

IronSpec: Increasing the Reliability of Formal Specifications

Paper #750

Abstract

The guarantees of formally verified systems are only as strong as their trusted specifications (specs). As observed by previous studies [19, 49], bugs in formal specs invalidate the assurances that proofs provide. Unfortunately, specs—by their very nature—cannot be *proven* correct. Currently, the only way to identify spec bugs is by careful, manual inspection.

In this paper we introduce IronSpec, a framework of automatic and manual techniques to increase the reliability of formal specifications. IronSpec draws inspiration from classical software testing practices, which we adapt to the realm of formal specs. IronSpec facilitates spec testing with automated sanity checking, a methodology for writing *Spec-Testing Proofs* (STPs), and automated spec mutation testing.

We evaluate IronSpec on 14 specs, including six specs of real-world verified codebases. Our results show that IronSpec is effective at flagging discrepancies between the spec and the developer’s intent, and has led to the discovery of *ten* specification bugs across all six real-world verified systems.

1 Introduction

Formal verification has emerged as a promising technique for increasing the robustness of complex systems by helping developers prove that their implementation meets a formal specification. As promising as this approach is, it has a fundamental Achilles’ heel: its guarantees of eliminating *all* bugs in the implementation rely on the *specification being correct*.

The crucial observation that the guarantees of a mechanized proof are only as strong as their specifications is not new and was first identified in 1985 [30]. Specifications (a.k.a. specs) are inherently trusted, rather than proven correct. Trust alone does not ensure that specs remain bug free. If a spec contains a bug, proving that the system meets this spec may be meaningless; the *proven* system could also contain a bug that is hidden by the buggy spec.

The correctness of specifications is the rock upon which the entire edifice of formal verification is built.

Despite the importance of writing correct specs, current best practices rely solely on manual inspection. Developers argue [22, 23] that because specs are typically small in comparison to the size of the corresponding proof and implementation, it is feasible to manually inspect specs thoroughly enough to ensure that they **capture the intended behavior** of the system. While expert developers are more likely to write correct specs, they are not infallible. As formal verification becomes widely adopted, more and more non-experts will write specs only, exacerbating the risk of introducing bugs. Thus, it is imperative that the process of writing specs be as robust as possible.

In fact, several studies [19, 29, 49], through extensive manual effort, have shown that formally verified systems—many of which were developed by experts in formal verification—contain critical bugs, which originate with problems and inconsistencies in their specs. For example, in January 2022, Notional Finance found a double-spending vulnerability in a deployed verified smart contract missed by manual inspection [31]. In this case, part of the spec was vacuous, causing it to be too weak, and thus the proof would still pass with literally *any* implementation.

Since a spec is a formal expression of a developer’s intent, *proving* the spec correct is ultimately impossible. Ensuring a spec matches a developer’s intent will always be best-effort. Whilst no approach can guarantee a bug-free spec, that does not mean attempts must exclusively rely on extensive manual effort and system expertise to resolve. Indeed, there are no structured or automated approaches for a developer to debug this complicated state space. To fill this gap, we propose a means to better handle this challenge.

Inspired by classic testing techniques [14, 21, 45], we introduce IronSpec, a spec testing framework. To enable testing specs, IronSpec adapts the automation of mutation testing [16, 28] and sanity checking along with a customized manual testing approach inspired by unit testing. Together, this framework introduces a systematic way to boost assurance that a spec captures the intended behavior of the system.

If there is a bug in a spec, it originates in the same manner

as any other type of bug; there is a disconnect between the intent of the developer and what is written. A spec is *incorrect*, if it is too weak and allows for any implementation to have undesired behavior, or if the spec is too strong, precluding some desired behavior. To identify spec bugs, we leverage the insight that spec bugs manifest themselves as a consequence of a *disconnect of intent* to search for and highlight such disconnects through structure and automation.

IronSpec aids in pinpointing where the intent of the developer diverged from the current spec by providing various tools that encapsulate this notion. Common cases where the intent of the developer deviated from the spec can be flagged with IronSpec’s Automatic Sanity Checker. If a system has a passing proof, IronSpec provides additional automation with spec mutation testing. Mutation testing can automatically identify cases where the behavior of the implementation differs from the spec, by using the proof to identify relevant mutations. The hints of potential intent disconnect provided by automation are bolstered by a manual methodology for writing *Spec-Testing Proofs (STPs)*. STPs are inspired by traditional unit testing and allow developers to test if their understanding of what behavior the spec should allow matches the current spec. STPs can further be used to investigate the hints provided by automation to either confirm the existence of a bug or absolve the disconnect as intended behavior.

We evaluate IronSpec to demonstrate its value over purely manual inspection by testing six specs produced in-house, two specs containing artificial bugs that were studied in [2], and six specs of open-source verified systems. We demonstrate the effectiveness of our approach by describing *ten* spec bugs found across a verified Distributed Validator Protocol [51], a verified SAT solver [3], a verified QBFT system [44], a formal spec of the Eth2.0 spec [11], daisy-nfsd [12], and a verified AWS Encryption SDK library [1].

Overall, this paper makes the following contributions:

- Introduces IronSpec, a spec testing framework that allows developers to pinpoint places where the current spec may have diverged from their original intent.
- Proposes an Automatic Sanity Checker, a testing methodology for writing *Spec-Testing Proofs (STPs)*, which are applicable to test specs even in the absence of a completed proof or implementation, and describes how to adapt mutation testing to specs to automatically identify divergences between the spec and the implementation.
- Demonstrates the effectiveness of IronSpec, by illustrating how we applied IronSpec to *six* real-world, verified systems leading to the discovery of *ten* spec bugs.

2 Manually Scrutinizing Specifications

Relying on manual inspection alone to ensure an intended specification is not practical. Fonseca et al. [19] performed a study aimed to challenge the assumption that just because a

system is verified, it is bug-free. In this study, the authors thoroughly examined three formally verified distributed systems, IronFleet [23], Verdi [52], and Chapar [38] and identified sixteen bugs across their specifications, verification tools, and their unverified shim layers. Two of these bugs were found to be in specifications. This study was chiefly manual and required close examination of the respective specifications to identify. The authors do introduce some basic automation, yet their techniques still rely predominantly on manual effort and expertise in the system. This work demonstrated the need for and acknowledges the lack of a more rigorous and automated approach to testing formal specifications. Similarly, Yang et al. [53] conducted a bug study of compilers and discovered two bugs within the verified compiler CompCert due to under-specification, and similarly observed that specifications are complex and lack scrutiny.

The concerning discovery of these previous works identifies the gap that this work aims to fill; to provide a means for developers to help automatically and methodically identify specification bugs across the spectrum of specifications.

Complicating this problem, specifications can take on different forms, making uniform debugging approaches difficult. In their simplest form, specifications can be in-line predicate assertions [26]; boolean functions that check the state of the system against some property. A more specific class of predicate assertions based on the Floyd-Hoare style logic [18, 25] are preconditions and postconditions, which establish invariants about the state of the program before and after a specific function. For more complex systems, rather than directly proving properties about the system, it can be easier [23, 52] to prove state machine refinement [34]. For refinement, the specification is an abstract state machine that encapsulates the desired behavior of the system.

To highlight the subtlety of trying to ensure a specification is correct manually, consider an incorrect specification for a simple `Sort` method found on line 3 of Specification 1. This `Sort` method takes a sequence of integers as input and promises to return a sorted sequence of integers in ascending order. The specification for this method is a single postcondition which ensures that the value at every index in the output sequence is less than or equal to the value at subsequent indices. At first glance, this may seem to be a correct specification for `Sort`—a mistake that many newcomers to verification make.

However, this specification is incorrect, as it neglects to mention any relationship between the input and output sequences. Given this buggy specification, a proof could still pass even with an incorrect implementation, erroneously giving the illusion of correctness. For example, if the input sequence was `[1, 6, 7, 2]` an incorrect implementation could arbitrarily return `[1, 42, 100]` or even the empty sequence `[]`. The *incorrect* implementation for `Sort` in Specification 1 is trivial and always returns an empty sequence. Yet, the proof for this method would still pass, as the trivial implementation satisfies

Specification 1 Incorrect Sort Spec

```
1  method Sort(input:seq<int>)
2    returns (out:seq<int>)
3    ensures forall i | 0 <= i < |output| - 1 ::
4      out[i] <= out[i+1]
5    { return []; }
```

Specification 2 Correct Sort Spec

```
1  method Sort(input:seq<int>)
2    returns (out:seq<int>)
3    ensures forall i | 0 <= i < |output| - 1 ::
4      out[i] <= out[i+1]
5    ensures multiset(input) == multiset(output)
6    { /* body omitted */ }
```

the incorrect, too-weak specification.

Manually identifying a spec bug, like that in Specification 1, can be challenging. In fact, a correct specification for `Sort` should also capture the relationship between the input and output by adding an additional post-condition to ensure that the multiset of the input is equal to the multiset of the output, see line 5 in Specification 2.

The opposite case, where a specification is too strong, can be equally as essential and challenging to manually identify. For example, if we replace line 5 in Specification 2 with `ensures input == output`, the specification becomes unnecessarily strong. Multisets do not take order into account, whereas sequences do, so the updated postcondition becomes overly strong by not including the multiset relationship. The only input and output pair that could satisfy this specification is if the sequences are identical and already in ascending order. Even if we have a correct implementation of `Sort`, proving that the implementation upholds the specification is impossible. To debug their inevitably failing proof, the developer must examine their implementation for bugs, check their proof for missing invariants *and* manually inspect their spec to make sure it captures the intended behavior. Having high confidence in the spec would make this scenario much more unlikely and would give the developer more time to focus on the proof itself, knowing they are proving the right thing.

3 How To Test A Specification

It is challenging to diagnose spec bugs because specs are trusted, and issues in specs can often be at odds with a developer's original understanding of the system. Complicating the problem, specs are often intended to be abstract, allowing different, correct implementations to meet the spec. Hence, we introduce IronSpec, a framework for testing specs to help gain confidence that a spec is correct. This work represents the first systematic effort to bridge the gap between the mature and extensive work in software testing and the lack of rigour for ensuring spec correctness.

IronSpec is inspired by the insight that the existence of a spec bug is inherently due to a disconnect between what the developer intended and what properties were captured in the spec. IronSpec provides tools to allow a tester to identify and test possible occurrences where the original intent of the developer may have diverged from the current spec. Some aspects of IronSpec only rely on the spec and do not rely on the existence of an implementation or a passing proof. If there is an implementation and a corresponding passing proof, however, IronSpec can leverage this to use the implementation as an additional reference point to help focus the investigation.

This section introduces and provides a high-level overview behind the ideas of why each testing component of IronSpec is useful in exposing disconnects between the developer and their spec. Section 4 discusses each in more detail.

3.1 Testing Specifications Without A Proof

Akin to test-driven development [6], it is desirable to test a spec without requiring a proof or corresponding implementation. If there is a bug in the spec when it comes time to write a proof, a developer may struggle and expend unnecessary manual effort in debugging in the wrong place. The Automatic Sanity Checker, and *Spec-Testing Proofs (STPs)*, provide two frames of reference for a tester to check their specs against, even in the absence of an implementation and proof.

Regardless of the context of the system, it is clearly never intended for a verified method to be permitted to return arbitrary values. If the spec is too weak, an incorrect implementation might be free to return *any* value, unconstrained by the spec. The Automatic Sanity Checker provides this type of generic, context-free hints about the existence of any such disconnects between the spec of a verified method and its input and output. The sanity checker also flags partial violations of these properties, which do not immediately indicate a bug but are hints to be investigated further. Because the specific disconnect that the Automatic Sanity checker searches for is only between the input, output, and spec, this technique can be used even in the absence of an implementation or proof. Section 4.1 describes the Automatic Sanity Checker in more detail.

The properties between the input, output, and spec that are checked by the Automatic Sanity Checker provide one reference point to gauge any decoupling of intent in the spec but ignore context specific to the spec. We address this problem by introducing a methodology for manually writing *Spec-Testing Proofs (STPs)*. STPs are proofs about spec context specific I/O. STPs help expose differences between the expected behaviors that a developer intends to be included in the spec and what is currently permitted. This testing methodology is useful when having passing or failing proofs, but can uniquely be applied to cases when only having a spec, regardless of the status of a proof. How to write STPs and interpret their results are explained in Section 4.2.

A STP is, by its definition, a proof. Being a proof is the key difference between STPs and standard unit tests. Instead of attempting to prove a general property, a STP demonstrates the validity of the spec for a specific, concretized value or a range of values. The testing methodology exploits the insight that crafting proofs for specific cases is often less challenging than producing a comprehensive proof and can frequently be proved by the verifier with minimal manual intervention. Rather than proving a property for *all* values, each STP is a proof for a *specific* value.

A consequence of STPs being proofs is that if the STP fails, it could be for various reasons. The failure could be due to a divergence between the expectation of the spec and the STP, indicating a bug; or it could be because there are too few strengthening assertions in the STP body for the verifier to prove the final postcondition of the STP. Distinguishing between a spec bug and the need to add proof to the body of the STP is impossible to immediately diagnose for every case due to the underlying undecidability of verifiers. If the tester is suspicious that the failure is indeed due to a disconnect, the appropriate next step is to write a concrete counterexample, proving the unintended behavior is allowed by the spec.

3.2 Testing Specifications With A Proof

Even when a system is verified with a passing proof, it is still possible for the system to contain bugs; thus testing a spec at this point is still very valuable. A too-weak spec could allow for a proof to pass with an incorrect implementation, falsely giving the illusion of correctness. Alternatively, even if the current implementation contains no bugs, a too-weak spec could allow for a different, incorrect implementation to exist, such that a proof would still pass with the same too-weak spec. Relying on a developer to arrive at a bug-free implementation given a buggy spec, goes against the very reason to verify systems in the first place; so its just as vital to identify spec bugs in this case too.

Both automatic sanity checking and writing STPs are applicable when testing a spec with a passing proof, but both the proof and implementation contain untapped information that can further assist testing. Like the spec, the implementation also captures the intent of the developer. Identifying the difference of intent between the behavior allowed by the spec and what is actually in the implementation, calls the developer's attention to potential disconnects. IronSpec can take advantage of the proof and implementation to test a spec with *mutation testing*. Mutation testing identifies cases when the spec is weaker than the current implementation. IronSpec uses the passing proof as a reference point to automatically distinguish cases where the existing implementation is weaker than the behavior allowed by the spec. Further details concerning how IronSpec adapts mutation testing to specs are described in Section 4.3

Departing from traditional mutation testing, IronSpec starts with a spec, implementation, and passing proof and then only mutates the spec. IronSpec relies on an existing passing proof to indicate whether a mutation should be killed, whereas traditional mutation testing relies on a test suite. A mutation is kept and considered *alive* if the original proof still passes with the mutated spec, indicating that the implementation also meets this different spec. The behavior allowed by the original spec but not the mutated spec serves as an example of a subset of behavior that may not be intended.

To reduce the chance of false positives only a subset of the generated mutations are eventually considered. Logically equivalent or weaker mutations than the original spec and mutations that trivially make the proof pass are ignored. The details for how specs are mutated and what constitutes valid mutations are expounded upon in Section 4.3.

Note that we deliberately mutate only specs and not implementations for two reasons. Firstly, specs are smaller than implementations, therefore reducing the number of mutations necessary to consider. Secondly, mutating only the spec rather than the implementation is advantageous for automation. Specs, being predicates, enable automatic filtering of irrelevant mutations. Assuming the proof passes given the original spec, any logically weaker spec mutation will still allow the proof to pass and does not provide any new relevant information. By automatically checking the relative logical strength of a mutated spec in relation to the original, weaker mutations can be identified and ignored. This automation is impossible when mutating the implementation, as determining relative logical strength is impossible in all cases. Logical relationships can be determined automatically for boolean functions, like specs, whereas not all imperative code shares this attribute. Implementation-based mutations would increase the manual burden on the tester, as many more false positives would be an unavoidable outcome that would require manual effort to sift through.

The existence of even a single alive spec mutation serves as a flag to the developer but could indicate various possibilities. An alive mutation is evidence that the spec is weaker in some way than the implementation. Understanding the implication of any such alive mutation cannot be automated and ultimately still relies on the developer's intuition to understand.

An alive mutation is only a hint of a gap between what is allowed by the spec and the current implementation. If a spec is correct but also weaker than the current implementation, there is the chance for alive mutations to be considered false positives and marked as intended behavior. In contrast, though rare, if both the spec and implementation are buggy, but the implementation is not weaker than the buggy spec, then no alive mutations will be found.

In certain cases, mutation testing is also useful in identifying too-strong specs. A spec in the Hoare-Logic style can also be considered incorrect by virtue of having a too strong precondition. IronSpec's mutation testing is still applicable in

Table 1: Automatic Sanity Checking Flags

Flag Severity	Condition
LOW	Post conditions only depend on a portion of the input
MED	Only part of the output is constrained by the postconditions
HIGH	None of the postconditions depend on any of the input
HIGH	None of the output is constrained by any of the postconditions

this case. If the spec mutation target is a precondition, rather than attempting to identify where the spec is disconnected from the implementation due to weakness, IronSpec reverses the criteria used to determine relevant mutations to consider mutations that are *weaker* than the original spec.

The automation of mutation testing does not provide complete coverage of spec testing but rather focuses the attention of the tester on a disconnect between the spec and the implementation. STPs can be used to help fill this gap. Focusing on writing STPs about the discrepancy hinted at by the alive mutation leads to a more efficient way of identifying bugs. STPs guided by the hint of alive mutations can allow a tester either to arrive at a counterexample, showing a bug in the spec, or absolve the alive mutation as intended behavior.

4 The IronSpec framework

IronSpec consists of three spec testing tools; an Automatic Sanity Checker, a methodology for writing *Spec-Testing Proofs* (STPs), and an automatic mutation testing framework. Each assists in identifying and flagging divergences between the developer’s intent and the existing spec.

4.1 Automatic Sanity Checker

The sanity checker examines the input, output, and spec of verified methods to identify cases where the spec may be weaker than intended. The properties that are checked and their severities are outlined in Table 1. Either of the HIGH severity flags signifies spec bugs, whereas the other severity levels indicate a cause for additional manual inspection. If a spec is weak in the ways indicated by the HIGH severity flags, an implementation could either return arbitrary output values or is not dependent on the input. Regardless of the particular functionality of the system, either case is a clear disconnect between the intent of a correct spec and the current spec.

The power of the Automatic Sanity Checker arises from exploiting the insight of the relationship between a spec and the I/O of its corresponding method. Both High-severity flags leverage the condition when the constraints between the I/O and spec are non-existent. If no postcondition depends on

Lemma 3 General Precondition STP

```

1 lemma PreconditionSTP(in:InType)
2   requires TestInputProperty(in)
3   ensures Precondition(in)
4   // or !Precondition(in)
```

Lemma 4 General Postcondition STP

```

1 lemma PostconditionSTP(in:InType,out:OutType)
2   requires TestInputProperty(in,out)
3   ensures Precondition(in)
4   ensures Postcondition(in,out)
5   // or !Postcondition(in,out)
```

any of the input values, then an obvious aspect of the spec is missing. The buggy sort spec in Specification 1 exemplifies this scenario. The spec is not constrained at all by the input, making the spec weak enough to allow for a proof for an incorrect implementation to pass. Similarly, if a method has an output not constrained by its postconditions, an incorrect implementation can return *any* output. The lower severity flags constitute partial violations of the general properties and do not immediately indicate bugs; rather, they signal a missing part of the spec that could be the source of a bug.

4.2 STP Methodology

The testing methodology outlines four classes of STPs to help guide developers in understanding the **Usefulness**, **Correctness**, and **Provability** of their specs, and if there is a bug in the spec, developers can prove its existence with a **Counterexample** STP. The methodology focuses on specs written following the Hoare-Logic style [48] but can be applied to any predicate-based spec. In all cases, the developer specifies context-specific *intended* valid or invalid input and output based on their understanding of what the spec should or should not permit. A passing STP shows that the intent of the developer matches the spec under investigation. How to write each class of STP is described in the following subsection.

4.2.1 Writing STPs

The construction of STPs share many similarities, but the results are interpreted differently. For all STPs, pre and post-conditions can be decoupled and tested individually. The general form for these STPs are found in Lemmas 3 and 4.

Usefulness STPs are concerned about whether the preconditions are weak enough to remain useful; the preconditions should accept all intended valid inputs. Usefulness STPs follow the general form of Lemma 3. The specific input values are defined as part of the precondition for this lemma as the `TestInputProperty`, and should be a value that the test writer *expects* to be a valid I/O allowed by the spec. The postcondition for a Usefulness STP should be the preconditions from the spec, i.e. `ensures Precondition(in)`.

Opposite to Usefulness STPs, Correctness STPs examine whether the postconditions are strong enough to reject all intended invalid outputs i.e. `ensures !Postcondition(in)`. Writing Correctness STPs is based on the general form of Lemma 4. To test if the postcondition is strong enough to reject buggy behavior, the test writer supplies an output value that is *expected* to be invalid and should not be allowed by the spec. To isolate testing the postcondition from the precondition, the test writer should also prove that the undesired output does not satisfy the spec as a result of an invalid input value (Line 3 in Lemma 4), ideally with a separate Usefulness STP validating the input.

Conversely to Usefulness and Correctness STPs, Provability STPs test whether the preconditions are strong enough and whether the postconditions are weak enough for the existence of a provable implementation. Provability STPs are most useful before having a passing proof, as a passing proof is evidence that the spec has this property. STPs for Provability are concerned with both preconditions and postconditions, thus follow from both Lemmas 3, and 4. Precondition STPs prompt the test writer to prove *expected* invalid input should not pass the precondition, i.e. `ensures !Precondition(in)`. Whereas postcondition STPs check that I/O *expected* to be permitted by the spec is allowed by the postconditions, i.e. `ensures Postcondition(in,out)`.

If suspicious of a spec bug, a test writer can also directly write a Counterexample STP. A passing Counterexample STP is concrete evidence of a bug in the spec. Counterexample STPs can take on two different forms but are still derivative of Lemma 4 if concerned with postconditions and Lemma 3 for preconditions. A Counterexample STP can either show that expected **valid I/O is rejected** by the spec or that expected **invalid I/O is accepted**.

4.2.2 Adding Proof Help To STPs

When a STP fails to verify, it could be due to a divergence between the expectations of the test writer and the current spec, indicating a spec bug, which can be confirmed with a counterexample STP, or it could be the result of the fundamental undecidability of the verifier. If it's the latter case, it is possible to circumvent this roadblock in some instances by adding additional proof to the STP body.

The process of proving an STP is no different than writing a proof for any lemma, but the specificity of the STP narrows the scope necessary to reason about. However, before spending the manual effort to add proof to the body of an STP, the first step is to negate the conclusion, i.e. the `ensures` of the STP. Negating the conclusion transforms a STP into a Counterexample STP. Thus, if the proof now passes there is a clear indication of a bug.

As an example of the process of writing an STP, consider the **Correctness** STP in Lemma 5 for the incorrect sort spec from Specification 1. This STP tests that the sort spec should

Lemma 5 Incorrect Sort Spec - Correctness STP Example

```

1 lemma CorrectnessSTPSort (
2   input:seq<int>, sorted:seq<int>)
3   requires input == [42, 1, 500]
4   requires sorted == [42, 500]
5   ensures !SortSpec(input,sorted)
6   { }
```

Table 2: Mutation Operators

Operator	Description
AOR	Arithmetic Operator Replacement
LOR	Logical Operator Replacement
ROR	Relational Operator Replacement
COI	Constant Operator Insertion
UOR	Unary Operator Replacement
ENO	Expression Negation Operator
VNOR	Variable Name Operator Replacement
SOR	Set Operator Replacement
HOR	Heap Operator Replacement

reject the case when the output sequence is sorted, but only contains a subset of the original input. Running a verifier on this STP would initially result in failing to automatically prove the postcondition. Before spending manual effort to prove this STP, the first step is always to negate the postcondition (i.e. `ensures SortSpec(input,sorted)`), transforming the Correctness STP into a Counterexample STP. The attempt to prove this counterexample would now pass and serves as a concrete example of where the spec has diverged from the test writer's understanding.

4.3 Mutation Testing

If the system has a passing proof, IronSpec can leverage the proof as a reference point for mutation testing. Mutation testing helps to identify any disconnect between the implementation and the spec. A mutation is *killed* if it fails any one of the three verification passes outlined in the following subsections. The existence of an *alive* mutation indicates that the original spec is weaker than the existing implementation. IronSpec uses the current implementation and proof to establish *alive* mutations, so any spec weakness is directly related to the current system implementation; this undue weakness exposes a risk of containing a spec bug.

4.3.1 Mutation Generation

We generate mutations inspired by the method-level mutation operators from MuJava [41, 42] and from [20], see Table 2. We further introduce an additional predicate-based mutation operator, Set Operator Replacement (SOR). SOR introduces mutations about set inclusion, for example, an expression, $e \in s$, would be mutated to become, $e \notin s$.

Lemma 6 *IsAtLeastAsWeak* Lemma

```
1 lemma IsAtLeastAsWeak(p:Params)
2   requires OriginalPredicate(p:Params)
3   ensures MutatedPredicate(p:Params)
```

The IronSpec prototype is based in Dafny, so all mutations are applied to expressions in the Dafny AST. For Dafny expressions that reason about the heap, we introduce the Heap Operator Replacement (HOR) mutation operator, which mutates expressions containing the Dafny keyword `old`.

Each generated mutated spec is the result of IronSpec applying only a single mutation operator at a time. The set of all mutated specs consists of all possible single-operator mutations for a given spec applied to each subexpression in the mutation target.

4.3.2 First Pass: Logical Redundancy

Not all mutations produced from the original spec are relevant. Depending on the mutation target, some mutations no longer provide any new information to the tester. For example, if the mutation target is part of a post-condition, then any mutation that is *at least as weak* as the original spec, does not provide any new information to the tester; a weaker spec allows for all the same set of behaviors and then some. The opposite case holds if the mutation target is part of the precondition.

Definition 4.1. A predicate, S' , *IsAtLeastAsWeak* as predicate S with parameters p iff $\forall p. S(p) \implies S'(p)$

Based on the mutation target, IronSpec tests the *IsAtLeastAsWeak* property by automatically formulating Lemma 6, which will pass iff the mutated spec is equivalent or strictly weaker than the original spec. Any spec that is not *AtLeastAsWeak* is kept and considered as a potential valid mutation. Conversely, if the mutation target requires the consideration of weaker mutations, the *IsAtLeastAsStrong* property is checked. The lemma for *IsAtLeastAsStrong* is similar to Lemma 6, however, the *requires* and *ensures* are reversed.

When IronSpec checks either the *IsAtLeastAsWeak* or *IsAtLeastAsStrong* lemma, the final lemmas are generated regarding the root high-level safety property. This is done to avoid false positives if the mutation target is a sub-predicate of the high-level safety property. A mutation to a sub-predicate may not cause the mutation target to become weaker or stronger but it may affect a higher-level predicate that calls the mutation target to become weaker or stronger.

As an example of these definitions, consider Figure 1, where each circle represents the set of behaviors allowed by each respective spec. Any behavior in the circle encapsulated by the Original Spec is still inside the set of allowable behaviors of Spec A, making Spec A strictly weaker than the Original Spec and thus would pass the *isAtLeastAsStrong* pass. Both Specs B and C are not at least as weak as the Original Spec and would survive the *isAtLeastAsWeak* pass.

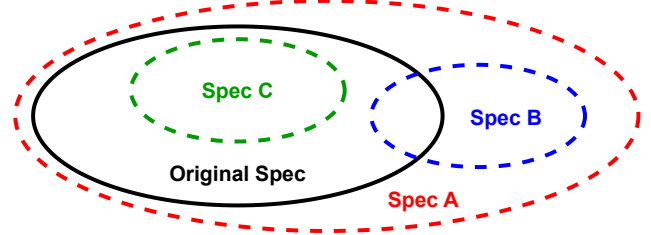


Figure 1: A mutated spec can either be strictly weaker (Spec A), strictly stronger (Spec C), logically equivalent, or partially stronger and weaker (Spec B) than the original spec.

4.3.3 Second Pass: Vacuity

IronSpec’s second pass aims to identify the mutations that cause vacuity [33]. A vacuous spec would allow for the system’s proof to pass trivially. Checking vacuity is more complicated than purely checking if the mutated predicate is itself vacuous, as the conditions of a caller in conjunction with the mutated predicate could result in an intermediary caller becoming vacuous. This is especially important with specs that are state machines where a mutation could cause a state transition to become *false*, removing that behavior from the spec. IronSpec automatically generates a lemma to check for vacuity by considering the full call path.

4.3.4 Third Pass: Full Proof

The final check is to see if the full proof will pass with the mutated spec. In this final pass, the system is re-verified with the addition of the mutated spec to ensure that no intermediary lemmas now fail. If the full proof passes, the mutation is considered *alive* and serves as a flag to the developer to re-examine the spec.

4.3.5 Hierarchical Classification of Alive Mutations

Rather than supplying the tester with a list of all alive mutations, IronSpec performs an additional pass to characterize the alive mutations, minimizing the output to the most relevant. To maximize the hint provided by an alive mutation, IronSpec evaluates the set of alive mutations to calculate a Direct-Acyclic Graph (DAG) indicating which mutations are weaker or stronger in relation to one another. The tester need only further concern themselves with the root of each connected component of this *mutation DAG*, as all children are weaker in each component. This hierarchical classification is inspired by previous research to classify and remove equivalent mutations [4, 20, 43, 47].

4.4 Using Alive Mutation As Hints For STPs

Human intuition is always the final oracle, thus, STPs are still needed to finish the investigation, but alive mutations

can be used as hints to write focused STPs. On their own, alive mutations only indicate a relative divergence between the spec and the implementation. The relative strength of an alive mutation can be used to shrink the state space necessary to test, focusing on the divergence between the two.

Armed with an alive mutation, a test writer can effectively exploit its hint by deviating from the standard guidelines of writing STPs and work **backward** from the spec difference. The spec difference is the set of behaviors allowed by the original spec, S , and not by the alive mutation, S' , ostensibly, $S - S'$. The spec difference embodies the fundamental insight IronSpec is based on; it captures a specific disconnect between the original spec and the implementation. The behavior allowed by this reduced expression is permitted by the original spec, but not by the more restrictive mutation. The spec difference uniquely presents the tester with this subset of behavior to determine if that particular disparity is intended.

Working backward allows the test writer to find concrete values that satisfy only the spec difference, achieving more concentrated STPs. This is useful because it removes the impetus of having the test writer identify values both valid in the original spec but also not valid in the mutated spec. The additional constraint of the spec difference increases the difficulty of manually identifying I/O that satisfies this more constrained space; the shift of working backward helps to alleviate this burden.

Driven by the insight that the actual semantic change between the mutation and the original spec is small, only a single mutation, the expression of the spec difference is minimal. Working backward allows the test writer to generate *any* I/O and check, using the verifier, that the I/O is accepted by the expression constituting the spec difference. After generating such values, the final decision relies on the tester to decide whether the I/O is intended. At this stage, the existence of an *unintended* value is a counterexample to the original spec.

5 Evaluation

We evaluate the effectiveness of the IronSpec prototype, by applying the Automated Sanity Checker, the STP Methodology, and the Automated Mutation Testing framework, to test *fourteen* different specifications written in Dafny [37]. Six of these specifications are produced in-house and include artificially introduced bugs, with an additional two specs containing artificial bugs described in [2]. Six of the specifications are of real-world open-source verified systems, which include: QBFT [44], DVT [51], TrueSat [3], Eth2.0 [11], daisy-nfsd [12], and an AWS Encryption SDK library [1].

When testing a spec, the ultimate oracle is the test writer, thus the final step is always to write a STP. When testing a spec, a tester could start with any aspect of IronSpec. In this section, we discuss the facets of IronSpec by highlighting their use in supplying the initial hint to discover of spec bugs.

Table 3: Spec bugs identified using the Automatic Sanity Checker. All bugs were confirmed with Counterexample STPs based on the initial hint of either MED or HIGH flags.

Bug	Specification	Method Name	Flag
TS1	TrueSat [3]	Formula Ctor	HIGH
TS2	TrueSat [3]	Start	MED
ETH1	Eth2.0 [11]	on_block	MED
AES1	AWS ESDK [1]	Encrypt	MED
AES2	AWS ESDK [1]	Decrypt	MED

5.1 Automatic Sanity Checking Evaluation

Applying the sanity checker to the six open-source verified systems, led to the discovery of five spec bugs across three specs, listed in Table 3.

Of the bugs identified, only TS1 was identified immediately with a HIGH severity flag, whereas the other four bugs were each discovered in less than an hour by writing STPs based on the hint of MED severity flags. The corresponding implementation for all five spec bugs appeared correct, but the specs were buggy, being too weak. To verify the existence of the spec weaknesses, we wrote buggy implementations as Counterexample STPs for each spec and demonstrated that the corresponding proof still passed.

Spec bugs TS2, ETH1, AES1, and AES2 were identified investigating the MED severity flag. We found that in these cases, the bug was a result of the output consisting of a complex datatype with many sub-fields and having postconditions concerning only a subset of these fields. This combination allows for a different implementation to update the remaining un-specified fields arbitrarily. A MED severity flag is not outright a spec bug because the un-specified fields may not be critical for safety, but if they are, it is a bug. The HIGH severity flag raised for TS1 was; "None of the postconditions depend on any of the input." This type of spec bug allows for a buggy implementation to completely ignore the input values when constructing the output.

The two bugs, AES1 and AES2, from the AWS Dafny Encryption SDK library (ESDK) [1] are both cases of spec weakness, confirmed by the authors. The Dafny ESDK is a verified SDK used as a reference to build ESDKs for other languages. These bugs exist for the high-level methods of Encrypt and Decrypt and are caused due to a combination of the postconditions under-constraining the output and because the postconditions of sub-methods are not exported. This underspecification allows for the proof of trivially incorrect implementations for Encrypt and Decrypt to pass.

Specs with output containing complex datatypes with many sub-fields are a critical point as a source of spec bugs. Judging from the results of the sanity checker this aspect can easily be overlooked. To avoid these types of spec bugs, it is vital to specify the values for all sub-fields of the output.

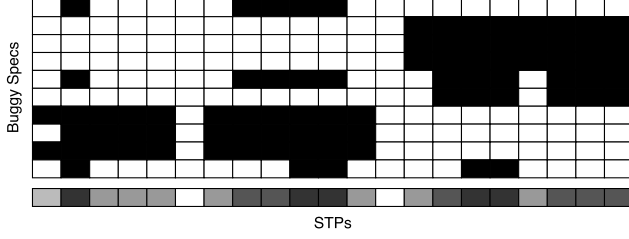


Figure 2: STP coverage for various buggy sort specs

5.2 STP Methodology Evaluation

In this section we describe our experience in writing Usefulness, Correctness, and Provability STPs for specs following the Hoare-Logic style, and in the cases of identifying a spec bug, counter-example STPs. We discuss the effectiveness of these STPs to expose differences in what behaviors the spec allows in contrast to a test writer’s expectations in the presence of artificially introduced bugs.

We also discuss a case study, where following the STP methodology we successfully discovered three spec bugs in a verified QBFT protocol. STPs were written for, and counterexamples derived, for all open source specs, but for the sake of brevity, the case study in this subsection focuses on the QBFT spec, where STPs acted as the initial hint of the existence of a spec bug.

5.2.1 STPs For Contrived Spec Bugs

We wrote STPs for five specs produced in-house that are in the Hoare-Logic style. These specs include methods for *Max*, *Sort*, *Binary Search*, *Token-with-revert-external* (*wre*), and *SimpleAuction-with-revert-external* (*wre*). The *Token* and *SimpleAuction* specs were modified from [10].

Even with the bias of constructing these examples ourselves, we wrote a comprehensive suite of STPs attempting to cover all behaviors. The bugs introduced into the spec vary but were comprised of an approximately equal split between too-strong and too-weak specs. A total of 70 STPs were written cumulative for all variations of the contrived specs. Only 30% of these STPs (21) needed additional proof help and of those, each STP needed, on average, an additional 1.7 lines of proof help. The STP suite was successful in all cases identifying introduced spec bugs—confirming the notion that a failing STP is a reliable flag suggesting a discrepancy between the intent encoded in the STPs and the spec.

Not all classes of STPs were useful in identifying all bugs because of the different natures of each type of STP. For example, consider Figure 2 that shows the coverage of 21 STPs for various bugs for a *Sort* spec. Each row corresponds to a different bug, each column matches a different STP, and a darkened cell is evidence of a specific STP identifying the bug. The bug was identified if a row contained a single darkened cell, whereas each column gives insight into the coverage of

Lemma 7 Simplified QBFT Provability Adversary STP

```

1 lemma AdversaryForwardMessageSTP (
2   a:Adversary,
3   a':Adversary,
4   inMsgs: set<Message>,
5   outMsgs: set<Message>)
6 requires validAdversaryConfig(a,a',inMsgs)
7 requires inMsgs == {ProposalMsg(CS1,block)}
8 requires outMsgs == {NewBlockMsg(CS1,block)}
9 ensures AdversaryNext(a, a', inMsgs, outMsgs)

```

a single STP in identifying different bugs. The bottom row is a heat map corresponding to the ratio of bugs identified by each STP. The results indicate that to take full advantage of the testing methodology and achieve adequate coverage to test both too-weak and too-strong specs, it is vital to write a comprehensive test suite comprising all three types of STPs.

5.2.2 STPs: QBFT Case Study

Writing STPs for the QBFT spec [50] led to the discovery of **three** spec bugs confirmed by the authors. The spec in this system consists of a single safety property *Blockchain Consistency* and the environment, which includes the high-level distributed system, the network, and the adversary which are all modeled as a state machine. Upon manual inspection of these STPs, we found that the adversary spec was *incomplete*, based on our understanding of what the adversary spec should be. The overall proof still passed even in the presence of these three bugs because they essentially cancel each other out; two making the spec weaker than it should be, and the other making the spec stronger.

The first inconsistency identified in the adversary spec was an example of the spec being too strong; limiting the actions of what an adversarial node *should* be able to do. The initial hint was provided by failing Provability STPs. The reason for writing Provability STPs was to initially answer if the adversary spec was too strong; which is answered by the general form of Provability STPs. An overly restricted adversary model would weaken and perhaps invalidate the guarantees of the overall proof. Following the guidelines for writing STPs in Section 4.2.2, negating the conclusion led to the discovery of a counterexample STP. Lemma 7 is a simplified example of such a failing Provability STP.

The failing simplified STP in Lemma 7 hints at the fact that the ability of the adversary to extract signed data structures is unnecessarily restricted. In the system model for QBFT, and other Byzantine fault-tolerant consensus protocols, a Byzantine node should be allowed to behave arbitrarily while not violating cryptographic assumptions. In this QBFT spec, adversaries are only able to extract and forward *CommitSeals* (CS) from a subset of received message types. The STP in Lemma 7 specifies the behavior of an adversary node receiving a *Proposal* message signed with a quorum

Lemma 8 Simplified QBFT Adversary Correctness STP

```
1 lemma AdversaryForgeMessageSTP (  
2   a:Adversary,  
3   a':Adversary,  
4   inMsgs: set<Message>,  
5   outMsgs: set<Message>)  
6 requires validAdversaryConfig(a,a',inMsgs)  
7 requires outMsgs == {ProposalMsg(CS2)}  
8           // forged msg  
9 ensures !AdversaryNext(a,a',inMsgs,outMsgs)
```

of CS1, and constructing and sending a *NewBlock* message containing the block and CS1 data structures copied from the Proposal message. The postcondition for this STP stipulates that this scenario constitutes a valid state transition from state *a* to state *a'*. After observing that this STP failed to immediately verify, negating the conclusion to, `!AdversaryNext(a, a', inMessages, outMessages)`, resulted in the proof for this transformed counterexample STP to pass.

While investigating the implications of the behavior in the failing STP, we modified the adversary spec, weakening it to allow an adversary to forward CS1 regardless of what message first contained it. After making this change, the full system's safety proof failed. To differentiate the proof failure from a now incomplete lemma, we constructed and proved a concrete counterexample resulting in a violation of the system's safety property, confirmed by the authors.

The second bug in the adversary spec is an example of the spec being too weak. This weakness is the reason why we can show a concrete counterexample to safety after addressing the first spec bug. The spec allows an adversarial node to send a Proposal message containing a block data structure with arbitrary values, including using the `CommitSeals` of honest nodes even if the adversary had not previously received such `CommitSeals` in previous messages. `CommitSeals` are only used to make a final decision of committing a block, but this weakness in the spec allows an adversary to propose new blocks containing `CommitSeals` as if from honest nodes. This lets an adversary send a message that appears to be signed by an honest node, violating the security assumptions made by the QBFT system model. The STP in Lemma 8 is a simplified version of the correctness STP used to discover this spec bug.

The third bug identified is related to the previous bug and is concerned with the underspecified spec of the function `getNewBlock()`. This function is empty-bodied and only contains the spec. Due to the underspecification of this function, a caller of this function, including an honest node can immediately send a message, such as a Proposal message, containing a full quorum of commit seals. If a buggy implementation is provided for this function, it too, could lead to a violation of the safety property.

5.2.3 STP Discussion

STPs enable fine-grain testing of specs and have been effective at helping to identify all ten spec bugs. By leveraging the insight that writing proofs for specific values is easier than a general proof, the manual effort required to write STPs remains minimal.

In the QBFT spec the presence of three spec bugs, two manifesting as a weakness in the spec and the other counteracting the first by overly restricting the adversary, makes manually or automatically identifying these bugs extremely difficult. Following the STP testing methodology, we efficiently, and without being experts in the system, identified these disconnects between what was written in the spec and our understanding of the intent of the spec.

5.3 Mutation Testing Evaluation

We applied IronSpec's automatic mutation testing to a set of six in-house specs, the two spec examples from [2], and the spec of six open-source verified codebases. The evaluation attempts to answer how prevalent are *alive* mutations in specs, and how useful are the provided hints in assisting in identifying any spec bugs.

All mutation testing experiments were performed on a cluster of 21 servers where each node was equipped with two Intel E5-2660 v2 10-core CPUs at 2.20 GHz and 256GB ECC Memory. In each experiment, one root node would create all mutations and send all subsequent verification requests in each stage of the mutation testing process to be processed in parallel at the other 20 nodes in the cluster using Dafny version 3.8.1. The results from running IronSpec can be found in Table 4 and are explained in the following subsections.

5.3.1 In-House Specifications

In addition to the five in-house specs introduced in Section 5.2.1, we applied mutation testing on a simple key-value store state machine spec. Identifying bugs in the in-house specs, served as a litmus test for the usefulness of mutations due to the biased nature of identifying self-introduced bugs.

The top half of Table 4 contains the experimental results of running mutation testing on the in-house specs. Each buggy spec was tested with a correct implementation (C) and an incorrect implementation (I). Mutations testing all in-house specs with a correct spec and a correct implementation resulted in no alive mutations.

Mutation testing identified relevant alive mutations, regardless of whether the implementation is correct. For all six incorrect specs, mutation testing resulted in helpful alive mutations except for Sort (C). All of the alive mutations were useful hints in manually identifying a weakness in the spec. The correct implementation for the sort spec was merge sort, and due to additional loop invariants in the implementation which added implicit strength to the spec.

Table 4: Results from running IronSpec’s automatic mutation testing. In-House, buggy, specs marked with "(C)" have a correct implementation, whereas "(I)" indicates an incorrect implementation. The Predicate Name is the specific mutation target within a spec. Spec LOC is the size of the mutation target, and Proof/Impl LOC is the size of the full end-to-end implementation and proof. Mutations are the total number of mutations generated, Alive mutations indicate the number of alive mutations after all three passes and hierarchy classifications.

	Specification	Predicate Name	Spec LOC	Proof/Impl LOC	# Mutations	# Alive Mutations	Time
In-House Specs	Max (C)	maxSpec	2	5	80	1	11.3s
	Max (I)			7		4	7.5s
	Sort (C)	sortSpec	1	55	50	0	4.5s
	Sort (I)			4		1	7.3s
	Binary Search (C)	searchSpec	4	31	170	1	10.4s
	Binary Search (I)			18		2	24.3s
	KV SM (C)	Query Op	4	187	37	7	21s
	KV SM (I)					7	28.8s
	Token-wre (C)	GInv	1	87	13	1	7.8s
	Token-wre (I)			91		1	7.8s
	SimpleAuction-wre (C)	GInv	9	181	187	3	15.25s
	SimpleAuction-wre (I)					3	15.5s
Open-Source Specs	Div	Div	3	14	50	3	3.5s
	NthHarmonic	NthHarmonic	1	4	11	2	3s
	QBFT	NetworkInit	3	15071	58	3	80 min
	QBFT	AdversaryNext	48		197	7	162 min
	QBFT	AdversaryInit	3		46	4	80 min
	Distributed Validator	AdversaryNext	23	24747	152	7	191 min
	daisy-nfsd	GETATTR	4	18	35	1	4.3 min
	daisy-nfsd	WRITE	7	54	119	3	4.6 min

5.3.2 Alive Mutations in Open Source Systems

The Div and NthHarmonic specs are simple buggy specs introduced by Abreu et.al. [2], where the authors proposed initial techniques to repair simple spec errors in Dafny. The alive mutations found for these specs coincide with the conclusions made by Abreu et.al. in demonstrating that these specs are buggy by being too weak.

QBFT Of the 58 generated mutants for the initial state for the network state machine spec, `NetworkInit`, three mutations remained *alive* as the roots of their respective components in the *mutation DAG*. Upon manual inspection of the surviving mutants, the spec differences all referenced an aspect of the Network’s state that was never mentioned elsewhere. Thus, any value for part of the state would be considered "safe". These mutations do not imply the existence of a bug, but neither are they strictly false positives; rather they are examples of spec bloat. These alive mutations should still serve as flags to the developer, forcing them to answer the question of whether this state is needed, and if so why are these parts of the state not referenced?

The alive mutations for the `AdversaryNext` and `AdversaryInit` predicates, both parts of the adversary state machine spec can be considered false positives. The alive

mutations were all *stronger* mutations, but it is always safe to restrict the actions allowed by an adversarial node. In some alive mutations, the mutation implied that the proof would still pass with no adversaries in the system, or only taking trivial actions. This observation led us to question and then to test with STPs if the adversary spec was initially more restricted than it should be, leading to the bugs discussed in Section 5.2.2.

DVT The Distributed Validator Technology Protocol (DVT) spec and proof [51] captures the behavior of an Ethereum Validator, where a group of nodes coordinates to perform the Ethereum validator duties. The DVT spec consists of the desired *non-slashable attestation* property, and the environment, with the latter defined as the high-level distributed system, an adversary, and the network. All aspects of the environment are modeled as state machine specs. The *non-slashable attestation* property ensures that the system avoids committing a slash-able offense and eventually produces valid attestations.

Applying mutation testing to the `AdversaryNext` predicate in the adversary spec resulted in seven *alive* mutations. One of the mutations was a false positive. Three of the mutations hinted towards a limitation of the messages allowed to be sent by an adversary, leading to a similar discovery as in the

first QBFT bug concerning the restricted nature of adversaries forwarding messages. The remaining three alive mutation roots were concerned with the creation of attestations. The weakness lies in the spec’s lack of specificity regarding the attestations an adversary can create. Armed with this observation, we show with a counter-example that this weakness could lead to a safety violation.

daisy-nfsd Applying the mutation framework to daisy-nfsd’s [12] top-level NFS API spec resulted in alive mutations in two different methods’ specs. These mutations hint at the same spec weakness; one that would allow for a different trivial implementation to always return an error. This bug was confirmed by the authors as a known issue in their spec.

5.3.3 Combining STPs With Mutation Testing Hints

For all alive mutations, writing STPs is always the final step in testing, but alive mutations can be used as hints, as described in Section 4.4. Consider the DVT spec bug from Section 5.3.2. The original predicate is non-trivial and consists of 22 lines including multiple quantified conjuncts. Working backward from the alive mutations and focusing only on the expression derived by the difference between the original spec and an alive mutation, resulted in shrinking the 22-line predicate into only a single conjunct, which becomes more tractable to write STPs for. The tradeoff of manual effort to calculate the simplified expression and writing STPs outweighs the effort needed to consider the entire spec.

5.3.4 Discussion

Mutation testing supplied the hints that led to the discovery of spec bugs in two verified codebases. These results exemplify the usefulness of adding automation to search for disconnects between the implementation and the spec. The insight of identifying tangible differences as potential areas of disconnected intent is a beneficial hint that can be leveraged to identify spec bugs, not just the ones spotlighted here. The results in Table 4 demonstrate that even with a small set of mutations, we were successful in identifying spec weaknesses.

The large increase in execution between the run time of different specs can be attested to the size of the varying sizes of the full system proof and the time that it takes to verify the entire proof with the mutations that survive the first two passes. The cost of running IronSpec on a large verified system is worth the execution time to debug a potential spec bug.

The IronSpec prototype takes the first steps to provide automation and structure to testing specs. The prototype is built to target Dafny specs, but the conceptual techniques are not tied to Dafny. So, while the evaluation demonstrated the effectiveness of IronSpec, the set of available verified systems in Dafny remains small. A broader repository of specs would introduce the opportunity to refine IronSpec across more variations of specs.

6 Related Work

Kemmerer [30] first identified the potential benefits of testing specifications. Since then, there have been several studies that have approached the abstract problem or aimed to test informal user requirements [9, 15, 32, 39, 40]. This line of work has since culminated in the study by Fonseca et al. [19].

Some works have begun to apply more structured approaches to increase reliability in formal methods. Kupferman [33] discussed the possible advantages of vacuity and coverage checks for temporal-logic model-checking tools. Inspired by vacuity testing, and the concept of *unit proofs* from [13], Priya et al. [49] performed a case study of some AWS verified libraries, uncovering some hidden bugs. Bernardi et al. [7] also identified formal specifications as a weak point in the verification process, and proposed to reuse specifications once correct, for smart contracts. Le Traon et al. [35] even discussed the notion of applying a mutation analysis to Eiffel contracts.

With verification becoming more commonplace and the discovery of spec bugs in verified systems, a few, mostly manual efforts have attempted to identify spec bugs. The 2022 Notional Finance bug found in verified code inspired Certora to investigate ways to introduce testing into the verification process [46]. Recently, Abreu et al. [2] proposed initial efforts in using the dynamic invariant inference tool *Daikon* [17] to aid in automatically repairing specifications. When faced with a failing proof, their prototype assumes that the implementation is correct, and uses the implementation to generate test cases for the spec. Any failing tests present an opportunity to attempt to fix the spec by suggesting strengthening or weakening modifications. Trusting the implementation to be correct lessens the guarantees of verification.

Testing and formal methods share a close relationship and a common goal. Often, rather than questioning specs, developers have relied on specs or other formal methods to assist in testing traditional software [7, 13, 24, 36]. Works concerned with the Oracle problem [5, 8] have often utilized specs thus. There has even been work to test verification tools [27].

7 Conclusion

The correctness of specifications is the rock upon which the entire edifice of formal verification is built. As formal verification becomes increasingly popular, it is imperative that the foundation be as solid as possible.

This work proposes IronSpec, a systematic framework of manual and automated approaches to aid developers in finding bugs in their specs. We show how IronSpec was used to identify a number of subtle bugs in the specs of open-source codebases, without requiring copious amounts of expertise on the proven system. We believe IronSpec is a necessary step forward towards writing correct software.

References

- [1] Aws encryption sdk for dafny. <https://github.com/aws/aws-encryption-sdk-dafny>, 2023.
- [2] A. Abreu, N. Macedo, and A. Mendes. Exploring automatic specification repair in dafny programs. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW)*, pages 105–112. IEEE, 2023.
- [3] C.-C. Andrici and Ș. Ciobăcă. Verifying the dpll algorithm in dafny. *arXiv preprint arXiv:1909.01743*, 2019.
- [4] D. Baldwin and F. Sayward. *Heuristics for determining equivalence of program mutations*. Yale University, Department of Computer Science, 1979.
- [5] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo. The oracle problem in software testing: A survey. *IEEE transactions on software engineering*, 41(5):507–525, 2014.
- [6] K. Beck. *Test driven development: By example*. Addison-Wesley Professional, 2022.
- [7] T. Bernardi, N. Dor, A. Fedotov, S. Grossman, N. Immerman, D. Jackson, A. Nutz, L. Oppenheim, O. Pistiner, N. Rinetzky, et al. Wip: Finding bugs automatically in smart contracts with parameterized invariants. *Retrieved July, 14:2020*, 2020.
- [8] M. Böhme, C. Cadar, and A. Roychoudhury. Fuzzing: Challenges and reflections. *IEEE Softw.*, 38(3):79–86, 2021.
- [9] M. Brockmeyer. Using modechart modules for testing formal specifications. In *Proceedings 4th IEEE International Symposium on High-Assurance Systems Engineering*, pages 20–26. IEEE, 1999.
- [10] F. Cassez, J. Fuller, and H. M. A. Quiles. Deductive verification of smart contracts with dafny. In *International Conference on Formal Methods for Industrial Critical Systems*, pages 50–66. Springer, 2022.
- [11] F. Cassez, J. Fuller, and R. Saltini. eth2.0-dafny. <https://github.com/Consensys/eth2.0-dafny/tree/master>, 2021.
- [12] T. Chajed, J. Tassarotti, M. Theng, M. F. Kaashoek, and N. Zeldovich. Verifying the {DaisyNFS} concurrent and crash-safe file system with sequential reasoning. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pages 447–463, 2022.
- [13] N. Chong, B. Cook, K. Kallas, K. Khazem, F. R. Monteiro, D. Schwartz-Narbonne, S. Tasiran, M. Tautschnig, and M. R. Tuttle. Code-level model checking in the software development workflow. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering: Software Engineering in Practice*, pages 11–20, 2020.
- [14] E. Daka and G. Fraser. A survey on unit testing practices and problems. In *2014 IEEE 25th International Symposium on Software Reliability Engineering*, pages 201–211. IEEE, 2014.
- [15] G. De Caso, V. Braberman, D. Garbervetsky, and S. Uchitel. Automated abstractions for contract validation. *IEEE Transactions on Software Engineering*, 38(1):141–162, 2010.
- [16] R. A. DeMillo, R. J. Lipton, and F. G. Sayward. Hints on test data selection: Help for the practicing programmer. *Computer*, 11(4):34–41, 1978.
- [17] M. D. Ernst, J. H. Perkins, P. J. Guo, S. McCamant, C. Pacheco, M. S. Tschantz, and C. Xiao. The daikon system for dynamic detection of likely invariants. *Science of computer programming*, 69(1-3):35–45, 2007.
- [18] R. W. Floyd. Assigning meanings to programs. *Program Verification: Fundamental Issues in Computer Science*, pages 65–81, 1993.
- [19] P. Fonseca, K. Zhang, X. Wang, and A. Krishnamurthy. An empirical study on the correctness of formally verified distributed systems. In *Proceedings of the Twelfth European Conference on Computer Systems*, pages 328–343, 2017.
- [20] R. Gheyi, M. Ribeiro, B. Souza, M. Guimarães, L. Fernandes, M. d’Amorim, V. Alves, L. Teixeira, and B. Fonseca. Identifying method-level mutation subsumption relations using z3. *Information and Software Technology*, 132:106496, 2021.
- [21] P. Hamill. *Unit test frameworks: tools for high-quality software development*. " O’Reilly Media, Inc.", 2004.
- [22] T. Hance, A. Lattuada, C. Hawblitzel, J. Howell, R. Johnson, and B. Parno. Storage systems are distributed systems (so verify them that way!). In *Proceedings of the 14th USENIX Conference on Operating Systems Design and Implementation*, pages 99–115, 2020.
- [23] C. Hawblitzel, J. Howell, M. Kapritsos, J. R. Lorch, B. Parno, M. L. Roberts, S. Setty, and B. Zill. Ironfleet: proving practical distributed systems correct. In *Proceedings of the 25th Symposium on Operating Systems Principles*, pages 1–17. ACM, 2015.
- [24] R. M. Hierons, K. Bogdanov, J. P. Bowen, R. Cleaveland, J. Derrick, J. Dick, M. Gheorghe, M. Harman, K. Kapoor,

- P. Krause, et al. Using formal specifications to support testing. *ACM Computing Surveys (CSUR)*, 41(2):1–76, 2009.
- [25] C. A. R. Hoare. An axiomatic basis for computer programming. *Communications of the ACM*, 12(10):576–580, 1969.
- [26] S. Igarashi, R. L. London, and D. C. Luckham. Automatic program verification i: A logical basis and its implementation. *Acta Informatica*, 4(2):145–182, 1975.
- [27] A. Irfan, S. Porncharoenwase, Z. Rakamarić, N. Rungta, and E. Torlak. Testing dafny (experience paper). In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 556–567, 2022.
- [28] Y. Jia and M. Harman. An analysis and survey of the development of mutation testing. *IEEE transactions on software engineering*, 37(5):649–678, 2010.
- [29] J. Kang, Y. Kim, C.-K. Hur, D. Dreyer, and V. Vafeiadis. Lightweight verification of separate compilation. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 178–190, 2016.
- [30] R. A. Kemmerer. Testing formal specifications to detect design errors. *IEEE transactions on software engineering*, (1):32–43, 1985.
- [31] U. Kirstein. Post-mortem analysis of the notional finance vulnerability — a tautological invariant, January 2022. [Online; posted 17-January-2022].
- [32] A. Knüppel, L. Schaer, and I. Schaefer. How much specification is enough? mutation analysis for software contracts. In *2021 IEEE/ACM 9th International Conference on Formal Methods in Software Engineering (FormalISE)*, pages 42–53. IEEE, 2021.
- [33] O. Kupferman. Sanity checks in formal verification. In *International Conference on Concurrency Theory*, pages 37–51. Springer, 2006.
- [34] L. Lamport. Specifying systems: the tla+ language and tools for hardware and software engineers. 2002.
- [35] Y. Le Traon, B. Baudry, and J.-M. Jézéquel. Design by contract to improve software vigilance. *IEEE Transactions on Software Engineering*, 32(8):571–586, 2006.
- [36] O. Legunsen, W. U. Hassan, X. Xu, G. Roşu, and D. Marinov. How good are the specs? a study of the bug-finding effectiveness of existing java api specifications. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, pages 602–613, 2016.
- [37] K. R. M. Leino. Dafny: An automatic program verifier for functional correctness. In *Proceedings of the 16th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning, LPAR’10*, pages 348–370, Berlin, Heidelberg, 2010. Springer-Verlag.
- [38] M. Lesani, C. J. Bell, and A. Chlipala. Chapar: certified causally consistent distributed key-value stores. *ACM SIGPLAN Notices*, 51(1):357–370, 2016.
- [39] M. Li and S. Liu. Reviewing formal specification for validation using animation and trace links. In *2014 21st Asia-Pacific Software Engineering Conference*, volume 1, pages 263–270. IEEE, 2014.
- [40] S. Liu, J. A. McDermid, and Y. Chen. A rigorous method for inspection of model-based formal specifications. *IEEE Transactions on Reliability*, 59(4):667–684, 2010.
- [41] Y.-S. Ma and J. Offutt. Description of mujava’s method-level mutation operators. *Update*, 2016.
- [42] Y.-S. Ma, J. Offutt, and Y. R. Kwon. Mujava: an automated class mutation system. *Software Testing, Verification and Reliability*, 15(2):97–133, 2005.
- [43] L. Madeyski, W. Orzeszyna, R. Torkar, and M. Jozala. Overcoming the equivalent mutant problem: A systematic literature review and a comparative experiment of second order mutation. *IEEE Transactions on Software Engineering*, 40(1):23–42, 2013.
- [44] H. Moniz. The istanbul bft consensus algorithm. *arXiv preprint arXiv:2002.03613*, 2020.
- [45] G. J. Myers, C. Sandler, and T. Badgett. *The art of software testing*. John Wiley & Sons, 2011.
- [46] S. Phipathananunth. Using mutations to analyze formal specifications. In *Companion Proceedings of the 2022 ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications: Software for Humanity*, pages 81–83, 2022.
- [47] A. V. Pizzoleto, F. C. Ferrari, J. Offutt, L. Fernandes, and M. Ribeiro. A systematic literature review of techniques and metrics to reduce the cost of mutation testing. *Journal of Systems and Software*, 157:110388, 2019.
- [48] V. R. Pratt. Semantical considerations on floyd-hoare logic. In *17th Annual Symposium on Foundations of Computer Science (sfcs 1976)*, pages 109–121. IEEE, 1976.
- [49] S. Priya, X. Zhou, Y. Su, Y. Vizel, Y. Bao, and A. Gurfinkel. Verifying verified code. *Innovations in Systems and Software Engineering*, 18(3):335–346, 2022.

- [50] R. Saltini. Qbft formal specification and verification. <https://github.com/Consensys/qbft-formal-spec-and-verification>, 2021.
- [51] R. Saltini. Formal verification of the distributed validator technology protocol. <https://github.com/Consensys/distributed-validator-formal-specs-and-verification>, 2023.
- [52] J. R. Wilcox, D. Woos, P. Panchekha, Z. Tatlock, X. Wang, M. D. Ernst, and T. Anderson. Verdi: a framework for implementing and formally verifying distributed systems. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 357–368, 2015.
- [53] X. Yang, Y. Chen, E. Eide, and J. Regehr. Finding and understanding bugs in c compilers. In *Proceedings of the 32nd ACM SIGPLAN conference on Programming language design and implementation*, pages 283–294, 2011.