# Proofster: Automated Formal Verification

| Arpan Agrawal | Emily First | Zhanna Kaufman | Tom Reichel |
|---|---|---|---|
| University of Illinois | University of Massachusetts | University of Massachusetts | University of Illinois |
| Urbana-Champaign, IL, USA | Amherst, MA, USA | Amherst, MA, USA | Urbana-Champaign, IL, USA |
| arpan2@illinois.edu | efirst@cs.umass.edu | zhannakaufma@cs.umass.edu | reichel3@illinois.edu |

| Shizhuo Zhang | Timothy Zhou | Alex Sanchez-Stern | Talia Ringer |
|---|---|---|---|
| University of Illinois | University of Illinois | University of Massachusetts | University of Illinois |
| Urbana-Champaign, IL, USA | Urbana-Champaign, IL, USA | Amherst, MA, USA | Urbana-Champaign, IL, USA |
| shizhuo2@illinois.edu | ttz2@illinois.edu | sanchezstern@cs.umass.edu | tringer@illinois.edu |

Yuriy Brun
University of Massachusetts
Amherst, MA, USA
brun@cs.umass.edu

*Abstract*—Formal verification is an effective but extremely work-intensive method of improving software quality. Verifying the correctness of software systems often requires significantly more effort than implementing them in the first place, despite the existence of proof assistants, such as Coq, aiding the process. Recent work has aimed to fully automate the synthesis of formal verification proofs, but little tool support exists for practitioners. This paper presents Proofster, a web-based tool aimed at assisting developers with the formal verification process via proof synthesis. Proofster inputs a Coq theorem specifying a property of a software system and attempts to automatically synthesize a formal proof of the correctness of that property. When it is unable to produce a proof, Proofster outputs the proof-space search tree its synthesis explored, which can guide the developer to provide a hint to enable Proofster to synthesize the proof. Proofster runs online at **https://proofster.cs.umass.edu/** and a video demonstrating Proofster is available at **https://youtu.be/xQAi66lRfwI/**.

## I. INTRODUCTION

Software bugs are so routine that the annual cost of operational software failures is more than $1.56 trillion [29], and software engineers spend 35–50% of their time validating and debugging software [37]. Formal verification is a promising method for building correct software systems. Proof assistants, such as Coq [51] and HOL4 [47], inherently support program verification and have had significant industrial impact. For example, Airbus France uses the Coq-verified CompCert C compiler [30] to ensure safety and improve performance of its aircraft [49], Chrome, Android, and Firefox use verified cryptographic libraries [11], [25], and Amazon Web Services applies formal verification to detect misconfigurations that can compromise cloud security [4].

Unfortunately, formal verification is challenging. Writing proofs in Coq is a painstaking exercise that requires deep expertise, as seen in the engineering processes behind several large proof developments [24], [53]. Even with the help of an Interactive Theorem Prover, the effort required to write proofs is often prohibitive. The Coq proof of the C compiler is more than three times that of the compiler code itself [30].

Meanwhile, it took 11 person-years to write the proofs required to verify the seL4 microkernel [35], which represents a tiny fraction of the functionality of a full kernel.

Recent work has aimed to simplify the process of writing proofs [5], [12], [13], [19], [20], [28], [23], [45], [46], [54]. Some formal verification can even be fully automated via proof synthesis. For example, CoqHammer [10] uses a set of precomputed mathematical facts to attempt to "hammer" out a proof. Meanwhile, ASTactic [54], Proverbot9001 [45], TacTok [13], Diva [12], and Passport [46] learn a predictive model from a corpus of existing proofs and use that model to guide a meta-heuristic search to synthesize a proof from scratch.

Unfortunately, relatively little tool support exists for practitioners to use these Coq proof-synthesis tools. For example, of the above-mentioned search-based tools, all but one have neither been integrated into IDEs nor built as stand-alone, graphical interfaces, making adoption difficult. Only Tactician [5] has a usable interface, by way of a plugin for Coq that can be integrated into Coq IDEs. But even then, the interface does not expose the features that help the user understand what the tool is doing under the hood, making debugging and explainability difficult.

In this paper, we present Proofster, a new graphical frontend for search-based proof-synthesis techniques that emphasizes explainability. Conceptually, Proofster can be straightforwardly extended to work with any proof-synthesis backend tool, and implements special features to support explainability for search-based backends. Here, we demonstrate Proofster with Proverbot9001 [45] as its backend.

Proofster's main contributions support the developer in two ways:

1) The developer can enter a theorem describing a software property they want proven, and Proofster uses its underlying backend to attempt to generate a proof. If successful, Proofster displays the Coq proof script, verifying that the

property is correct. PꞦoofster uses the Alectryon library to render literate Coq code [39], which is interactive and easy to read, even when one does not have immediate access to a proof assistant to step through the synthesized proof. The developer can explore the context throughout the proof to better understand why the property is verifiably correct.

2) If the synthesis is unsuccessful, PꞦoofster uses the D3.js library [6] to allow the developer to interactively explore the search tree it used in trying to synthesize a proof, and understand the relevant context. The developer can then identify the most promising search-path, augment it, and have PꞦoofster attempt to synthesize a proof again, using that information.

A live PꞦoofster deployment is available at https://proofster.cs.umass.edu/.

## II. PꞦOOFSTER

PꞦoofster is a frontend tool that interfaces with Coq-based proof synthesis tools. Section II-A discusses how proof engineers interactively write proofs in Coq and how machine-learning-guided proof synthesis tools automatically generate proofs. Section II-B then describes the PꞦoofster implementation and Section II-C illustrates, with examples, how a proof engineer can use PꞦoofster to construct proofs.

### A. Proofs and proof synthesis in Coq

When using the Coq proof assistant, a developer begins by specifying a theorem to prove. This theorem is a type definition in Coq's internal language, Gallina. A proof of that theorem is a term of that type. However, writing that *proof term* directly is difficult, and so Coq provides an interactive environment for reasoning through a proof at a higher level, via a *proof script*.

The developer can use Coq's Ltac language to construct a proof script, a sequence of *tactics* which Coq uses to guide its internal search for a Gallina-based proof term. The theorem prover is called *interactive*, because the developer can specify a tactic to try, have the theorem prover execute the tactic to update the *proof state* (the set of goals that need to be proven, and the known facts), and use that proof state to decide on the next tactic. This interactive process continues until no goals remain, meaning the theorem is proven.

The burden is on the developer to come up with the sequence of tactics. To ease this burden, recent work has created search-based, machine-learning-guided proof-synthesis tools that perform automatic proof-script generation. Most of these tools train a predictive model on a corpus of human-written proof scripts. This model uses a partially written proof script and the theorem being proven to predict a ranked list of the most likely next tactics that should come in the proof script.

The tools differ in how they model the proof scripts when making predictions. For example, ASTactic considers only the current proof state (and ignores the current, partial proof script) [54]. TacTok is a collection of two models — Tac and Tok — both of which encode both the proof state and the partial proof script. Tac works at the tactic granularity, whereas Tok

works at the token granularity; the two prove complementary sets of theorems [13]. These tools model abstract syntax trees using TreeLSTM [50] and proof-script sequences using bidirectional LSTM [38], whereas Proverbot9001, which also models proof state and partial proof script, uses a sequence model [45]. Passport further enhances the model by encoding identifier information for the names of theorems, datatypes, functions, type constructors, and local variables [46]. GamePad, meanwhile, uses its own RNN-based tree encoder and targets only synthetic lemmas [23]. Finally, Diva observes that the variability inherent in machine learning — small perturbations in the learning process, such as hyperparameters, the order in which the training data is seen, and the encoded richness of the training data — leads to diversity in the sets of theorems the learned models can prove. Using the theorem prover's unique ability to serve as an oracle for correctness, Diva uses this diversity to significantly increase its proving power [12].

Armed with a predictive model, these search-based tools search through the space of possible proof scripts. They use the model to predict the likely next proof steps, and the theorem prover to compute the new proof states or errors resulting from these steps. They prune search paths unlikely to be successful or that repeat an already explored state; Proverbot9001, in particular, also prunes states that would explore a subgoal for which a solution was already found. This search through the space of proof scripts represents a set of potential partial proof scripts that aim to make progress toward the goal of proving the theorem. We call the set of explored search paths, together, the *search tree*.

### B. The PꞦoofster implementation

PꞦoofster is implemented as a Flask app and uses BeautifulSoup to create the results page with the synthesized proof and the search graph. PꞦoofster allows the developer to enter a theorem into a text box (or select one from several examples, as a demonstration). PꞦoofster then passes the developer-specified theorem to its proof-synthesis backend and retrieves the search tree, and, if the backend is successful, the synthesized proof. PꞦoofster then uses Alectryon to render the proof as an interactive, literate Coq object. Hovering over a tactic displays the context and goals at that stage of the proof.

PꞦoofster uses the the D3.js library display the search tree and allow the developer to interact with it. Subtrees can be collapsed and expanded to see the tactics tried by the proof synthesis model. This information can also be helpful to developers to provide hints to PꞦoofster in the case where PꞦoofster fails to prove the theorem initially.

PꞦoofster is deployed on AWS and is publicly available at https://proofster.cs.umass.edu/. PꞦoofster is open-source, and is publicly available at https://github.com/UCSD-PL/proverbot9001/tree/demowebtool.

Next, we illustrate PꞦoofster's two use cases using examples.

### C. Using PꞦoofster

Supposed a developer has written a function, `max_elem_list`, that takes a list of natural numbers and returns its largest
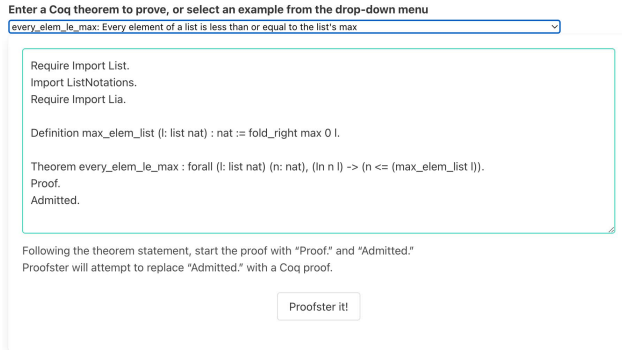
# Proofster

Enter a Coq theorem to prove, or select an example from the drop-down menu

every_elem_le_max: Every element of a list is less than or equal to the list's max

```
Require Import List.
Import ListNotations.
Require Import Lia.

Definition max_elem_list (l: list nat) : nat := fold_right max 0 l.

Theorem every_elem_le_max : forall (l: list nat) (n: nat), (In n l) -> (n <= (max_elem_list l)).
Proof.
Admitted.
```

Following the theorem statement, start the proof with "Proof." and "Admitted."
Proofster will attempt to replace "Admitted." with a Coq proof.

Proofster it!

Fig. 1. A Proofster screenshot of the developer asking to prove the theorem **every_elem_le_max** about the function **max_elem_list**.

# Proofster

```
Require Import List.
Import ListNotations.
Require Import Lia.

Definition max_elem_list (l: list nat) : nat := fold_right max 0 l.

Theorem every_elem_le_max : forall (l: list nat) (n: nat), (In n l) → (n ≤
(max_elem_list l)).
Proof.
induction l.
intros.
simpl.
destruct n.
eauto.
destruct H.
simpl.
intros.
destruct H.
rewrite H.
intuition.
rewrite IHl.
intuition.
eauto.
Qed.
```

```
a : nat
l : list nat
IHl : forall n : nat, In n l → n ≤ max_elem_list l

forall n : nat,
In n (a :: l) → n ≤ max_elem_list (a :: l)
```
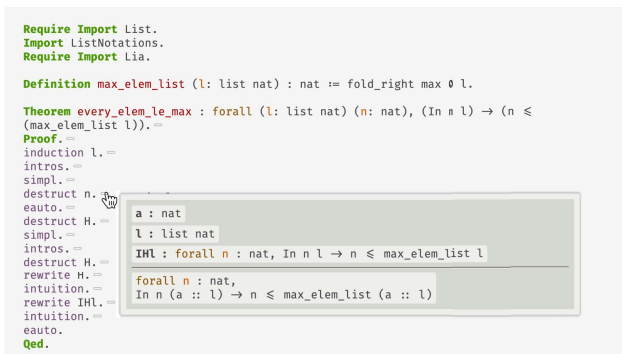
Fig. 2. When Proofster executes the query from Figure 1, it produces a complete proof for the theorem **every_elem_le_max**. Hovering over a tactic in the proof shows the proof state at that point in the proof, which allows the developer to explore and understand how the proof verifies the property.

element. The developer would like to verify this function's correctness by formally proving the property that each element of the list is less than or equal to the result of executing **max_elem_list** on that list.

The developer decides to use Proofster to prove the above property, in Coq. She heads over to https://proofster.cs.umass.edu/ and enters some basic imports, the definition of the **max_elem_list** function, and the theorem **every_elem_le_max**. She does not enter the proof of the theorem, but only starts it with **Proof.** and **Admitted.** to tell Proofster to generate a proof for that theorem. (Proofster will replace **Admitted.** with the proof.)

Figure 1 shows a Proofster screenshot with the developer's inputs. Clicking "Proofster it!" tells Proofster to run its backend to attempt to generate a proof. It succeeds, and Proofster displays the full proof (partial screenshot in Figure 2).

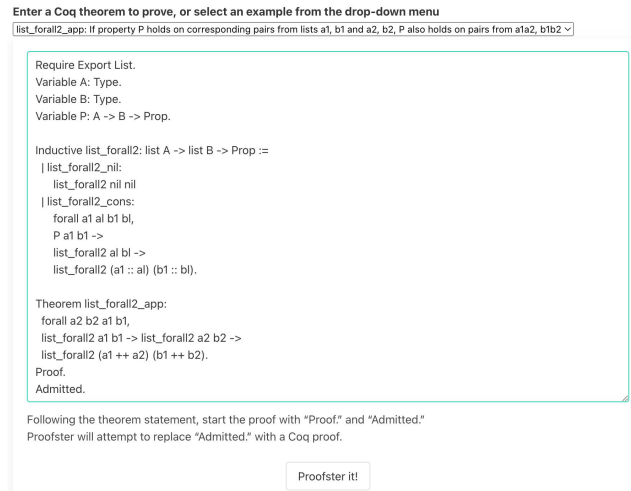The backend will not always be able to produce a proof fully

# Proofster

Enter a Coq theorem to prove, or select an example from the drop-down menu

list_forall2_app: If property P holds on corresponding pairs from lists a1 and a2, b2, P also holds on pairs from a1a2, b1b2

```
Require Export List.
Variable A: Type.
Variable B: Type.
Variable P: A -> B -> Prop.

Inductive list_forall2: list A -> list B -> Prop :=
  | list_forall2_nil:
    list_forall2 nil nil
  | list_forall2_cons:
    forall a1 al b1 bl,
    P a1 b1 ->
    list_forall2 al bl ->
    list_forall2 (a1 :: al) (b1 :: bl).

Theorem list_forall2_app:
  forall a2 b2 a1 b1,
  list_forall2 a1 b1 -> list_forall2 a2 b2 ->
  list_forall2 (a1 ++ a2) (b1 ++ b2).
Proof.
Admitted.
```

Following the theorem statement, start the proof with "Proof." and "Admitted."
Proofster will attempt to replace "Admitted." with a Coq proof.

Proofster it!

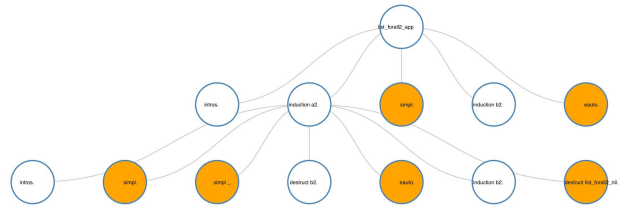Fig. 3. A Proofster screenshot of the developer asking to prove the theorem **list_forall2**.

Fig. 4. When Proofster executes the query from Figure 3, it is not able to generate a complete proof, but displays its search tree, instead. (Image has been rotated for space.)

automatically. Suppose the developer wants to verify another property. Given two lists, let proposition **P** be a proposition on two elements, and let theorem **list_forall2** say that proposition **P** holds for every pair formed by zipping the two lists together. Suppose the developer wants to then prove another property, captured by theorem **list_forall2_app**, which states that for all lists **a1**, **a2**, **b1**, **b2**, if **list_forall2** holds for **a1**, **b1** and for **a1**, **b1**, then it also holds for the pair of lists formed by appending **a1** and **a2**, and appending **b1** and **b2**.

Figure 3 shows the query the developers submits to Proofster to prove this theorem. However, Proofster's backend fails to automatically synthesize a proof for this theorem. Instead of a proof, Proofster displays the search tree for the developer to investigate (Figure 4). She sees that Proofster tried a few forms of induction on the input lists and gets an idea: perhaps inducting over terms of the *relation* between lists **list_forall2 a1 b1**, rather than over the lists directly, will result in a more informative inductive hypothesis. The developer returns to the query page and suggests a hint for Proofster: **induction l**, which inducts over the first unnamed hypothesis (here,

# Pꓤoofster

```coq
Require Export List.
Variable A: Type.
Variable B: Type.
Variable P: A → B → Prop.

Inductive list_forall2: list A → list B → Prop :=
  | list_forall2_nil:
      list_forall2 nil nil
  | list_forall2_cons:
      forall a1 al b1 bl,
      P a1 b1 →
      list_forall2 al bl →
      list_forall2 (a1 :: al) (b1 :: bl).

Theorem list_forall2_app:
  forall a2 b2 a1 b1,
  list_forall2 a1 b1 → list_forall2 a2 b2 →
  list_forall2 (a1 ++ a2) (b1 ++ b2).
Proof.
induction 1.
simpl.
intros.
eauto.
intros.
econstructor.
eauto.
eauto.
Qed.
```

Fig. 5. The succesful result of running the query in Figure 3, modified by adding `induction 1` before `Admitted`.

the term of type `list_forall2 a1 b1`), something Pꓤoofster had failed to try. She then admits the rest and queries Pꓤoofster. Armed with this hint, Pꓤoofster synthesizes the correct proof (Figure 5).

### D. Evaluation Plan

We plan to evaluate Pꓤoofster by soliciting feedback from developers, and by using it in a proof engineering graduate class. Pꓤoofster's backends have been thoroughly evaluated on a benchmark of 68K Coq theorems from 122 open-source projects. ASTactic can fully automatically prove 12.3% of the theorems [54], Passport 12.7% [46], TacTok 12.9 [13], Proverbot9001 [45] 19.2%, and Diva 21.7% [12]. Together with CoqHammer, these tools can prove more than 33% of the theorems.

### III. Related Work

The Pꓤoofster web interface provides an environment to interactively explore both the synthesized proof, and the synthesis search process. It uses the Alectryon [39] library to render literate Coq code, which is interactive and easy to read, even when one does not have immediate access to a proof assistant to step through the synthesized proof. jsCoq [15] and PeaCoq [44] also allow you to interact with formal proofs via web interfaces, but neither synthesize proofs. Tactician tactic-learning Coq plugin can be accessed through a web demonstration of two examples using jsCoq [5]. Section 7.1 of "QED at Large" [42] provides a thorough survey of user interfaces for formal proofs.

Automatically synthesizing proofs from scratch is a promising direction in easing formal verification [5], [10], [13], [12], [23], [26], [45], [46], [54]. For the Coq proof assistant, these methods have been able to prove as many as $\frac{1}{3}$ of the theorems [12] in a large benchmark of correctness properties of software systems [54]. However, these efforts have not yet

directly addressed usability and adoption, which is Pꓤoofster's goal. Such tools could potentially prove mathematical theorems [27] or nonfunctional software properties, such as privacy [9]. For software properties such as fairness [3], [7], [8], [14], [22] and safety [52], complementary approaches provide high-confidence, probabilistic guarantees based on statistical tests and confidence bounds [2], [16], [21], [31], [52].

Proof repair is an important open problem in formal verification [41], [43], which Pꓤoofster may aid by allowing the user to provide hints based on information gained from failed proof-synthesis attempts. This problem is related to automated program repair, which aims to patch defects in systems [18], e.g., using tests and bug reports [33], [17] or inferred constraints on program behavior [1]. In automated program repair, a major challenge is that the tests used to validate the generated patches only partially describe the expected system behavior, and thus the patches can overfit to those tests, failing to correctly repair the program while appearing to do so [34], [48], [40]. Among other methods, extracting test oracles from natural language specifications [32] or using bug reports to help localize defects [33] can help.

### IV. Contributions and Future Work

We have presented Pꓤoofster, a web-based tool aimed at assisting developers with the formal verification process via proof synthesis. Pꓤoofster uses a proof synthesis backend to attempt to automatically generate Coq proofs for user-supplied theorems. The user can use Pꓤoofster to explore the proof state at various stages of the synthesized proof, as well as the search tree generated during synthesis. When Pꓤoofster fails to produce a proof, the user can provide hints as partial proofs, helping Pꓤoofster try again. While our implementation currently works with a specific backend [45], its design is general and aims to work with any Coq proof-synthesis tool, e.g., Diva [12], among others. But reifying that ability is left to future work. Similarly, Pꓤoofster works specifically with proofs for the Coq proof assistant, but, in theory, can be made to work with proof-synthesis tools for other proof assistants, e.g., Thor [26] for the Isabelle/HOL proof assistant [36], among others. Finally, while Pꓤoofster's web-based interface makes it accessible to a broad set of users, we are currently building a version as a Coq plugin, that would integrate it into the IDEs more commonly used by proof engineers.

### References

[1] Afsoon Afzal, Manish Motwani, Kathryn T. Stolee, Yuriy Brun, and Claire Le Goues. SOSRepair: Expressive semantic search for real-world program repair. *IEEE TSE*, 47(10):2162–2181, October 2021.

[2] Aws Albarghouthi, Loris D'Antoni, Samuel Drews, and Aditya Nori. FairSquare: Probabilistic verification for program fairness. In *OOPSLA*, 2017.

[3] Rico Angell, Brittany Johnson, Yuriy Brun, and Alexandra Meliou. Themis: Automatically testing software for discrimination. In *ESEC/FSE Demo*, pages 871–875, November 2018.

[4] AWS Provable Security. https://aws.amazon.com/security/provable-security.

[5] Lasse Blaauwbroek, Josef Urban, and Herman Geuvers. Tactic learning and proving for the Coq proof assistant. In *International Conference on Logic for Programming, Artificial Intelligence, and Reasoning (LPAR)*, volume 73, pages 138–150, 2020.

[6] Mike Bostock. D3.js — Data-driven documents, 2012.

[7] Yuriy Brun and Alexandra Meliou. Software fairness. In *ESEC/FSE NIER*, pages 754–759, November 2018.

[8] Zhenpeng Chen, Jie M. Zhang, Max Hort, Federica Sarro, and Mark Harman. Fairness testing: A comprehensive survey and analysis of trends. *CoRR*, abs/2207.10223, 2022.

[9] Véronique Cortier, Niklas Grimm, Joseph Lallemand, and Matteo Maffei. A type system for privacy properties. In *CCS*, pages 409–423, 2017.

[10] Łukasz Czajka and Cezary Kaliszyk. Hammer for Coq: Automation for dependent type theory. *Journal of Automated Reasoning*, 61(1-4):423–453, 2018.

[11] Andres Erbsen, Jade Philipoom, Jason Gross, Robert Sloan, and Adam Chlipala. Simple high-level code for cryptographic arithmetic — with proofs, without compromises. In *IEEE S&P*, pages 1202–1219, 2019.

[12] Emily First and Yuriy Brun. Diversity-driven automated formal verification. In *ICSE*, 2022.

[13] Emily First, Yuriy Brun, and Arjun Guha. TacTok: Semantics-aware proof synthesis. *PACMPL (OOPSLA)*, 4:231:1–231:31, November 2020.

[14] Sainyam Galhotra, Yuriy Brun, and Alexandra Meliou. Fairness testing: Testing software for discrimination. In *ESEC/FSE*, pages 498–510, September 2017.

[15] Emilio Jesús Gallego Arias, Benoît Pin, and Pierre Jouvelot. jsCoq: Towards hybrid theorem proving interfaces. In *Workshop on User Interfaces for Theorem Provers*, pages 15–27, 2017.

[16] Stephen Giguere, Blossom Metevier, Yuriy Brun, Bruno Castro da Silva, Philip S. Thomas, and Scott Niekum. Fairness guarantees under demographic shift. In *ICLR*, April 2022.

[17] Alberto Goffi, Alessandra Gorla, Michael D. Ernst, and Mauro Pezzè. Automatic generation of oracles for exceptional behaviors. In *ISSTA*, pages 213–224, July 2016.

[18] Claire Le Goues, Michael Pradel, and Abhik Roychoudhury. Automated program repair. *CACM*, 62(12):56–65, November 2019.

[19] Vincent J. Hellendoorn, Premkumar T. Devanbu, and Mohammad Amin Alipour. On the naturalness of proofs. In *ESEC/FSE NIER*, pages 724–728, 2018.

[20] Jónathan Heras and Ekaterina Komendantskaya. Recycling proof patterns in Coq: Case studies. *Mathematics in Computer Science*, 8(1), 2014.

[21] Austin Hoag, James E. Kostas, Bruno Castro da Silva, Philip S. Thomas, and Yuriy Brun. Seldonian toolkit: Building software with safe and fair machine learning. In *ICSE Demo*, May 2023.

[22] Max Hort, Zhenpeng Chen, Jie M. Zhang, Federica Sarro, and Mark Harman. Bias mitigation for machine learning classifiers: A comprehensive survey. *CoRR*, abs/2207.07068, 2022.

[23] Daniel Huang, Prafulla Dhariwal, Dawn Song, and Ilya Sutskever. GamePad: A learning environment for theorem proving. In *ICLR*, 2019.

[24] Jonathan Jacky, Stefani Banerian, Michael D. Ernst, Calvin Loncaric, Stuart Pernsteiner, Zachary Tatlock, and Emina Torlak. Automatic formal verification for EPICS. In *International Conference on Accelerator and Large Experimental Physics Control Systems (ICALEPCS)*, 2017.

[25] Kevin Jacobs and Benjamin Beurdouche. Performance improvements via formally-verified cryptography in Firefox. blog.mozilla.org/security/2020/07/06/performance-improvements-via-formally-verified-cryptography-in-firefox/, 2020.

[26] Albert Jiang, Konrad Czechowski, Mateja Jamnik, Piotr Milos, Szymon Tworkowski, Wenda Li, and Yuhuai Tony Wu. Thor: Wielding hammers to integrate language models and automated theorem provers. In *NeurIPS*, New Orleans, LA, USA, 2022.

[27] Albert Qiaochu Jiang, Wenda Li, Jesse Michael Han, and Yuhuai Wu. LISA: Language models of ISAbelle proofs. In *AITP*, pages 17.1–17.3, September 2021.

[28] Ekaterina Komendantskaya, Jónathan Heras, and Gudmund Grov. Machine learning in proof general: Interfacing interfaces. In *UITP*, 2012.

[29] Herb Krasner. The cost of poor software quality in the US: A 2020 report. https://www.it-cisq.org/pdf/CPSQ-2020-report.pdf, 2020.

[30] Xavier Leroy. Formal verification of a realistic compiler. *CACM*, 52(7):107–115, 2009.

[31] Blossom Metevier, Stephen Giguere, Sarah Brockman, Ari Kobren, Yuriy Brun, Emma Brunskill, and Philip S. Thomas. Offline contextual bandits with high probability fairness guarantees. In *NeurIPS*, pages 14893–14904, December 2019.

[32] Manish Motwani and Yuriy Brun. Automatically generating precise oracles from structured natural language specifications. In *ICSE*, pages 188–199, May 2019.

[33] Manish Motwani and Yuriy Brun. Better automatic program repair by using bug reports and tests together. In *ICSE*, May 2023.

[34] Manish Motwani, Mauricio Soto, Yuriy Brun, René Just, and Claire Le Goues. Quality of automated program repair on real-world defects. *IEEE TSE*, 48(2):637–661, February 2022.

[35] Toby Murray, Daniel Matichuk, Matthew Brassil, Peter Gammie, Timothy Bourke, Sean Seefried, Corey Lewis, Xin Gao, and Gerwin Klein. seL4: From general purpose to a proof of information flow enforcement. In *IEEE S&P*, pages 415–429, 2013.

[36] Tobias Nipkow, Lawrence C Paulson, and Markus Wenzel. *Isabelle/HOL: A proof assistant for higher-order logic*, volume 2283. Springer Science & Business Media, 2002.

[37] Devon H. O'Dell. The debugging mindset: Understanding the psychology of learning strategies leads to effective problem-solving skills. *Queue*, 15(1):71–90, February 2017.

[38] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *NAACL-HLT*, pages 2227–2237, 2018.

[39] Clément Pit-Claudel. Untangling mechanized proofs. In *SLE*, pages 155–174, 2020.

[40] Zichao Qi, Fan Long, Sara Achour, and Martin Rinard. An analysis of patch plausibility and correctness for generate-and-validate patch generation systems. In *ISSTA*, pages 24–36, 2015.

[41] Talia Ringer. *Proof Repair*. PhD thesis, University of Washington, 2021.

[42] Talia Ringer, Karl Palmskog, Ilya Sergey, Milos Gligoric, and Zachary Tatlock. QED at large: A survey of engineering of formally verified software. *Foundations and Trends® in Programming Languages*, 5(2-3):102–281, 2019.

[43] Talia Ringer, RanDair Porter, Nathaniel Yazdani, John Leo, and Dan Grossman. Proof repair across type equivalences. In *PLDI*, pages 112–127, June 2021.

[44] Valentin Robert. *Front-end tooling for building and maintaining dependently-typed functional programs*. PhD thesis, UC San Diego, 2018.

[45] Alex Sanchez-Stern, Yousef Alhessi, Lawrence Saul, and Sorin Lerner. Generating correctness proofs with neural networks. In *MAPL*, 2020.

[46] Alex Sanchez-Stern, Emily First, Timothy Zhou, Zhanna Kaufman, Yuriy Brun, and Talia Ringer. Passport: Improving automated formal verification using identifiers. *ACM TOPLAS*, 2023.

[47] Konrad Slind and Michael Norrish. A brief overview of HOL4. In *TPHOLs*, 2008.

[48] Edward K. Smith, Earl Barr, Claire Le Goues, and Yuriy Brun. Is the cure worse than the disease? Overfitting in automated program repair. In *ESEC/FSE*, pages 532–543, September 2015.

[49] Jean Souyris. Industrial use of CompCert on a safety-critical software product. projects.laas.fr/IFSE/FMF/J3/slides/P05_Jean_Souyiris.pdf, 2014.

[50] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations from tree-structured long short-term memory networks. In *ACL*, pages 1556–1566, 2015.

[51] The Coq Development Team. Coq, v.8.7. https://coq.inria.fr, 2017.

[52] Philip S. Thomas, Bruno Castro da Silva, Andrew G. Barto, Stephen Giguere, Yuriy Brun, and Emma Brunskill. Preventing undesirable behavior of intelligent machines. *Science*, 366(6468):999–1004, 22 November 2019.

[53] James R. Wilcox, Doug Woos, Pavel Panchekha, Zachary Tatlock, Xi Wang, Michael D. Ernst, and Thomas Anderson. Verdi: A framework for implementing and formally verifying distributed systems. In *PLDI*, pages 357–368, 2015.

[54] Kaiyu Yang and Jia Deng. Learning to prove theorems via interacting with proof assistants. In *ICML*, 2019.