

A Review on “AUTOVERUS: Automated Proof Generation for Rust Code”

Gijung Im

FCAI Lab
Yonsei University

Jan 22, 2026



- **Title:** AUTOVERUS: Automated Proof Generation for Rust Code
- **Source:** OOPSLA 2025
- **Authors:** Microsoft Research (Weidong Cui, et al.), UIUC, and Columbia Univ.
- **Presenter:** Gijung Im



Table of Contents

- 1 Introduction
 - Motivation
 - Background: Verus
 - Unique Challenges
- 2 Overview
 - Architecture
 - Procedure
- 3 Experiment & Evaluation
 - Evaluation Setup
 - Experimental Results
- 4 Conclusion



Motivation

1. Generative AI in Software Engineering:
 - LLMs excel at code generation (e.g., GitHub Copilot)
 - but lag in generating formal proofs.
2. The Need for Correctness:
 - Developers often distrust AI-generated code.
3. Goal:
 - ↪ To automatically generate correctness proofs for Rust code, allowing developers to enjoy both **productivity** (AI) and **reliability** (Verification).



Verus: Verified Rust for Low-level Systems

- A static verification tool for Rust that uses an SMT solver (Z3).
- Allows developers to write specifications and proofs directly in Rust-like syntax.
- Key Philosophy:
 1. *Implementation and Proof in the same language.*
No need to learn separate languages like Rocq or Dafny.
 2. *Zero-overhead.*
All “ghost code” (specifications/proofs) are erased at compile time, leaving only optimized Rust executables.



The Mode System (spec vs. proof vs. exec)

```

1 spec fn is_digit(c: u8) -> bool {
2   c >= 48 && c <= 57
3 }
4
5 spec fn cnt_dig(seq: Seq<u8>) -> int
6   decreases seq.len(),
7 {
8   if seq.len() == 0 {
9     0
10  } else {
11    cnt_dig(seq.drop_last()) +
12    if is_digit(seq.last()) {
13      1 as int
14    } else {
15      0 as int
16    }
17  }
18 }

```

(a) The spec functions.

```

19 fn count_digits(text: &Vec<u8>) -> (ret: usize)
20   ensures ret == cnt_dig(text@),
21 {
22   let mut count = 0;
23   let mut i = 0;
24   while i < text.len()
25     invariant
26       i <= text.len(),
27       count <= i,
28       count == cnt_dig(text@.subrange(0, i as int)),
29   {
30     if text[i] >= 48 && text[i] <= 57
31       {count += 1;}
32     i += 1;
33     assert(text@.subrange(0, i - 1 as int)
34       == text@.subrange(0, i as int).drop_last());
35   }
36   assert(text@ == text@.subrange(0, i as int));
37   count
38 }

```

(b) The implementation.

Fig. 1. A Rust function, in gray background, with Verus annotations. The dark-yellow background highlights the proof annotation needed by Verus to prove the specifications highlighted in the light-yellow.



Challenges in Verus Proof Generation

1. *Data Scarcity.*

- Verus is young (≤ 3 years old), #GitHub repos ≤ 10 .
- Unlike Rocq or Isabelle, there is no massive dataset for training.

2. *Syntax Subtleties.*

- Verus introduces specific syntax extensions—e.g., the @ operator, `int` vs. `u64`.
- LLMs often confuse *ghost code* rules with standard Rust rules.

3. *Feedback Ambiguity.*

- Unlike Interactive Theorem Provers (ITP) that show proof state, SMT solvers essentially return either “Pass” or “Fail”.
- It is difficult to measure progress when verification fails.



Workflow

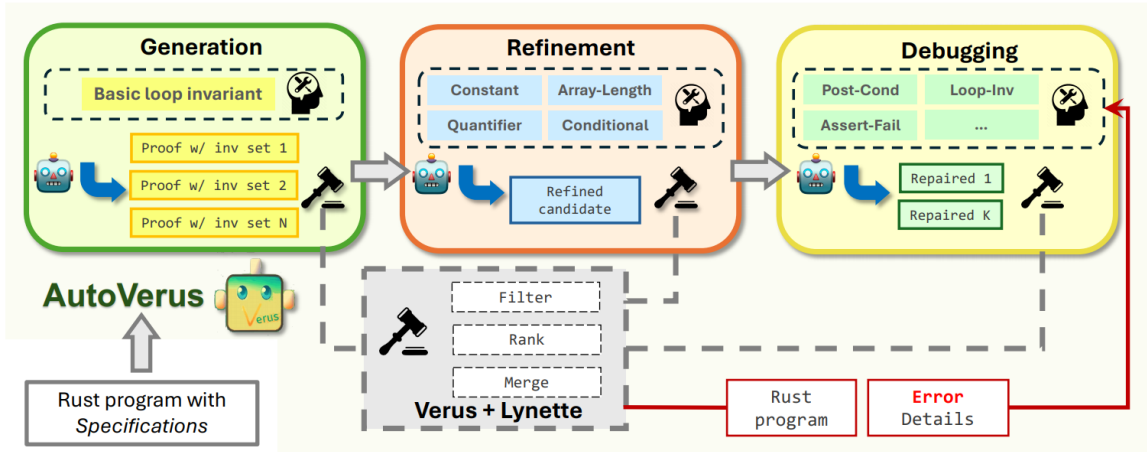


Fig. 2. The workflow of AUTOVERUS

Phase I–Generation (Preliminary Proof)

- **Objective:** Generate initial loop invariants.
- **Post-Processing.**

```
1 def is_correct(p):  
2     return verus_correct(p) or verus_correct(houdini(p))  
3     # The Houdini algorithm deletes the unprovable proof annotations from p.  
4  
5 def GenPrelimProof(program):  
6     pCandidates = LoopInvAgent(program) # proof candidates  
7     safeCandidates = { p ∈ pCandidates | is_safe(program, p) }  
8     Sort(safeCandidates, comparison=(V,E)-rank) # V: verified, E: erroneous  
9     mergeP = program # merged proof  
10    for p in safeCandidates:  
11        mergeP = merge(mergeP, p)  
12        if is_correct(mergeP):  
13            return mergeP  
14    return safeCandidates[0] # first element
```

Listing 1: Preliminary proof generation with post-processing



Example—Diffy Benchmark

```

1 invariant
2   N > 0,
3   i <= N as usize,
4   sum.len() == 1,
5   sum[0] <= N + i,
6   forall|k: int|0<=k<N ==> b[k]==1,
7   b.len() == N,
8 // (V, E) = (4, 2)

```

(a) Loop invariants of best-so-far P_B .

```

1 invariant
2   N > 0,
3   i <= N as usize,
4   sum.len() == 1,
5   sum[0] <= N + i,
6   forall|k: int|0<=k<N ==> b[k]==1,
7   N < 1000,
8 // (V, E) = (3, 2)

```

(b) Loop invariants of a lower-ranked P_i .

```

1 i = 0;
2
3 while (i < N as usize)
4   invariant
5     N > 0,
6     i <= N as usize,
7     sum.len() == 1,
8     sum[0] <= N + i,
9     forall|k: int|0<=k<N ==> b[k]==1,
10    b.len() == N,
11    N < 1000,
12 {
13   sum.set(0, sum[0] + b[i]);
14   i = i + 1;
15 }
16 // (V, E) = (5, 0)
17 // Verified

```

(c) The merged program P' which is *verified*.

Fig. 5. An example of a merged program for a problem from Diffy benchmark.



Phase II–Refinement (Generic Fixes)

- **Objective:** Correct common mistakes before consulting the verifier.
- **4 Specialized Agents.**
 - *Constant-Propagation:* Propagates preconditions into every loop.
 - *Array-Length:* Adds invariants about every array/container sizes (e.g., `len()`).
 - *Quantifier:* Checks whether quantifier-related invariants are used correctly.
 - *Conditional-Loop-Invariant:* Adjusts invariants for specific iterations of a loop.

```
1 def RefineProof(program):
2     best = program
3     # from the simplest (i.e., Constant-Propagation) to the more complex ones
4     for refineAgent in refineAgentList:
5         p = refineAgent(best)
6         if is_correct(p):
7             return p
8         if accept_refine(p):
9             best = p
10    return best
```

Listing 2: Proof Refinement



Phase III-Debugging (Error-Driven)

- **Objective:** Iteratively fix verification errors reported by Verus.

```
1 def DebugProof(program):  
2     for _ in range(MaxIter):  
3         if is_correct(program):  
4             return program  
5     errors = verus_compile(program)  
6     e = select_error(errors)  
7     repair_agent = select_repair_agent(e)  
8     repairedCandidates = repair_agent(program, e)  
9     for p in repairedCandidates:  
10        if is_correct(p):  
11            return p  
12        if accept_repair(p):  
13            program = p  
14            break  
15 return program
```

Listing 3: Debug Proof Generation

repair_precond_error

repair_precond_veclen

repair_postcond_error

repair_invfail_front

repair_invfail_end

repair_assertion_error

repair_arithmetic_flow

repair_mismatch_type

repair_assertion_error
_with_proof_func

repair_default

Fig. 6. AUTOVERUS repair agent list



Fibonacci Sequence (Iteration 1 & 2)

```

1 fn fibonacci(n: usize) -> (ret: Vec<i32>)
2   requires
3     fibo_fits_i32(n as int),
4     n >= 2,
5   ensures
6     forall|i: int| 2 <= i < n ==>
7       #[trigger] ret@[i] == fibo(i),
8     ret@.len() == n,
9 {
10  let mut fib = Vec::new();
11  fib.push(0);
12  fib.push(1);
13  let mut i = 2;
14  +1 // Assert the invariant right before the loop
15  +1 assert(fib.len() == 2);
16  while i < n
17    invariant
18      0 <= 2 <= i <= n,
19      fib.len() == i,
20      fibo_fits_i32(n as int),
21      forall|j: int| 0 <= j < i ==>
22        #[trigger] fib[j] == fibo(j),

```

```

23 -1 fib.len() == n;
24 {
25 +2 // Assert the addition won't overflow an i32
26 +2 assert(fib[i-1] as int + fib[i-2] as int
27 +2         < 0x8000_0000) by {
28 +2   assert(fib[i - 1] == fibo(i - 1));
29 +2   assert(fib[i - 2] == fibo(i - 2));
30 +2   assert(fibo(i-1)+fibo(i-2)==fibo(i));
31 +2   lemma_fibo_monotonic(i, n);
32 +2 };
33   let next_fib = fib[i - 1] + fib[i - 2];
34   fib.push(next_fib);
35   i += 1;
36 }
37 fib
38 }
39 +2 proof fn lemma_fibo_monotonic(..) {..}

```

(a) Repair the “the loop invariant not satisfied before loop” and “arithmetic overflow” errors.



Fibonacci Sequence (Iteration 3)

```
1 spec fn fibo(n: int) -> nat
2   decreases n,
3 {
4   if n <= 0 { 0 }
5   else if n == 1 { 1 }
6   else { fibo(n - 2) + fibo(n - 1) }
7 }
8 spec fn fibo_fits_i32(n: int) -> bool {
9   fibo(n) < 0x8000_0000
10 }
```

(b) The spec functions.

```
1 proof fn lemma_fibo_monotonic(n:int, m:int)
2   requires n <= m,
3   ensures fibo(n) <= fibo(m),
4   decreases m - n
5 {
6   if n < m {
7     lemma_fibo_monotonic(n, m - 1);
8 +   assert(fibo(n) <= fibo(m - 1));
9 +   assert(fibo(m - 1) <= fibo(m));
10  }
11 }
```

(c) Repair the post-condition not satisfied error.



Fibonacci Sequence (Iteration 4 & 5)

```
1 proof fn lemma_fibo_monotonic(n:int, m:int)
2   ...
3 {
4   if n < m {
5     lemma_fibo_monotonic(n, m - 1);
6     assert(fibo(n) <= fibo(m - 1));
7 -   assert(fibo(m - 1) <= fibo(m));
8 +   if m > 1 { // Avoiding the case m == 1 to
9               // prevent negative indexing in fibo(m-2)
10 +    assert(fibo(m-2)+fibo(m-1)==fibo(m));
11 +    assert(fibo(m - 1) <= fibo(m));
12   }
13 }
```

(d) Repair the assertion failure error.



Verus-Bench

- A curated suit of $150 (= 78 + 38 + 11 + 23)$ non-trivial proof tasks.
 - *MBPP*: 78 tasks (Python \rightarrow Rust)
 - *Diffy*: 38 tasks (C benchmarks)
 - *CloverBench*: 11 tasks
 - *Misc*: 23 tasks (Algorithms)
- Excludes trivial tasks that Verus can verify without annotations.
- *Comparison Baseline*: GPT-4o with sophisticated prompting (repeated invocations).

Table 1. Summary of Verus-Bench

Benchmark Sources	CloverBench	Diffy	MBPP	Misc	Total
# of Proof Tasks	11	38	78	23	150
Executable LOC	175	951	1,333	390	2,849
Specification LOC	80	265	700	207	1,252



Results

- **Success Rate:** AutoVerus = 91.3% vs. Baseline = 44.7%.
- **Efficiency:**
 - Solved $> 50\%$ of tasks within 30[sec] or 3 LLM calls.
 - Baseline required significantly more time and calls.
- **Cost:** Total cost for the entire benchmark was only $\leq \$37$.



Summar & Significance

Summary.

- *Synergy of Techniques*: Combines LLM creativity with formal rigor.
- *Workflow Mimicry*: Imitated the “Draft-Refine-Debug” loop used by human experts.
- *Overcoming Data Scarcity*: Demonstrated that In-Context Learning + Static Analysis can solve problems in few resources languages without fine-tuning.

Significance.

- *Democratizing Verification*: Lowers the barrier to entry for verifying Rust systems.
- *Practicality*: Works on standard hardware with low cost and high speed.



Future Work

Future Work.

1. Scaling to large-scale system verification with complex dependencies.
2. Automating specification generation.



Thank you for listening!

