 23:30 Building a Large Proof Repair Dataset

color	6743	113300	<a href="https://github.com/fblanqui/color">https://github.com/fblanqui/color</a>
CompCert	7835	156983	<a href="https://github.com/AbsInt/CompCert">https://github.com/AbsInt/CompCert</a>
CompCertM	1499	72820	<a href="https://github.com/snu-sf/CompCertM">https://github.com/snu-sf/CompCertM</a>
CompCertR	7817	161208	<a href="https://github.com/snu-sf/CompCertR">https://github.com/snu-sf/CompCertR</a>
concat	514	6417	<a href="https://github.com/coq-contribs/concat">https://github.com/coq-contribs/concat</a>
ConCert	1581	33188	<a href="https://github.com/AU-COBRA/ConCert">https://github.com/AU-COBRA/ConCert</a>
constructive-geometry	116	902	<a href="https://github.com/coq-contribs/constructive-geometry">https://github.com/coq-contribs/constructive-geometry</a>
constructive-ltl	634	7621	<a href="https://github.com/jwiegley/constructive-ltl">https://github.com/jwiegley/constructive-ltl</a>
containers	1417	12466	<a href="https://github.com/coq-contribs/containers">https://github.com/coq-contribs/containers</a>
coq_real	843	24386	<a href="https://github.com/roglo/coq_real">https://github.com/roglo/coq_real</a>
coq-bitset	147	3755	<a href="https://github.com/artart78/coq-bitset">https://github.com/artart78/coq-bitset</a>
coq-ceres	122	1263	<a href="https://github.com/Lysxia/coq-ceres">https://github.com/Lysxia/coq-ceres</a>
Coq-Combi	3520	37470	<a href="https://github.com/math-comp/Coq-Combi">https://github.com/math-comp/Coq-Combi</a>
coq-compile	29	2330	<a href="https://github.com/coq-ext-lib/coq-compile">https://github.com/coq-ext-lib/coq-compile</a>
coq-cunit	0	7	<a href="https://github.com/clarus/coq-cunit">https://github.com/clarus/coq-cunit</a>
coqeal	1609	17503	<a href="https://github.com/coq-community/coqeal">https://github.com/coq-community/coqeal</a>
coq-error-handlers	0	10	<a href="https://github.com/coq-error-handlers">https://github.com/coq-error-handlers</a>
coq-ext-lib	297	4842	<a href="https://github.com/coq-community/coq-ext-lib">https://github.com/coq-community/coq-ext-lib</a>
Coq-Flow-Equivalence	376	8278	<a href="https://github.com/GaloisInc/Coq-Flow-Equivalence">https://github.com/GaloisInc/Coq-Flow-Equivalence</a>
coq-forcing	72	1190	<a href="https://github.com/CoqHott/coq-forcing">https://github.com/CoqHott/coq-forcing</a>
coq-forcing	72	1190	<a href="https://github.com/ppedrot/coq-forcing">https://github.com/ppedrot/coq-forcing</a>
coq-function-ninjas	0	3	<a href="https://github.com/coq-function-ninjas">https://github.com/coq-function-ninjas</a>
coq-guarded-computational-type-theory	205	3824	<a href="https://github.com/jonsterling/coq-guarded-computational-type-theory">https://github.com/jonsterling/coq-guarded-computational-type-theory</a>
coq-haskell	481	5512	<a href="https://github.com/jwiegley/coq-haskell">https://github.com/jwiegley/coq-haskell</a>
coq-http	27	795	<a href="https://github.com/liyishuai/coq-http">https://github.com/liyishuai/coq-http</a>
coq-http2	47	943	<a href="https://github.com/liyishuai/coq-http2">https://github.com/liyishuai/coq-http2</a>
coq-iterable	0	15	<a href="https://github.com/clarus/coq-iterable">https://github.com/clarus/coq-iterable</a>

coq-library-complexity	2438	43607	<a href="https://github.com/uds-psl/">https://github.com/uds-psl/</a>
coq-library-undecidability	10281	153181	<a href="https://github.com/uds-psl/coq-library-undecidability">https://github.com/uds-psl/coq-library-undecidability</a>
coq-library-undecidability	10281	153181	<a href="https://github.com/uds-psl/coq-library-undecidability">https://github.com/uds-psl/coq-library-undecidability</a>
coq-list-plus	0	42	<a href="https://github.com/clarus/coq-list-plus">https://github.com/clarus/coq-list-plus</a>
coq-list-string	9	207	<a href="https://github.com/clarus/coq-list-string">https://github.com/clarus/coq-list-string</a>
coqoban	3	455	<a href="https://github.com/coq-community/coqoban">https://github.com/coq-community/coqoban</a>
coq-performance-tests	170	2536	<a href="https://github.com/coq-community/coq-performance-tests">https://github.com/coq-community/coq-performance-tests</a>
Coq-Polyhedra	1086	11773	<a href="https://github.com/Coq-Polyhedra/Coq-Polyhedra">https://github.com/Coq-Polyhedra/Coq-Polyhedra</a>
coq-procrastination	60	506	<a href="https://github.com/Armael/coq-procrastination">https://github.com/Armael/coq-procrastination</a>
coq-record-update	6	243	<a href="https://github.com/tchajed/coq-record-update">https://github.com/tchajed/coq-record-update</a>
coqrel	270	1763	<a href="https://github.com/CertiKOS/coqrel">https://github.com/CertiKOS/coqrel</a>
coq-robot	1505	13575	<a href="https://github.com/affeldt-aist/coq-robot">https://github.com/affeldt-aist/coq-robot</a>
coq-simple-io	0	267	<a href="https://github.com/Lysxia/coq-simple-io">https://github.com/Lysxia/coq-simple-io</a>
coqtail-math	2515	35411	<a href="https://github.com/coq-community/coqtail-math">https://github.com/coq-community/coqtail-math</a>
coqutil	797	9516	<a href="https://github.com/mit-plv/coqutil">https://github.com/mit-plv/coqutil</a>
coq-utils	298	2820	<a href="https://github.com/arthuraa/coq-utils">https://github.com/arthuraa/coq-utils</a>
Core-Erlang-Formalization	550	12170	<a href="https://github.com/harp-project/Core-Erlang-Formalization">https://github.com/harp-project/Core-Erlang-Formalization</a>
corespec	1510	21682	<a href="https://github.com/sweirich/corespec">https://github.com/sweirich/corespec</a>
cours-de-coq	92	792	<a href="https://github.com/coq-contribs/cours-de-coq">https://github.com/coq-contribs/cours-de-coq</a>
cpdt-japanese	431	3152	<a href="https://github.com/cpdt-japanese/cpdt-japanese">https://github.com/cpdt-japanese/cpdt-japanese</a>
cps	436	9211	<a href="https://github.com/takanuva/cps">https://github.com/takanuva/cps</a>
crellvm	722	22052	<a href="https://github.com/snu-sf/crellvm">https://github.com/snu-sf/crellvm</a>
cryptis	448	5751	<a href="https://github.com/arthuraa/cryptis">https://github.com/arthuraa/cryptis</a>
cspec	884	11566	<a href="https://github.com/mit-pdos/cspec">https://github.com/mit-pdos/cspec</a>
ct	161	1948	<a href="https://github.com/relrod/ct">https://github.com/relrod/ct</a>

## 23:32 Building a Large Proof Repair Dataset

ctlctl	28	383	<a href="https://github.com/coq-contribs/ctlctl">https://github.com/coq-contribs/ctlctl</a>
dblib	244	1994	<a href="https://github.com/coq-community/dblib">https://github.com/coq-community/dblib</a>
DeepSpecDB	2719	86898	<a href="https://github.com/PrincetonUniversity/DeepSpecDB">https://github.com/PrincetonUniversity/DeepSpecDB</a>
demos	69	477	<a href="https://github.com/coq-contribs/demos">https://github.com/coq-contribs/demos</a>
dep-map	90	1843	<a href="https://github.com/coq-contribs/dep-map">https://github.com/coq-contribs/dep-map</a>
deriving	67	1567	<a href="https://github.com/arthuraa/deriving">https://github.com/arthuraa/deriving</a>
dez	1523	20350	<a href="https://github.com/Tuplanolla/dez">https://github.com/Tuplanolla/dez</a>
dictionaries	67	845	<a href="https://github.com/coq-contribs/dictionaries">https://github.com/coq-contribs/dictionaries</a>
disel	844	12582	<a href="https://github.com/DistributedComponents/disel">https://github.com/DistributedComponents/disel</a>
distributed-reference-counting	954	24025	<a href="https://github.com/coq-contribs/distributed-reference-counting">https://github.com/coq-contribs/distributed-reference-counting</a>
domains	1738	58184	<a href="https://github.com/robdockins/domains">https://github.com/robdockins/domains</a>
domain-theory	51	812	<a href="https://github.com/coq-contribs/domain-theory">https://github.com/coq-contribs/domain-theory</a>
dot-iris	1590	12473	<a href="https://github.com/Blaisorblade/dot-iris">https://github.com/Blaisorblade/dot-iris</a>
engine-bench	94	1264	<a href="https://github.com/mit-plv/engine-bench">https://github.com/mit-plv/engine-bench</a>
ett-to-wtt	433	13925	<a href="https://github.com/TheoWinterhalter/ett-to-wtt">https://github.com/TheoWinterhalter/ett-to-wtt</a>
euler-formula	162	6757	<a href="https://github.com/coq-contribs/euler-formula">https://github.com/coq-contribs/euler-formula</a>
event-struct	1039	11142	<a href="https://github.com/Event-Structures/event-struct">https://github.com/Event-Structures/event-struct</a>
exceptions	3	94	<a href="https://github.com/coq-contribs/exceptions">https://github.com/coq-contribs/exceptions</a>
exploring-robust-property-preservation	295	6518	<a href="https://github.com/secure-compilation/exploring-robust-property-preservation">https://github.com/secure-compilation/exploring-robust-property-preservation</a>
extructures	357	2986	<a href="https://github.com/arthuraa/extructures">https://github.com/arthuraa/extructures</a>
fcs1-pcm	2118	15806	<a href="https://github.com/imdea-software/fcs1-pcm">https://github.com/imdea-software/fcs1-pcm</a>
fermat4	130	841	<a href="https://github.com/coq-contribs/fermat4">https://github.com/coq-contribs/fermat4</a>
fiat	5932	84544	<a href="https://github.com/mit-plv/fiat">https://github.com/mit-plv/fiat</a>
finmap	870	5058	<a href="https://github.com/math-comp/finmap">https://github.com/math-comp/finmap</a>
finmap	870	5058	<a href="https://github.com/math-comp/finmap">https://github.com/math-comp/finmap</a>
float	780	10192	<a href="https://github.com/coq-contribs/float">https://github.com/coq-contribs/float</a>

FormalML	7026	135311	<a href="https://github.com/IBM/FormalML">https://github.com/IBM/FormalML</a>
formal-type-theory	67	10680	<a href="https://github.com/TheoWinterhalter/formal-type-theory">https://github.com/TheoWinterhalter/formal-type-theory</a>
frap	1870	23058	<a href="https://github.com/achlipala/frap">https://github.com/achlipala/frap</a>
free-groups	33	484	<a href="https://github.com/coq-contribs/free-groups">https://github.com/coq-contribs/free-groups</a>
fssec-model	153	3255	<a href="https://github.com/coq-contribs/fssec-model">https://github.com/coq-contribs/fssec-model</a>
functional-algebra	278	1286	<a href="https://github.com/llee454/functional-algebra">https://github.com/llee454/functional-algebra</a>
functions-in-zfc	632	2437	<a href="https://github.com/coq-contribs/functions-in-zfc">https://github.com/coq-contribs/functions-in-zfc</a>
fundamental-arithmetics	152	1941	<a href="https://github.com/coq-contribs/fundamental-arithmetics">https://github.com/coq-contribs/fundamental-arithmetics</a>
general-type-theories	546	10172	<a href="https://github.com/peterlefanulumsdaine/general-type-theories">https://github.com/peterlefanulumsdaine/general-type-theories</a>
generic-environments	269	2818	<a href="https://github.com/coq-community/generic-environments">https://github.com/coq-community/generic-environments</a>
GeoCoq	4087	122527	<a href="https://github.com/GeoCoq/GeoCoq">https://github.com/GeoCoq/GeoCoq</a>
goedel	103	8548	<a href="https://github.com/coq-community/goedel">https://github.com/coq-community/goedel</a>
GraphCoQL	220	2513	<a href="https://github.com/imfd/GraphCoQL">https://github.com/imfd/GraphCoQL</a>
graphs	174	4460	<a href="https://github.com/coq-contribs/graphs">https://github.com/coq-contribs/graphs</a>
graph-theory	2272	29435	<a href="https://github.com/coq-community/graph-theory">https://github.com/coq-community/graph-theory</a>
groups	14	134	<a href="https://github.com/coq-contribs/groups">https://github.com/coq-contribs/groups</a>
group-theory	105	1099	<a href="https://github.com/coq-contribs/group-theory">https://github.com/coq-contribs/group-theory</a>
hahn	1458	9836	<a href="https://github.com/vafeiadis/hahn">https://github.com/vafeiadis/hahn</a>
hanoi	577	10514	<a href="https://github.com/thery/hanoi">https://github.com/thery/hanoi</a>
hardware	186	1643	<a href="https://github.com/coq-contribs/hardware">https://github.com/coq-contribs/hardware</a>
hedges	110	2310	<a href="https://github.com/coq-contribs/hedges">https://github.com/coq-contribs/hedges</a>
helix	2046	56461	<a href="https://github.com/vzaliva/helix">https://github.com/vzaliva/helix</a>
higman-cf	30	279	<a href="https://github.com/coq-contribs/higman-cf">https://github.com/coq-contribs/higman-cf</a>
higman-s	49	913	<a href="https://github.com/coq-contribs/higman-s">https://github.com/coq-contribs/higman-s</a>
hoare-tut	25	346	<a href="https://github.com/coq-community/hoare-tut">https://github.com/coq-community/hoare-tut</a>
HoTT	4974	62461	<a href="https://github.com/HoTT/HoTT">https://github.com/HoTT/HoTT</a>

## 23:34 Building a Large Proof Repair Dataset

HoTT-categories	380	4559	<a href="https://github.com/CategoricalData/HoTT-categories">https://github.com/CategoricalData/HoTT-categories</a>
htt	410	4143	<a href="https://github.com/imdea-software/htt">https://github.com/imdea-software/htt</a>
huffman	285	3910	<a href="https://github.com/coq-community/huffman">https://github.com/coq-community/huffman</a>
hybrid	524	7099	<a href="https://github.com/Eelis/hybrid">https://github.com/Eelis/hybrid</a>
idxassoc	58	713	<a href="https://github.com/coq-contribs/idxassoc">https://github.com/coq-contribs/idxassoc</a>
ieee754	15	239	<a href="https://github.com/coq-contribs/ieee754">https://github.com/coq-contribs/ieee754</a>
IFC	756	11435	<a href="https://github.com/QuickChick/IFC">https://github.com/QuickChick/IFC</a>
imm	1444	24624	<a href="https://github.com/weakmemory/imm">https://github.com/weakmemory/imm</a>
infinitary-rewriting-coq	262	4207	<a href="https://github.com/martijnvermaat/infinitary-rewriting-coq">https://github.com/martijnvermaat/infinitary-rewriting-coq</a>
infotheo	3335	45155	<a href="https://github.com/affeldt-aist/infotheo">https://github.com/affeldt-aist/infotheo</a>
InfSeqExt	145	1943	<a href="https://github.com/DistributedComponents/InfSeqExt">https://github.com/DistributedComponents/InfSeqExt</a>
InteractionTrees	1822	45795	<a href="https://github.com/DeepSpec/InteractionTrees">https://github.com/DeepSpec/InteractionTrees</a>
int-map	494	7816	<a href="https://github.com/coq-contribs/int-map">https://github.com/coq-contribs/int-map</a>
ipc	479	9102	<a href="https://github.com/coq-contribs/ipc">https://github.com/coq-contribs/ipc</a>
iris-coq	4299	32669	<a href="https://github.com/izgzhen/iris-coq">https://github.com/izgzhen/iris-coq</a>
iron	1519	30137	<a href="https://github.com/discus-lang/iron">https://github.com/discus-lang/iron</a>
izf	122	827	<a href="https://github.com/coq-contribs/izf">https://github.com/coq-contribs/izf</a>
jordan-curve-theorem	663	42492	<a href="https://github.com/coq-contribs/jordan-curve-theorem">https://github.com/coq-contribs/jordan-curve-theorem</a>
kami	2092	41182	<a href="https://github.com/mit-plv/kami">https://github.com/mit-plv/kami</a>
Kami	2120	37732	<a href="https://github.com/sifive/Kami">https://github.com/sifive/Kami</a>
katamaran	854	14929	<a href="https://github.com/katamaran-project/katamaran">https://github.com/katamaran-project/katamaran</a>
Ktheory	40	505	<a href="https://github.com/DanGrayson/Ktheory">https://github.com/DanGrayson/Ktheory</a>
lambda	90	941	<a href="https://github.com/coq-contribs/lambda">https://github.com/coq-contribs/lambda</a>
largecatmodules	467	9644	<a href="https://github.com/UniMath/largecatmodules">https://github.com/UniMath/largecatmodules</a>
lazy-pcf	83	1587	<a href="https://github.com/coq-contribs/lazy-pcf">https://github.com/coq-contribs/lazy-pcf</a>
lemma-overloading	840	6476	<a href="https://github.com/coq-community/lemma-overloading">https://github.com/coq-community/lemma-overloading</a>
lesniewski-mereology	140	1396	<a href="https://github.com/coq-contribs/lesniewski-mereology">https://github.com/coq-contribs/lesniewski-mereology</a>

lin-alg	543	10797	<a href="https://github.com/coq-contribs/lin-alg">https://github.com/coq-contribs/lin-alg</a>
linearscan	304	4786	<a href="https://github.com/jwiegle/linearscan">https://github.com/jwiegle/linearscan</a>
llvmtwin-coq	675	34538	<a href="https://github.com/snu-sf/llvmtwin-coq">https://github.com/snu-sf/llvmtwin-coq</a>
LOGIC	1729	25450	<a href="https://github.com/QinxiangCao/LOGIC">https://github.com/QinxiangCao/LOGIC</a>
loopring-protocol2-verification	155	3718	<a href="https://github.com/sec-bit/loopring-protocol2-verification">https://github.com/sec-bit/loopring-protocol2-verification</a>
lvc	3782	60531	<a href="https://github.com/sigurdschneider/lvc">https://github.com/sigurdschneider/lvc</a>
maple-mode	20	154	<a href="https://github.com/coq-contribs/maple-mode">https://github.com/coq-contribs/maple-mode</a>
MapleS	7629	153467	<a href="https://github.com/namefanwjcom/MapleS">https://github.com/namefanwjcom/MapleS</a>
markov	31	951	<a href="https://github.com/coq-contribs/markov">https://github.com/coq-contribs/markov</a>
math-classes	1741	17370	<a href="https://github.com/coq-community/math-classes">https://github.com/coq-community/math-classes</a>
math-comp	14263	104712	<a href="https://github.com/math-comp/math-comp">https://github.com/math-comp/math-comp</a>
maths	38	386	<a href="https://github.com/coq-contribs/math">https://github.com/coq-contribs/math</a>
memory-safe-language	72	1325	<a href="https://github.com/arthuraa/memory-safe-language">https://github.com/arthuraa/memory-safe-language</a>
metacoq	8885	161078	<a href="https://github.com/MetaCoq/metacoq">https://github.com/MetaCoq/metacoq</a>
metalib	796	8026	<a href="https://github.com/plclub/metalib">https://github.com/plclub/metalib</a>
micro-policies-coq	569	13727	<a href="https://github.com/micro-policies/micro-policies-coq">https://github.com/micro-policies/micro-policies-coq</a>
mini-compiler	4	86	<a href="https://github.com/coq-contribs/mini-compiler">https://github.com/coq-contribs/mini-compiler</a>
miniml	12	681	<a href="https://github.com/coq-contribs/miniml">https://github.com/coq-contribs/miniml</a>
mirror-core	1380	36458	<a href="https://github.com/gmalecha/mirror-core">https://github.com/gmalecha/mirror-core</a>
mirror-shard	636	12465	<a href="https://github.com/gmalecha/mirror-shard">https://github.com/gmalecha/mirror-shard</a>
mod-red	271	6204	<a href="https://github.com/coq-contribs/mod-red">https://github.com/coq-contribs/mod-red</a>
monae	1661	16336	<a href="https://github.com/affeldt-aist/monae">https://github.com/affeldt-aist/monae</a>
multinomials	1008	7573	<a href="https://github.com/math-comp/multinomials">https://github.com/math-comp/multinomials</a>
oeuf	10471	186512	<a href="https://github.com/uwplse/oeuf">https://github.com/uwplse/oeuf</a>
orb-stab	40	1590	<a href="https://github.com/coq-contribs/orb-stab">https://github.com/coq-contribs/orb-stab</a>
otway-rees	21	613	<a href="https://github.com/coq-contribs/otway-rees">https://github.com/coq-contribs/otway-rees</a>
paco	1519	23749	<a href="https://github.com/snu-sf/paco">https://github.com/snu-sf/paco</a>
paramcoq-iff	285	6527	<a href="https://github.com/aa755/paramcoq-iff">https://github.com/aa755/paramcoq-iff</a>
param-pi	72	3853	<a href="https://github.com/coq-contribs/param-pi">https://github.com/coq-contribs/param-pi</a>

## 23:36 Building a Large Proof Repair Dataset

parsing-parses	274	3512	<a href="https://github.com/JasonGross/parsing-parses">https://github.com/JasonGross/parsing-parses</a>
perennial	5284	137211	<a href="https://github.com/mit-pdos/perennial">https://github.com/mit-pdos/perennial</a>
pigeons	6	34	<a href="https://github.com/llee454/pigeons">https://github.com/llee454/pigeons</a>
pipcore	1637	130325	<a href="https://github.com/2xs/pipcore">https://github.com/2xs/pipcore</a>
pnf	356	2857	<a href="https://github.com/ilyasergey/pnf">https://github.com/ilyasergey/pnf</a>
pocklington	265	6546	<a href="https://github.com/coq-community/pocklington">https://github.com/coq-community/pocklington</a>
PolTac	309	2247	<a href="https://github.com/thery/PolTac">https://github.com/thery/PolTac</a>
probchain	246	4801	<a href="https://github.com/certichain/probchain">https://github.com/certichain/probchain</a>
ProcKami	54	2055	<a href="https://github.com/sifive/ProcKami">https://github.com/sifive/ProcKami</a>
promising-coq	835	29947	<a href="https://github.com/snu-sf/promising-coq">https://github.com/snu-sf/promising-coq</a>
proofs	174	2558	<a href="https://github.com/stepchowfun/proofs">https://github.com/stepchowfun/proofs</a>
propcalc	66	629	<a href="https://github.com/coq-contribs/propcalc">https://github.com/coq-contribs/propcalc</a>
pts	350	4513	<a href="https://github.com/coq-contribs/pts">https://github.com/coq-contribs/pts</a>
puiseuxth	910	20553	<a href="https://github.com/roglo/puiseuxth">https://github.com/roglo/puiseuxth</a>
qarith	60	1151	<a href="https://github.com/coq-contribs/qarith">https://github.com/coq-contribs/qarith</a>
qarith-stern-brocot	1131	16216	<a href="https://github.com/coq-community/qarith-stern-brocot">https://github.com/coq-community/qarith-stern-brocot</a>
quadcopter	1061	34482	<a href="https://github.com/dricketts/quadcopter">https://github.com/dricketts/quadcopter</a>
QuantumLib	1482	20473	<a href="https://github.com/inQWIRE/QuantumLib">https://github.com/inQWIRE/QuantumLib</a>
QuickChick	1211	13104	<a href="https://github.com/QuickChick/QuickChick">https://github.com/QuickChick/QuickChick</a>
quicksort-complexity	591	8858	<a href="https://github.com/coq-contribs/quicksort-complexity">https://github.com/coq-contribs/quicksort-complexity</a>
QWIRE	1167	18932	<a href="https://github.com/inQWIRE/QWIRE">https://github.com/inQWIRE/QWIRE</a>
railroad-crossing	100	1108	<a href="https://github.com/coq-contribs/railroad-crossing">https://github.com/coq-contribs/railroad-crossing</a>
ramsey	8	55	<a href="https://github.com/coq-contribs/ramsey">https://github.com/coq-contribs/ramsey</a>
regex	105	1556	<a href="https://github.com/coq-contribs/regex">https://github.com/coq-contribs/regex</a>
regex-reexamined-coq	476	7609	<a href="https://github.com/awalterschulze/regex-reexamined-coq">https://github.com/awalterschulze/regex-reexamined-coq</a>
rem	8	113	<a href="https://github.com/coq-contribs/rem">https://github.com/coq-contribs/rem</a>
rewriter	1436	13713	<a href="https://github.com/mit-plv/rewriter">https://github.com/mit-plv/rewriter</a>
rezk_completion	350	4261	<a href="https://github.com/benediktahrens/rezk_completion">https://github.com/benediktahrens/rezk_completion</a>



riscv-coq	209	3367	<a href="https://github.com/samuelgruetter/riscv-coq">https://github.com/samuelgruetter/riscv-coq</a>
rsa	113	1194	<a href="https://github.com/coq-contribs/rsa">https://github.com/coq-contribs/rsa</a>
ruler-compass-geometry	332	4376	<a href="https://github.com/coq-contribs/ruler-compass-geometry">https://github.com/coq-contribs/ruler-compass-geometry</a>
SCEV-coq	23	475	<a href="https://github.com/bollu/SCEV-coq">https://github.com/bollu/SCEV-coq</a>
schroeder	14	276	<a href="https://github.com/coq-contribs/schroeder">https://github.com/coq-contribs/schroeder</a>
search-trees	50	593	<a href="https://github.com/coq-contribs/search-trees">https://github.com/coq-contribs/search-trees</a>
SeLoC	381	7840	<a href="https://github.com/co-dan/SeLoC">https://github.com/co-dan/SeLoC</a>
Set-Theory	2703	66081	<a href="https://github.com/choukh/Set-Theory">https://github.com/choukh/Set-Theory</a>
shuffle	40	430	<a href="https://github.com/coq-contribs/shuffle">https://github.com/coq-contribs/shuffle</a>
silveroak	1519	26477	<a href="https://github.com/project-oak/silveroak">https://github.com/project-oak/silveroak</a>
sirtt	179	4645	<a href="https://github.com/TheoWinterhalter/sirtt">https://github.com/TheoWinterhalter/sirtt</a>
smc	685	25108	<a href="https://github.com/coq-contribs/smc">https://github.com/coq-contribs/smc</a>
SQIR	1713	37680	<a href="https://github.com/inQWIRE/SQIR">https://github.com/inQWIRE/SQIR</a>
SquiggleEq	1302	16922	<a href="https://github.com/aa755/SquiggleEq">https://github.com/aa755/SquiggleEq</a>
ssprove	1298	21179	<a href="https://github.com/SSProve/ssprove">https://github.com/SSProve/ssprove</a>
ssrbit	344	2966	<a href="https://github.com/ejgallego/ssrbit">https://github.com/ejgallego/ssrbit</a>
StdLibKami	72	4505	<a href="https://github.com/sifive/StdLibKami">https://github.com/sifive/StdLibKami</a>
StructTact	287	3563	<a href="https://github.com/uwplse/StructTact">https://github.com/uwplse/StructTact</a>
subst	417	3650	<a href="https://github.com/coq-contribs/subst">https://github.com/coq-contribs/subst</a>
sudoku	250	4015	<a href="https://github.com/coq-community/sudoku">https://github.com/coq-community/sudoku</a>
system	1	112	<a href="https://github.com/coq-concurrency/system">https://github.com/coq-concurrency/system</a>
tarjan	459	5187	<a href="https://github.com/math-comp/tarjan">https://github.com/math-comp/tarjan</a>
tarski-geometry	136	2702	<a href="https://github.com/coq-contribs/tarski-geometry">https://github.com/coq-contribs/tarski-geometry</a>
three-gap	81	1091	<a href="https://github.com/coq-contribs/three-gap">https://github.com/coq-contribs/three-gap</a>
tlc	2314	24176	<a href="https://github.com/charguer/tlc">https://github.com/charguer/tlc</a>
topology	616	16461	<a href="https://github.com/coq-community/topology">https://github.com/coq-community/topology</a>
tortoise-hare-algorithm	4	82	<a href="https://github.com/coq-contribs/tortoise-hare-algorithm">https://github.com/coq-contribs/tortoise-hare-algorithm</a>
toychain	264	4320	<a href="https://github.com/certichain/toychain">https://github.com/certichain/toychain</a>
transfer	201	1726	<a href="https://github.com/Zimmi48/transfer">https://github.com/Zimmi48/transfer</a>

## 23:38 Building a Large Proof Repair Dataset

traversable-fincontainer	71	871	<a href="https://github.com/coq-contribs/traversable-fincontainer">https://github.com/coq-contribs/traversable-fincontainer</a>
tree-automata	834	24988	<a href="https://github.com/coq-contribs/tree-automata">https://github.com/coq-contribs/tree-automata</a>
twoSquare	202	1769	<a href="https://github.com/theytwoSquare">https://github.com/theytwoSquare</a>
TypeTheory	1923	30575	<a href="https://github.com/UniMath/TypeTheory">https://github.com/UniMath/TypeTheory</a>
UnifySL	1193	19223	<a href="https://github.com/QinxiangCao/UnifySL">https://github.com/QinxiangCao/UnifySL</a>
UnifySL	1193	19223	<a href="https://github.com/QinxiangCao/UnifySL">https://github.com/QinxiangCao/UnifySL</a>
UniMath	15201	256571	<a href="https://github.com/UniMath/UniMath">https://github.com/UniMath/UniMath</a>
univalent_parametricity	399	5121	<a href="https://github.com/CoqHott/univalent_parametricity">https://github.com/CoqHott/univalent_parametricity</a>
Valuations	184	5079	<a href="https://github.com/FFaissolle/Valuations">https://github.com/FFaissolle/Valuations</a>
vellvm-legacy	2484	40111	<a href="https://github.com/vellvm/vellvm-legacy">https://github.com/vellvm/vellvm-legacy</a>
velus	3338	58852	<a href="https://github.com/INRIA/velus">https://github.com/INRIA/velus</a>
verdi	618	12534	<a href="https://github.com/uwplse/verdi">https://github.com/uwplse/verdi</a>
verdi-chord	990	16420	<a href="https://github.com/DistributedComponents/verdi-chord">https://github.com/DistributedComponents/verdi-chord</a>
verdi-raft	2272	42557	<a href="https://github.com/uwplse/verdi-raft">https://github.com/uwplse/verdi-raft</a>
Verified-FEC	994	20957	<a href="https://github.com/verified-network-toolchain/Verified-FEC">https://github.com/verified-network-toolchain/Verified-FEC</a>
verified-ifc	884	15495	<a href="https://github.com/micro-policies/verified-ifc">https://github.com/micro-policies/verified-ifc</a>
VeriGHC	118	1522	<a href="https://github.com/trommler/VeriGHC">https://github.com/trommler/VeriGHC</a>
VST	27777	570848	<a href="https://github.com/PrincetonUniversity/VST">https://github.com/PrincetonUniversity/VST</a>
WasmCert-Coq	574	13803	<a href="https://github.com/WasmCert/WasmCert-Coq">https://github.com/WasmCert/WasmCert-Coq</a>
weakestmoToImm	1282	33541	<a href="https://github.com/weakmemory/weakestmoToImm">https://github.com/weakmemory/weakestmoToImm</a>
weak-up-to	150	1720	<a href="https://github.com/coq-contribs/weak-up-to">https://github.com/coq-contribs/weak-up-to</a>
when-good-components-go-bad	1030	23295	<a href="https://github.com/secure-compilation/when-good-components-go-bad">https://github.com/secure-compilation/when-good-components-go-bad</a>
yalla	1019	30165	<a href="https://github.com/olaure01/yalla">https://github.com/olaure01/yalla</a>
ynot	885	9996	<a href="https://github.com/ynot-harvard/ynot">https://github.com/ynot-harvard/ynot</a>
zchinese	43	781	<a href="https://github.com/coq-contribs/zchinese">https://github.com/coq-contribs/zchinese</a>
zf	213	2977	<a href="https://github.com/coq-contribs/zf">https://github.com/coq-contribs/zf</a>
zfc	241	2450	<a href="https://github.com/coq-contribs/zfc">https://github.com/coq-contribs/zfc</a>

zorns-lemma	183	4482	<a href="https://github.com/coq-community/zorns-lemma">https://github.com/coq-community/zorns-lemma</a>
zsearch-trees	48	639	<a href="https://github.com/coq-contribs/zsearch-trees">https://github.com/coq-contribs/zsearch-trees</a>
<b>Total</b>	<b>344,316</b>	<b>6,199,931</b>	

■ **Table 2** The list of projects accepted for inclusion into either the repair or pretraining datasets. Note that the proof and sentence counts are estimated based upon the heuristic parsing described in **TODO: TODO**.

Project	Proof Count	Sentence Count	Repository URL
Abel	429	4063	<a href="https://github.com/math-comp/Abel">https://github.com/math-comp/Abel</a>
additions	103	1425	<a href="https://github.com/coq-contribs/additions">https://github.com/coq-contribs/additions</a>
ails	225	4605	<a href="https://github.com/coq-contribs/ails">https://github.com/coq-contribs/ails</a>
algebra	559	6609	<a href="https://github.com/coq-contribs/algebra">https://github.com/coq-contribs/algebra</a>
amm11262	28	1154	<a href="https://github.com/coq-contribs/amm11262">https://github.com/coq-contribs/amm11262</a>
analysis	4663	39651	<a href="https://github.com/math-comp/analysis">https://github.com/math-comp/analysis</a>
angles	93	2013	<a href="https://github.com/coq-contribs/angles">https://github.com/coq-contribs/angles</a>
area-method	779	10603	<a href="https://github.com/coq-contribs/area-method">https://github.com/coq-contribs/area-method</a>
asn1fpcq	194	2982	<a href="https://github.com/digamma-ai/asn1fpcq">https://github.com/digamma-ai/asn1fpcq</a>
atbr	819	11089	<a href="https://github.com/coq-community/atbr">https://github.com/coq-community/atbr</a>
automata	344	5833	<a href="https://github.com/coq-contribs/automata">https://github.com/coq-contribs/automata</a>
axiomatic-abp	399	4494	<a href="https://github.com/coq-contribs/axiomatic-abp">https://github.com/coq-contribs/axiomatic-abp</a>
bbv	653	7343	<a href="https://github.com/mit-plv/bbv">https://github.com/mit-plv/bbv</a>
bdds	232	15510	<a href="https://github.com/coq-contribs/bdds">https://github.com/coq-contribs/bdds</a>
bellantonicook	498	7515	<a href="https://github.com/davidnowak/bellantonicook">https://github.com/davidnowak/bellantonicook</a>
bigenough	5	43	<a href="https://github.com/math-comp/bigenough">https://github.com/math-comp/bigenough</a>
bignums	638	9191	<a href="https://github.com/coq-community/bignums">https://github.com/coq-community/bignums</a>
bits	436	4895	<a href="https://github.com/coq-community/bits">https://github.com/coq-community/bits</a>
buchberger	753	9705	<a href="https://github.com/coq-community/buchberger">https://github.com/coq-community/buchberger</a>
Categories	519	7757	<a href="https://github.com/amintimany/Categories">https://github.com/amintimany/Categories</a>
category-theory	1377	15827	<a href="https://github.com/jwiegley/category-theory">https://github.com/jwiegley/category-theory</a>

## 23:40 Building a Large Proof Repair Dataset

celsius	211	4470	<a href="https://github.com/clementbladeau/celsius">https://github.com/clementbladeau/celsius</a>
cerise	1026	35029	<a href="https://github.com/logsem/cerise">https://github.com/logsem/cerise</a>
certicoq	4350	133553	<a href="https://github.com/CertiCoq/certicoq">https://github.com/CertiCoq/certicoq</a>
CertiGraph	3522	106982	<a href="https://github.com/CertiGraph/CertiGraph">https://github.com/CertiGraph/CertiGraph</a>
cgraphs	812	16559	<a href="https://github.com/julesjacobs/cgraphs">https://github.com/julesjacobs/cgraphs</a>
ChargeCore	299	3750	<a href="https://github.com/jesper-bengtson/ChargeCore">https://github.com/jesper-bengtson/ChargeCore</a>
checker	4	55	<a href="https://github.com/coq-contribs/checker">https://github.com/coq-contribs/checker</a>
chinese	137	1769	<a href="https://github.com/coq-contribs/chinese">https://github.com/coq-contribs/chinese</a>
circuits	220	2723	<a href="https://github.com/coq-contribs/circuits">https://github.com/coq-contribs/circuits</a>
ClassicalReal	1303	48201	<a href="https://github.com/QinxiangCao/ClassicalReal">https://github.com/QinxiangCao/ClassicalReal</a>
cls-coq	611	20968	<a href="https://github.com/combinators/cls-coq">https://github.com/combinators/cls-coq</a>
CompCert	7835	156983	<a href="https://github.com/AbsInt/CompCert">https://github.com/AbsInt/CompCert</a>
CompCertR	7817	161208	<a href="https://github.com/snu-sf/CompCertR">https://github.com/snu-sf/CompCertR</a>
concat	514	6417	<a href="https://github.com/coq-contribs/concat">https://github.com/coq-contribs/concat</a>
ConCert	1581	33188	<a href="https://github.com/AU-COBRA/ConCert">https://github.com/AU-COBRA/ConCert</a>
constructive-geometry	116	902	<a href="https://github.com/coq-contribs/constructive-geometry">https://github.com/coq-contribs/constructive-geometry</a>
constructive-ltl	634	7621	<a href="https://github.com/jwieglej/constructive-ltl">https://github.com/jwieglej/constructive-ltl</a>
containers	1417	12466	<a href="https://github.com/coq-contribs/containers">https://github.com/coq-contribs/containers</a>
coq-ceres	122	1263	<a href="https://github.com/Lysxia/coq-ceres">https://github.com/Lysxia/coq-ceres</a>
coq-cunit	0	7	<a href="https://github.com/clarus/coq-cunit">https://github.com/clarus/coq-cunit</a>
coqeal	1609	17503	<a href="https://github.com/coq-community/coqeal">https://github.com/coq-community/coqeal</a>
coq-error-handlers	0	10	<a href="https://github.com/clarus/coq-error-handlers">https://github.com/clarus/coq-error-handlers</a>
coq-ext-lib	297	4842	<a href="https://github.com/coq-community/coq-ext-lib">https://github.com/coq-community/coq-ext-lib</a>
Coq-Flow-Equivalence	376	8278	<a href="https://github.com/GaloisInc/Coq-Flow-Equivalence">https://github.com/GaloisInc/Coq-Flow-Equivalence</a>
coq-function-ninjas	0	3	<a href="https://github.com/clarus/coq-function-ninjas">https://github.com/clarus/coq-function-ninjas</a>
coq-haskell	481	5512	<a href="https://github.com/jwieglej/coq-haskell1">https://github.com/jwieglej/coq-haskell1</a>
coq-http	27	795	<a href="https://github.com/liyishuai/coq-http">https://github.com/liyishuai/coq-http</a>

coq-iterable	0	15	<a href="https://github.com/clarus/coq-iterable">https://github.com/clarus/coq-iterable</a>
coq-list-plus	0	42	<a href="https://github.com/clarus/coq-list-plus">https://github.com/clarus/coq-list-plus</a>
coq-list-string	9	207	<a href="https://github.com/clarus/coq-list-string">https://github.com/clarus/coq-list-string</a>
coqoban	3	455	<a href="https://github.com/coq-community/coqoban">https://github.com/coq-community/coqoban</a>
coq-performance-tests	170	2536	<a href="https://github.com/coq-community/coq-performance-tests">https://github.com/coq-community/coq-performance-tests</a>
coq-procrastination	60	506	<a href="https://github.com/Armael/coq-procrastination">https://github.com/Armael/coq-procrastination</a>
coqrel	270	1763	<a href="https://github.com/CertiKOS/coqrel">https://github.com/CertiKOS/coqrel</a>
coq-robot	1505	13575	<a href="https://github.com/affeldt-aist/coq-robot">https://github.com/affeldt-aist/coq-robot</a>
coq-simple-io	0	267	<a href="https://github.com/Lysxia/coq-simple-io">https://github.com/Lysxia/coq-simple-io</a>
coqtail-math	2515	35411	<a href="https://github.com/coq-community/coqtail-math">https://github.com/coq-community/coqtail-math</a>
Core-Erlang-Formalization	550	12170	<a href="https://github.com/harp-project/Core-Erlang-Formalization">https://github.com/harp-project/Core-Erlang-Formalization</a>
cours-de-coq	92	792	<a href="https://github.com/coq-contribs/cours-de-coq">https://github.com/coq-contribs/cours-de-coq</a>
cpdt-japanese	431	3152	<a href="https://github.com/cpdt-japanese/cpdt-japanese">https://github.com/cpdt-japanese/cpdt-japanese</a>
cryptis	448	5751	<a href="https://github.com/arthuraa/cryptis">https://github.com/arthuraa/cryptis</a>
ctlctl	28	383	<a href="https://github.com/coq-contribs/ctlctl">https://github.com/coq-contribs/ctlctl</a>
dblib	244	1994	<a href="https://github.com/coq-community/dblib">https://github.com/coq-community/dblib</a>
demos	69	477	<a href="https://github.com/coq-contribs/demos">https://github.com/coq-contribs/demos</a>
dep-map	90	1843	<a href="https://github.com/coq-contribs/dep-map">https://github.com/coq-contribs/dep-map</a>
deriving	67	1567	<a href="https://github.com/arthuraa/deriving">https://github.com/arthuraa/deriving</a>
dictionaries	67	845	<a href="https://github.com/coq-contribs/dictionaries">https://github.com/coq-contribs/dictionaries</a>
diesel	844	12582	<a href="https://github.com/DistributedComponents/diesel">https://github.com/DistributedComponents/diesel</a>
distributed-reference-counting	954	24025	<a href="https://github.com/coq-contribs/distributed-reference-counting">https://github.com/coq-contribs/distributed-reference-counting</a>
domain-theory	51	812	<a href="https://github.com/coq-contribs/domain-theory">https://github.com/coq-contribs/domain-theory</a>
dot-iris	1590	12473	<a href="https://github.com/Blaisorblade/dot-iris">https://github.com/Blaisorblade/dot-iris</a>
euler-formula	162	6757	<a href="https://github.com/coq-contribs/euler-formula">https://github.com/coq-contribs/euler-formula</a>

## 23:42 Building a Large Proof Repair Dataset

event-struct	1039	11142	<a href="https://github.com/Event-Structures/event-struct">https://github.com/Event-Structures/event-struct</a>
exceptions	3	94	<a href="https://github.com/coq-contribs/exceptions">https://github.com/coq-contribs/exceptions</a>
fcs1-pcm	2118	15806	<a href="https://github.com/imdea-software/fcs1-pcm">https://github.com/imdea-software/fcs1-pcm</a>
fermat4	130	841	<a href="https://github.com/coq-contribs/fermat4">https://github.com/coq-contribs/fermat4</a>
finmap	870	5058	<a href="https://github.com/math-comp/finmap">https://github.com/math-comp/finmap</a>
float	780	10192	<a href="https://github.com/coq-contribs/float">https://github.com/coq-contribs/float</a>
FormalML	7026	135311	<a href="https://github.com/IBM/FormalML">https://github.com/IBM/FormalML</a>
frap	1870	23058	<a href="https://github.com/achlipala/frap">https://github.com/achlipala/frap</a>
free-groups	33	484	<a href="https://github.com/coq-contribs/free-groups">https://github.com/coq-contribs/free-groups</a>
fssec-model	153	3255	<a href="https://github.com/coq-contribs/fssec-model">https://github.com/coq-contribs/fssec-model</a>
functions-in-zfc	632	2437	<a href="https://github.com/coq-contribs/functions-in-zfc">https://github.com/coq-contribs/functions-in-zfc</a>
fundamental-arithmetics	152	1941	<a href="https://github.com/coq-contribs/fundamental-arithmetics">https://github.com/coq-contribs/fundamental-arithmetics</a>
generic-environments	269	2818	<a href="https://github.com/coq-community/generic-environments">https://github.com/coq-community/generic-environments</a>
GeoCoq	4087	122527	<a href="https://github.com/GeoCoq/GeoCoq">https://github.com/GeoCoq/GeoCoq</a>
goedel	103	8548	<a href="https://github.com/coq-community/goedel">https://github.com/coq-community/goedel</a>
GraphCoQL	220	2513	<a href="https://github.com/imfd/GraphCoQL">https://github.com/imfd/GraphCoQL</a>
graphs	174	4460	<a href="https://github.com/coq-contribs/graphs">https://github.com/coq-contribs/graphs</a>
graph-theory	2272	29435	<a href="https://github.com/coq-community/graph-theory">https://github.com/coq-community/graph-theory</a>
groups	14	134	<a href="https://github.com/coq-contribs/groups">https://github.com/coq-contribs/groups</a>
group-theory	105	1099	<a href="https://github.com/coq-contribs/group-theory">https://github.com/coq-contribs/group-theory</a>
hardware	186	1643	<a href="https://github.com/coq-contribs/hardware">https://github.com/coq-contribs/hardware</a>
hedges	110	2310	<a href="https://github.com/coq-contribs/hedges">https://github.com/coq-contribs/hedges</a>
higman-cf	30	279	<a href="https://github.com/coq-contribs/higman-cf">https://github.com/coq-contribs/higman-cf</a>
higman-s	49	913	<a href="https://github.com/coq-contribs/higman-s">https://github.com/coq-contribs/higman-s</a>
hoare-tut	25	346	<a href="https://github.com/coq-community/hoare-tut">https://github.com/coq-community/hoare-tut</a>
HoTT	4974	62461	<a href="https://github.com/HoTT/HoTT">https://github.com/HoTT/HoTT</a>
htt	410	4143	<a href="https://github.com/imdea-software/htt">https://github.com/imdea-software/htt</a>

huffman	285	3910	<a href="https://github.com/coq-community/huffman">https://github.com/coq-community/huffman</a>
idxassoc	58	713	<a href="https://github.com/coq-contribs/idxassoc">https://github.com/coq-contribs/idxassoc</a>
ieee754	15	239	<a href="https://github.com/coq-contribs/ieee754">https://github.com/coq-contribs/ieee754</a>
infotheo	3335	45155	<a href="https://github.com/affeldt-aist/infotheo">https://github.com/affeldt-aist/infotheo</a>
InfSeqExt	145	1943	<a href="https://github.com/DistributedComponents/InfSeqExt">https://github.com/DistributedComponents/InfSeqExt</a>
InteractionTrees	1822	45795	<a href="https://github.com/DeepSpec/InteractionTrees">https://github.com/DeepSpec/InteractionTrees</a>
int-map	494	7816	<a href="https://github.com/coq-contribs/int-map">https://github.com/coq-contribs/int-map</a>
ipc	479	9102	<a href="https://github.com/coq-contribs/ipc">https://github.com/coq-contribs/ipc</a>
izf	122	827	<a href="https://github.com/coq-contribs/izf">https://github.com/coq-contribs/izf</a>
jordan-curve-theorem	663	42492	<a href="https://github.com/coq-contribs/jordan-curve-theorem">https://github.com/coq-contribs/jordan-curve-theorem</a>
katamaran	854	14929	<a href="https://github.com/katamaran-project/katamaran">https://github.com/katamaran-project/katamaran</a>
lambda	90	941	<a href="https://github.com/coq-contribs/lambda">https://github.com/coq-contribs/lambda</a>
lazy-pcf	83	1587	<a href="https://github.com/coq-contribs/lazy-pcf">https://github.com/coq-contribs/lazy-pcf</a>
lemma-overloading	840	6476	<a href="https://github.com/coq-community/lemma-overloading">https://github.com/coq-community/lemma-overloading</a>
lesniewski-mereology	140	1396	<a href="https://github.com/coq-contribs/lesniewski-mereology">https://github.com/coq-contribs/lesniewski-mereology</a>
MapleS	7629	153467	<a href="https://github.com/namefanwjcom/MapleS">https://github.com/namefanwjcom/MapleS</a>
markov	31	951	<a href="https://github.com/coq-contribs/markov">https://github.com/coq-contribs/markov</a>
math-classes	1741	17370	<a href="https://github.com/coq-community/math-classes">https://github.com/coq-community/math-classes</a>
math-comp	14263	104712	<a href="https://github.com/math-comp/math-comp">https://github.com/math-comp/math-comp</a>
maths	38	386	<a href="https://github.com/coq-contribs/math">https://github.com/coq-contribs/math</a>
mini-compiler	4	86	<a href="https://github.com/coq-contribs/mini-compiler">https://github.com/coq-contribs/mini-compiler</a>
miniml	12	681	<a href="https://github.com/coq-contribs/miniml">https://github.com/coq-contribs/miniml</a>
mod-red	271	6204	<a href="https://github.com/coq-contribs/mod-red">https://github.com/coq-contribs/mod-red</a>
monae	1661	16336	<a href="https://github.com/affeldt-aist/monae">https://github.com/affeldt-aist/monae</a>
multinomials	1008	7573	<a href="https://github.com/math-comp/multinomials">https://github.com/math-comp/multinomials</a>
otway-rees	21	613	<a href="https://github.com/coq-contribs/otway-rees">https://github.com/coq-contribs/otway-rees</a>
paco	1519	23749	<a href="https://github.com/snu-sf/paco">https://github.com/snu-sf/paco</a>

## 23:44 Building a Large Proof Repair Dataset

param-pi	72	3853	<a href="https://github.com/coq-contribs/param-pi">https://github.com/coq-contribs/param-pi</a>
pigeons	6	34	<a href="https://github.com/llee454/pigeons">https://github.com/llee454/pigeons</a>
pipcore	1637	130325	<a href="https://github.com/2xs/pipcore">https://github.com/2xs/pipcore</a>
pnf	356	2857	<a href="https://github.com/ilyasergey/pnf">https://github.com/ilyasergey/pnf</a>
pocklington	265	6546	<a href="https://github.com/coq-community/pocklington">https://github.com/coq-community/pocklington</a>
proofs	174	2558	<a href="https://github.com/stepchowfun/proofs">https://github.com/stepchowfun/proofs</a>
propcalc	66	629	<a href="https://github.com/coq-contribs/propcalc">https://github.com/coq-contribs/propcalc</a>
pts	350	4513	<a href="https://github.com/coq-contribs/pts">https://github.com/coq-contribs/pts</a>
qarith	60	1151	<a href="https://github.com/coq-contribs/qarith">https://github.com/coq-contribs/qarith</a>
qarith-stern-brocot	1131	16216	<a href="https://github.com/coq-community/qarith-stern-brocot">https://github.com/coq-community/qarith-stern-brocot</a>
QuantumLib	1482	20473	<a href="https://github.com/inQWIRE/QuantumLib">https://github.com/inQWIRE/QuantumLib</a>
QuickChick	1211	13104	<a href="https://github.com/QuickChick/QuickChick">https://github.com/QuickChick/QuickChick</a>
quicksort-complexity	591	8858	<a href="https://github.com/coq-contribs/quicksort-complexity">https://github.com/coq-contribs/quicksort-complexity</a>
QWIRE	1167	18932	<a href="https://github.com/inQWIRE/QWIRE">https://github.com/inQWIRE/QWIRE</a>
railroad-crossing	100	1108	<a href="https://github.com/coq-contribs/railroad-crossing">https://github.com/coq-contribs/railroad-crossing</a>
ramsey	8	55	<a href="https://github.com/coq-contribs/ramsey">https://github.com/coq-contribs/ramsey</a>
regex	105	1556	<a href="https://github.com/coq-contribs/regex">https://github.com/coq-contribs/regex</a>
regex-reexamined-coq	476	7609	<a href="https://github.com/awalterschulze/regex-reexamined-coq">https://github.com/awalterschulze/regex-reexamined-coq</a>
rem	8	113	<a href="https://github.com/coq-contribs/rem">https://github.com/coq-contribs/rem</a>
rsa	113	1194	<a href="https://github.com/coq-contribs/rsa">https://github.com/coq-contribs/rsa</a>
ruler-compass-geometry	332	4376	<a href="https://github.com/coq-contribs/ruler-compass-geometry">https://github.com/coq-contribs/ruler-compass-geometry</a>
SCEV-coq	23	475	<a href="https://github.com/bollu/SCEV-coq">https://github.com/bollu/SCEV-coq</a>
schroeder	14	276	<a href="https://github.com/coq-contribs/schroeder">https://github.com/coq-contribs/schroeder</a>
search-trees	50	593	<a href="https://github.com/coq-contribs/search-trees">https://github.com/coq-contribs/search-trees</a>
SeLoC	381	7840	<a href="https://github.com/co-dan/SeLoC">https://github.com/co-dan/SeLoC</a>
Set-Theory	2703	66081	<a href="https://github.com/choukh/Set-Theory">https://github.com/choukh/Set-Theory</a>
shuffle	40	430	<a href="https://github.com/coq-contribs/shuffle">https://github.com/coq-contribs/shuffle</a>
sirtt	179	4645	<a href="https://github.com/TheoWinterhalter/sirtt">https://github.com/TheoWinterhalter/sirtt</a>



smc	685	25108	<a href="https://github.com/coq-contribs/smc">https://github.com/coq-contribs/smc</a>
SQIR	1713	37680	<a href="https://github.com/inQWIRE/SQIR">https://github.com/inQWIRE/SQIR</a>
SquiggleEq	1302	16922	<a href="https://github.com/aa755/SquiggleEq">https://github.com/aa755/SquiggleEq</a>
ssprove	1298	21179	<a href="https://github.com/SSProve/ssprove">https://github.com/SSProve/ssprove</a>
subst	417	3650	<a href="https://github.com/coq-contribs/subst">https://github.com/coq-contribs/subst</a>
sudoku	250	4015	<a href="https://github.com/coq-community/sudoku">https://github.com/coq-community/sudoku</a>
tarjan	459	5187	<a href="https://github.com/math-comp/tarjan">https://github.com/math-comp/tarjan</a>
tarski-geometry	136	2702	<a href="https://github.com/coq-contribs/tarski-geometry">https://github.com/coq-contribs/tarski-geometry</a>
three-gap	81	1091	<a href="https://github.com/coq-contribs/three-gap">https://github.com/coq-contribs/three-gap</a>
tlc	2314	24176	<a href="https://github.com/charguer/tlc">https://github.com/charguer/tlc</a>
topology	616	16461	<a href="https://github.com/coq-community/topology">https://github.com/coq-community/topology</a>
tortoise-hare-algorithm	4	82	<a href="https://github.com/coq-contribs/tortoise-hare-algorithm">https://github.com/coq-contribs/tortoise-hare-algorithm</a>
toychain	264	4320	<a href="https://github.com/certichain/toychain">https://github.com/certichain/toychain</a>
transfer	201	1726	<a href="https://github.com/Zimmi48/transfer">https://github.com/Zimmi48/transfer</a>
traversable-fincontainer	71	871	<a href="https://github.com/coq-contribs/traversable-fincontainer">https://github.com/coq-contribs/traversable-fincontainer</a>
tree-automata	834	24988	<a href="https://github.com/coq-contribs/tree-automata">https://github.com/coq-contribs/tree-automata</a>
UnifySL	1193	19223	<a href="https://github.com/QinxiangCao/UnifySL">https://github.com/QinxiangCao/UnifySL</a>
UnifySL	1193	19223	<a href="https://github.com/QinxiangCao/UnifySL">https://github.com/QinxiangCao/UnifySL</a>
UniMath	15201	256571	<a href="https://github.com/UniMath/UniMath">https://github.com/UniMath/UniMath</a>
univalent_parametricity	399	5121	<a href="https://github.com/CoqHott/univalent_parametricity">https://github.com/CoqHott/univalent_parametricity</a>
velus	3338	58852	<a href="https://github.com/INRIA/velus">https://github.com/INRIA/velus</a>
verdi-chord	990	16420	<a href="https://github.com/DistributedComponents/verdi-chord">https://github.com/DistributedComponents/verdi-chord</a>
weak-up-to	150	1720	<a href="https://github.com/coq-contribs/weak-up-to">https://github.com/coq-contribs/weak-up-to</a>
zchinese	43	781	<a href="https://github.com/coq-contribs/zchinese">https://github.com/coq-contribs/zchinese</a>
zf	213	2977	<a href="https://github.com/coq-contribs/zf">https://github.com/coq-contribs/zf</a>
zfc	241	2450	<a href="https://github.com/coq-contribs/zfc">https://github.com/coq-contribs/zfc</a>
zorns-lemma	183	4482	<a href="https://github.com/coq-community/zorns-lemma">https://github.com/coq-community/zorns-lemma</a>

## 23:46 Building a Large Proof Repair Dataset

zsearch-trees	48	639	<a href="https://github.com/coq-contribs/zsearch-trees">https://github.com/coq-contribs/zsearch-trees</a>
<b>Total</b>	<b>344,316</b>	<b>6,199,931</b>	

■ **Table 3** The list of projects considered but rejected for inclusion into either the repair or pretraining datasets. Projects were rejected if they did not build under Coq 8.10.2 using either the main branch or an obviously named tag that should support Coq 8.10.2. Note that the proof and sentence counts are estimated based upon the heuristic parsing described in **TODO: TODO**.

Project	Proof Count	Sentence Count	Repository URL
AML-Formalization	1532	43978	<a href="https://github.com/harp-project/AML-Formalization">https://github.com/harp-project/AML-Formalization</a>
bedrock	4423	66779	<a href="https://github.com/mit-plv/bedrock">https://github.com/mit-plv/bedrock</a>
bedrock2	1233	32668	<a href="https://github.com/mit-plv/bedrock2">https://github.com/mit-plv/bedrock2</a>
bedrock-mirror-shard	1726	18771	<a href="https://github.com/gmalecha/bedrock-mirror-shard">https://github.com/gmalecha/bedrock-mirror-shard</a>
cage	456	7425	<a href="https://github.com/gstew5/cage">https://github.com/gstew5/cage</a>
cecoa	1281	33151	<a href="https://github.com/davidnowak/cecoa">https://github.com/davidnowak/cecoa</a>
color	6743	113300	<a href="https://github.com/fblanqui/color">https://github.com/fblanqui/color</a>
CompCertM	1499	72820	<a href="https://github.com/snu-sf/CompCertM">https://github.com/snu-sf/CompCertM</a>
Coq-Combi	3520	37470	<a href="https://github.com/math-comp/Coq-Combi">https://github.com/math-comp/Coq-Combi</a>
coq-compile	29	2330	<a href="https://github.com/coq-ext-lib/coq-compile">https://github.com/coq-ext-lib/coq-compile</a>
coq-forcing	72	1190	<a href="https://github.com/CoqHott/coq-forcing">https://github.com/CoqHott/coq-forcing</a>
coq-guarded-computational-type-theory	205	3824	<a href="https://github.com/jonsterling/coq-guarded-computational-type-theory">https://github.com/jonsterling/coq-guarded-computational-type-theory</a>
coq-http2	47	943	<a href="https://github.com/liyishuai/coq-http2">https://github.com/liyishuai/coq-http2</a>
coqutil	797	9516	<a href="https://github.com/mit-plv/coqutil">https://github.com/mit-plv/coqutil</a>
crellvm	722	22052	<a href="https://github.com/snu-sf/crellvm">https://github.com/snu-sf/crellvm</a>
DeepSpecDB	2719	86898	<a href="https://github.com/PrincetonUniversity/DeepSpecDB">https://github.com/PrincetonUniversity/DeepSpecDB</a>
engine-bench	94	1264	<a href="https://github.com/mit-plv/engine-bench">https://github.com/mit-plv/engine-bench</a>
fiat	5932	84544	<a href="https://github.com/mit-plv/fiat">https://github.com/mit-plv/fiat</a>
formal-type-theory	67	10680	<a href="https://github.com/TheoWinterhalter/formal-type-theory">https://github.com/TheoWinterhalter/formal-type-theory</a>
functional-algebra	278	1286	<a href="https://github.com/llee454/functional-algebra">https://github.com/llee454/functional-algebra</a>

general-type-theories	546	10172	<a href="https://github.com/peterlefanulumsdaine/general-type-theories">https://github.com/peterlefanulumsdaine/general-type-theories</a>
HoTT-categories	380	4559	<a href="https://github.com/CategoricalData/HoTT-categories">https://github.com/CategoricalData/HoTT-categories</a>
hybrid	524	7099	<a href="https://github.com/Eelis/hybrid">https://github.com/Eelis/hybrid</a>
IFC	756	11435	<a href="https://github.com/QuickChick/IFC">https://github.com/QuickChick/IFC</a>
infinitary-rewriting-coq	262	4207	<a href="https://github.com/martijnvermaat/infinitary-rewriting-coq">https://github.com/martijnvermaat/infinitary-rewriting-coq</a>
iris-coq	4299	32669	<a href="https://github.com/izgzhen/iris-coq">https://github.com/izgzhen/iris-coq</a>
iron	1519	30137	<a href="https://github.com/discus-lang/iron">https://github.com/discus-lang/iron</a>
kami	2092	41182	<a href="https://github.com/mit-plv/kami">https://github.com/mit-plv/kami</a>
lin-alg	543	10797	<a href="https://github.com/coq-contribs/lin-alg">https://github.com/coq-contribs/lin-alg</a>
linearscan	304	4786	<a href="https://github.com/jwiegle/linearscan">https://github.com/jwiegle/linearscan</a>
llvmtwin-coq	675	34538	<a href="https://github.com/snu-sf/llvmtwin-coq">https://github.com/snu-sf/llvmtwin-coq</a>
maple-mode	20	154	<a href="https://github.com/coq-contribs/maple-mode">https://github.com/coq-contribs/maple-mode</a>
metalib	796	8026	<a href="https://github.com/plclub/metalib">https://github.com/plclub/metalib</a>
mirror-core	1380	36458	<a href="https://github.com/gmalecha/mirror-core">https://github.com/gmalecha/mirror-core</a>
mirror-shard	636	12465	<a href="https://github.com/gmalecha/mirror-shard">https://github.com/gmalecha/mirror-shard</a>
orb-stab	40	1590	<a href="https://github.com/coq-contribs/orb-stab">https://github.com/coq-contribs/orb-stab</a>
paramcoq-iff	285	6527	<a href="https://github.com/aa755/paramcoq-iff">https://github.com/aa755/paramcoq-iff</a>
parsing-parses	274	3512	<a href="https://github.com/JasonGross/parsing-parses">https://github.com/JasonGross/parsing-parses</a>
perennial	5284	137211	<a href="https://github.com/mit-pdos/perennial">https://github.com/mit-pdos/perennial</a>
probchain	246	4801	<a href="https://github.com/certichain/probchain">https://github.com/certichain/probchain</a>
promising-coq	835	29947	<a href="https://github.com/snu-sf/promising-coq">https://github.com/snu-sf/promising-coq</a>
quadcopter	1061	34482	<a href="https://github.com/dricketts/quadcopter">https://github.com/dricketts/quadcopter</a>
rewriter	1436	13713	<a href="https://github.com/mit-plv/rewriter">https://github.com/mit-plv/rewriter</a>
rezk_completion	350	4261	<a href="https://github.com/benediktahrens/rezk_completion">https://github.com/benediktahrens/rezk_completion</a>
ssrbit	344	2966	<a href="https://github.com/ejgallego/ssrbit">https://github.com/ejgallego/ssrbit</a>
system	1	112	<a href="https://github.com/coq-concurrency/system">https://github.com/coq-concurrency/system</a>
yalla	1019	30165	<a href="https://github.com/olaure01/yalla">https://github.com/olaure01/yalla</a>
<b>Total</b>	<b>344,316</b>	<b>6,199,931</b>	

## 23:48 Building a Large Proof Repair Dataset

■ **Table 4** The list of projects considered but skipped due to time restrictions or build complications for inclusion into either the repair or pretraining datasets. Note that the proof and sentence counts are estimated based upon the heuristic parsing described in **TODO: TODO**.

Project	Proof Count	Sentence Count	Repository URL
argosy	304	3614	<a href="https://github.com/mit-pdos/argosy">https://github.com/mit-pdos/argosy</a>
banach_tarski	479	8608	<a href="https://github.com/roglo/banach_tarski">https://github.com/roglo/banach_tarski</a>
ceramist	427	8223	<a href="https://github.com/verse-lab/ceramist">https://github.com/verse-lab/ceramist</a>
cheerios	83	1078	<a href="https://github.com/uwplse/cheerios">https://github.com/uwplse/cheerios</a>
coq_real	843	24386	<a href="https://github.com/roglo/coq_real">https://github.com/roglo/coq_real</a>
coq-bitset	147	3755	<a href="https://github.com/artart78/coq-bitset">https://github.com/artart78/coq-bitset</a>
coq-forcing	72	1190	<a href="https://github.com/CoqHott/coq-forcing">https://github.com/CoqHott/coq-forcing</a>
coq-library-complexity	2438	43607	<a href="https://github.com/uds-psl/coq-library-complexity">https://github.com/uds-psl/coq-library-complexity</a>
coq-library-undecidability	10281	153181	<a href="https://github.com/uds-psl/coq-library-undecidability">https://github.com/uds-psl/coq-library-undecidability</a>
Coq-Polyhedra	1086	11773	<a href="https://github.com/Coq-Polyhedra/Coq-Polyhedra">https://github.com/Coq-Polyhedra/Coq-Polyhedra</a>
coq-record-update	6	243	<a href="https://github.com/tchajed/coq-record-update">https://github.com/tchajed/coq-record-update</a>
coq-utils	298	2820	<a href="https://github.com/arthuraa/coq-utils">https://github.com/arthuraa/coq-utils</a>
corespec	1510	21682	<a href="https://github.com/sweirich/corespec">https://github.com/sweirich/corespec</a>
cps	436	9211	<a href="https://github.com/takanuva/cps">https://github.com/takanuva/cps</a>
cspec	884	11566	<a href="https://github.com/mit-pdos/cspec">https://github.com/mit-pdos/cspec</a>
ct	161	1948	<a href="https://github.com/relrod/ct">https://github.com/relrod/ct</a>
dez	1523	20350	<a href="https://github.com/Tuplanolla/dez">https://github.com/Tuplanolla/dez</a>
domains	1738	58184	<a href="https://github.com/robdockins/domains">https://github.com/robdockins/domains</a>
ett-to-wtt	433	13925	<a href="https://github.com/TheoWinterhalter/ett-to-wtt">https://github.com/TheoWinterhalter/ett-to-wtt</a>
exploring-robust-property-preservation	295	6518	<a href="https://github.com/secure-compilation/exploring-robust-property-preservation">https://github.com/secure-compilation/exploring-robust-property-preservation</a>
extructures	357	2986	<a href="https://github.com/arthuraa/extructures">https://github.com/arthuraa/extructures</a>
hahn	1458	9836	<a href="https://github.com/vafeiadis/hahn">https://github.com/vafeiadis/hahn</a>
hanoi	577	10514	<a href="https://github.com/thery/hanoi">https://github.com/thery/hanoi</a>
helix	2046	56461	<a href="https://github.com/vzaliva/helix">https://github.com/vzaliva/helix</a>
imm	1444	24624	<a href="https://github.com/weakmemory/imm">https://github.com/weakmemory/imm</a>

Kami	2120	37732	<a href="https://github.com/sifive/Kami">https://github.com/sifive/Kami</a>
Ktheory	40	505	<a href="https://github.com/DanGrayson/Ktheory">https://github.com/DanGrayson/Ktheory</a>
largecatmodules	467	9644	<a href="https://github.com/UniMath/largecatmodules">https://github.com/UniMath/largecatmodules</a>
LOGIC	1729	25450	<a href="https://github.com/QinxiangCao/LOGIC">https://github.com/QinxiangCao/LOGIC</a>
loopring-protocol2-verification	155	3718	<a href="https://github.com/sec-bit/loopring-protocol2-verification">https://github.com/sec-bit/loopring-protocol2-verification</a>
lvc	3782	60531	<a href="https://github.com/sigurdschneider/lvc">https://github.com/sigurdschneider/lvc</a>
memory-safe-language	72	1325	<a href="https://github.com/arthuraa/memory-safe-language">https://github.com/arthuraa/memory-safe-language</a>
metacoq	8885	161078	<a href="https://github.com/MetaCoq/metacoq">https://github.com/MetaCoq/metacoq</a>
micro-policies-coq	569	13727	<a href="https://github.com/micro-policies/micro-policies-coq">https://github.com/micro-policies/micro-policies-coq</a>
oeuf	10471	186512	<a href="https://github.com/uwplse/oeuf">https://github.com/uwplse/oeuf</a>
PolTac	309	2247	<a href="https://github.com/thery/PolTac">https://github.com/thery/PolTac</a>
ProcKami	54	2055	<a href="https://github.com/sifive/ProcKami">https://github.com/sifive/ProcKami</a>
puiseuxth	910	20553	<a href="https://github.com/roglo/puiseuxth">https://github.com/roglo/puiseuxth</a>
riscv-coq	209	3367	<a href="https://github.com/samuelgruetter/riscv-coq">https://github.com/samuelgruetter/riscv-coq</a>
silveroak	1519	26477	<a href="https://github.com/project-oak/silveroak">https://github.com/project-oak/silveroak</a>
StdLibKami	72	4505	<a href="https://github.com/sifive/StdLibKami">https://github.com/sifive/StdLibKami</a>
StructTact	287	3563	<a href="https://github.com/uwplse/StructTact">https://github.com/uwplse/StructTact</a>
twoSquare	202	1769	<a href="https://github.com/thery/twoSquare">https://github.com/thery/twoSquare</a>
TypeTheory	1923	30575	<a href="https://github.com/UniMath/TypeTheory">https://github.com/UniMath/TypeTheory</a>
Valuations	184	5079	<a href="https://github.com/FFaissolle/Valuations">https://github.com/FFaissolle/Valuations</a>
vellvm-legacy	2484	40111	<a href="https://github.com/vellvm/vellvm-legacy">https://github.com/vellvm/vellvm-legacy</a>
verdi	618	12534	<a href="https://github.com/uwplse/verdi">https://github.com/uwplse/verdi</a>
verdi-raft	2272	42557	<a href="https://github.com/uwplse/verdi-raft">https://github.com/uwplse/verdi-raft</a>
Verified-FEC	994	20957	<a href="https://github.com/verified-network-toolchain/Verified-FEC">https://github.com/verified-network-toolchain/Verified-FEC</a>
verified-ifc	884	15495	<a href="https://github.com/micro-policies/verified-ifc">https://github.com/micro-policies/verified-ifc</a>
VeriGHC	118	1522	<a href="https://github.com/trommler/VeriGHC">https://github.com/trommler/VeriGHC</a>
VST	27777	570848	<a href="https://github.com/PrincetonUniversity/VST">https://github.com/PrincetonUniversity/VST</a>
WasmCert-Coq	574	13803	<a href="https://github.com/WasmCert/WasmCert-Coq">https://github.com/WasmCert/WasmCert-Coq</a>

## 23:50 Building a Large Proof Repair Dataset

weakestmoToImm	1282	33541	<a href="https://github.com/weakmemory/weakestmoToImm">https://github.com/weakmemory/weakestmoToImm</a>
when-good-components-go-bad	1030	23295	<a href="https://github.com/secure-compilation/when-good-components-go-bad">https://github.com/secure-compilation/when-good-components-go-bad</a>
ynot	885	9996	<a href="https://github.com/ynot-harvard/ynot">https://github.com/ynot-harvard/ynot</a>
<b>Total</b>	<b>344,316</b>	<b>6,199,931</b>	

■ **Table 5** The list of projects considered for inclusion into either the repair or pretraining datasets that have opam files and thus easily-inferred dependencies. Note that the proof and sentence counts are estimated based upon the heuristic parsing described in **TODO: TODO**.

Project	Proof Count	Sentence Count	Repository URL
Abel	429	4063	<a href="https://github.com/math-comp/Abel">https://github.com/math-comp/Abel</a>
analysis	4663	39651	<a href="https://github.com/math-comp/analysis">https://github.com/math-comp/analysis</a>
atbr	819	11089	<a href="https://github.com/coq-community/atbr">https://github.com/coq-community/atbr</a>
bellantonicook	498	7515	<a href="https://github.com/davidnowak/bellantonicook">https://github.com/davidnowak/bellantonicook</a>
bigenough	5	43	<a href="https://github.com/math-comp/bigenough">https://github.com/math-comp/bigenough</a>
bits	436	4895	<a href="https://github.com/coq-community/bits">https://github.com/coq-community/bits</a>
celsius	211	4470	<a href="https://github.com/clementbladeau/celsius">https://github.com/clementbladeau/celsius</a>
cerise	1026	35029	<a href="https://github.com/logsem/cerise">https://github.com/logsem/cerise</a>
certicoq	4350	133553	<a href="https://github.com/CertiCoq/certicoq">https://github.com/CertiCoq/certicoq</a>
CertiGraph	3522	106982	<a href="https://github.com/CertiGraph/CertiGraph">https://github.com/CertiGraph/CertiGraph</a>
cgraphs	812	16559	<a href="https://github.com/julesjacobs/cgraphs">https://github.com/julesjacobs/cgraphs</a>
ConCert	1581	33188	<a href="https://github.com/AU-COBRA/ConCert">https://github.com/AU-COBRA/ConCert</a>
coq-ceres	122	1263	<a href="https://github.com/Lysxia/coq-ceres">https://github.com/Lysxia/coq-ceres</a>
coqeal	1609	17503	<a href="https://github.com/coq-community/coqeal">https://github.com/coq-community/coqeal</a>
coq-ext-lib	297	4842	<a href="https://github.com/coq-community/coq-ext-lib">https://github.com/coq-community/coq-ext-lib</a>
coq-haskell	481	5512	<a href="https://github.com/jwieglej/coq-haskell">https://github.com/jwieglej/coq-haskell</a>
coq-http	27	795	<a href="https://github.com/liyishuai/coq-http">https://github.com/liyishuai/coq-http</a>
coq-list-string	9	207	<a href="https://github.com/clarus/coq-list-string">https://github.com/clarus/coq-list-string</a>
coqoban	3	455	<a href="https://github.com/coq-community/coqoban">https://github.com/coq-community/coqoban</a>
coq-procrastination	60	506	<a href="https://github.com/Armael/coq-procrastination">https://github.com/Armael/coq-procrastination</a>
coq-robot	1505	13575	<a href="https://github.com/affeldt-aist/coq-robot">https://github.com/affeldt-aist/coq-robot</a>

coq-simple-io	0	267	<a href="https://github.com/Lysxia/coq-simple-io">https://github.com/Lysxia/coq-simple-io</a>
coqtail-math	2515	35411	<a href="https://github.com/coq-community/coqtail-math">https://github.com/coq-community/coqtail-math</a>
dblib	244	1994	<a href="https://github.com/coq-community/dblib">https://github.com/coq-community/dblib</a>
deriving	67	1567	<a href="https://github.com/arthuraa/deriving">https://github.com/arthuraa/deriving</a>
dot-iris	1590	12473	<a href="https://github.com/Blaisorblade/dot-iris">https://github.com/Blaisorblade/dot-iris</a>
event-struct	1039	11142	<a href="https://github.com/Event-Structures/event-struct">https://github.com/Event-Structures/event-struct</a>
fcs1-pcm	2118	15806	<a href="https://github.com/imdea-software/fcs1-pcm">https://github.com/imdea-software/fcs1-pcm</a>
finmap	870	5058	<a href="https://github.com/math-comp/finmap">https://github.com/math-comp/finmap</a>
FormalML	7026	135311	<a href="https://github.com/IBM/FormalML">https://github.com/IBM/FormalML</a>
generic-environments	269	2818	<a href="https://github.com/coq-community/generic-environments">https://github.com/coq-community/generic-environments</a>
goedel	103	8548	<a href="https://github.com/coq-community/goedel">https://github.com/coq-community/goedel</a>
graph-theory	2272	29435	<a href="https://github.com/coq-community/graph-theory">https://github.com/coq-community/graph-theory</a>
hoare-tut	25	346	<a href="https://github.com/coq-community/hoare-tut">https://github.com/coq-community/hoare-tut</a>
HoTT	4974	62461	<a href="https://github.com/HoTT/HoTT">https://github.com/HoTT/HoTT</a>
htt	410	4143	<a href="https://github.com/imdea-software/htt">https://github.com/imdea-software/htt</a>
huffman	285	3910	<a href="https://github.com/coq-community/huffman">https://github.com/coq-community/huffman</a>
infotheo	3335	45155	<a href="https://github.com/affeldt-aist/infotheo">https://github.com/affeldt-aist/infotheo</a>
InfSeqExt	145	1943	<a href="https://github.com/DistributedComponents/InfSeqExt">https://github.com/DistributedComponents/InfSeqExt</a>
InteractionTrees	1822	45795	<a href="https://github.com/DeepSpec/InteractionTrees">https://github.com/DeepSpec/InteractionTrees</a>
katamaran	854	14929	<a href="https://github.com/katamaran-project/katamaran">https://github.com/katamaran-project/katamaran</a>
lemma-overloading	840	6476	<a href="https://github.com/coq-community/lemma-overloading">https://github.com/coq-community/lemma-overloading</a>
math-classes	1741	17370	<a href="https://github.com/coq-community/math-classes">https://github.com/coq-community/math-classes</a>
math-comp	14263	104712	<a href="https://github.com/math-comp/math-comp">https://github.com/math-comp/math-comp</a>
monae	1661	16336	<a href="https://github.com/affeldt-aist/monae">https://github.com/affeldt-aist/monae</a>
multinomials	1008	7573	<a href="https://github.com/math-comp/multinomials">https://github.com/math-comp/multinomials</a>
pocklington	265	6546	<a href="https://github.com/coq-community/pocklington">https://github.com/coq-community/pocklington</a>

## 23:52 Building a Large Proof Repair Dataset

qarith-stern-brocot	1131	16216	<a href="https://github.com/coq-community/qarith-stern-brocot">https://github.com/coq-community/qarith-stern-brocot</a>
QuickChick	1211	13104	<a href="https://github.com/QuickChick/QuickChick">https://github.com/QuickChick/QuickChick</a>
quicksort-complexity	591	8858	<a href="https://github.com/coq-contribs/quicksort-complexity">https://github.com/coq-contribs/quicksort-complexity</a>
SeLoC	381	7840	<a href="https://github.com/co-dan/SeLoC">https://github.com/co-dan/SeLoC</a>
ssprove	1298	21179	<a href="https://github.com/SSProve/ssprove">https://github.com/SSProve/ssprove</a>
sudoku	250	4015	<a href="https://github.com/coq-community/sudoku">https://github.com/coq-community/sudoku</a>
tarjan	459	5187	<a href="https://github.com/math-comp/tarjan">https://github.com/math-comp/tarjan</a>
tlc	2314	24176	<a href="https://github.com/charguer/tlc">https://github.com/charguer/tlc</a>
topology	616	16461	<a href="https://github.com/coq-community/topology">https://github.com/coq-community/topology</a>
toychain	264	4320	<a href="https://github.com/certichain/toychain">https://github.com/certichain/toychain</a>
verdi-chord	990	16420	<a href="https://github.com/DistributedComponents/verdi-chord">https://github.com/DistributedComponents/verdi-chord</a>
<b>Total</b>	<b>344,316</b>	<b>6,199,931</b>	

■ **Table 6** The list of projects considered for inclusion into either the repair or pretraining datasets that do not have opam files and thus do not have easily-inferred dependencies. Note that the proof and sentence counts are estimated based upon the heuristic parsing described in **TODO: TODO**.

Project	Proof Count	Sentence Count	Repository URL
additions	103	1425	<a href="https://github.com/coq-contribs/additions">https://github.com/coq-contribs/additions</a>
ails	225	4605	<a href="https://github.com/coq-contribs/ails">https://github.com/coq-contribs/ails</a>
algebra	559	6609	<a href="https://github.com/coq-contribs/algebra">https://github.com/coq-contribs/algebra</a>
AML-Formalization	1532	43978	<a href="https://github.com/harp-project/AML-Formalization">https://github.com/harp-project/AML-Formalization</a>
amm11262	28	1154	<a href="https://github.com/coq-contribs/amm11262">https://github.com/coq-contribs/amm11262</a>
angles	93	2013	<a href="https://github.com/coq-contribs/angles">https://github.com/coq-contribs/angles</a>
area-method	779	10603	<a href="https://github.com/coq-contribs/area-method">https://github.com/coq-contribs/area-method</a>
argosy	304	3614	<a href="https://github.com/mit-pdos/argosy">https://github.com/mit-pdos/argosy</a>
asn1fpcoq	194	2982	<a href="https://github.com/digamma-ai/asn1fpcoq">https://github.com/digamma-ai/asn1fpcoq</a>
automata	344	5833	<a href="https://github.com/coq-contribs/automata">https://github.com/coq-contribs/automata</a>
axiomatic-abp	399	4494	<a href="https://github.com/coq-contribs/axiomatic-abp">https://github.com/coq-contribs/axiomatic-abp</a>
banach_tarski	479	8608	<a href="https://github.com/roglo/banach_tarski">https://github.com/roglo/banach_tarski</a>



bbv	653	7343	<a href="https://github.com/mit-plv/bbv">https://github.com/mit-plv/bbv</a>
bdds	232	15510	<a href="https://github.com/coq-contribs/bdds">https://github.com/coq-contribs/bdds</a>
bedrock	4423	66779	<a href="https://github.com/mit-plv/bedrock">https://github.com/mit-plv/bedrock</a>
bedrock2	1233	32668	<a href="https://github.com/mit-plv/bedrock2">https://github.com/mit-plv/bedrock2</a>
bedrock-mirror-shard	1726	18771	<a href="https://github.com/gmalecha/bedrock-mirror-shard">https://github.com/gmalecha/bedrock-mirror-shard</a>
bignums	638	9191	<a href="https://github.com/coq-community/bignums">https://github.com/coq-community/bignums</a>
buchberger	753	9705	<a href="https://github.com/coq-community/buchberger">https://github.com/coq-community/buchberger</a>
cage	456	7425	<a href="https://github.com/gstew5/cage">https://github.com/gstew5/cage</a>
Categories	519	7757	<a href="https://github.com/amintimany/Categories">https://github.com/amintimany/Categories</a>
category-theory	1377	15827	<a href="https://github.com/jwiegley/category-theory">https://github.com/jwiegley/category-theory</a>
cecoa	1281	33151	<a href="https://github.com/davidnowak/cecoa">https://github.com/davidnowak/cecoa</a>
ceramist	427	8223	<a href="https://github.com/verse-lab/ceramist">https://github.com/verse-lab/ceramist</a>
ChargeCore	299	3750	<a href="https://github.com/jesper-bengtson/ChargeCore">https://github.com/jesper-bengtson/ChargeCore</a>
checker	4	55	<a href="https://github.com/coq-contribs/checker">https://github.com/coq-contribs/checker</a>
cheerios	83	1078	<a href="https://github.com/uwplse/cheerios">https://github.com/uwplse/cheerios</a>
chinese	137	1769	<a href="https://github.com/coq-contribs/chinese">https://github.com/coq-contribs/chinese</a>
circuits	220	2723	<a href="https://github.com/coq-contribs/circuits">https://github.com/coq-contribs/circuits</a>
ClassicalReal	1303	48201	<a href="https://github.com/QinxiangCao/ClassicalReal">https://github.com/QinxiangCao/ClassicalReal</a>
cls-coq	611	20968	<a href="https://github.com/combinators/cls-coq">https://github.com/combinators/cls-coq</a>
color	6743	113300	<a href="https://github.com/fblanqui/color">https://github.com/fblanqui/color</a>
CompCert	7835	156983	<a href="https://github.com/AbsInt/CompCert">https://github.com/AbsInt/CompCert</a>
CompCertM	1499	72820	<a href="https://github.com/snu-sf/CompCertM">https://github.com/snu-sf/CompCertM</a>
CompCertR	7817	161208	<a href="https://github.com/snu-sf/CompCertR">https://github.com/snu-sf/CompCertR</a>
concat	514	6417	<a href="https://github.com/coq-contribs/concat">https://github.com/coq-contribs/concat</a>
constructive-geometry	116	902	<a href="https://github.com/coq-contribs/constructive-geometry">https://github.com/coq-contribs/constructive-geometry</a>
constructive-ltl	634	7621	<a href="https://github.com/jwiegley/constructive-ltl">https://github.com/jwiegley/constructive-ltl</a>
containers	1417	12466	<a href="https://github.com/coq-contribs/containers">https://github.com/coq-contribs/containers</a>
coq_real	843	24386	<a href="https://github.com/roglo/coq_real">https://github.com/roglo/coq_real</a>

## 23:54 Building a Large Proof Repair Dataset

coq-bitset	147	3755	<a href="https://github.com/artart78/coq-bitset">https://github.com/artart78/coq-bitset</a>
Coq-Combi	3520	37470	<a href="https://github.com/math-comp/Coq-Combi">https://github.com/math-comp/Coq-Combi</a>
coq-compile	29	2330	<a href="https://github.com/coq-ext-lib/coq-compile">https://github.com/coq-ext-lib/coq-compile</a>
coq-cunit	0	7	<a href="https://github.com/clarus/coq-cunit">https://github.com/clarus/coq-cunit</a>
coq-error-handlers	0	10	<a href="https://github.com/clarus/coq-error-handlers">https://github.com/clarus/coq-error-handlers</a>
Coq-Flow-Equivalence	376	8278	<a href="https://github.com/GaloisInc/Coq-Flow-Equivalence">https://github.com/GaloisInc/Coq-Flow-Equivalence</a>
coq-forcing	72	1190	<a href="https://github.com/CoqHott/coq-forcing">https://github.com/CoqHott/coq-forcing</a>
coq-function-ninjas	0	3	<a href="https://github.com/clarus/coq-function-ninjas">https://github.com/clarus/coq-function-ninjas</a>
coq-guarded-computational-type-theory	205	3824	<a href="https://github.com/jonsterling/coq-guarded-computational-type-theory">https://github.com/jonsterling/coq-guarded-computational-type-theory</a>
coq-http2	47	943	<a href="https://github.com/liyishuai/coq-http2">https://github.com/liyishuai/coq-http2</a>
coq-iterable	0	15	<a href="https://github.com/clarus/coq-iterable">https://github.com/clarus/coq-iterable</a>
coq-library-complexity	2438	43607	<a href="https://github.com/uds-psl/coq-library-complexity">https://github.com/uds-psl/coq-library-complexity</a>
coq-library-undecidability	10281	153181	<a href="https://github.com/uds-psl/coq-library-undecidability">https://github.com/uds-psl/coq-library-undecidability</a>
coq-list-plus	0	42	<a href="https://github.com/clarus/coq-list-plus">https://github.com/clarus/coq-list-plus</a>
coq-performance-tests	170	2536	<a href="https://github.com/coq-community/coq-performance-tests">https://github.com/coq-community/coq-performance-tests</a>
Coq-Polyhedra	1086	11773	<a href="https://github.com/Coq-Polyhedra/Coq-Polyhedra">https://github.com/Coq-Polyhedra/Coq-Polyhedra</a>
coq-record-update	6	243	<a href="https://github.com/tchajed/coq-record-update">https://github.com/tchajed/coq-record-update</a>
coqrel	270	1763	<a href="https://github.com/CertiKOS/coqrel">https://github.com/CertiKOS/coqrel</a>
coqutil	797	9516	<a href="https://github.com/mit-plv/coqutil">https://github.com/mit-plv/coqutil</a>
coq-utils	298	2820	<a href="https://github.com/arthuraa/coq-utils">https://github.com/arthuraa/coq-utils</a>
Core-Erlang-Formalization	550	12170	<a href="https://github.com/harp-project/Core-Erlang-Formalization">https://github.com/harp-project/Core-Erlang-Formalization</a>
corespec	1510	21682	<a href="https://github.com/sweirich/corespec">https://github.com/sweirich/corespec</a>
cours-de-coq	92	792	<a href="https://github.com/coq-contribs/cours-de-coq">https://github.com/coq-contribs/cours-de-coq</a>
cpdt-japanese	431	3152	<a href="https://github.com/cpdt-japanese/cpdt-japanese">https://github.com/cpdt-japanese/cpdt-japanese</a>
cps	436	9211	<a href="https://github.com/takanuva/cps">https://github.com/takanuva/cps</a>

crellvm	722	22052	<a href="https://github.com/snu-sf/crellvm">https://github.com/snu-sf/crellvm</a>
cryptis	448	5751	<a href="https://github.com/arthuraa/cryptis">https://github.com/arthuraa/cryptis</a>
cspec	884	11566	<a href="https://github.com/mit-pdos/cspec">https://github.com/mit-pdos/cspec</a>
ct	161	1948	<a href="https://github.com/relrod/ct">https://github.com/relrod/ct</a>
ctlctl	28	383	<a href="https://github.com/coq-contribs/ctlctl">https://github.com/coq-contribs/ctlctl</a>
DeepSpecDB	2719	86898	<a href="https://github.com/PrincetonUniversity/DeepSpecDB">https://github.com/PrincetonUniversity/DeepSpecDB</a>
demos	69	477	<a href="https://github.com/coq-contribs/demos">https://github.com/coq-contribs/demos</a>
dep-map	90	1843	<a href="https://github.com/coq-contribs/dep-map">https://github.com/coq-contribs/dep-map</a>
dez	1523	20350	<a href="https://github.com/Tuplanolla/dez">https://github.com/Tuplanolla/dez</a>
dictionaries	67	845	<a href="https://github.com/coq-contribs/dictionaries">https://github.com/coq-contribs/dictionaries</a>
diesel	844	12582	<a href="https://github.com/DistributedComponents/diesel">https://github.com/DistributedComponents/diesel</a>
distributed-reference-counting	954	24025	<a href="https://github.com/coq-contribs/distributed-reference-counting">https://github.com/coq-contribs/distributed-reference-counting</a>
domains	1738	58184	<a href="https://github.com/robdockins/domains">https://github.com/robdockins/domains</a>
domain-theory	51	812	<a href="https://github.com/coq-contribs/domain-theory">https://github.com/coq-contribs/domain-theory</a>
engine-bench	94	1264	<a href="https://github.com/mit-plv/engine-bench">https://github.com/mit-plv/engine-bench</a>
ett-to-wtt	433	13925	<a href="https://github.com/TheoWinterhalter/ett-to-wtt">https://github.com/TheoWinterhalter/ett-to-wtt</a>
euler-formula	162	6757	<a href="https://github.com/coq-contribs/euler-formula">https://github.com/coq-contribs/euler-formula</a>
exceptions	3	94	<a href="https://github.com/coq-contribs/exceptions">https://github.com/coq-contribs/exceptions</a>
exploring-robust-property-preservation	295	6518	<a href="https://github.com/secure-compilation/exploring-robust-property-preservation">https://github.com/secure-compilation/exploring-robust-property-preservation</a>
extructures	357	2986	<a href="https://github.com/arthuraa/extructures">https://github.com/arthuraa/extructures</a>
fermat4	130	841	<a href="https://github.com/coq-contribs/fermat4">https://github.com/coq-contribs/fermat4</a>
fiat	5932	84544	<a href="https://github.com/mit-plv/fiat">https://github.com/mit-plv/fiat</a>
float	780	10192	<a href="https://github.com/coq-contribs/float">https://github.com/coq-contribs/float</a>
formal-type-theory	67	10680	<a href="https://github.com/TheoWinterhalter/formal-type-theory">https://github.com/TheoWinterhalter/formal-type-theory</a>
frap	1870	23058	<a href="https://github.com/achlipala/frap">https://github.com/achlipala/frap</a>
free-groups	33	484	<a href="https://github.com/coq-contribs/free-groups">https://github.com/coq-contribs/free-groups</a>

## 23:56 Building a Large Proof Repair Dataset

fssec-model	153	3255	<a href="https://github.com/coq-contribs/fssec-model">https://github.com/coq-contribs/fssec-model</a>
functional-algebra	278	1286	<a href="https://github.com/llee454/functional-algebra">https://github.com/llee454/functional-algebra</a>
functions-in-zfc	632	2437	<a href="https://github.com/coq-contribs/functions-in-zfc">https://github.com/coq-contribs/functions-in-zfc</a>
fundamental-arithmetics	152	1941	<a href="https://github.com/coq-contribs/fundamental-arithmetics">https://github.com/coq-contribs/fundamental-arithmetics</a>
general-type-theories	546	10172	<a href="https://github.com/peterlefanulumsdaine/general-type-theories">https://github.com/peterlefanulumsdaine/general-type-theories</a>
GeoCoq	4087	122527	<a href="https://github.com/GeoCoq/GeoCoq">https://github.com/GeoCoq/GeoCoq</a>
GraphCoQL	220	2513	<a href="https://github.com/imfd/GraphCoQL">https://github.com/imfd/GraphCoQL</a>
graphs	174	4460	<a href="https://github.com/coq-contribs/graphs">https://github.com/coq-contribs/graphs</a>
groups	14	134	<a href="https://github.com/coq-contribs/groups">https://github.com/coq-contribs/groups</a>
group-theory	105	1099	<a href="https://github.com/coq-contribs/group-theory">https://github.com/coq-contribs/group-theory</a>
hahn	1458	9836	<a href="https://github.com/vafeiadis/hahn">https://github.com/vafeiadis/hahn</a>
hanoi	577	10514	<a href="https://github.com/they/hanoi">https://github.com/they/hanoi</a>
hardware	186	1643	<a href="https://github.com/coq-contribs/hardware">https://github.com/coq-contribs/hardware</a>
hedges	110	2310	<a href="https://github.com/coq-contribs/hedges">https://github.com/coq-contribs/hedges</a>
helix	2046	56461	<a href="https://github.com/vzaliva/helix">https://github.com/vzaliva/helix</a>
higman-cf	30	279	<a href="https://github.com/coq-contribs/higman-cf">https://github.com/coq-contribs/higman-cf</a>
higman-s	49	913	<a href="https://github.com/coq-contribs/higman-s">https://github.com/coq-contribs/higman-s</a>
HoTT-categories	380	4559	<a href="https://github.com/CategoricalData/HoTT-categories">https://github.com/CategoricalData/HoTT-categories</a>
hybrid	524	7099	<a href="https://github.com/Eelis/hybrid">https://github.com/Eelis/hybrid</a>
idxassoc	58	713	<a href="https://github.com/coq-contribs/idxassoc">https://github.com/coq-contribs/idxassoc</a>
ieee754	15	239	<a href="https://github.com/coq-contribs/ieee754">https://github.com/coq-contribs/ieee754</a>
IFC	756	11435	<a href="https://github.com/QuickChick/IFC">https://github.com/QuickChick/IFC</a>
imm	1444	24624	<a href="https://github.com/weakmemory/imm">https://github.com/weakmemory/imm</a>
infinitary-rewriting-coq	262	4207	<a href="https://github.com/martijnvermaat/infinitary-rewriting-coq">https://github.com/martijnvermaat/infinitary-rewriting-coq</a>
int-map	494	7816	<a href="https://github.com/coq-contribs/int-map">https://github.com/coq-contribs/int-map</a>
ipc	479	9102	<a href="https://github.com/coq-contribs/ipc">https://github.com/coq-contribs/ipc</a>
iris-coq	4299	32669	<a href="https://github.com/izgzhen/iris-coq">https://github.com/izgzhen/iris-coq</a>

iron	1519	30137	<a href="https://github.com/discus-lang/iron">https://github.com/discus-lang/iron</a>
izf	122	827	<a href="https://github.com/coq-contribs/izf">https://github.com/coq-contribs/izf</a>
jordan-curve-theorem	663	42492	<a href="https://github.com/coq-contribs/jordan-curve-theorem">https://github.com/coq-contribs/jordan-curve-theorem</a>
kami	2092	41182	<a href="https://github.com/mit-plv/kami">https://github.com/mit-plv/kami</a>
Kami	2120	37732	<a href="https://github.com/sifive/Kami">https://github.com/sifive/Kami</a>
Ktheory	40	505	<a href="https://github.com/DanGrayson/Ktheory">https://github.com/DanGrayson/Ktheory</a>
lambda	90	941	<a href="https://github.com/coq-contribs/lambda">https://github.com/coq-contribs/lambda</a>
largecatmodules	467	9644	<a href="https://github.com/UniMath/largecatmodules">https://github.com/UniMath/largecatmodules</a>
lazy-pcf	83	1587	<a href="https://github.com/coq-contribs/lazy-pcf">https://github.com/coq-contribs/lazy-pcf</a>
lesniewski-mereology	140	1396	<a href="https://github.com/coq-contribs/lesniewski-mereology">https://github.com/coq-contribs/lesniewski-mereology</a>
lin-alg	543	10797	<a href="https://github.com/coq-contribs/lin-alg">https://github.com/coq-contribs/lin-alg</a>
linearscan	304	4786	<a href="https://github.com/jwiegle/linearscan">https://github.com/jwiegle/linearscan</a>
llvmtwin-coq	675	34538	<a href="https://github.com/snu-sf/llvmtwin-coq">https://github.com/snu-sf/llvmtwin-coq</a>
LOGIC	1729	25450	<a href="https://github.com/QinxiangCao/LOGIC">https://github.com/QinxiangCao/LOGIC</a>
loopring-protocol2-verification	155	3718	<a href="https://github.com/sec-bit/loopring-protocol2-verification">https://github.com/sec-bit/loopring-protocol2-verification</a>
lvc	3782	60531	<a href="https://github.com/sigurdschneider/lvc">https://github.com/sigurdschneider/lvc</a>
maple-mode	20	154	<a href="https://github.com/coq-contribs/maple-mode">https://github.com/coq-contribs/maple-mode</a>
MapleS	7629	153467	<a href="https://github.com/namefanwjcom/MapleS">https://github.com/namefanwjcom/MapleS</a>
markov	31	951	<a href="https://github.com/coq-contribs/markov">https://github.com/coq-contribs/markov</a>
maths	38	386	<a href="https://github.com/coq-contribs/maths">https://github.com/coq-contribs/maths</a>
memory-safe-language	72	1325	<a href="https://github.com/arthuraa/memory-safe-language">https://github.com/arthuraa/memory-safe-language</a>
metacoq	8885	161078	<a href="https://github.com/MetaCoq/metacoq">https://github.com/MetaCoq/metacoq</a>
metalib	796	8026	<a href="https://github.com/plclub/metalib">https://github.com/plclub/metalib</a>
micro-policies-coq	569	13727	<a href="https://github.com/micro-policies/micro-policies-coq">https://github.com/micro-policies/micro-policies-coq</a>
mini-compiler	4	86	<a href="https://github.com/coq-contribs/mini-compiler">https://github.com/coq-contribs/mini-compiler</a>
miniml	12	681	<a href="https://github.com/coq-contribs/miniml">https://github.com/coq-contribs/miniml</a>
mirror-core	1380	36458	<a href="https://github.com/gmalecha/mirror-core">https://github.com/gmalecha/mirror-core</a>
mirror-shard	636	12465	<a href="https://github.com/gmalecha/mirror-shard">https://github.com/gmalecha/mirror-shard</a>

## 23:58 Building a Large Proof Repair Dataset

mod-red	271	6204	<a href="https://github.com/coq-contribs/mod-red">https://github.com/coq-contribs/mod-red</a>
oeuf	10471	186512	<a href="https://github.com/uwplse/oeuf">https://github.com/uwplse/oeuf</a>
orb-stab	40	1590	<a href="https://github.com/coq-contribs/orb-stab">https://github.com/coq-contribs/orb-stab</a>
otway-rees	21	613	<a href="https://github.com/coq-contribs/otway-rees">https://github.com/coq-contribs/otway-rees</a>
paco	1519	23749	<a href="https://github.com/snu-sf/paco">https://github.com/snu-sf/paco</a>
paramcoq-iff	285	6527	<a href="https://github.com/aa755/paramcoq-iff">https://github.com/aa755/paramcoq-iff</a>
param-pi	72	3853	<a href="https://github.com/coq-contribs/param-pi">https://github.com/coq-contribs/param-pi</a>
parsing-parses	274	3512	<a href="https://github.com/JasonGross/parsing-parses">https://github.com/JasonGross/parsing-parses</a>
perennial	5284	137211	<a href="https://github.com/mit-pdos/perennial">https://github.com/mit-pdos/perennial</a>
pigeons	6	34	<a href="https://github.com/llee454/pigeons">https://github.com/llee454/pigeons</a>
pipcore	1637	130325	<a href="https://github.com/2xs/pipcore">https://github.com/2xs/pipcore</a>
pnf	356	2857	<a href="https://github.com/ilyasergey/pnf">https://github.com/ilyasergey/pnf</a>
PolTac	309	2247	<a href="https://github.com/thery/PolTac">https://github.com/thery/PolTac</a>
probchain	246	4801	<a href="https://github.com/certichain/probchain">https://github.com/certichain/probchain</a>
ProcKami	54	2055	<a href="https://github.com/sifive/ProcKami">https://github.com/sifive/ProcKami</a>
promising-coq	835	29947	<a href="https://github.com/snu-sf/promising-coq">https://github.com/snu-sf/promising-coq</a>
proofs	174	2558	<a href="https://github.com/stepchowfun/proofs">https://github.com/stepchowfun/proofs</a>
propcalc	66	629	<a href="https://github.com/coq-contribs/propcalc">https://github.com/coq-contribs/propcalc</a>
pts	350	4513	<a href="https://github.com/coq-contribs/pts">https://github.com/coq-contribs/pts</a>
puiseuxth	910	20553	<a href="https://github.com/roglo/puiseuxth">https://github.com/roglo/puiseuxth</a>
qarith	60	1151	<a href="https://github.com/coq-contribs/qarith">https://github.com/coq-contribs/qarith</a>
quadcopter	1061	34482	<a href="https://github.com/drocketts/quadcopter">https://github.com/drocketts/quadcopter</a>
QuantumLib	1482	20473	<a href="https://github.com/inQWIRE/QuantumLib">https://github.com/inQWIRE/QuantumLib</a>
QWIRE	1167	18932	<a href="https://github.com/inQWIRE/QWIRE">https://github.com/inQWIRE/QWIRE</a>
railroad-crossing	100	1108	<a href="https://github.com/coq-contribs/railroad-crossing">https://github.com/coq-contribs/railroad-crossing</a>
ramsey	8	55	<a href="https://github.com/coq-contribs/ramsey">https://github.com/coq-contribs/ramsey</a>
regex	105	1556	<a href="https://github.com/coq-contribs/regex">https://github.com/coq-contribs/regex</a>
regex-reexamined-coq	476	7609	<a href="https://github.com/awalterschulze/regex-reexamined-coq">https://github.com/awalterschulze/regex-reexamined-coq</a>
rem	8	113	<a href="https://github.com/coq-contribs/rem">https://github.com/coq-contribs/rem</a>
rewriter	1436	13713	<a href="https://github.com/mit-plv/rewriter">https://github.com/mit-plv/rewriter</a>

rezk_completion	350	4261	<a href="https://github.com/benediktahrens/rezk_completion">https://github.com/benediktahrens/rezk_completion</a>
riscv-coq	209	3367	<a href="https://github.com/samuelgruetter/riscv-coq">https://github.com/samuelgruetter/riscv-coq</a>
rsa	113	1194	<a href="https://github.com/coq-contribs/rsa">https://github.com/coq-contribs/rsa</a>
ruler-compass-geometry	332	4376	<a href="https://github.com/coq-contribs/ruler-compass-geometry">https://github.com/coq-contribs/ruler-compass-geometry</a>
SCEV-coq	23	475	<a href="https://github.com/bollu/SCEV-coq">https://github.com/bollu/SCEV-coq</a>
schroeder	14	276	<a href="https://github.com/coq-contribs/schroeder">https://github.com/coq-contribs/schroeder</a>
search-trees	50	593	<a href="https://github.com/coq-contribs/search-trees">https://github.com/coq-contribs/search-trees</a>
Set-Theory	2703	66081	<a href="https://github.com/choukh/Set-Theory">https://github.com/choukh/Set-Theory</a>
shuffle	40	430	<a href="https://github.com/coq-contribs/shuffle">https://github.com/coq-contribs/shuffle</a>
silveroak	1519	26477	<a href="https://github.com/project-oak/silveroak">https://github.com/project-oak/silveroak</a>
sirtt	179	4645	<a href="https://github.com/TheoWinterhalter/sirtt">https://github.com/TheoWinterhalter/sirtt</a>
smc	685	25108	<a href="https://github.com/coq-contribs/smc">https://github.com/coq-contribs/smc</a>
SQIR	1713	37680	<a href="https://github.com/inQWIRE/SQIR">https://github.com/inQWIRE/SQIR</a>
SquiggleEq	1302	16922	<a href="https://github.com/aa755/SquiggleEq">https://github.com/aa755/SquiggleEq</a>
ssrbit	344	2966	<a href="https://github.com/ejgallego/ssrbit">https://github.com/ejgallego/ssrbit</a>
StdLibKami	72	4505	<a href="https://github.com/sifive/StdLibKami">https://github.com/sifive/StdLibKami</a>
StructTact	287	3563	<a href="https://github.com/uwplse/StructTact">https://github.com/uwplse/StructTact</a>
subst	417	3650	<a href="https://github.com/coq-contribs/subst">https://github.com/coq-contribs/subst</a>
system	1	112	<a href="https://github.com/coq-concurrency/system">https://github.com/coq-concurrency/system</a>
tarski-geometry	136	2702	<a href="https://github.com/coq-contribs/tarski-geometry">https://github.com/coq-contribs/tarski-geometry</a>
three-gap	81	1091	<a href="https://github.com/coq-contribs/three-gap">https://github.com/coq-contribs/three-gap</a>
tortoise-hare-algorithm	4	82	<a href="https://github.com/coq-contribs/tortoise-hare-algorithm">https://github.com/coq-contribs/tortoise-hare-algorithm</a>
transfer	201	1726	<a href="https://github.com/Zimmi48/transfer">https://github.com/Zimmi48/transfer</a>
traversable-fincontainer	71	871	<a href="https://github.com/coq-contribs/traversable-fincontainer">https://github.com/coq-contribs/traversable-fincontainer</a>
tree-automata	834	24988	<a href="https://github.com/coq-contribs/tree-automata">https://github.com/coq-contribs/tree-automata</a>
twoSquare	202	1769	<a href="https://github.com/thery/twoSquare">https://github.com/thery/twoSquare</a>
TypeTheory	1923	30575	<a href="https://github.com/UniMath/TypeTheory">https://github.com/UniMath/TypeTheory</a>

## 23:60 Building a Large Proof Repair Dataset

UnifySL	1193	19223	<a href="https://github.com/QinxiangCao/UnifySL">https://github.com/QinxiangCao/UnifySL</a>
UnifySL	1193	19223	<a href="https://github.com/QinxiangCao/UnifySL">https://github.com/QinxiangCao/UnifySL</a>
UniMath	15201	256571	<a href="https://github.com/UniMath/UniMath">https://github.com/UniMath/UniMath</a>
univalent_parametricity	399	5121	<a href="https://github.com/CoqHott/univalent_parametricity">https://github.com/CoqHott/univalent_parametricity</a>
Valuations	184	5079	<a href="https://github.com/FFaissole/Valuations">https://github.com/FFaissole/Valuations</a>
vellvm-legacy	2484	40111	<a href="https://github.com/vellvm/vellvm-legacy">https://github.com/vellvm/vellvm-legacy</a>
velus	3338	58852	<a href="https://github.com/INRIA/velus">https://github.com/INRIA/velus</a>
verdi	618	12534	<a href="https://github.com/uwplse/verdi">https://github.com/uwplse/verdi</a>
verdi-raft	2272	42557	<a href="https://github.com/uwplse/verdi-raft">https://github.com/uwplse/verdi-raft</a>
Verified-FEC	994	20957	<a href="https://github.com/verified-network-toolchain/Verified-FEC">https://github.com/verified-network-toolchain/Verified-FEC</a>
verified-ifc	884	15495	<a href="https://github.com/micro-policies/verified-ifc">https://github.com/micro-policies/verified-ifc</a>
VeriGHC	118	1522	<a href="https://github.com/trommler/VeriGHC">https://github.com/trommler/VeriGHC</a>
VST	27777	570848	<a href="https://github.com/PrincetonUniversity/VST">https://github.com/PrincetonUniversity/VST</a>
WasmCert-Coq	574	13803	<a href="https://github.com/WasmCert/WasmCert-Coq">https://github.com/WasmCert/WasmCert-Coq</a>
weakestmoToImm	1282	33541	<a href="https://github.com/weakmemory/weakestmoToImm">https://github.com/weakmemory/weakestmoToImm</a>
weak-up-to	150	1720	<a href="https://github.com/coq-contribs/weak-up-to">https://github.com/coq-contribs/weak-up-to</a>
when-good-components-go-bad	1030	23295	<a href="https://github.com/secure-compilation/when-good-components-go-bad">https://github.com/secure-compilation/when-good-components-go-bad</a>
yalla	1019	30165	<a href="https://github.com/olaure01/yalla">https://github.com/olaure01/yalla</a>
ynot	885	9996	<a href="https://github.com/ynot-harvard/ynot">https://github.com/ynot-harvard/ynot</a>
zchinese	43	781	<a href="https://github.com/coq-contribs/zchinese">https://github.com/coq-contribs/zchinese</a>
zf	213	2977	<a href="https://github.com/coq-contribs/zf">https://github.com/coq-contribs/zf</a>
zfc	241	2450	<a href="https://github.com/coq-contribs/zfc">https://github.com/coq-contribs/zfc</a>
zorns-lemma	183	4482	<a href="https://github.com/coq-community/zorns-lemma">https://github.com/coq-community/zorns-lemma</a>
zsearch-trees	48	639	<a href="https://github.com/coq-contribs/zsearch-trees">https://github.com/coq-contribs/zsearch-trees</a>
<b>Total</b>	<b>344,316</b>	<b>6,199,931</b>	