

# Leveraging Large Language Models to Boost Dafny’s Developers Productivity

Álvaro Silva  
Independent Researcher  
Porto, Portugal  
amfpsilva@hotmail.com

Alexandra Mendes  
HASLab / INESC TEC & Faculty of  
Engineering, University of Porto  
Porto, Portugal  
alexandra@archimendes.com

João F. Ferreira  
INESC-ID & IST, University of Lisbon  
Lisbon, Portugal  
joao@joaoff.com

## ABSTRACT

This research idea paper proposes leveraging Large Language Models (LLMs) to enhance the productivity of Dafny developers. Although the use of verification-aware languages, such as Dafny, has increased considerably in the last decade, these are still not widely adopted. Often the cost of using such languages is too high, due to the level of expertise required from the developers and challenges that they often face when trying to prove a program correct. Even though Dafny automates a lot of the verification process, sometimes there are steps that are too complex for Dafny to perform on its own. One such case is that of missing lemmas, i.e. Dafny is unable to prove a result without being given further help in the form of a theorem that can assist it in the proof of the step.

In this paper, we describe preliminary work on a new Dafny plugin that leverages LLMs to assist developers by generating suggestions for relevant lemmas that Dafny is unable to discover and use. Moreover, for the lemmas that cannot be proved automatically, the plugin also attempts to provide accompanying calculational proofs. We also discuss ideas for future work by describing a research agenda on using LLMs to increase the adoption of verification-aware languages in general, by increasing developers productivity and by reducing the level of expertise required for crafting formal specifications and proving program properties.

## CCS CONCEPTS

• **Software and its engineering** → **Software notations and tools; Software verification and validation.**

## KEYWORDS

Verification-Aware Languages, Dafny, Large Language Models, Generative AI, Software Productivity, Software Verification, Lemma Inference, Proof Inference, Automated Program Repair, Code Summarization

## 1 INTRODUCTION

As software is becoming increasingly pervasive in our daily lives, it is more important than ever to ensure that it does not contain bugs. Particularly in critical systems, software testing is not enough as stronger assurances are needed. These assurances can be achieved through software verification, which mathematically ensures that software behaves exactly as intended, with respect to a formal specification. Several verification-aware programming languages and systems, where logical constructs such as pre and postconditions, invariants, and assertions provide assurances about the correctness

of the program, are available and enable verification alongside code development. One such language is Dafny [21].

Dafny is celebrated for its advanced deductive verification support that offers a sophisticated backend, encompassing a compiler capable of producing executable binaries and a verifier tasked with meticulously validating code conformity to specified requirements. While Dafny stands as a state-of-the-art tool in software verification, its application demands a profound grasp of concepts often encountered exclusively during the formalization of specifications, which is not commonplace in conventional programming languages. This complexity can hinder software development productivity in Dafny. In addition, and although Dafny has considerable adoption in industry (e.g. at Amazon Web Services<sup>1</sup> and at Consensys [5–7]), it is not widely used, most likely due to the cost in effort, time, and, as stated above, the need for expert knowledge when using Dafny.

Large Language Models (LLMs) have recently demonstrated extraordinary capabilities, exhibiting proficiency in diverse tasks such as engaging in conversations, retrieving and summarizing extensive information, and even generating and explaining text and code [2]. Their utility is further underscored by their application as code suggestion tools, exemplified by GitHub Copilot [28].

In this research idea paper, we describe how we plan to explore LLMs’ capabilities to boost Dafny’s developers productivity and to support its wider adoption. First, we describe some of the challenges identified in the use of Dafny and how we intend to address them. Second, we showcase preliminary work on prompting GPT-4 [2, 4] to support the automated inference of lemmas and their proofs. Finally, we discuss implications of these results, challenges in using these tools within the context of Dafny, and the next steps of our research agenda. Our proposed solution encompasses the integration of the proposed features into the Dafny plugin for VS Code, seamlessly integrating the presented ideas to assist developers in using Dafny to write specifications and their corresponding implementations.

Our overarching goal is to enhance the accessibility of Dafny for new users and cultivate increased autonomy and productivity in Dafny’s use.

## 2 SYNERGIES BETWEEN DAFNY AND LLMs

LLMs and Dafny verification capabilities can complement each other, as LLMs are inherently creative but not always reliable, whilst Dafny’s verification engine acts as a stringent gatekeeper, ensuring that artifacts generated by LLMs undergo a rigorous verification process. This interplay makes their combination promising.

<sup>1</sup><https://github.com/aws/aws-encryption-sdk-dafny>

There are several challenges that Dafny users face when using the language. An important challenge is the common case that non-trivial Dafny developments entail substantial effort in writing lemmas and proofs. For example, in Cassez’s verification of the Incremental Merkle Tree Algorithm [5], almost 90% of the lines of code are proofs and function definitions used in the proofs. Moreover, an experience report on using Dafny at the VerifyThis 2021 verification competition shows that the interpretability of Dafny’s error messages is challenging as these are often not informative [9].

In this section, we describe four challenges that we intend to tackle using LLMs.

## 2.1 Predicate and Lemma Inference

Even though Dafny automates a lot of the verification process, sometimes there are steps that are too complex for Dafny to perform on its own. One such case is that of missing lemmas, i.e. Dafny is unable to prove a result without being given further help in the form of a theorem that can assist it in the proof of the step. The challenge lies not only in proving the lemmas and theorems but also in determining the specifications of lemmas and predicates capable of solving the problem at hand.

Various approaches have been explored in the literature to address these challenges across a diverse range of verification tools. For instance, inferring inductive invariants in TLA+ [34], inference of lemmas in Coq [33, 35] and in symbolic-Heap separation logic [37].

Our objective is to address lemma inference and general predicate inference using LLMs. Our approach involves using prompting to infer lemmas and predicates, and fine-tuning LLMs to infer them, particularly if the zero-shot or few-shot learning approaches yield unsatisfactory results. Section 3 provides an illustrative example of lemma inference using prompting and GPT-4 [2].

## 2.2 Proof Inference

Dafny supports calculational proofs (aka verified calculations), which are proofs by stepwise formula manipulation. As pointed out by Leino and Polikarpova [23], calculational proofs are praised for their rigor, readability, and elegance. Indeed, it has been shown that calculational proofs can greatly improve on traditional verbose proofs in natural language [10–12].

Writing calculational proofs is challenging and any method that can help infer these proofs or steps of these proofs can boost Dafny’s developers productivity. Moreover, proof inference can complement the lemma inference described in subsection 2.1, as when lemmas are added to the code, it is often the case that a proof needs to be provided by the user. LLMs can assist in this process by providing the full proof or by giving hints to the user regarding the proof steps. This also applies to lemmas inferred by the LLM, as proofs will also need to be provided. In addition, even though Dafny can often prove the lemmas on its own, it has been shown that, even in those cases, adding proof steps that are not absolutely necessary can reduce considerably the verification time [5].

To infer proofs, we plan to use models fine-tuned on proof data, but we will also explore few-shot prompted and even zero-shot prompted approaches. Section 3 provides an illustrative example of few-shot prompted proof generation.

## 2.3 Automated Repair

Most programmers make mistakes when writing code. Automated Program Repair (APR) can help with this by supporting developers with automatic fix of software bugs. There are several existing works that successfully applied LLMs for automated program repair [19, 29, 40, 41]. However, as far as we are aware, there are no previous developments on APR for Dafny leveraging LLMs. In addition, in the context of Dafny, when a bug is detected, i.e. the program fails to verify, the issue can be due to an incorrect specification or an incorrect implementation. Although most current research assumes that the specification is correct, focusing on repairing the implementation, it is known that many issues with software stem from incorrect specifications [24]. Previous work on specification repair in Dafny relies on Daikon [8] for dynamic invariant inference which is then used for generating weakening and strengthening candidate fixes (and their combination) [1].

We are not aware of any work on Dafny’s specification or program repair using LLMs. Our goal is to explore existing approaches that use LLMs for automated repair and adapt/improve them to be used for Dafny. We also intend to tackle proof repair and to do so in at least two contexts: when the code changes and existing proofs become invalid and when the user or the LLM suggest a proof that does not verify. We will follow a continuous feedback loop between LLMs and the Dafny’s verifier, where the LLM produces proofs or fixes to proofs and, should these not verify, the feedback produced by Dafny is fed into the LLM to produce another suggestion.

## 2.4 Summarization and Natural Language Specs

Code summarization, the task of generating summaries that accurately describe the functionalities of the code, can reduce developers’ efforts in interpreting the goals of a program or snippet of code. In the context of Dafny, code summarization can also potentially help with understanding the intentions set out in specifications. Further, as detailed above, Dafny’s error messages can be challenging to understand; code summarization has the potential to enhance error messages with further information about the error and about what is the mismatch between the specifications and the code. As developers spend around 58% percent of their time on program comprehension activities [42], features that assist them in code comprehension seem valuable to enhance their productivity.

We intend to explore LLMs’ capabilities in code summarization to enhance error messages in Dafny, to provide explanations of specs/code, and also complement and update specs/code comments which are helpful in code understanding activities. Previous works have shown the enormous potential of LLMs in code summarization tasks [16, 20, 39]. Another feature that we intend to explore is the use of LLMs to support the translation of natural language specifications into Dafny, assisting developers in this task that requires knowledge that is not commonplace.

## 3 PRELIMINARY EXPERIMENTS

This section describes preliminary experiments on using LLMs to infer lemmas and calculational proofs. In particular we used the latest version of GPT-4 Turbo, gpt-4-1106-preview, trained on data up to April 2023, enabling support for contexts with 128,000 tokens.

```

method CoincidenceCount(a: array<int>,
                        b: array<int>)
  returns (c: nat)
  requires forall i,j ::
    0<=i<j<a.Length ==> a[i]<=a[j]
  requires forall i,j ::
    0<=i<j<b.Length ==> b[i]<=b[j]
  ensures c == |multiset(a[..)] * multiset(b[..)]| {
    c := 0;
    var m, n := 0, 0;
    while m < a.Length &&
          n < b.Length
    invariant 0 <= m <= a.Length && 0 <= n <= b.Length
    // Fails to prove the following invariant
    invariant c + |multiset(a[m..]) * multiset(b[n..])|
              == |multiset(a[..)] * multiset(b[..)]|
    decreases a.Length - m + b.Length - n {
      if {
        case a[m] == b[n] => {
          /* Suggested by GPT-4: */
          LemmaIntersectionAfterIncrease_mn(a, b, m, n);
          c, m, n := c + 1, m + 1, n + 1;}
        case a[m] < b[n] => {
          /* Suggested by GPT-4: */
          LemmaIntersectionAfterIncrease_m(a, b, m, n);
          m := m + 1;}
        case b[n] < a[m] => {
          /* Suggested by GPT-4: */
          LemmaIntersectionAfterIncrease_m(a, b, m, n);
          n := n + 1; }}}

```

(a) Method that computes how many elements two sorted arrays have in common. The second invariant cannot be verified automatically. GPT-4 extends the code with calls to lemmas that enable the verification.

You are a software expert specializing in formal methods using the Dafny programming language. You receive the following program where a loop invariant could not be proven. The verifier error message is inside `// VERIFIER ERROR ... //`. Your task is to create lemmas and insert them into the code to facilitate verification.

(b) Prompt used to generate lemmas.

```

lemma LemmaIntersectionAfterIncrease_m
  (a: array<int>, b: array<int>, m: nat, n: nat)
  requires 0 <= m < a.Length && 0 <= n <= b.Length
  requires forall i,j :: 0<=i<j<a.Length ==> a[i]<=a[j]
  requires forall i,j :: 0<=i<j<b.Length ==> b[i]<=b[j]
  ensures |multiset(a[m..]) * multiset(b[n..])| ==
         |multiset(a[m+1..]) * multiset(b[n..])|

lemma LemmaIntersectionAfterIncrease_n
  (a: array<int>, b: array<int>, m: nat, n: nat)
  requires 0 <= m <= a.Length && 0 <= n < b.Length
  requires forall i,j :: 0<=i<j<a.Length ==> a[i]<=a[j]
  requires forall i,j :: 0<=i<j<b.Length ==> b[i]<=b[j]
  ensures |multiset(a[m..]) * multiset(b[n..])| ==
         |multiset(a[m..]) * multiset(b[n+1..])|

lemma LemmaIntersectionAfterIncrease_mn
  (a: array<int>, b: array<int>, m: nat, n: nat)
  requires 0 <= m < a.Length && 0 <= n < b.Length
  requires forall i,j :: 0<=i<j<a.Length ==> a[i]<=a[j]
  requires forall i,j :: 0<=i<j<b.Length ==> b[i]<=b[j]
  //ensures |multiset(a[m..]) * multiset(b[n..])| ==
           |multiset(a[m+1..]) * multiset(b[n+1..])|
  ensures |multiset(a[m..]) * multiset(b[n..])| ==
         |multiset(a[m+1..]) * multiset(b[n+1..])| + 1

```

(c) Lemmas generated by GPT-4. The commented assertion in gray had to be replaced with the assertion in blue.

Figure 1: GPT-4 suggests lemmas to complete method verification.

### 3.1 Lemma Inference

We use an example taken from Leino's book [22] and shown in Figure 1a. Given two sorted arrays, the method `CoincidenceCount` computes how many elements they have in common. The postcondition is expressed in terms of multisets (note that `*` denotes multiset intersection and the vertical-bar brackets denote the cardinality of a multiset). Since the two arrays are sorted, the algorithm can be efficiently implemented using two indices (`m` and `n`) to keep track of how many elements of `a` and `b` have been processed. Dafny's verifier is able to prove the method postcondition from the loop invariants. However, it cannot automatically prove the second invariant (identified by the red comment), because it is not able to automatically prove relevant properties about multisets. The typical approach to solve this is to annotate the program with lemmas that provide enough information for the proof to be completed.

To determine whether GPT-4 can assist us by inferring useful lemmas, we used the code shown in Figure 1a and the prompt shown

in Figure 1b<sup>2</sup>. GPT-4 inferred the lemmas shown in Figure 1c and placed them in the code as shown in Figure 1a (in blue). Only a small correction was required to make the program verify using these lemmas as axioms (highlighted in blue). However, note that these lemmas cannot be proved since they require more information in the precondition (e.g. to prove `LemmaIntersectionAfterIncrease_mn`, it is required to add `a[m] == b[n]` as a precondition; similar annotations are needed for the other lemmas).

### 3.2 Proof Inference

Our experiments with prompt engineering to infer proofs for the three lemmas shown in Figure 1c were not as successful. In general, there were many syntactic errors and we had to provide many examples of proofs to generate plausible solutions. The best result

<sup>2</sup>The code used was the one shown in black, with the additional comment above the invariant that cannot be proved: `// VERIFIER_ERROR loop invariant violation. This invariant could not be proved to be maintained by the loop //`.

```

type pos = x | 1 <= x witness 1

ghost predicate IsFactor(p: pos, x: pos) {
  exists q :: p * q == x }

ghost function Factors(x: pos): set<pos> {
  set p: pos | p <= x && IsFactor(p, x) }

lemma Factor0(p: pos, y: pos, x: pos)
  requires exists a :: x == p*a
  requires exists b :: y == p*b
  ensures IsFactor(p, y + x) {
  /* Proof Suggested by GPT-4: */
  var a: pos := x / p; // Witness for x == p*a
  var b: pos := y / p; // Witness for y == p*b

  calc {
    p * (a + b);
    == { arithmetic }
    p * a + p * b;
    == { definition of a and b }
    x + y; }}

```

Figure 2: GPT-4-generated proof for lemma `Factor0`.

we obtained was when we asked to prove one of the lemmas, but gave a complete proof of one of the other similar lemmas.

To test whether GPT-4 would be able to help Dafny developers when they need to prove less domain-specific results, we also attempted to infer proofs for statements that are more widely known. An example is shown in Figure 2, where a proof for the lemma `Factor0` was generated by GPT-4. The lemma is a version of a basic property of number theory:  $p$  is a factor of any linear combination  $p * a + p * b$ . The predicate is a variation of the lemmas required to prove correctness of Euclid’s algorithm as implemented in Dafny’s integration tests.<sup>3</sup>

The best attempt by GPT-4 to prove the lemma is shown in blue in Figure 2. To achieve this, we designed a prompt incorporating extensive Dafny code and examples; in particular we provided all the contents of the integration test file for calculations.<sup>4</sup> Even though the hints are correct and make sense, they are not accepted by Dafny. Nevertheless, once we comment them, the proof is accepted. Moreover, in the version of Dafny that we used, 4.3.0, we had to change the variable definitions to:

```
var a :| x == p*a; var b :| y == p*b;
```

## 4 RELATED WORK

Clover [36] is the only work that we are aware of on using LLMs in the context of Dafny. Clover consists on using a checker that performs consistency checks among code, docstrings, and formal

annotations. The main idea is to reduce correctness checking to a problem of consistency checking. However, the reduction from correctness to consistency is not mathematically complete. Our goals are different and more fine-grained: we aim to prove correctness and to contribute with features that continuously focus on parts of the program, not necessarily the program as a whole, to enhance developers productivity in achieving program correctness.

Regarding lemma inference, we are not aware of any work focused on Dafny nor work that uses LLMs, but there have been efforts to synthesize lemmas for Coq [33, 35], and for symbolic-Heap separation logic [37]. `ADTIND` automates proofs by induction over algebraic data types, where lemmas are synthesized by term enumeration guided by user-specified templates [45].

Recent work on neural methods to automate proof synthesis is related to our goals of inferring calculational proofs. Given a partial proof and the proof state, neural theorem provers use neural networks to predict the next likely proof steps, which are then evaluated by a proof assistant to return new proof states or errors. Several neural theorem provers have been proposed for Coq [3, 13, 14, 17, 31, 32, 38, 43], for Lean [44] and Isabelle [18, 27]. `Baldur` uses LLMs to repair proofs as part of its proof synthesis approach [15]. Finally, regarding computer support for calculational proofs, there have been efforts to create structure editors that also support verified calculations [25, 26] but none using LLMs.

## 5 CONCLUSION

The lemma inference experiment yielded promising results, demonstrating GPT-4’s ability to infer lemmas with minor errors. However, throughout the experiments, several responses contained syntax errors that proved challenging to rectify through reprompting. The program struggled to complete lemma proofs or correct minor errors in lemma specifications.

On the other hand, the proof inference experiment revealed that only through a meticulously crafted prompt with rich examples could GPT-4 successfully complete the proofs. This underscores the need for improved LLMs that reduce the need for extensive manual prompt curation.

In order to achieve these improvements, our next steps will focus on prompting engineering and fine-tuning models with well-crafted datasets containing relevant Dafny code, with a specific emphasis on lemma inference, proof inference, and suggestions that only contain correct syntax.

In terms of dataset creation, we plan to contribute to and extend `CloverBench` [36], the dataset used by `Clover`. At the time of writing, `CloverBench` contains only 60 small CS textbook examples. Many more, and more complex, examples will be required to ensure the practicality of our ideas. Moreover, we plan to create additional and specialized datasets for specific tasks; for example, we intend to create a Dafny dataset akin to Reichel et al.’s dataset for Coq [30], to assist developers in proof repair.

Finally, to increase the adoption of our ideas, we started their integration into Dafny’s `VSCode` plugin.

<sup>3</sup>Permalink: <https://github.com/dafny-lang/dafny/blob/eae8fc97aabc13ea665b486060e831188628d42b/Source/IntegrationTests/TestFiles/LitTests/LitTest/dafny4/gcd.dfy>. See also: <https://leino.science/papers/krml279.html>

<sup>4</sup>Permalink: <https://github.com/dafny-lang/dafny/blob/eae8fc97aabc13ea665b486060e831188628d42b/Source/IntegrationTests/TestFiles/LitTests/LitTest/dafny2/Calculations.dfy>

## REFERENCES

- [1] Alexandre Abreu, Nuno Macedo, and Alexandra Mendes. 2023. Exploring Automatic Specification Repair in Dafny Programs. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW)*. 105–112. <https://doi.org/10.1109/ASEW60602.2023.00019>
- [2] Open AI. 2023. GPT-4 Technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [3] Lasse Blaauwbroek, Josef Urban, and Herman Geuvers. 2020. The Tactician. In *Intelligent Computer Mathematics*. 271–277.
- [4] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712* (2023).
- [5] Franck Cassez. 2021. Verification of the incremental Merkle tree algorithm with Dafny. In *24th International Symposium on Formal Methods (FM) (LNCS, Vol. 13047)*. Springer, 445–462.
- [6] Franck Cassez, Joanne Fuller, Milad K Ghale, David J Pearce, and Horacio MA Quiles. 2023. Formal and Executable Semantics of the Ethereum Virtual Machine in Dafny. In *International Symposium on Formal Methods*. Springer, 571–583.
- [7] Franck Cassez, Joanne Fuller, and Horacio Mijail Antón Quiles. 2022. Deductive verification of smart contracts with Dafny. In *International Conference on Formal Methods for Industrial Critical Systems*. Springer, 50–66.
- [8] Michael D Ernst, Jeff H Perkins, Philip J Guo, Stephen McCamant, Carlos Pacheco, Matthew S Tschantz, and Chen Xiao. 2007. The Daikon system for dynamic detection of likely invariants. *Sci. Comput. Program.* 69, 1-3 (2007), 35–45.
- [9] Marie Farrell, Conor Reynolds, and Rosemary Monahan. 2021. Using dafny to solve the VerifyThis 2021 challenges. In *Proceedings of the 23rd ACM International Workshop on Formal Techniques for Java-like Programs*. 32–38.
- [10] João F. Ferreira, Alexandra Mendes, Roland Backhouse, and Luis S Barbosa. 2009. Which mathematics for the information society?. In *Teaching Formal Methods: Second International Conference, TFM 2009, Eindhoven, The Netherlands, November 2-6, 2009. Proceedings 2*. Springer, 39–56.
- [11] João F. Ferreira, Alexandra Mendes, Alcino Cunha, Carlos Baquero, and Paulo Silva. 2011. Logic training through algorithmic problem solving. In *Tools for Teaching Logic: Third International Congress, TICTTL 2011, Salamanca, Spain, June 1-4, 2011. Proceedings*. Springer, 62–69.
- [12] João F. Ferreira and Alexandra Mendes. 2009. Students' feedback on teaching mathematics through the calculational method. In *2009 39th IEEE Frontiers in Education Conference*. IEEE, 1–6.
- [13] Emily First and Yuriy Brun. 2022. Diversity-Driven Automated Formal Verification. In *ICSE (22-27)*. Pittsburgh, PA, USA, 749–761. <https://doi.org/10.1145/3510003.3510138>
- [14] Emily First, Yuriy Brun, and Arjun Guha. 2020. TacTok: Semantics-Aware Proof Synthesis. *PACMPL OOPSLA 4* (November 2020), 231:1–231:31. <https://doi.org/10.1145/3428299>
- [15] Emily First, Markus Rabe, Talia Ringer, and Yuriy Brun. 2023. Baldur: Whole-Proof Generation and Repair with Large Language Models. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (, San Francisco, CA, USA.) (ESEC/FSE 2023)*. Association for Computing Machinery, New York, NY, USA, 1229–1241. <https://doi.org/10.1145/3611643.3616243>
- [16] Shuzheng Gao, Cuiyun Gao, Yulan He, Jichuan Zeng, Lunyui Nie, Xin Xia, and Michael Lyu. 2023. Code Structure-Guided Transformer for Source Code Summarization. *ACM Trans. Softw. Eng. Methodol.* 32, 1, Article 23 (feb 2023), 32 pages. <https://doi.org/10.1145/3522674>
- [17] Daniel Huang, Prafulla Dhariwal, Dawn Song, and Ilya Sutskever. 2019. GamePad: A Learning Environment for Theorem Proving. In *ICLR*. <https://openreview.net/forum?id=r1xwKoR9Y7>
- [18] Albert Jiang, Konrad Czechowski, Mateja Jamnik, Piotr Milos, Szymon Tworkowski, Wenda Li, and Yuhuai Tony Wu. 2022. Thor: Wielding Hammers to Integrate Language Models and Automated Theorem Provers. In *NeurIPS*.
- [19] Matthew Jin, Syed Shahriar, Michele Tufano, Xin Shi, Shuai Lu, Neel Sundaresan, and Alexey Svyatkovskiy. 2023. Inferfix: End-to-end program repair with llms. *arXiv preprint arXiv:2303.07263* (2023).
- [20] Junaed Younus Khan and Gias Uddin. 2023. Automatic Code Documentation Generation Using GPT-3. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering (Rochester, MI, USA) (ASE '22)*. Association for Computing Machinery, New York, NY, USA, Article 174, 6 pages. <https://doi.org/10.1145/3551349.3559548>
- [21] K Rustan M Leino. 2010. Dafny: An automatic program verifier for functional correctness. In *International conference on logic for programming artificial intelligence and reasoning*. Springer, 348–370.
- [22] K Rustan M Leino. 2023. *Program Proofs*. MIT Press.
- [23] K Rustan M Leino and Nadia Polikarpova. 2013. Verified calculations. In *Working Conference on Verified Software: Theories, Tools, and Experiments*. Springer, 170–190.
- [24] Nancy Leveson. 2020. Are you sure your software will not kill anyone? *Commun. ACM* 63, 2 (2020), 25–28.
- [25] Alexandra Mendes, Roland Backhouse, and João F. Ferreira. 2014. Structure editing of handwritten mathematics: Improving the computer support for the calculational method. In *Proceedings of the Ninth ACM International Conference on Interactive Tabletops and Surfaces*. 139–148.
- [26] Alexandra Mendes and João F. Ferreira. 2018. Towards Verified Handwritten Calculational Proofs: (Short Paper). In *Interactive Theorem Proving: 9th International Conference, ITP 2018, Held as Part of the Federated Logic Conference, FloC 2018, Oxford, UK, July 9-12, 2018, Proceedings 9*. Springer, 432–440.
- [27] Maciej Mikula, Szymon Antoniak, Szymon Tworkowski, Albert Qiaochu Jiang, Jin Peng Zhou, Christian Szegedy, Lukasz Kuciński, Piotr Miłoś, and Yuhuai Wu. 2023. Magnushammer: A Transformer-based Approach to Premise Selection. *arXiv:2303.04488*
- [28] Nhan Nguyen and Sarah Nadi. 2022. An empirical evaluation of GitHub copilot's code suggestions. In *Proceedings of the 19th International Conference on Mining Software Repositories*. 1–5.
- [29] Julian Aron Prenner, Hlib Babii, and Romain Robbes. 2022. Can OpenAI's Codex Fix Bugs? An Evaluation on QuixBugs. In *Proceedings of the Third International Workshop on Automated Program Repair (Pittsburgh, Pennsylvania) (APR '22)*. Association for Computing Machinery, New York, NY, USA, 69–75. <https://doi.org/10.1145/3524459.3527351>
- [30] Tom Reichel, R Henderson, Andrew Touchet, Andrew Gardner, and Talia Ringer. 2023. Proof repair infrastructure for supervised models: Building a large proof repair dataset. In *14th International Conference on Interactive Theorem Proving (ITP 2023)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik.
- [31] Alex Sanchez-Stern, Yousef Alhessi, Lawrence Saul, and Sorin Lerner. 2020. Generating correctness proofs with neural networks. In *Proceedings of the 4th ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*. 1–10.
- [32] Alex Sanchez-Stern, Emily First, Timothy Zhou, Zhanna Kaufman, Yuriy Brun, and Talia Ringer. 2023. Passport: Improving Automated Formal Verification Using Identifiers. *ACM TOPLAS* 45, 2, Article 12 (June 2023), 12:1-12:30 pages. <https://doi.org/10.1145/3593374>
- [33] John Sarracino, Tobias Kappé, Ryan Doenges, and Greg Morrisett. [n. d.]. Metaprogramming for Practical and Extensible SMT-Powered Coq Proof Automation. ([n. d.]).
- [34] William Schultz, Ian Dardik, and Stavros Tripakis. 2022. Plain and Simple Inductive Invariant Inference for Distributed Protocols in TLA+. In *2022 Formal Methods in Computer-Aided Design (FMCAD)*. IEEE, 273–283.
- [35] Aishwarya Sivaraman, Alex Sanchez-Stern, Brettton Chen, Sorin Lerner, and Todd Millstein. 2022. Data-driven lemma synthesis for interactive proofs. *Proceedings of the ACM on Programming Languages* 6, OOPSLA2 (2022), 505–531.
- [36] Chuyue Sun, Ying Sheng, Oded Padon, and Clark Barrett. 2023. Clover: Closed-Loop Verifiable Code Generation. *arXiv preprint arXiv:2310.17807* (2023).
- [37] Quang-Trung Ta, Ton Chanh Le, Siau-Cheng Khoo, and Wei-Ngan Chin. 2017. Automated lemma synthesis in symbolic-heap separation logic. *Proceedings of the ACM on Programming Languages* 2, POPL (2017), 1–29.
- [38] Amitayush Thakur, Yeming Wen, and Swarat Chaudhuri. 2023. A Language-Agent Approach to Formal Theorem-Proving. *arXiv:2310.04353*
- [39] Ruyun Wang, Hanwen Zhang, Guoliang Lu, Lei Lyu, and Chen Lyu. 2020. Fret: Functional Reinforced Transformer With BERT for Code Summarization. *IEEE Access* 8 (2020), 135591–135604. <https://doi.org/10.1109/ACCESS.2020.3011744>
- [40] Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. 2023. Automated Program Repair in the Era of Large Pre-Trained Language Models. In *Proceedings of the 45th International Conference on Software Engineering (Melbourne, Victoria, Australia) (ICSE '23)*. IEEE Press, 1482–1494. <https://doi.org/10.1109/ICSE48619.2023.00129>
- [41] Chunqiu Steven Xia and Lingming Zhang. 2022. Less Training, More Repairing Please: Revisiting Automated Program Repair via Zero-Shot Learning. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (, Singapore, Singapore.) (ESEC/FSE 2022)*. Association for Computing Machinery, New York, NY, USA, 959–971. <https://doi.org/10.1145/3540250.3549101>
- [42] Xin Xia, Lingfeng Bao, David Lo, Zhenchang Xing, Ahmed E. Hassan, and Shanying Li. 2018. Measuring Program Comprehension: A Large-Scale Field Study with Professionals. *IEEE Transactions on Software Engineering* 44, 10 (2018), 951–976. <https://doi.org/10.1109/TSE.2017.2734091>
- [43] Kaiyu Yang and Jia Deng. 2019. Learning to prove theorems via interacting with proof assistants. In *ICML*. 6984–6994.
- [44] Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. 2023. LeanDojo: Theorem Proving with Retrieval-Augmented Language Models. In *Neural Information Processing Systems (NeurIPS)*.
- [45] Weikun Yang, Grigory Fedyukovich, and Aarti Gupta. 2019. Lemma synthesis for automating induction over algebraic data types. In *Principles and Practice of Constraint Programming: 25th International Conference, CP 2019, Stamford, CT, USA, September 30–October 4, 2019, Proceedings 25*. Springer, 600–617.