

FMitF: Track I: A Pathway for Combining Formal, Static, and Dynamic Analysis of Real-World Embedded Systems

1 Problem Statement, Overview, and Objectives

The core problem we aim to address in this proposal is that *use of formal modeling, advanced static analysis, and advanced dynamic analysis* for verification and validation, especially on critical embedded systems, is prohibitively difficult and lacks sufficient synergy for cost-effective application. This is true even of systems built in an academic research context: that is, unless the research is primarily *about* methods for verifying and testing systems, rather than work on an embedded system for its own sake, these methods are hard to apply. Furthermore, even when use of these techniques to ensure correctness, reliability, or security *is* a focus of the project, such use is almost always limited to one type of effort—model checking, theorem proving, or automated test generation. A major cause for this difficulty is the *lack of synergy* between these related efforts, the failure of effort in one context to transfer to another context. In short:

- Learning to use a formal modeling language and tool provides help in discovering defects in a high-level model or protocol, but seldom helps with implementation-level bugs not thus modeled.
- Many static analysis tools are primarily “bug detectors” (e.g., Coverity or CodeSonar), whose output is essentially limited to a list of possible problems; devising test inputs to reject false positives is hard.
- More powerful static analysis tools, such as FRAMA-C [47], provide proofs of correctness for limited aspects of a system, and a rich specification and annotation language. However, there is little or no connection between this annotation and either formal modeling or state-of-the-art test generation.
- There are a large variety of automated test generation tools; however, effort spent learning one tool only partially transfers to another tool. Many tools (e.g., AFL [79]) focus on finding *crashes*, and do not leverage other types of specification.

Consider the case of an embedded systems engineer working on a custom, low-energy consumption, communication protocol for use in a network of low-power sensors and actuators. If she builds a formal model of the protocol, she will likely find that this extensive effort provides no help, other than an improved concept of the system, in proving the correctness of her *implementation*. If the engineer begins instead by building an automated test generation harness she will find that, despite having spent considerable time expressing pre-conditions and post-conditions for various functions in the implementation, the work must be duplicated when she decides to try to formally prove the correctness of core functionality. Had she begun with proofs, again, logically related (or even equivalent) information would have to be re-expressed, in a different language, to perform test generation. Effort spent in using one tool almost never carries over to another approach. There is simply not enough time or energy available to make use of the available technology. In practice, *no advanced correctness technology may be used at all*. Since it is hard to predict which one(s) will have the greatest payoff, or even work at all, so perhaps it is best to just focus on manual testing.

Proposed Solution: While allowing efforts from any form of formal or automated verification or validation attempt to carry over to other forms (e.g., formal models to code annotations for static analysis, code annotations to test harnesses, test harnesses to formal models, formal models to test harnesses, etc.) is the ideal goal, simply making it possible to follow *one* critical path to combine methods is feasible and desirable.

Which path is most important to realize? Our approach is based in the reality of the embedded systems domain, where, while formal modeling is sometimes used, there is, in real-world efforts, *always* an implementation. The most basic obstacle to the adoption of formal methods in embedded systems is that if there is only the usual informal design effort or the adaptation of a legacy implementation, formal methods are often inapplicable. By focusing on *adding annotations to implementation code*, and exploiting those annotations to enable analyses, we promise to *always* give embedded systems engineers a reasonable payoff.

This project proposes to make it possible to introduce specifications into implementation code that can be directly checked using sophisticated automated test generation strategies, including symbolic execution, fuzzing, and model checking. Furthermore, these specifications can be directly exported to form the basis for formal models using, e.g., timed automata. In the long run, to benefit those developers who are more

open to formal methods already, we hope that these annotations can also be *imported* from a timed automata representation, but we begin where most embedded systems developers are, now, not where we hope they may be, someday. We additionally focus on using frameworks/front-ends allowing application of multiple approaches. Our commitment is to enabling a *maximum diversity of analysis methods* with a *minimum of specification and tool-learning effort*, to make formal methods attractive in terms of cost-benefit ratio.

This project is specifically focused on embedded and networked systems but we anticipate that our solution will generalize to other application domains with similar characteristics and challenges. We expect developers to learn new tools, but not new programming paradigms or languages. The proposed contribution to embedded systems design is not a radical reworking of development methods, which, like many formal methods efforts in the past, would be unlikely to achieve widespread adoption, but the introduction of an *advanced form of unit testing*, that works with, e.g., legacy code, with more powerful methods for specification and checking of correctness. This will modify development, in that design-for-testability and design-for-verifiability will become second nature. This project is therefore based on the following core ideas:

- 1. The primary obstacle to adoption of formal methods approaches in embedded systems development is not a lack of relevant methods and tools.**
- 2. In particular, there are methods and tools that apply to the *implementation* of embedded systems in C and C++; every embedded software system requires an implementation.**
- 3. However, learning and using any one of these tools may or may not “pay off” and the effort spent is only of limited application when applying another tool.**
- 4. Therefore, to improve embedded systems development via formal methods we need:**
 - (a) an *implementation-focused common framework* for applying methods and tools and**
 - (b) a focus on *practically-inspired* improvements to methods and tools.**

Field Application Domains: The proposed research is motivated by applications in embedded and networked systems, in particular wireless sensor networks and multi-robot systems. A wireless sensor network (WSN) consists of multiple, often many, sensor nodes that communicate via a wireless network with servers and/or other sensor nodes. In a WSN, each sensor monitors some physical quantities, such as air temperature, humidity, and CO₂ concentration level, and exchanges its measurement data with other entities for specific objectives, such as for monitoring and forecasting wildfires. Each sensor node typically has some limited computing power (an embedded microprocessor) and limited energy provided by a battery and/or a renewable energy source. Furthermore, wireless communication in a WSN is often subject to frequent packet drops and other failures. Consequently, some of the major challenges for the reliability and correctness of WSNs include timing and communication uncertainty, and computing and energy constraints. In addition, testing and verifying code for sensor nodes faces another significant hurdle due to the nature of interacting directly with the physical world, which involves another level of uncertainty and is not always easily replicable in a software testing or verification tool. A multi-robot system (MRS) consists of multiple robots, possibly of different types such as terrestrial robots and aerial drones, which coordinate to perform certain tasks, such as package delivery, disaster rescue, and surveillance. An MRS carries similar major challenges as a WSN. Due to the complexity of their dynamics and interactions with the physical world, testing and verifying software for robots, especially cooperative robots in a MRS, is particularly challenging.

In summary, this project focuses on the following major challenges of the field application domains of embedded and networked systems: (1) timing and communication uncertainty, (2) complex dynamic behaviors, and (3) interactions with the physical world. They will be addressed by a combination of approaches, from both the formal method and software testing disciplines and the field application disciplines. Case studies in real-world WSNs and MRS, which are representative of embedded and networked systems and their challenges, will be used for validating and demonstrating our proposed solution.

1.1 PI Qualifications

See the collaboration plan for an extensive examination of PI Qualifications; in brief, PI Groce has extensive history with formal methods and testing tool development, and practical application to real-world embedded

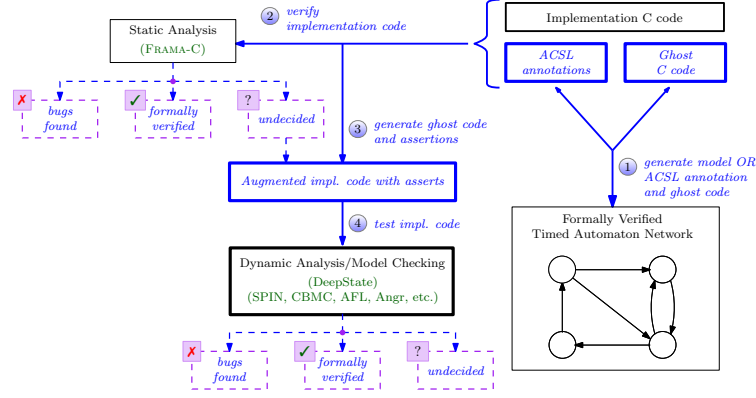


Figure 1: Overview of the proposed research.

systems. Co-PI Nghiem is an expert in control and autonomy for robotics, including use of formal methods, and Co-PI Flikkema has extensive experience with deployed real-world embedded systems and networks.

1.2 Intellectual Merit

The aim of this proposal is to (1) identify a set of principles for the analysis (formal, static, and dynamic) of implementations of embedded and networked systems; (2) match these theoretical principles with tools usable by engineers developing such systems; and (3) enable the synergistic use of enhanced versions of these tools in real applications through a common framework with minimal duplication of effort and maximal extraction of information from shared annotations. In the first case study, we will analyze networks of wireless sensor nodes deployed in the Southwest Experimental Garden Array (SEGA) [77, 25] – a distributed facility for examining climatic, genetic, and environmental factors in plant ecology – and in the Distributed Sensing & Computing Over Sparse Environments (DISCOVER) Platform – an NSF CCRI project that develops a large-scale and diverse testbed for wireless sensor networks and multi-robot systems. The second case study will use the framework to formally verify and dynamically test implementations of autonomous control and coordination of multiple autonomous terrestrial and aerial robots on the DISCOVER platform.

Figure 1 shows the overall concept. The core open research problems addressed are represented by two sets of arrows. First, an engineering design, expressed as C code, is provided, and annotated with correctness properties, information about the expected environment (constraints on sensor values, etc.), and hints to guide heuristic application of tools ranging from fuzzers to symbolic execution engines to model checkers. While they are not the focus of this project, code to apply advanced static analysis or timed automata (TA) model skeletons can also be automatically generated:

1. A (generated) TA model can be used to check high-level properties of the design, ignoring many low-level implementation details. However, this step is often skipped in practice.
2. The implementation code with the ACSL annotations and ghost code can be checked by a static analysis tool, such as FRAMA-C. In some cases this will verify the code, and in other cases a definite bug will be found; but often the result will be “undecided” and further analysis required to see if a bug is spurious or real. Again, this step can be skipped, though it is likely low-cost and beneficial.
3. Finally, the focus of our efforts is a multi-pronged attempt to refute correctness (or increase our confidence in it) via dynamic analysis—automated test generation—and implementation-level model checking. Ghost code, additional assertions, and needed test-harnesses are *automatically generated* from annotated code.
4. The augmented implementation code is then analyzed using the DeepState [29] framework, which serves as a front-end to highly scalable fuzzers, as well as to symbolic execution tools and model checkers.

Our focus is on providing a unified specification method that can be applied to source code itself, and on enabling and enhancing the dynamic analysis and model checking aspects of the approach. We believe these approaches are currently difficult to apply, but likely to dramatically improve the ability to detect bugs.

Principles: Assuming that an algorithm for an embedded and networked system, such as a communication protocol, can in theory be described as (probabilistic) timed automata [3], which satisfies temporal logic formulas [13], and implemented as a set of imperative programs, we ask:

- Given a set of annotated programs, how can we best automatically find bugs in those programs (and, in some circumstances, for some properties, prove correctness), based on an annotated specification?
- Can we generate a skeleton of a timed automata model from annotated programs, in order to facilitate adoption of design-level analysis by traditional embedded systems developers?

Note that these problems differ considerably from the more studied, but more limited, synthesis problem. We are not assuming that system development will involve first producing a formal model, then using that model to automatically generate an implementation; rather, we consider the typical real-world scenario, where modeling is a separate activity, either undertaken after implementation due to concerns about reliability, or an activity during design that only indirectly informs the implementation. That is, the more studied problem is producing a runtime semantics for a model; we address the problem of reconciling a runtime semantics with a model semantics, without unrealistic burden on engineers. Any effort to increase the adoption of formal methods and automated test generation approaches is likely to be successful only to the extent that it enters embedded systems engineering practice via this existing pathway. Engineers at Google have referred to this approach as “meeting developers where they are” [61].

Tools: We will focus on C code, using DeepState [29] as a front-end for dynamic and implementation-level model checking approaches, and UPPAAL [8] and PRISM [50] for the analysis of protocols; FRAMA-C will provide a powerful static analysis framework, and we will adopt its ACSL language developed for FRAMA-C as a basis for our specification language. The primary open research questions here are numerous, and include: (1) how to extend existing specification languages to support timing, interrupts, and uncertainty; (2) how to assign the same meaning to a specification construct in various contexts, ranging from fuzzing to symbolic execution to explicit-state model checking to bounded model checking in the “dynamic” DeepState world, and including a static context for FRAMA-C and a modeling context for timed automata; (3) how to minimize annotation burden while allowing developers to include information that can be exploited by those methods: e.g., to make intelligent use of pre-conditions in fuzzing, to automatically derive loop bounds in bounded model checking, and to restrict branching factors and store state in explicit-state model checking; (4) how to handle intra-program parallelism; (5) how to ensure that the methods are sufficiently automatic and behave in ways engineers will expect. Our focus will be on *practical* solutions, guided by domain experts, rather than on purely theoretical approaches that do not scale. Practical solutions here require fundamental contributions to system and specification design and dynamic analysis and model checking technologies.

Field applications: This project will contribute significantly to the field domains of embedded and networked systems, especially in safety-critical applications where correctness and reliability are vital. The methods and tools developed by our team will enable domain experts without formal training in formal methods and software testing, to test and verify their code to not only reduce bugs but also gain confidence in their systems.

2 Background and Preliminary Research

A Foundation for Implementation-Level Specification: Ideally, the development of critical embedded systems should rely on a combination of formal methods to achieve an appropriate degree of guarantee: automatic static analysis to ensure the absence of some runtime errors, deductive verification to prove functional correctness of aspects of the code, and runtime verification for parts of the code that cannot be (or are not yet) proved using deductive verification, or that generated *warnings* from static analysis requiring confirmation of a problem or refutation of the feasibility of an error. A core question is how to represent a specification that is intimately tied to real implementation code.

This project will therefore make use of the ideas developed in the ongoing work on the FRAMA-C (<https://frama-c.com>) [47] tool as a foundation for a specification and annotation language. FRAMA-C is a widely-used source code analysis platform that aims at enabling verification of industrial-size programs,

supports combinations of different approaches, by providing its users with a collection of *plugins* for analyses. Moreover, collaborative verification across cooperating plugins is enabled by their integration on top of a shared kernel, and, most critically for our purposes, their compliance to a common specification language: ACSL [6]. ACSL, the ANSI/ISO C Specification Language, is based on the notion of a contract as in JML [52]: ACSL allows users to specify functional properties of programs through pre/post-conditions. Many built-in predicates and logic functions are provided, to handle, for example, pointer validity and separation. Using ACSL/FRAMA-C means that while focusing on dynamic analysis and model checking, our approach automatically provides an additional benefit for embedded systems engineers: access to the current set of powerful FRAMA-C static analyses. FRAMA-C provides both abstract interpretation [23] based analysis plugins and deductive verification plugins based on a weakest precondition calculus; the latter have recently been improved to make proof without interacting with a theorem prover easier for engineers [9].

FRAMA-C was designed as a static analysis platform, but has been extended with limited plugins for dynamic analysis. One of these plugins is E-ACSL, which supports runtime assertion checking [20]. In FRAMA-C, E-ACSL is both the name of the assertion language and the name of a plugin that generates C code to check these assertions at runtime. E-ACSL is a subset of ACSL. The plugin E-ACSL is used to translate a subset of FRAMA-C assertions into executable C code. However, the E-ACSL plugin does not support all the specification constructs we need, or assist developers in the most difficult part of dynamic analysis: constructing a set of tests that exercise the checks. The only assistance provided by FRAMA-C for this is very limited in capability. Rather than “re-inventing the wheel” and offering a solution that lacks a strong static aspect we therefore extend ACSL and E-ACSL.

Dynamic Analysis with DeepState: While FRAMA-C provides powerful tools for static detection of program faults and generation of runtime checks for properties that cannot be discharged by formal proof or sound static analysis, it provides only limited, and difficult-to-scale, ability to generate program inputs to exercise runtime checks, limited to one tool, PathCrawler [75], that aims to produce a unit test for a single function, using concolic testing (dynamic symbolic execution [28]). In cases where this fails to scale, PathCrawler will fail. Furthermore, PathCrawler is tuned to the problem of testing a single function, not producing more complex scenario-based tests of a set of functions that must coordinate state changes. Finally, PathCrawler is not an open source, extensible system, may be costly to acquire and use, and is arguably impossible to extend.

The limitation of dynamic analysis tools to PathCrawler is a major weakness of FRAMA-C from the perspective of a user. Scalability of symbolic-execution-based test generation methods is extremely difficult to predict, and producing complete and exhaustive preconditions that allow a function to be tested entirely in isolation is often either too time-consuming or essentially impossible, because the actual environment is only represented by the set of states reachable using a set of coordinating functions or a library. These problems are pressing, for several reasons. First, full formal proof of correctness is, at present, impractical for most realistic systems. The actual work of fault detection and validation of software still relies, fundamentally, on effective testing. Moreover, modeling and even static approaches often must rest on a basis of numerous un-examined assumptions about the behavior of hardware systems and low-level system behavior. Only actual concrete inputs can be executed on real hardware, and satisfy regulatory requirements on code coverage such as those imposed on civilian avionics by DO-178B, etc. [63]. Furthermore, only testing can prove faults are not spurious, the result of imprecise abstraction or weak assumptions.

Most developers do not know how to use symbolic execution tools; developers seldom even know how to use less challenging tools such as gray-box fuzzers, even relatively push-button ones such as AFL [79]. Even those developers whose primary focus is critical security infrastructure such as OpenSSL are often not users, much less expert users, of such tools. Furthermore, different tools find different faults, have different scalability limitations, and even have different show-stopping bugs that prevent them from being applied to specific testing problems. DeepState [29, 30] addresses these problems. First, developers *do*, usually, know how to use unit testing frameworks, such as JUnit [26] or Google Test [1]. DeepState makes it possible to write parameterized unit tests [72] in a GoogleTest-like framework, and automatically produce tests using symbolic execution tools [66, 68, 65, 55], or fuzzers like AFL [79] or libFuzzer [64] (as well as Eclipse [17],

Angora [16], and Honggfuzz [2]). DeepState targets the same space as property-based testing tools such as QuickCheck [18], ScalaCheck [57], Hypothesis [54], and TSTL [34, 43], but for C/C++ unit tests. DeepState is, most importantly, the first tool to provide a front-end that can make use of a growing variety of back-ends for test generation. Developers who write tests using DeepState can expect that DeepState will let them, without rewriting their tests, make use of new symbolic execution or fuzzing advances. The harness/test definition remains the same, but the method(s) used to generate tests may change over time. Most property-based tools only provide random testing, and symbolic execution tools such as Pex [71, 73] or KLEE [15] offer only a single back-end. DeepState has already been used to test (and find bugs in) an ext3-like file system [69, 31] and a widely used compression library (<https://github.com/Blosc/c-blosc2/issues/95>), and is being considered as a basis for automatic testing for in NASA’s open source flight software framework FPrime [11, 56]. Although only released in early 2018, DeepState is already one of the most popular property-based testing and fuzzing projects on GitHub, and has been used internally by both startups and well-established companies, and in security audits by Trail of Bits. There have even been informal discussions of integrating DeepState, once matured, into a future release of the GoogleTest [1] platform. PI Groce is at present the lead developer for DeepState.

3 Research Plan

3.1 DeepState and Automated Test Generation

Applying DeepState to real embedded systems requires us to meet many challenges:

1. The *specification* of correctness must be translated into an executable form. To some extent, the existence of the E-ACSL executable subset of ACSL, and libraries for runtime checking of properties satisfies this condition. DeepState can support any C/C++ executable method of checking for correctness. However, some executable specifications need to be modified to be efficiently handled when the DeepState back-end is a symbolic execution tool. DeepState’s nature as a test generation tool means that it supports constructs, such as Minimum, Maximum, and Pump, not usually available in executable specifications. Tailoring E-ACSL usage for DeepState therefore requires a custom effort, including extending the semantics of executable specifications and optimizing the implementation for symbolic execution and fuzzing. Finally, because our domain critically involves timing, we need to implement DeepState handling of (and E-ACSL representations for) deadlines, and specification of function-level deadlines including arbitrary, specified, “runtimes” for code that operates via simulation rather than real hardware (or in symbolic execution). Similar, but in some ways even more complex, challenges are posed by the ubiquity of *interrupts* in embedded code, a problem addressed by very little previous work in fuzzing [67].
2. The *assumptions* that control which tests are considered valid must be translated in the same way; normally, E-ACSL simply translates these into further assertions (as pre-conditions to check at runtime), but in DeepState, we need to distinguish between ASSUME failures (invalid tests) and ASSERT failures (bugs).
3. The inputs to a function must be translated into code controlling the input values that DeepState provides, including ranges and types. When input types are simple, this process is straightforward; however, when functions take, e.g., arbitrarily sized arrays, linked lists, or other complex structures, this becomes a problem of constructing a test harness that (1) makes fuzzing and symbolic execution scalable but (2) uses large enough structures to expose subtle bugs. Moreover, because DeepState supports strategies for input generation, such as forking concrete states for values too complex for symbolic execution using the Pump construct, the translation must determine when such strategies are appropriate, and apply them.
4. In many cases, checking a single function may not be an effective way to detect faults; only a sequence of API calls can expose a problem in a system (e.g., that a function produce a state that causes another function to violate an invariant). ACSL annotations provide enough information for a fully-automated translation to a harness enabling dynamic analysis in the case of proving properties of a single function, but not for groups of functions. Moreover, even in cases where the violation of a specification can, in theory, be discovered without calling multiple functions, the state space may be too large to explore with a fuzzer or symbolic execution tool. In such cases, exploring only states produced by valid call sequences has two

```

void update_state(struct state_t *s, uint64_t bv) {
    ASSUME(valid_state(s));
    ASSUME(valid_bv(bv));
    ...
}
void process_both_sensor_readings(struct state_t *s) {
    ASSUME(valid_state(s));
    unit64_t s1_bv = acquire_s1(), s2_bv = acquire_s2();
    update_state(s, s1_bv); update_state(s, s2_bv);
}
void process_one_sensor_reading(struct state_t *s) {
    ASSUME(valid_state(s));
    unit64_t s1_bv = acquire_s1();
    update_state(s, s1_bv);
}

struct state_t *NewState() {
    return DeepState_Malloc(sizeof(struct state_t));
}
TEST(SensorReading, UpdateNeverSlow) {
    struct state_t *s = NewState();
    DeepState_Timeout(
        [&]{update_state(s, DeepState_UInt64());},
        MAX_EXPECTED_UPDATE_TIME);
}
TEST(SensorReading, AvoidCrashes) {
    struct state_t *s = NewState();
    for(int i = 0; i < TEST_LENGTH; i++) {
        OneOf(
            [&]{process_both_sensor_readings(s);},
            [&]{process_one_sensor_reading(s);});
    }
}

```

Figure 2: Sensor reading code and DeepState test harness

benefits: first, the space itself may be much smaller, and easier to explore, than the full set of possible input values. Second, errors in this part of the input space are more important. Even if a precondition is not sufficiently restrictive to guarantee correct behavior, if the “bad” inputs are never, in practice, generated by the functions that modify system state, the fault may not matter. In cases where constructing a sufficiently exact precondition is difficult for engineers, such “in-use” verification may be the only avenue to system assurance. We propose to let users annotate *sets* of functions to be tested as an API-call-sequence group, extending recent work exploring this concept [10, 62].

5. Finally, DeepState and, in fact, general-purpose fuzzers such as AFL, have, to date, been exclusively (to our knowledge) used in what might be deemed conventional environments. As recently noted, “the tight coupling between hardware and firmware and the diversity found in embedded systems makes it hard to perform dynamic analysis on firmware” and existing mainstream fuzzing tools offer almost no support to embedded developers for simulation and emulation [21].

These goals require significant advances in three areas of dynamic analysis: first, a complete and principled approach to the problem of handling pre-conditions/assumption semantics, and second, an investigation of how to let fuzzers take advantage of the significant additional structure provided by property-based testing, including such assumptions. Consider the code in Figure 2. This defines two different tests of software that reads sensor values and incorporates them into a system state. The two tests check two different properties: `UpdateNeverSlow` ensures that updating the sensor is never too slow. It is checked, potentially, over *all* valid inputs, not just ones produced by the actual sensor reading code in `acquire_s1` and `acquire_s2`. The second test, `AvoidCrashes` starts the system up in some valid state, and repeatedly either reads both sensors or only sensor one. There is no explicit property, only the expectation that the system will not crash; tests can be executed using LLVM sanitizers to check for integer overflow and other undefined behavior. Generating such harnesses automatically from ACSL specifications is a significant challenge, but our research agenda also includes solving problems that would appear even for manual harnesses. For example, what is the proper semantics of the `ASSUME` in `update_state`? It depends on the test. In `UpdateNeverSlow`, a fuzzer will often generate an input value that violates the (possibly complex) requirements on valid states and sensor readings. These invalid inputs should not be flagged as bugs (the default behavior of E-ACSL), but instead the test should be abandoned without indicating that it failed. However, in `AvoidCrashes`, since we are not directly generating state values, that is, `update_state` is not an *entry point* for the test, assumption violations should result in failed tests. We aim to synthesize code to make assumptions automatically take on the proper semantics during test execution (including symbolic execution using constraint solvers).

This point about preconditions/`ASSUME` brings up a second point. Preconditions, when they have an `ASSUME` semantics, are fundamentally different than other branches in code. A fuzzer will attempt to explore the behavior of branches in `valid_state` and `valid_bv` just as it explores branches in `update_state` or

acquire. However, it is often possible to enumerate a vast number of paths that differentiate only invalid inputs, and so produce very little real testing. A classic example is “testing” a file system by producing a huge variety of unmountable file system images, rather than actually executing POSIX operations [35, 36]. DeepState knows which branches are pre-conditions, and so can help avoid this problem. In some fuzzers, this means prioritizing inputs to mutate based on whether they execute any code other than validity checks; but in fuzzers, such as Angora [16] and Eclipser [17], that use lightweight constraint-solving to cover branches, the process can be more sophisticated. We have begun discussions with the Eclipser team, and they confirm that identifying precondition code and devising suitable heuristics to handle it (e.g., never solve for a negation of a passed check) should improve performance. Fuzzing of individual functions or sets of functions is a highly promising area: most fuzzing is applied at the whole-program level, where input generation can simply be too hard. By focusing on a middle-ground between unit testing and whole-program fuzzing—using fuzzer technology to drive property-driven testing—the problem is made tractable. Prioritizing paths that include more than just input validation is an explicit goal of, e.g., AFLFast [12], but it must work with an implicit definition based on path frequencies, while we have access to ground truth. Given the complexity of state validity checks, there may be hard-to-reach—but uninteresting—ways to create invalid input; AFLFast will *prioritize* such paths, while we will (correctly) avoid them.

This effort also connects to a second fuzzing research thrust: making specification elements that do not correspond to simple code coverage visible to a fuzzer. In this example, consider the `DeepState.Timeout` check (note that this itself is functionality we will develop as part of handling timing constraints in FRAMA-C and DeepState). Unless we break down the timing analysis explicitly using a set of conditional branches, coverage-driven fuzzers cannot distinguish an execution that is very slow (close to violating the constraint) from one that has the minimum execution time possible. We propose to make timing of such specified events visible to a fuzzer, by modifying coverage bit-vectors to incorporate bucketing of execution time. Once we add such novel coverage measures, and introduce distinctions between coverage classes (as with preconditions), we will research how to balance competing priorities in more complex notions of coverage. In addition to implicit execution properties such as timing, this can apply to coverage of data structures, for fuzzing data-driven code such as machine-learning algorithms, where much behavior is implicit—e.g., the route taken through a forest of decision trees. In general we aim to extend the work [12, 53, 60, 80, 5], that prioritizes certain program paths in an intelligent way, by exploiting our extended ACSL/E-ACSL.

Finally, these elements must be tied to the problem of applying fuzzing and related methods in embedded-relevant execution environments. We plan to investigate multiple potential solutions, initially focusing on integrating fuzzing the approach taken by the just-released HALucinator (<https://github.com/embedded-sec/halucinator> tool [21] for virtualizing firmware via the Hardware Abstraction Layer, which is likely to add emulator-specific notions of coverage and path relevance.

3.2 DeepState and SAT/SMT-Based Bounded Model Checking

While automated test generation by fuzzing or binary-level symbolic execution can be highly effective as a means for finding bugs in code, other approaches are also needed to handle the kinds of code especially common in embedded contexts. In particular, embedded software often includes a large number of functions that perform complex low-level bit operations, especially for interacting with hardware and “parsing” network packets (from traditional wireless or RF-derived signals) communicating in very low-level protocols. Fuzzing or binary symbolic analysis often has trouble finding exact bit-values; it is well known that, e.g., inverting even non-cryptographic hashes is hard for either approach. Direct translation to bounded SAT (from source, not compiled code) on the other hand, often easily handles such input generation problems.

CBMC, the C Bounded Model Checker [49] is a well-known tool that analyzes C programs using a translation to SAT or SMT queries based on a bounded unrolling of loops. CBMC is an actively developed project, and has been used extensively in real-world development for years, including in automotive/embedded code development at Bosch and General Electric [70], in analysis of Amazon Web Services infrastructure [22], and in the analysis of flight software systems at NASA’s Jet Propulsion Laboratory [36]. Using CBMC requires writing custom test harnesses using CBMC’s API for expressing nondeterminism, and running the

tool with a specified bound on loop executions, in addition to other complex configuration options.

We propose to allow CBMC to be used as a backend for verification by DeepState, with a seamless interface, just as DeepState currently supports symbolic analysis engines such as angr and Manticore. It is notoriously hard to guess when a SAT/SMT based approach to code analysis will work well and when it will fail to scale; using a DeepState harness will allow users to try CBMC at “no cost.”

Moreover, because choosing loop unwinding bounds imposes a serious burden on embedded engineers, we will investigate their automatic determinations. One approach is to instrument fuzzer or symbolic-execution engine generated tests to record iterations of loops, and then use the maximum bound observed. Additionally, for small functions (the most likely targets for DeepState-CBMC: complex but compact bit-manipulation code), the mutation-based approach proposed by Groce et. al [37] may work. Finally, in some cases CBMC may be able to find interesting bugs for cases where the loop unrollings are limited, but cannot scale to larger depth limits. Using the same instrumentation that we use to estimate loop bounds, we will use the ability to guide fuzzers by alternative “coverage” to focus fuzzer runs on executions with more loop iterations than the bound explored by CBMC. This will offer engineers a true partnership between verification methods.

3.3 DeepState and Explicit-State Model Checking

Just as some functions are best analyzed using bounded model checking, some dynamic analysis problems are best handled by explicit-state model checking that actually executes C code, like a fuzzer, but with the capability to store states and backtrack, in order to exhaustively explore a state space, using either actual comparison of stored states or comparison of abstractions of states to guide exploration. This approach is particularly attractive for exploring sequences of API calls; this kind of test generation was used in efforts that uncovered dozens of errors in file systems at NASA/JPL [36].

The SPIN model checker [46] offers execution of C code with backtracking [44, 45]. DeepState’s `OneOf` construct has a semantics that can be matched with the SPIN nondeterministic choice, which in part inspired the DeepState construct [33, 32]. However, integrating SPIN as a back-end for DeepState is even more challenging than integrating CBMC. With CBMC, the mapping from DeepState to CBMC semantics may be performed by changing included headers so that CBMC-specific constructs have differing implementations (but not semantics); SPIN however executes C code in the context of a PROMELA model, which requires rewriting a DeepState model to embed test choices inside SPIN’s constructs. This also means “lifting” DeepState API calls to the PROMELA level outside the C code, and bridging between nondeterminism visible to SPIN and determinism within C code; PI Groce’s previous work [33] can serve as a foundation. A more fundamental problem is that while CBMC and DeepState can share a semantics for, e.g., `DeepState_Int64()`, a PROMELA model with a branching factor of, e.g., 2^{64} will not work. Solutions range from using results from fuzzing to choose a limited range, to translating “flat” bit-value selection into a sequence of choices with a larger range but bias towards certain values, to using SPIN to control a seed and deterministically choosing random values [33], a hybrid approach.

3.4 DeepState and Timed Automata Model Skeletons

As noted above, one of our core assumptions is that timed automata can model the underlying protocols in many embedded systems. However, writing timed automata models using UPPAAL [8] and PRISM [50] is at present a skill only a small number of embedded engineers have mastered. In order to encourage more engineers to make use of these powerful formalisms, we propose 1) to enable DeepState to generate *traces* of the annotations related to timing that are covered during a run and 2) to build a tool to combine and reconcile these traces into a skeleton model for UPPAAL [8] or PRISM [50]. The structure of code (function locations of DeepState annotations) will be used to form the structure of the model. Additional annotations for, e.g. probabilities, may need to be added if not present in the code annotations, though DeepState already has a primitive support for expressing probabilities that we plan to extend.

3.5 Case Studies

The above briefly introduces a number of problems that we know in advance must be dealt with in order to enable a pathway for combining formal, static, and dynamic analysis. At heart, however, we aim to allow

case studies to prioritize our efforts, and are certain that other challenges will arise during these efforts. The studies informing this research is the embedded software of wireless sensor nodes used in the Southwest Experimental Garden Array (SEGA) [19, 27, 7] and of wireless sensor nodes and mobile robots in the Distributed Sensing & Computing Over Sparse Environments (DISCOVER) Platform.

3.5.1 Communication Protocol for Wireless Sensor Nodes in SEGA

Overview: SEGA is a large collection of operational wireless sensor/actuator networks for monitoring and control of ecological systems, located at 17 sites in the states of Arizona and California. Currently, SEGA consists of 138 wireless nodes and is planned to expand to a total of 154 nodes at 21 sites in the coming years. As a genetics-based climate change research platform, SEGA allows scientists to quantify the ecological and evolutionary responses of species to changing climate conditions. Multiple long-term and large-scale scientific experiments are conducted at SEGA sites. The SEGA project was led by NAU, and the group of co-PI Flikkema developed and are maintaining the wireless nodes, including their embedded software. The SEGA nodes use a multi-processor architecture, in which a central processor provides OS-level services while plug-in satellite processors handle transducer sampling, actuation, and related computational tasks. In addition to allowing true parallelism, this architecture enables hardware-level improvements in energy efficiency, since each satellite can be optimized for its specific task. The nodes synchronously interact with neighbors in a multi-hop, self-organizing/healing network, implemented by a custom communication protocol designed by the group of co-PI Flikkema. The nodes use a custom time-triggered RTOS tightly integrated with a time/frequency-hopped PHY/MAC protocol. This approach minimizes communication energy cost.

Challenges: Because timing is critical and is determined by the embedded system hardware and software, most testing has occurred at the network level, with extensive in-lab testing with small networks and instrumented field tests. However, it has been found in long-term deployments that occasionally the networking fails and nodes become isolated—we think due to a complex set of subtle bugs rooted in different levels of timing abstraction. When such a failure occurs, it often spreads from one node to others, causing nodes to seek to rejoin and expend high levels of energy for radio operation and eventually deplete their energy sources. Eventually, subnets, or sometimes the entire site, are disabled and humans must visit the site to reboot it. Such failures could affect or even destroy (e.g., via over-watering), long-running scientific experiments. We aim to use SEGA (in particular the protocol in question and its implementation) as our case study. This will enable us to apply our approach in a practical setting, and ensure that what we produce is actually usable by engineers of real systems. SEGA is an ideal case study for several reasons. First, the above mentioned network problem enables exploring how to design, prove, and test time-critical systems in a way that does no harm: human life is not affected in this application, and data is not lost since all sensed information is logged as a local back-up. On the other hand, reliable operation is important. Finally, this application uses common data structures for task control blocks, and the operating system at each node schedules and dispatches both periodic and pseudo-randomly scheduled tasks.

Plan: Following our proposed workflow, we will first annotate the implementation with specifications of correctness properties. We may model the protocol itself as a timed automaton in UPPAAL or PRISM, in order to ensure that there is not a subtle flaw in the protocol itself, and to model our expectations of behavior in the real system (and to better understand needed specifications). Either of these steps may expose the source of the mysterious networking failures. We will use DeepState, driven by harnesses automatically generated by our tools, to generate tests of the implementation components in question, using fuzzing at first, followed by CBMC and SPIN model checking once prototype back-ends are available. For the purpose of verification and testing, the uncertainty inherent in the physical environment and in wireless communication will be modeled by probabilistic timed automata. The above workflow will be conducted by an Embedded System Engineering student, who is familiar with the SEGA IoT system but does not have expertise in software verification and testing, using the software tools developed in this project. Feedback from the engineer in this case study will inform us how to develop and improve the theory and tools for practical usage.

3.5.2 Embedded Software of Wireless Sensor Nodes and Robots in DISCOVER

Overview: DISCOVER is a cyberinfrastructure testbed for remote, rural, and sparsely populated areas. The project is funded by NSF and led by NAU, whose team includes co-PI Flikkema (PI of DISCOVER) and co-PI Nghiem (co-PI of DISCOVER in charge of robotics). DISCOVER consists of a fabric of highly configurable Internet-of-Things (IoT) sensor nodes, autonomous and highly capable terrestrial robots and drones, and a heterogeneous wireless network. DISCOVER sites will be located at the campuses of NAU, Navajo Technical University, and Clemson University, as well as several remote sites. The platform will enable focused research in many domains, including data science and machine learning, heterogeneous networked services, distributed computing and AI, control, autonomous robots, and in-network computation, among many others. We will use DISCOVER for two case studies: one on embedded software for wireless sensor nodes and the other on distributed coordination in multi-robot systems. Coordinated operation of multiple autonomous robots has many important real-world applications [14, 78], e.g., in rescue, security, or disaster response missions. In such applications, each robot is autonomous but has the capability to coordinate efficiently and safely with other robots to complete a shared mission, often in a distributed manner. Such coordination is essential in real-world applications where the environment is constantly and unexpectedly changing. One of the most critical challenges of this application is to guarantee the safety of a coordination plan, which is typically implemented in C code on the embedded computers of the robots and usually involves wireless inter-robot communication, sensing, and actuation.

Challenges: As a community research platform, DISCOVER will allow users, who are researchers in relevant field domains, to develop software code and experiments for the DISCOVER stationary nodes (i.e., sensor nodes) and mobile nodes (i.e., terrestrial robots and drones). A critical step of the process supported by DISCOVER is automatic verification and testing of the embedded code submitted by users for live experiments. We expect that submitted code is developed by researchers who are not trained in computer science and who do not usually apply best practices in software engineering. We also expect that submitted code will have a wide spectrum of code quality and may be malicious, either by chance or intentionally. Automatic functional testing of user code will therefore be of critical importance for the operation and sustainability of DISCOVER. Another challenge is the fact that the code to be tested and verified is for embedded systems that have highly complex physical dynamics (in the case of robots) and interactions with the physical world and with other physical systems. Such physical aspects cannot be easily described in code for the purpose of software verification and testing. Therefore, sophisticated software-based and hardware-in-the-loop simulations, including embedded hardware emulators and robot simulators, are necessary, for which several options will be developed by the DISCOVER team.

Plan: For the wireless sensor nodes, where a major part of their operation is related to data storage and communication, our planned approach will be similar to that of the SEGA case study. For the robotic nodes, we will focus on multi-robot coordination since other user errors, such as unauthorized hardware access and unauthorized maneuvers of robots, can be prevented or mitigated through a combination of techniques such as containerization, access rights, and geofencing. Our approach will be briefly described below. First, we will model a coordination plan for multiple robots as a (potentially very complex) network of timed automata. Performance specifications will be expressed in temporal logics, e.g., the Signal Temporal Logic (STL) [24], and checked against the model using verification and testing tools such as UPPAAL or S-TaLiRo [4]. While we do not expect actual user code to be accompanied by formal models, in our case study, this step ensures that the original coordination plan has no subtle flaws, and helps us determine properties that need formulation at the implementation level. An implementation of the algorithm in C code, distributed among the robots, will be developed by a robotics/control student. The implementation will be annotated with a specification in our extended ACSL/E-ACSL. We will then use DeepState harnesses to generate tests of the implementation components using fuzzing, symbolic execution, and both bounded SAT/SMT based and explicit-state model checking. Finally, we will determine if a timed automata skeleton extracted from the implementation code corresponds to and would help create a full specification such as we developed before beginning implementation. The very different nature and complexity of this study, compared to stationary

sensor nodes, will ensure that our methods and tools work in a variety of kinds of real systems. To overcome the challenge stemming from the complex physical dynamics and interactions of the robots, we will utilize a sophisticated robot simulation environment, based on the Robot Operating System (ROS) [59], with a rich set of predefined scenarios, developed by the DISCOVER team (specifically by the group of co-PI Nghiem). An interface between the robot simulation software and the tools developed in this project will be created to enable seamless verification and testing of the robotic code.

3.6 Work Plan

The project will be organized into two phases, described by work packages. In the first phase, T3.1 will be conducted along with and inform T1.1 (see Figure 3). In the second phase, the focus will be on the application of tools in T1.2 in tandem with T2. Tasks related to case studies (tasks T4.1, and T4.2) will help refine the developed tools especially in the final phases of the project.

Work Package 1 (WP1): This work package concerns the development of and use of ACSL and E-ACSL extensions.

- T1.1: This task will consider needed extensions for handling real-world embedded systems. In particular, there will be a focus on a study of the formal semantics of timed automaton networks defined in UPPAAL and PRISM, to determine the extent to which shared semantics can be assigned making it possible to carry implementation annotations into such formal models. In addition, this task will include initial consultation with engineers from Galois (<https://galois.com>) to discuss needs for their customers and tools. Galois is a key player in the space of annotations and tools for critical low-level system verification.
- T1.2: This task will take feedback from applications of tools to generate tests and proofs (T2) into account, to add annotations that are focused on heuristic guidance for tools, not correctness per se.

One Ph.D. student will conduct this work, which will last for the entire duration of the project.

Evaluation: Evaluation of WP1 will be determined by ability of embedded engineers to agree that the key properties, including those related to timed automata models, to be checked are (1) all representable by the annotations (2) easy to construct (3) easy to read when produced by others and (4) maintainable after introduction. In addition to our own case studies, discussion with Galois engineers will inform our evaluation.

Work Package 2 (WP2): This work package covers automatically translation of ACSL/E-ACSL-annotated code into a DeepState test harness (Section 3.1), development of back-ends for CBMC and SPIN, and improvements to fuzzers:

- T2.1: This task will optimize the implementation of symbolic execution and fuzzing in DeepState, so that ACSL/E-ACSL annotations and extensions from WP1 can be used effectively.
- T2.2: This task will develop DeepState back-ends for CBMC and SPIN, inform annotations needed to handle loop bounds, memory tracking and matching, and make use of feedback from fuzzing.

The execution of this work package will also span the entire duration of the project. Because the tasks in this package are also based on developing verification and test generation tools (thus formal methods expertise), the same Ph.D. student will work on WP1 and WP2. We separate the WPs primarily to emphasize that specification extensions and tool support are somewhat orthogonal concerns, and evaluated differently.

Evaluation: Evaluation of WP2 will be determined by the application of DeepState harnesses to generate tests for realistic systems. We will use benchmarks and simple examples to some extent, but primarily rely on our connection to case studies. For test generation, we will use coverage and faults detected as measures [48].

Work Package 3 (WP3): This work package will focus on the field applications described in Section 3.5, as both a way to inform the methodology and tool developments. WP3 includes the following case studies:

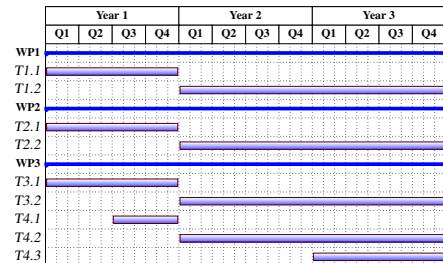


Figure 3: Project schedule.

Wireless sensor network (WSN) case studies on SEGA and DISCOVER. This is divided into two tasks:

- T3.1: In this task, the wireless sensor node systems will be studied thoroughly to extract the key requirements and characteristics of the embedded system implementations. Timed automaton models of the communication protocol in each system, at different levels of abstraction, may be developed and formally verified in UPPAAL and/or PRISM, to inform task T1.1. The system information and models resulting from this task will inform the semantics design and method developments in WP1 and WP2. As time allows we will extend this work to include sensing elements.
- T3.2: This task will apply the tools developed in WP1 and WP2 to the WSN systems, to detect and fix bugs in the communication protocol implementations; in particular, the bugs that cause the intermittent failures in SEGA. It will also provide feedback to the other work packages to refine and improve our tools.

Multi-robot system case study on DISCOVER. This study is divided into three tasks:

- T4.1: In this task, a standard multi-robot coordination algorithm will be modeled as a network of timed automata. Using our insights into the robotics application, we will express its performance specifications, particularly its safety requirements, in temporal logics and formally verify or test them in tools like UPPAAL, PRISM, or S-TaLiRo. This task will extend the developed semantics and methods to applications beyond communication protocols, to identify further needed runtime extensions and semantic connections between timed automata theory, implementation annotations, and runtime checks.
- T4.2: This task will apply the tools developed in WP1 and WP2, and the robot simulation environment of the DISCOVER platform, to the coordinated multi-robot system, in order to validate the implementation code, detect and fix possible bugs, and improve the tools developed in this project.
- T4.3: This task will aim to apply the DeepState-trace driven route to produce timed automata skeletons.

As the tasks in this work package are conducted in tandem with WP1 and WP2, to form a feedback loop with the developments in other work packages, it will last for the entire duration of the project. We expect that groups of undergraduate students, in collaboration with an embedded systems Ph.D. student and the Ph.D. students in WP1 and WP2, will perform the work. Close collaboration with the DISCOVER team, led by Dr. Flikkema and Dr. Nghiem, is expected.

Evaluation: In essence, this task is the evaluation aspect of our project, which forms one of the major thrusts of the project. The successful application of WP1 and WP2 tools to the case studies is essentially the driving factor in determining our success in the project. The measure of success is: (1) faults detected and corrected; (2) functionality proven correct using CBMC, symbolic execution engines, or SPIN; (3) coverage and other measures of generated tests; and (4) reported usability and value by engineers, particularly students. For T4.3, evaluation will be based on comparison of extracted skeletons with independently developed full models.

4 Contributions to Formal Methods and the Field

The contributions to formal methods proposed include: (1) Fundamental contributions to integrating formal specification languages developed for use in static analysis and theorem proving with dynamic analysis, producing a common semantics for formal, static, and dynamic checking of correctness; handling of timing and interrupts are notable examples of problems to be addressed in this effort; (2) Enhanced ability of fuzzing and other test generation methods to make use of information from formal specifications, and integrate feedback about, e.g., specification coverage into test generation heuristics; (3) Common semantics and a framework for fuzzing, symbolic execution, and model checking. (4) Approaches to using feedback from fuzzing to guide bounded or explicit-state model checking; (5) Translations from implementation-level specification to (probabilistic) timed automata models.

The contributions to the field include: (1) New development and design methods that focus on implementation-level specification as a guiding principle for embedded systems; (4) Tactics and strategies for incorporating the above methods into legacy efforts, where existing code bases require additional specification and annotation; (5) Best-practices for using formal, static, and dynamic tools in debugging legacy systems problems.

Our *evaluation* of the degree to which these contributions have been realized is described in the work plan above, integrated with description of case study efforts.

5 Related Work

A fundamental goal of this project is to reduce both user effort and the opportunity for user effort by allowing minimizing (ideally to one) the number of times a user must specify an aspect of system correctness. The principle that important information should have a “single point of truth” is widely accepted in software engineering, even in such foundational early advances as avoiding repeated magic numbers by the use of named constants. Such a principle can be extended to specification and definition of test harnesses. Early work emphasizing this goal of both reducing work and chance of error in specification and test generation included the effort by Groce and Joshi to use a single harness for both model-checking and random testing, in the verification of the Mars Science Laboratory’s file system [33, 35, 36]. In later work, Groce and Erwig extended this idea to propose development of a single language with a unified semantics for a wide variety of dynamic test generation tools [32]; this approach is essentially realized in the DeepState [29] system. Indeed, FRAMA-C and ACSL [6] and DeepState are both arguably limited instantiations of this goal: providing a single language, interface, and semantics that is applied to a variety of methods (static or dynamic) for checking that a specification holds. This project aims to further extend this goal by extending it to include a formal timed-automata model and to connect the primarily static and dynamic approaches.

The implementation and verification of distributed systems, and code extraction from automata modeling in the proof assistant COQ [74, 76, 58] is a topic of some previous work. Such proposals require that the developers master COQ, and start from the modeling activity to generate code. They are therefore not applicable in the context of the verification of legacy embedded C code, the common case in the real world.

Testing real-time systems modeled by networks of timed automata was investigated by the authors of the tool UPPAAL [42, 41, 51] and implemented in the tools UPPAAL-TRON (<http://people.cs.aau.dk/~marius/tron/index.html>) and UPPAAL-COVER (<http://www.hessel.nu/CoVer/index.php>). These tools generate tests, either offline or online, for conformance testing of a real-time system with respect to its model and an environment model, both as timed automaton networks. In both cases, the real-time system is considered a black-box with an input/output interface through which the test generator or monitor can change the system inputs and observe the system outputs. The actual implementation code is not considered and is in fact hidden from the testing tools. While this approach is general, it has several drawbacks. It requires a centralized input/output interface accessible to the testing tools. Such an interface is not always available in all systems, especially in large-scale distributed systems like the sensor/actuator networks considered in our case study. Furthermore, by considering only the (timed) input/output behavior of a system, this approach may not be able to test internal system behaviors and therefore miss opportunities for a better test coverage.

6 Broader Impacts

Improving Software System Reliability: A key element of our approach is to focus on realistically deployable techniques. We aim for early integration with NASA’s FPrime [11, 56] open source flight software architecture and platform; PI Groce is already in discussion with engineers at NASA’s Jet Propulsion Laboratory, and engaged in producing tests for the FPrime autocoder using DeepState. This integration will allow our methods to be applied to CubeSat missions (and other flight software systems), leading to improved reliability for low-budget space-based scientific efforts. We expect, in the long run, that our approaches will lead to more reliable and robust development in many embedded and cyberphysical systems domains. Our engagement with interested Galois engineers ensures the applicability of our methods will be as wide as possible, and impact tools outside our initial scope.

Education and Outreach: The proposed research yields several opportunities for enhancing CS education, recruiting new CS majors, and retaining CS students, particularly members of underrepresented groups. In addition to the activities discussed at length in the Broadening Participation in Computing plan, PI Groce will work with the NAU Student ACM Chapter to present a series of “excursions in testing” that introduce automated testing to students, using DeepState to find bugs in real world code, including code from media player libraries. The work of Guzdial [40] has shown that media computation is an effective way to both recruit and retain female and underrepresented students in computer science. Groce is also teaching a class

on automated testing of embedded systems. Co-PI Nghiem has developed a course on autonomous vehicles, based on the F1/10 platform (funded by NSF). To prepare students for addressing safety in autonomous driving, future offerings of the course will incorporate the methods and tools developed in this project.

Broadening Participation in Computing (BPC): The goal of the BPC component of this project is to *increase the number of females who are involved or choose careers in computing, at NAU and in the local community of Flagstaff, Arizona.* Our plan carefully integrates active learning experiences designed for female students at both the undergraduate and middle school/junior high levels. **Undergraduate Education Experience** - We will reach female students in two degree programs at the 2nd-year level: Computer Science and Electrical and Computer Engineering. In CS, we will target CS 200 Introduction to Computer Organization; in ECE, we will target EE 215 Microprocessors. We will integrate a new project in which teams of female (and possibly male, due to the current lack of females in ECE and CS) students imagine and create exciting and meaningful one-day active learning experiences and projects for female student teams in grades 7-9. We will provide full support to these teams, especially female students, and design the project so that female students will take leadership roles to gain confidence. In both courses, we will bring in expert speakers to facilitate development of students' understanding how to design these projects so they are marker events in the students' lives. **Outreach to grades 7-9** - As noted above, the undergraduate teams will develop active learning and design project "Build Events" for girls in grades 7-9. We will recruit female undergraduates who have taken CS 200/EE 215 to become mentors in the one-day events for the grade 7-9 students. We will schedule these events as part of the annual Flagstaff Festival of Science, and plan them for Saturdays to avoid conflicts with school schedules, maximizing participation. The Flagstaff Festival of Science (www.scifest.org), now in its 32nd year and enjoying wide financial and participatory support in the community, holds over 100 events for all ages over a 10-day period in the Fall, and is an ideal venue. **Facilities and Support** - By scheduling the grade 7-9 Build Events on Saturdays, we will be able to use the educational laboratories of the School of Informatics, Computing & Cyber Systems (SICCS) for the Flagstaff Festival of Science events. We have requested \$2,000 for each in years 2 and 3 for materials (primarily embedded development boards) for these experiences. **Assessment** - We will conduct focused feedback sessions and administer short surveys of the participants to aid continuous improvement of the activities over the course of the project.

7 Results From Prior NSF Support

PI Groce: The most relevant prior NSF support for PI Groce is CCF- CCF-2129446, "Feedback-Driven Mutation Testing for Any Language," with a total budget of \$500,000 from 9/2021 until 8/2024, a collaborative proposal with Claire Le Goues of Carnegie Mellon University. **Intellectual Merit:** This project focuses on a synergistic approach for allowing developers to improve testing by using mutation testing to identify weaknesses in tests and to generate tests. Though in its first year, this project has already resulted in three publications [39, 38?]. **Broader Impact:** Work from this project has already resulted in the reporting and correction of multiple bugs in software systems, including production compilers for smart contracts.

Co-PI Flikkema and co-PI Nghiem are PI and co-PI, respectively, on the Distributed Sensing & Computing Over Sparse Environments (DISCOVER) Platform funded by an NSF CCRI grant (2120485), with a total budget of \$1,366,513 from 10/2021 until 9/2024. **Intellectual Merit:** DISCOVER is a cyberinfrastructure testbed for remote, rural, and sparsely populated areas, consisting of a fabric of highly configurable and fixed and mobile Internet-of-Things nodes and a heterogeneous wireless network. It will enable focused research on new breeds of algorithms that address a range of challenges around prioritization and optimization of computation and communication in remote, less populous and rural areas. **Broader Impact:** DISCOVER will provide a research platform for investigation of distributed computing, networking, security, control and coordination solutions in a heterogeneous configurable cyber-physical system infrastructure that will provide critical services for areas and populations at increasing risk of being underserved. The education and outreach impacts of this project include training and research opportunities for undergraduate students, engaging underrepresented minority students, and developing hands-on research experiments for K-12 students.

References

- [1] Google Test. <https://github.com/google/googletest>, 2008.
- [2] honggfuzz. <https://github.com/google/honggfuzz>, 2010.
- [3] Rajeev Alur and David L. Dill. A theory of timed automata. *Theor Comput Sci*, 126(2):183 – 235, 1994. ISSN 0304-3975. doi:[https://doi.org/10.1016/0304-3975\(94\)90010-8](https://doi.org/10.1016/0304-3975(94)90010-8).
- [4] Yashwanth Annpureddy, Che Liu, Georgios Fainekos, and Sriram Sankaranarayanan. S-taliro: A tool for temporal logic falsification for hybrid systems. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 254–257. Springer, 2011.
- [5] Cornelius Aschermann, Sergej Schumilo, Tim Blazytko, Robert Gawlik, and Thorsten Holz. REDQUEEN: fuzzing with input-to-state correspondence. In *NDSS (Network and Distributed Security Symposium)*, 2019.
- [6] Patrick Baudin, Jean C. Filliâtre, Thierry Hubert, Claude Marché, Benjamin Monate, Yannick Moy, and Virgile Prevosto. *ACSL: ANSI/ISO C Specification Language*, February 2011. <http://frama-c.cea.fr/acsl.html>.
- [7] David M Bell, Eric J Ward, A Christopher Oishi, Ram Oren, Paul G Flikkema, and James S Clark. A state-space modeling approach to estimating canopy conductance and associated uncertainties from sap flux density data. *Tree Physiology*, 2015.
- [8] Johan Bengtsson, Kim Larsen, Fredrik Larsson, Paul Pettersson, and Wang Yi. Uppaal—a tool suite for automatic verification of real-time systems. In *International Hybrid Systems Workshop*, pages 232–243. Springer, 1995.
- [9] Allan Blanchard, Frédéric Loulergue, and Nikolai Kosmatov. Towards Full Proof Automation in Frama-C using Auto-Active Verification. In *Nasa Formal Methods*, LNCS, pages 88–105. Springer, 2019. doi:10.1007/978-3-030-20652-9_6.
- [10] Lionel Blatter, Nikolai Kosmatov, Pascale Le Gall, Virgile Prevosto, and Guillaume Petiot. Static and dynamic verification of relational properties on self-composed C code, 2018.
- [11] Robert Bocchino, Timothy Canham, Garth Watney, Leonard Reder, and Jeffrey Levison. F prime: An open-source framework for small-scale flight software systems. In *Small Satellite Conference*, 2018.
- [12] Marcel Böhme, Van-Thuan Pham, and Abhik Roychoudhury. Coverage-based greybox fuzzing as Markov chain. *IEEE Transactions on Software Engineering*, 45(5):489–506, 2017.
- [13] Patricia Bouyer, François Laroussinie, Nicolas Markey, Joël Ouaknine, and James Worrell. Timed temporal logics. In *Models, Algorithms, Logics and Tools - Essays Dedicated to Kim Guldstrand Larsen on the Occasion of His 60th Birthday*, volume 10460 of LNCS, pages 211–230. Springer, 2017. doi:10.1007/978-3-319-63121-9_11.
- [14] W. Burgard, M. Moors, C. Stachniss, and F. E. Schneider. Coordinated multi-robot exploration. *IEEE Transactions on Robotics*, 21(3):376–386, June 2005. doi:10.1109/TRO.2004.839232.
- [15] Cristian Cadar, Daniel Dunbar, and Dawson Engler. KLEE: unassisted and automatic generation of high-coverage tests for complex systems programs. In *OSDI*, 2008.
- [16] Peng Chen and Hao Chen. Angora: Efficient fuzzing by principled search. In *2018 IEEE Symposium on Security and Privacy (SP)*, pages 711–725, 2018.
- [17] Jaeseung Choi, Joonun Jang, Choongwoo Han, and Sang Kil Cha. Grey-box concolic testing on binary code. In *International Conference on Software Engineering*, pages 736–747, 2019.
- [18] Koen Claessen and John Hughes. QuickCheck: a lightweight tool for random testing of Haskell programs. In *International Conference on Functional Programming*, pages 268–279, 2000.
- [19] J.S. Clark, P.K. Agarwal, D. M. Bell, P.G. Flikkema, A. Gelfand, X. Nguyen, E. Ward, and J. Yang. Inferential ecosystem models, from network data to prediction. *Ecological Applications*, 21(5), July 2011. In press.
- [20] Lori A. Clarke and David S. Rosenblum. A historical perspective on runtime assertion checking in software development. *SIGSOFT Softw. Eng. Notes*, 31(3):25–37, May 2006. doi:10.1145/1127878.1127900.

- [21] Abraham A Clements, Eric Gustafson, Tobias Scharnowski, Paul Grosen, David Fritz, Christopher Kruegel, Giovanni Vigna, Saurabh Bagchi, and Mathias Payer. {HALucinator}: Firmware re-hosting through abstraction layer emulation. In *29th USENIX Security Symposium (USENIX Security 20)*, pages 1201–1218, 2020.
- [22] Byron Cook, Kareem Khazem, Daniel Kroening, Serdar Tasiran, Michael Tautschnig, and Mark R Tuttle. Model checking boot code from aws data centers. *Formal Methods in System Design*, pages 1–19, 2020.
- [23] Patrick Cousot and Radhia Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proc. of the 4th ACM Symposium on Principles of Programming Languages (POPL 1977)*, pages 238–252. ACM, 1977.
- [24] Alexandre Donzé and Oded Maler. Robust satisfaction of temporal logic over real-valued signals. In *International Conference on Formal Modeling and Analysis of Timed Systems*, pages 92–106. Springer, 2010.
- [25] P.G. Flikkema, K.R. Yamamoto, S. Boegli, C. Porter, and P. Heinrich. Towards cyber-eco systems: Networked sensing, inference and control for distributed ecological experiments. In *2012 IEEE International Conference on Green Computing and Communications (GreenCom)*, pages 372–381, Nov 2012. doi:10.1109/GreenCom.2012.61.
- [26] Eric Gamma and Kent Beck. JUnit 5. <http://junit.org/junit5/>.
- [27] Souparno Ghosh, David M. Bell, James S Clark, Alan E. Gelfand, and Paul G. Flikkema. Process modeling for soil moisture using sensor network data. *Statistical Methodology*, 17:99–112, 2014.
- [28] Patrice Godefroid, Nils Klarlund, and Koushik Sen. DART: directed automated random testing. In *Programming Language Design and Implementation*, pages 213–223, 2005.
- [29] Peter Goodman and Alex Groce. DeepState: Symbolic unit testing for C and C++. In *NDSS Workshop on Binary Analysis Research*, 2018.
- [30] Peter Goodman, Gustavo Greco, and Alex Groce. Tutorial: DeepState: Bringing vulnerability detection tools into the development cycle. In *IEEE Cybersecurity Development Conference (SECDEV)*, 2018.
- [31] Alex Groce. Test harness for testfs. <https://github.com/agroce/testfs>, 2018.
- [32] Alex Groce and Martin Erwig. Finding common ground: Choose, assert, and assume. In *International Workshop on Dynamic Analysis*, pages 12–17, 2012.
- [33] Alex Groce and Rajeev Joshi. Random testing and model checking: Building a common framework for nondeterministic exploration. In *Workshop on Dynamic Analysis*, pages 22–28, 2008.
- [34] Alex Groce and Jervis Pinto. A little language for testing. In *NASA Formal Methods Symposium*, pages 204–218, 2015.
- [35] Alex Groce, Gerard Holzmann, Rajeev Joshi, and Ru-Gang Xu. Putting flight software through the paces with testing, model checking, and constraint-solving. In *Workshop on Constraints in Formal Verification*, pages 1–15, 2008.
- [36] Alex Groce, Klaus Havelund, Gerard Holzmann, Rajeev Joshi, and Ru-Gang Xu. Establishing flight software reliability: Testing, model checking, constraint-solving, monitoring and learning. *Annals of Mathematics and Artificial Intelligence*, 70(4):315–349, 2014.
- [37] Alex Groce, Iftexhar Ahmed, Carlos Jensen, Paul E McKenney, and Josie Holmes. How verified (or tested) is my code? falsification-driven verification and testing. *Automated Software Engineering Journal*, 2018. Accepted for publication.
- [38] Alex Groce, Kush Jain, , Rijnard van Tonder, Goutamkumar Tulajappa Kalburgi, and Claire Le Goues. Looking for lacunae in bitcoin core’s fuzzing efforts. In *ACM/IEEE International Conference on Software Engineering*, 2022. accepted for publication.
- [39] Alex Groce, Rijnard van Tonder, Goutamkumar Tulajappa Kalburgi, and Claire Le Goues. Making no-fuss compiler fuzzing effective. In *ACM SIGPLAN International Conference on Compiler Construction*, 2022. accepted for publication.
- [40] Mark Guzdial. A media computation course for non-majors. In *Proceedings of the 8th Annual*

- Conference on Innovation and Technology in Computer Science Education*, ITiCSE '03, pages 104–108, New York, NY, USA, 2003. ACM. ISBN 1-58113-672-2. doi:10.1145/961511.961542. URL <http://doi.acm.org/10.1145/961511.961542>.
- [41] Anders Hessel and Paul Pettersson. CoVer - a real-time test case generation tool. In *19th IFIP International Conference on Testing of Communicating Systems and 7th International Workshop on Formal Approaches to Testing of Software*, 2007.
 - [42] Anders Hessel, Kim Larsen, Marius Mikucionis, Brian Nielsen, Paul Pettersson, and Arne Skou. Testing Real-Time systems using UPPAAL. In *Formal Methods and Testing*, pages 77–117. 2008.
 - [43] Josie Holmes, Alex Groce, Jervis Pinto, Pranjal Mittal, Pooria Azimi, Kevin Kellar, and James O'Brien. TSTL: the template scripting testing language. *International Journal on Software Tools for Technology Transfer*, 20(1):57–78, February 2018.
 - [44] Gerard Holzmann and Rajeev Joshi. Model-driven software verification. In *SPIN Workshop on Model Checking of Software*, pages 76–91, 2004.
 - [45] Gerard Holzmann, Rajeev Joshi, and Alex Groce. Model driven code checking. *Automated Software Engineering*, 15(3–4):283–297, 2008.
 - [46] Gerard J. Holzmann. *The SPIN Model Checker: Primer and Reference Manual*. Addison-Wesley Professional, 2003.
 - [47] Florent Kirchner, Nikolai Kosmatov, Virgile Prevosto, Julien Signoles, and Boris Yakobowski. Frama-C: A software analysis perspective. *Formal Asp. Comput.*, 27(3):573–609, 2015. doi:10.1007/s00165-014-0326-7.
 - [48] George Klees, Andrew Ruef, Benji Cooper, Shiyi Wei, and Michael Hicks. Evaluating fuzz testing. *arXiv preprint arXiv:1808.09700*, 2018.
 - [49] Daniel Kroening, Edmund M. Clarke, and Flavio Lerda. A tool for checking ANSI-C programs. In *Tools and Algorithms for the Construction and Analysis of Systems*, pages 168–176, 2004.
 - [50] M. Kwiatkowska, G. Norman, and D. Parker. PRISM 4.0: Verification of probabilistic real-time systems. In *Proc. 23rd International Conference on Computer Aided Verification (CAV'11)*, volume 6806 of *LNCS*, pages 585–591. Springer, 2011.
 - [51] Kim G. Larsen, Marius Mikucionis, and Brian Nielsen. Testing real-time embedded software using UPPAAL-TRON: an industrial case study. In *the 5th ACM international conference on Embedded software*, pages 299 – 306. ACM Press New York, NY, USA, September 18–22 2005.
 - [52] Gary T. Leavens, Albert L. Baker, and Clyde Ruby. JML: a Java Modeling Language. In *Formal Underpinnings of Java Workshop (at OOPSLA '98)*, October 1998. <http://www-dse.doc.ic.ac.uk/~sue/oopsla/cfp.html>.
 - [53] Caroline Lemieux and Koushik Sen. FairFuzz: a targeted mutation strategy for increasing greybox fuzz testing coverage. In *International Conference on Automated Software Engineering*, pages 475–485, 2018.
 - [54] David R. MacIver. Hypothesis: Test faster, fix more. <http://hypothesis.works/>, March 2013.
 - [55] Mark Mossberg, Felipe Manzano, Eric Hennenfent, Alex Groce, Gustavo Grieco, Josselin Feist, Trent Brunson, and Artem Dinaburg. Manticore: A user-friendly symbolic execution framework for binaries and smart contracts. In *34th IEEE/ACM International Conference on Automated Software Engineering, ASE 2019, San Diego, CA, USA, November 11-15, 2019*, pages 1186–1189. IEEE, 2019. doi:10.1109/ASE.2019.00133. URL <https://doi.org/10.1109/ASE.2019.00133>.
 - [56] NASA. F prime: A flight-proven, multi-platform, open-source flight software framework. <https://github.com/nasa/fprime>, 2018.
 - [57] Rickard Nilsson, Shane Auckland, Mark Sumner, and Sanjiv Sahayam. ScalaCheck user guide. <https://github.com/rickynils/scalacheck/blob/master/doc/UserGuide.md>, September 2016.
 - [58] Christine Paulin-Mohring. Modelisation of timed automata in coq. In *Theoretical Aspects of Computer Software, 4th International Symposium, TACS 2001, Sendai, Japan, October 29-31, 2001, Proceedings*,

- volume 2215 of *LNCS*, pages 298–315. Springer, 2001. doi:10.1007/3-540-45500-0_15.
- [59] Morgan Quigley, Brian Gerkey, Ken Conley, Josh Faust, Tully Foote, Jeremy Leibs, Eric Berger, Rob Wheeler, and Andrew Ng. Ros: an open-source robot operating system. In *Proc. of the IEEE Intl. Conf. on Robotics and Automation (ICRA) Workshop on Open Source Robotics*, Kobe, Japan, May 2009.
 - [60] Sanjay Rawat, Vivek Jain, Ashish Kumar, Lucian Cojocar, Cristiano Giuffrida, and Herbert Bos. VUzzer: application-aware evolutionary fuzzing. In *NDSS (Network and Distributed Security Symposium)*, 2017.
 - [61] Alastair Reid, Luke Church, Shaked Flur, Sarah de Haas, Maritza Johnson, and Ben Laurie. Towards making formal methods normal: meeting developers where they are, 2020.
 - [62] Virgile Robles, Nikolai Kosmatov, Virgile Prevosto, Louis Rilling, and Pascale Le Gall. MetAcsl: Specification and verification of high-level properties. In Tomáš Vojnar and Lijun Zhang, editors, *Tools and Algorithms for the Construction and Analysis of Systems*, pages 358–364, Cham, 2019. Springer International Publishing. ISBN 978-3-030-17462-0.
 - [63] RTCA Special Committee 167. Software considerations in airborne systems and equipment certification. Technical Report DO-178-B, RTCA, Inc., 1992.
 - [64] Kostya Serebryany. Continuous fuzzing with libfuzzer and addresssanitizer. In *Cybersecurity, Development (SecDev)*, IEEE, pages 157–157. IEEE, 2016.
 - [65] Yan Shoshitaishvili, Ruoyu Wang, Christophe Hauser, Christopher Kruegel, and Giovanni Vigna. Firmalice - automatic detection of authentication bypass vulnerabilities in binary firmware. In *NDSS*, 2015.
 - [66] Yan Shoshitaishvili, Ruoyu Wang, Christopher Salls, Nick Stephens, Mario Polino, Audrey Dutcher, John Grosen, Siji Feng, Christophe Hauser, Christopher Kruegel, and Giovanni Vigna. Sok: (state of) the art of war: Offensive techniques in binary analysis. In *IEEE Symposium on Security and Privacy*, 2016.
 - [67] Dokyung Song, Felicitas Hetzelt, Dipanjan Das, Chad Spensky, Yeoul Na, Stijn Volckaert, Giovanni Vigna, Christopher Kruegel, Jean-Pierre Seifert, and Michael Franz. Periscope: An effective probing and fuzzing framework for the hardware-OS boundary. In *NDSS*, 2019.
 - [68] Nick Stephens, John Grosen, Christopher Salls, Audrey Dutcher, Ruoyu Wang, Jacopo Corbetta, Yan Shoshitaishvili, Christopher Kruegel, and Giovanni Vigna. Driller: Augmenting fuzzing through selective symbolic execution. In *Network and Distributed System Security Symposium*, 2016.
 - [69] Jack Sun, Daniel Fryer, Ashvin Goel, and Angela Demke Brown. Using declarative invariants for protecting file-system integrity. In *Proceedings of the 6th Workshop on Programming Languages and Operating Systems*, page 6. ACM, 2011.
 - [70] Andreas Tiemeyer, Tom Melham, Daniel Kroening, and John O’Leary. Crest: Hardware formal verification with ansi-c reference specifications, 2019.
 - [71] Nikolai Tillmann and Jonathan De Halleux. Pex—white box test generation for .NET. In *Tests and Proofs*, pages 134–153, 2008.
 - [72] Nikolai Tillmann and Wolfram Schulte. Parameterized unit tests. In *ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 253–262, 2005.
 - [73] Nikolai Tillmann and Wolfram Schulte. Parameterized unit tests with Unit Meister. In *Proceedings of the 10th European Software Engineering Conference Held Jointly with 13th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 241–244, 2005.
 - [74] James R. Wilcox, Doug Woos, Pavel Panchekha, Zachary Tatlock, Xi Wang, Michael D. Ernst, and Thomas Anderson. Verdi: A framework for implementing and formally verifying distributed systems. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 357–368, New York, NY, USA, 2015. ACM. doi:10.1145/2737924.2737958.
 - [75] N. Williams, B. Marre, P. Mouy, and M. Roger. PathCrawler: automatic generation of path tests by combining static and dynamic analysis. In *EDCC*, 2005.
 - [76] Doug Woos, James R. Wilcox, Steve Anton, Zachary Tatlock, Michael D. Ernst, and Thomas Anderson. Planning for change in a formal verification of the raft consensus protocol. In *Proceedings of the 5th*

- ACM SIGPLAN Conference on Certified Programs and Proofs (CPP)*, pages 154–165, New York, NY, USA, 2016. ACM. doi:10.1145/2854065.2854081.
- [77] K. Yamamoto, Y. He, P. Heinrich, A. Orange, B. Ruggeri, H. Wilberger, and P.G. Flikkema. WiSARDNet field-to-desktop: Building a wireless cyberinfrastructure for environmental monitoring. In Charles van Riper III, Brian F. Wakeling, and Thomas D. Sisk, editors, *The Colorado Plateau IV: Shaping Conservation Through Science and Management*, pages 101–108. The University of Arizona Press, 2010.
 - [78] Zhi Yan, Nicolas Jouandeau, and Arab Ali Cherif. A survey and analysis of multi-robot coordination. *International Journal of Advanced Robotic Systems*, 10(12):399, 2013. doi:10.5772/57313.
 - [79] Michal Zalewski. american fuzzy lop (2.35b). <http://lcamtuf.coredump.cx/af1/>, November 2014.
 - [80] Lei Zhao, Yue Duan, Heng Yin, and Jifeng Xuan. Send hardest problems my way: Probabilistic path prioritization for hybrid fuzzing. In *NDSS (Network and Distributed Security Symposium)*, 2019.