

# Extending Source Code Generators for Evidence-based Software Certification

Ewen Denney

USRA/RIACS, NASA Ames  
edenney@email.arc.nasa.gov

Bernd Fischer

ECS, University of Southampton  
B.Fischer@ecs.soton.ac.uk

**Abstract**—Automated code generation offers many advantages over manual software development but treating generators as trusted black boxes raise problems for certification. Traditional process-oriented approaches to certification thus require that the generator be verified to the same level of assurance as the generated code, but this is infeasible for realistic generators. However, generators can be extended to support an *evidence-based* approach to certification. By careful design of the trusted kernel, assurance of the generator itself is not required.

In this paper, we describe several related extensions to two in-house code generators to provide two forms of evidence along with the code: safety proofs and safety explanations. We also describe how additionally provided links are used to trace between the code and the safety artifacts.

**Keywords:** automated code generation, safety, certification, qualification, evidence-based, user interfaces, theorem provers, traceability

## I. INTRODUCTION

Automated code generation is an enabling technology for model-based software development and has significant potential to improve the entire software development process. It promises many benefits, including reduced turn-around times, increased programmer productivity, and elimination of manual coding errors. However, the key to realizing these benefits is of course generator correctness—nothing is gained from replacing manual coding errors with automatic coding errors.

Consequently, a wide variety of techniques have been investigated to provide evidence that the generated code is correct. The existing approaches broadly fall into three different categories. In *certified* code generation, the code generator itself is certified (or *qualified*), using any technology that is appropriate or required by a certification authority. This category ranges from the systematic construction of generator test suites [1] over the application of compiler verification techniques [2] to the extraction of the code generator from a correctness proof in a logical framework like Isabelle [3] or Coq [4]. It also includes all process-oriented certification approaches, in particular code generator qualification as mandated by DO-178B [5]. In *certifying* code generation, the code generator simultaneously derives code and certificates. The best example for this approach is deductive program synthesis based on the proofs-as-programs principle, using an off-the-shelf theorem prover [6]. In *certifiable* code generation, the code generator is extended by a (separate) certification component that derives

a certificate for the generated code after the fact, using hints (e.g., loop invariants) provided by the generator, or by exploiting the idiomatic structure of the generated code. This category includes the different approaches to proof-carrying code (PCC) [7], [8], [9] as well as our own previous work [10], [11], [12].

In this paper, we present an integrated certifiable code generation system that combines program verification, proof checking, tracing, and explanation generation to support the evidence-based safety certification of automatically generated code. Following our previous work, we focus on the Hoare-style certification of specific safety properties (similar to the different PCC approaches) rather than showing full correctness of the generated programs. The evidence constructed by our system thus consists primarily of safety *proofs* but since certification is a social as much as a technical process, proofs in isolation from the program are not sufficient, and our system also supports *explanations* and *links* as equally important aspects of the evidence. We can thus consider the combination of proofs, explanations, and links as explicit *certificates*, i.e., independently checkable evidence of the claimed safety properties. The independent checking of the proofs can be automated by a separate proof checking tool that ensures that the formal proofs are solutions for the correct tasks associated with the given safety policy, and that their individual steps are correct. We have used this overall approach and the described tools to certify a variety of safety properties for code generated by the AUTOBAYES [13] and AUTOFILTER [14] systems. However, we concentrate on the certification extensions to the generators here, and omit details of the code generation process.

In the next section, we briefly provide the logical background of our safety certification approach. The following two sections then describe the use of proofs and explanations as evidence. Section 5 describes an interactive certification assistant, and Section 6 concludes.

## II. SOURCE-LEVEL SAFETY CERTIFICATION

**Safety Certification.** Software safety certification demonstrates that a program does not violate certain conditions during its execution. A *safety property* is an exact characterization of these conditions based on the operational semantics of the

language. A *safety policy* is a set of Hoare rules designed to show that safe programs satisfy the safety property of interest. We focus on source-level certification because (i) high-level domain-specific policies such as *frame safety* [15] can be formulated only on the source code level, and (ii) we are extending a source code generator. Source-level certification is complementary to object-level approaches like PCC and to ensure that compilation does not compromise the demonstrated safety policy, source-level certification should be followed by object-level certification.

For each notion of safety the appropriate safety property and corresponding policy must be formulated. This is usually straightforward; in particular, a safety policy can be constructed systematically from a safety property by instantiating a generic rule set that is derived from the standard rules of the Hoare calculus [10]. The basic idea is to extend the standard environment of program variables with a “shadow” environment of safety variables which record safety information related to the corresponding program variables. The rules are then responsible for maintaining this environment and producing the appropriate safety obligations. This is done using a family of *safety substitutions* that are added to the normal substitutions, and a family of *safety predicates* that are added to the calculated weakest preconditions (WPCs). Safety certification then starts with the postcondition *true* and computes the weakest safety precondition (WSPC), i.e., the WPC together with all applied safety predicates and safety substitutions. If the program is safe then the WSPC and all intermediate proof (i.e., safety and verification) obligations will be provable without any assumptions.

As an example, consider initialization safety, which ensures that each variable or individual array element has been explicitly assigned a value before it is used. Here, the safety environment consists of shadow variables  $x_{\text{init}}$  that contain the value `INIT` after the variable  $x$  has been assigned a value. Arrays are represented by shadow arrays to capture the status of the individual elements. The rules of the policy can be formulated in a “backwards” style and then used to compute the WSPCs. For example, the *for*-rule shown in Figure 1 says that for an arbitrary postcondition,  $Q$ , if  $c$  has WSPC  $P$  for the postcondition  $I[i+1/i]$ , and if the two intermediate obligations are true, then the WSPC of the loop is as shown. Since the *for*-statement assigns a value to the loop variable,  $i$ , it also affects the value of the corresponding shadow variable,  $i_{\text{init}}$  (reflected in the first intermediate obligation). The rule also applies the safety predicate  $\text{safe}_{\text{init}}$  to the immediate subexpressions  $e_1$  and  $e_2$  of the *for*-statement. Since the initialization safety property defines an expression to be safe if all corresponding shadow variables have the value `INIT`,  $\text{safe}_{\text{init}}(x[i])$  for example simply translates to  $i_{\text{init}} = \text{INIT} \wedge x_{\text{init}}[i] = \text{INIT}$ .

**Logical Annotations.** The *for*-rule, with its explicit loop invariant, highlights the central role logical annotations (i.e., pre- and postconditions and loop invariants) play in Hoare-style techniques. Fortunately, only relatively simple annotations are required, even for fully automated program proofs of the different safety properties. This is a consequence of the highly idiomatic structure of the automatically generated code and the restriction to specific safety properties. In our

*certifiable code generation* approach [11], the code generator itself is extended in such a way that it produces the necessary annotations together with the code. This is achieved by embedding annotation templates into the code templates, which are instantiated and refined in parallel by the generator. The generated logical annotations are then propagated throughout the code.

**Generating Obligations.** The annotated code is processed by a verification condition generator (VCG), which applies the rules of the safety policy in order to generate the safety obligations. As usual, the VCG works backwards through the code. At each statement, the safety predicates are added and the safety substitutions are applied. The VCG has been designed to be “correct-by-inspection”, i.e., to be sufficiently simple that it is straightforward to see that it correctly implements the rules of the logic. Hence, the VCG does not implement any optimizations or apply any simplifications; in particular, it does not actually apply the substitutions but maintains explicit formal substitution terms.

### III. PROOFS AS EVIDENCE

#### A. Proof Construction

**Simplification.** Since the VCG does not apply any optimizations and simplifications, the generated obligations tend to be large and to overwhelm current automated theorem provers (ATPs). They need to be simplified aggressively, therefore, before they can be submitted to an ATP with any hope of success. Our system thus includes several rewrite-based simplifiers. We focus on rewrite-based simplifications rather than decision procedures because rewriting is easier to certify: each individual rewrite step  $T \rightsquigarrow S$  could be traced and checked independently, e.g., by using an ATP to prove that  $S \Rightarrow T$  holds. However, this rewrite checking is not yet implemented.

**Processing Obligations.** The simplified safety obligations are exported as a number of individual proof obligations using TPTP first-order logic syntax [16]. For provers that do not accept the TPTP syntax, the appropriate (trusted) TPTP2X-converter is used before invoking the theorem prover. A small script then adds the axioms of the domain theory; parts of the domain theory are generated dynamically in order to facilitate reasoning with (small) integers. The completed proof task is processed by the ATP and the proof is stored as part of the certificate.

**Results.** We have evaluated several state-of-the-art ATPs on more than 25,000 proof tasks generated by our system. As expected, the unsimplified tasks prove to be too difficult for the provers, and only about two-thirds of the “out-of-the-box” tasks could be proven. After aggressive simplification, however, most of the provers could solve almost all emerging tasks. More details of the evaluation can be found in [17].

#### B. Proof Checking

Safety certification remains a challenging task for ATPs: the longest proof found during our previous experiments involved more than 8000 inference steps. Consequently, simple “correct-by-inspection” theorem provers like `leanTAP` [18]

$$\frac{P \{c\} \ I[i + 1/i] \ I[\text{INIT}/i_{\text{ini}}] \wedge e_1 \leq i \leq e_2 \Rightarrow P \ I[e_2 + 1/i] \Rightarrow Q}{I[e_1/i] \wedge e_1 \leq e_2 \wedge \text{safe}_{\text{init}}(e_1) \wedge \text{safe}_{\text{init}}(e_2) \ \{\mathbf{for} \ i := e_1 \ \mathbf{to} \ e_2 \ \mathbf{inv} \ I \ \mathbf{do} \ c\} \ Q}$$

Fig. 1. Hoare rule for *for*-loops

are not powerful enough. Instead, we need to employ high-performance ATPs, which use complicated calculi, elaborate data structures, and optimized implementations. This makes formal verification of their correctness infeasible [19]. Moreover, since most ATPs are under continuous development, single versions are never subjected to enough validation so that “tool pedigree” arguments remain weak.<sup>1</sup> In fact, despite the soundness checks applied in the CADE ATP System Competition (CASC), [21], there have been several unsoundnesses in participating ATPs, which have been detected only afterwards [22], [23], [24].

If the ATPs generate evidence in the form of sufficiently detailed proofs, they can be independently verified by a *proof checker*. Its function is to ensure that the ATP’s output really is a proof in the logical system in use. Techniques include the syntactic validation of Otter proof steps by Ivy [19], higher-order proof term reconstruction in Isabelle [25], higher-order proof step checking in HOL [26], reducing proof checking to type checking as in Coq [4], and semantic derivation verification [27], which has been used in this work. Here, the required semantic properties of each proof step are encoded in one or more *proof check obligations*, which are then discharged by trusted ATPs. If all obligations are discharged, the proof output of the original ATP is verified. This approach is tractable because the correctness proof for each individual step in the original proof is substantially easier than the original proof itself, and thus within reach of the trusted ATP. For certification purposes, all proofs found by the trusted ATP constitute evidence, and become part of the certificate constructed by the certification system.

**Semantic Derivation Verification.** The proofs produced by an ATP can be considered abstractly as *derivations*, i.e., directed acyclic graphs (DAG), whose leaf nodes are formulae (possibly derived) from the input problem (i.e., the original VC and any added axioms and lemmas), whose interior nodes are formulae inferred from parent formulae, and whose unique root nodes are the final derived formulae. Derivation verification then involves three notionally distinct phases. First, it is necessary to check the overall structure of the ATP output: specifically, that the derivation is a well-formed DAG. Second, it is necessary to check that each leaf node is a formula that occurs in, or is derived from, the input problem. This ensures that the ATP solves the original problem. Third, it is necessary to check that each inferred formula has the required semantic relationship (typically an implication from the premises of the applied inference rule to its conclusion) to its parents. This ensures that the proof is correct.

**Theorem Obligations.** For each application of an in-

ference rule that derives a logical consequence, a *theorem obligation* is formed to show that the inferred formula is indeed a logical consequence of the parent formulae. If the inference rule implements any theory (e.g., paramodulation implements most of equality theory), then the corresponding axioms of the theory are added as axioms to the obligation. The obligation is then handed to the trusted ATP system. If the trusted system finds a proof, the inference step is correct.

For inference rules that introduce branches in the original proof search, it is necessary for the checker to discharge multiple theorem obligations. For example, *explicit splitting* as implemented in SPASS [28] takes a problem  $S \cup \{L \vee R\}$  in clausal normal form (CNF), where  $L$  and  $R$  do not share any variables, replaces it by two subproblems  $S \cup \{L\}$  and  $S \cup \{R\}$  and searches for refutations of both subproblems. Obviously, both refutations must be checked to assure that a proof of the original problem has been found. In addition, to verify the splitting step’s role in overall proof, a theorem obligation to prove  $\neg(L \vee R)$  from  $\{\neg L, \neg R\}$  must be discharged.

**Leaf Theorem Obligations.** The leaf formulae of a derivation must occur in or be derived from the original input problem—otherwise, the ATP solves a different problem. To verify this, *leaf theorem obligations* to prove each leaf formula from the input formulae must be discharged. This makes the technique robust to some of the preprocessing inferences that are performed by ATP systems, e.g., factoring and simplification of the input. If the input problem is in first-order form (including quantifiers), and the derivation is a CNF refutation, the leaf clauses may have been formed with the use of Skolemization. Such leaf clauses are not logical consequences of the original input formulae. Skolemization steps can be incompletely verified by discharging a theorem obligation to prove the parent formula from the Skolemized formula. Although this is an incomplete verification step (i.e., unsound Skolemization steps can pass this check), it catches some simple errors and thus provides additional assurance.

**Experimental Evaluation.** As a practical test and evaluation of the proof checking approach described here, we scrutinized the proofs for 109 safety obligations generated from the certification of programs generated by the AUTOBAYES and AUTOFILTER code generators [17]. These obligations are also included as “difficult” problems in the TPTP problem library [16], the standard corpus for testing and evaluating ATPs.

The original proofs were found by the resolution provers EP (Version 0.82) [29] and by SPASS (Version 2.1). The proofs output by EP include details of the CNF-conversion, and the subsequent CNF-refutation while SPASS omits the CNF-conversion. Both systems are based on the superposition calculus, but differ in the specific inference rules used. Additionally, the systems have quite different control heuristics. As a result the proofs produced by the two systems have quite different characteristics.

<sup>1</sup>The notable exception is Otter [20], which has been essentially unchanged since 1996. However, our previous experiments have shown that its performance is not sufficient for discharging the safety obligations we generate [17].

The proof checking was done using the GDV system [27]. For the EP proofs, GDV was configured to check all aspects of each proof. For the SPASS proofs, GDV was configured to check only selected aspects of each proof: leaves were not verified because SPASS does not document the CNF-conversion, all inferred formulae and splitting steps were semantically verified, and the derivation was checked structurally, with the exception of structural aspects specific to splitting steps that were too time-consuming for the full set of proofs. We used Otter 3.3 [20] as trusted ATP for discharging theorem obligations. The experiments were run on Linux-based PCs with 2.8GHz and 1GB RAM, with a 10s CPU time limit for each discharge.

EP can solve 48 of the 109 problems, with 46 of the proofs fully verified. Both failure cases were caused by Otter’s inability to discharge obligations arising from steps in the CNF-conversion. In particular, the obligations to verify the step that negates the conjecture, which entails proving the negation of the negation from the original, could not be discharged. Most of the proofs induce less than 10 theorem obligations and only one proof induces more than 100 obligations. Most obligations were discharged quickly, with only three of the 590 obligations requiring more than 0.3s. SPASS can solve 83 of the 109 problems, which includes the 48 problems solved by EP, but the proofs are obviously different. All 83 of the SPASS proofs passed the verification checks chosen. Again, most of the proofs require less than 10 obligations to be discharged, but SPASS produces some very large proofs that consequently induce a very large number of obligations: 18 proof induces more than 100 theorem obligations and the largest proof resulted in 3493 obligations. At the same time, all 19737 SPASS obligations were discharged in less than 0.1s. More details can be found in [30].

**Results.** In the absence of proof checking, the applied ATP must become part of the trusted infrastructure, which substantially increases its size and complexity. Proof checkers are much smaller and can provide additional assurance, namely that the proofs correctly solve the original problem. We have applied semantic derivation verification to successfully check the safety proofs found by SPASS and EP. Our results indicate that the approach is feasible for the proofs found in this application domain, despite the substantial computational costs incurred by proof checking.

#### IV. EXPLANATIONS AS EVIDENCE

Although formal proofs can be an effective way of demonstrating safety and correctness, certification traditionally requires documentary evidence either that the software development complies with some process (e.g., DO-178B [5]), or that the artifacts are safe.

Treating a prover as a trusted black-box, or even checking its proofs, however, does not help in understanding why code is safe and is therefore difficult to reconcile with traditional approaches. Although proofs generated by an ATP can be verbalized, they are still difficult to understand and, more significantly, to relate to the actual program. We claim, however, that it is unnecessary to render actual proof steps; the

verification conditions alone provide sufficient insight into the safety of a program, can be related to the corresponding parts of the program, and can be rendered as comprehensible text. Based on this assumption, we have developed two related techniques based on extensions to the underlying logic: explanation of the VCs, which is useful for debugging and tracing; and explanation of program safety.

##### A. Explaining VCs

In practice, many things can—and typically do—go wrong with program verification: the program may be incorrect or unsafe, the annotations may be incorrect or incomplete, the simplifier may be too weak or counter-productive, the domain theory may be incomplete, and the ATP may run out of resources. In each of these cases, users are confronted only with failed VCs, but are left without any information about the causes of the failure. They must thus analyze the VCs by interpreting their constituent parts, and relating them through the applied Hoare rules and simplifications to the corresponding source code locations. Unfortunately, VCs are a detailed low-level representation of both the underlying information and the process used to derive it, so this is often difficult to achieve.

In this paper we sketch an implemented technique that helps users to trace, analyze, and understand VCs. The idea is to systematically extend the Hoare rules by “semantic mark-up” so that we can use the calculus itself to build up *explanations* of the VCs. This mark-up takes the form of *semantic labels* that are attached to the meta-variables used in the Hoare rules, so that the VCG then produces labeled versions of the VCs. The labels are maintained through the different processing steps, and are then extracted from the final VCs and rendered as natural language explanations. The main feature of VCs that we consider here is their *structure*. More domain-specific mark-up can be used to explain the *purpose* of VCs.

Figure 2 shows two different versions of a small example program to illustrate the process. Figure 2(a) shows the original annotations required (before propagation) to certify the program as initialization safe while Figure 2(b) shows the result of the propagation phase. Note that the propagation step already introduces some labels; for example, in line 7 the sub-formulas  $x_{\text{init}} = \text{INIT}$  and  $y_{\text{init}} = \text{INIT}$  are labeled with their original locations (i.e., lines 3 and 5). We use the notation  $[t]^{\text{lab}}$  to denote a term  $t$  that is labeled with a label  $\text{lab}$ ; the labels can also have internal structure.

Although the example is deliberately kept very simple, it already illustrates several of the difficulties that arise in interpreting VCs, in particular the combination of information from throughout the program into a single VC. The explanations become particularly complicated when the substitutions arising from the *assign*- and *array-update*-rules of the Hoare calculus are taken into account because the non-local effects of substitution applications need explaining: the sub-formulas from the annotations are no longer preserved intact and so need to be traced to their respective origins. For larger programs the overall structure quickly becomes complex. Figure 3 shows the automatically generated structural explanation for an example

<pre> 1  var i, x, y, z; 2  x := 1; 3  post x<sub>init</sub> = INIT 4  y := 2; 5  post y<sub>init</sub> = INIT 6  z := x + y; 7  post z<sub>init</sub> = INIT 8  for i := 0 to 2 9    inv true do 10   z := z * z; </pre>	<pre> var i, x, y, z; x := 1; post x<sub>init</sub> = INIT y := 2; post <math>\lceil x_{init} = \text{INIT}^{\text{orig}(3)} \wedge y_{init} = \text{INIT}</math> z := x + y; post <math>\lceil x_{init} = \text{INIT}^{\text{orig}(3)} \wedge \lceil y_{init} = \text{INIT}^{\text{orig}(5)} \wedge z_{init} = \text{INIT}</math> for i := 0 to 2   inv <math>\lceil x_{init} = \text{INIT}^{\text{orig}(3)} \wedge \lceil y_{init} = \text{INIT}^{\text{orig}(5)} \wedge \lceil z_{init} = \text{INIT}^{\text{orig}(7)} \wedge \text{true}</math> do   z := z * z; </pre>
(a)	(b)

Fig. 2. (a) Code with actual annotations (b) Code with annotations after propagation

VC derived using the full annotations in Figure 2(b). For presentation purposes we did not simplify the VCs as this would obscure their structure and complicate their understanding even further.

### B. Labeled Rules

Modified Hoare rules concisely capture the semantic mark-up (i.e., label types and positions) required for any given feature of the VCs that is to be explained. Labels can be added in four places: to the “incoming” postcondition of a recursive VCG call in the premise of an inference rule, to the WSPC, to a generated VC, or to a substitution. The labels are not dependent on the specific safety property and could contain additional embedded labels for more detailed or property-specific explanations.

We restrict our attention here to the *for*-rule shown in Figure 1, which we extend with semantic labels to give the rule in Figure 4. Note that we omit location information from all labels shown here in order to keep the presentation clear. The WSPC comprises the safety predicates and the invariant, which has to be established in the entry form (i.e., at the lower bound of the loop) and is thus labeled with  $\lceil \text{test\_inv}$ . In the premise, individual sub-formulas of both the exit-condition  $I \wedge \neg b \Rightarrow Q$  and the step-condition  $I \wedge b \Rightarrow P$  are labeled appropriately. In the triple for the step condition,  $P \{c\} I[i+1/i]$ , the incoming postcondition  $I[i+1/i]$  must be labeled with its purpose  $\lceil \text{test\_inv\_iter}$  for the recursive call; moreover, all emerging VCs must be marked up with the secondary purpose  $\lceil \text{pres\_inv}$ , meaning that they contribute to the preservation of the invariant. We indicate this by labeling the entire triple. In addition, the substitutions (more precisely, the right-hand sides of the individual replacements) are marked-up to record their type and the origin of the substituted expressions.

Note how the same formula  $I$  is used in four different roles and consequently labeled in four different ways. This contextual knowledge is only available at the point of rule application and can not be easily recovered by a post-hoc analysis of the generated VCs.

### C. Explanation Generation

The generation of the actual textual explanations is independent of the particular feature which is to be explained and

proceeds in two phases. First, there is a rewrite-based normalization of the VCs and corresponding labels. The rewrite rules used for the normalization of unlabeled VCs (Section III-A) are not label-aware and cannot be applied “as is” to the labeled case because (i) the labeling changes the term structure and thus the applicability of the rules and (ii) the labels need special handling. We have therefore defined a set of label-aware rewrite rules (omitted here) that are used together with additional unlabeled rules to simplify the labeled VCs.

The normalization is then followed by a rendering phase that extracts and further normalizes the final label structure and, using feature-specific explanation templates, turns it into natural language text.

### D. Explaining Program Safety

In contrast to Section IV-A, where individual VCs are rendered to give a *problem-centric* explanation of the verification, we can use the same underlying information (along with some more information about the program), to give a *program-specific* explanation.

Figure 5 gives an example program and the corresponding initialization safety explanation provided by the system. The program needs an invariant (not given here) in order to prove its safety. The explanation is only generated if the theorem prover successfully proves all the corresponding verification conditions. Note that we currently perform no symbolic evaluation during the rendering. The safety of the final assignment (line 14) is proven using the invariant but the explanation simply indicates where this is used (see [31] for more details).

## V. CERTIFICATION ASSISTANT

The previous two sections have discussed two important forms of evidence. As we have argued above, it is crucial for certification to relate this evidence to the program under consideration. We have therefore built a *certification assistant* that provides access to the auxiliary artifacts that are produced during the certification. This includes the intermediate stages in the processing chain (generated axioms, clausal normal form etc.), prover log files, and actual proofs, depending on the required level of evidence. These artifacts can support, or in the absence of a proof collectively serve as, the certificate, and can be inspected as raw text files, or using third-party tools,

The purpose of this proof obligation is to show that the loop invariant at line 9 under the substitution originating from line 10 is still true after each loop iteration; it is also used to show the preservation of the loop invariant at line 9. Hence, given

- the postcondition at line 3 propagated into the invariant at line 9,
- the postcondition at line 5 propagated into the invariant at line 9,
- the postcondition at line 7 propagated into the invariant at line 9,
- the invariant at line 9,
- the loop bounds at line 10,

show that the loop invariant at line 9 under the substitution originating from line 11 is still true after each iteration to line 11.

Fig. 3. Explanation automatically generated for the VC  $0 \leq i \leq 2 \wedge x_{\text{init}} = \text{INIT} \wedge y_{\text{init}} = \text{INIT} \wedge z_{\text{init}} = \text{INIT} \Rightarrow \text{INIT} = \text{INIT}$  derived from Figure 2(b)

$$\frac{\begin{array}{l} \lceil P \{c\} \lceil I[i+1]^{\text{sub}/i} \rceil^{\text{est\_inv\_iter}} \rceil^{\text{pres\_inv}} \\ \lceil I[\lceil \text{INIT} \rceil^{\text{sub\_safety}/i_{\text{init}}} \rceil^{\text{ass\_inv}} \wedge \lceil e_1 \leq i \leq e_2 \rceil^{\text{ass\_bounds}} \Rightarrow P \rceil^{\text{pres\_inv}} \\ \lceil I[e_2+1]^{\text{sub}/i} \rceil^{\text{ass\_inv\_exit}} \Rightarrow Q \end{array}}{\left( \begin{array}{l} \lceil I[e_1]^{\text{sub}/i} \rceil^{\text{est\_inv}} \wedge \lceil e_1 \leq e_2 \rceil^{\text{est\_bounds}} \\ \wedge \lceil \text{safe}_{\text{init}}(e_1) \rceil^{\text{safety}} \wedge \lceil \text{safe}_{\text{init}}(e_2) \rceil^{\text{safety}} \end{array} \right) \{ \text{for } i := e_1 \text{ to } e_2 \text{ inv } I \text{ do } c \} Q$$

Fig. 4. Hoare rule for *for*-loops with semantic markup for initialization safety

	Safety Explanation for Initialization of Variables
1 <b>var</b> a[0:9];	<i>The assignment b:=1 at line 6 is safe.</i>
2 <b>var</b> b;	<i>The assignment c:=2 at line 7 is safe.</i>
3 <b>var</b> c;	<i>The assignment d:=b*b+c*c at line 8 is safe; the term b is initialized from b:=1 at line 6; the term c is initialized from c:=2 at line 7.</i>
4 <b>var</b> d;	
5 <b>var</b> x;	<i>The loop index i ranges from 0 to 9 and is initialized at line 9.</i>
6    b:=1;	<i>The conditional expression i&lt;5 appears at line 10; the loop index i ranges from 0 to 9 and is initialized at line 9.</i>
7    c:=2;	
8    d:=b*b+c*c;	<i>The assignment a[d+i]:=i at line 11 is safe (if the condition i&lt;5 at line 10 is true); the term d is initialized from d:=b*b+c*c at line 8; the term b is initialized from b:=1 at line 6; the term c is initialized from c:=2 at line 7; the loop index i ranges from 0 to 9 and is initialized at line 9.</i>
9 <b>for</b> i:=0 to 9	
10 <b>if</b> i<5	
11       a[d+i]:=i;	<i>The assignment a[2*d-1-i]:=i at line 13 is safe (if the condition i&lt;5 at line 10 is false); the term d is initialized from d:=b*b+c*c at line 8; the term b is initialized from b:=1 at line 6; the term c is initialized from c:=2 at line 7; the loop index i ranges from 0 to 9 and is initialized at line 9.</i>
12 <b>else</b>	
13       a[2*d-1-i]:=i;	<i>The assignment x:=a[a[5]] at line 14 is safe; using the invariant for the loop at line 9 and the postcondition i:=9+1 after the loop.</i>
14    x:=a[a[5]];	<i>[Certified by e-setheo on Mon Mar 15 18:02:24 PST 2004 for init property.]</i>

Fig. 5. Left: Example program (annotations omitted) Right: Auto-generated explanation for *init* safety property

e.g., the GDV derivation verifier [27] and the proof visualizer from the TPTP tool suite [16].

The assistant also provides some limited functionality for creating proofs: it allows a (TPTP-compliant) prover to be chosen and invoked for selected VCs, and for the resulting proofs to be checked. We will concentrate here, however, on the assistant's use in tracing the VCs.

As discussed above, manually tracing VCs back to their source is quite difficult as the verification process is inherently complex and a single VC can depend on a variety of information distributed throughout the program.

Section IV-A described the mark-up for explanations. Since this includes location information, it can also be used to trace between the VCs and the source code. The VCG adds the

appropriate information to the formulas it constructs as it processes a statement at a given source code location. We currently use simple line numbers as locations rather than individual subterm positions [32].

Figure 6 shows how the tracing information can be used to support the certification process. A click on the source link associated with each verification condition prompts the certification assistant to highlight in boldface all affected lines of the code. A further click on the verification condition link itself displays the formula and explanation, which can then be interpreted in the context of the relevant program fragments. This helps domain experts assess whether the safety policy is actually violated when a proof attempt has failed, which parts of the program are affected, and eventually how the violation

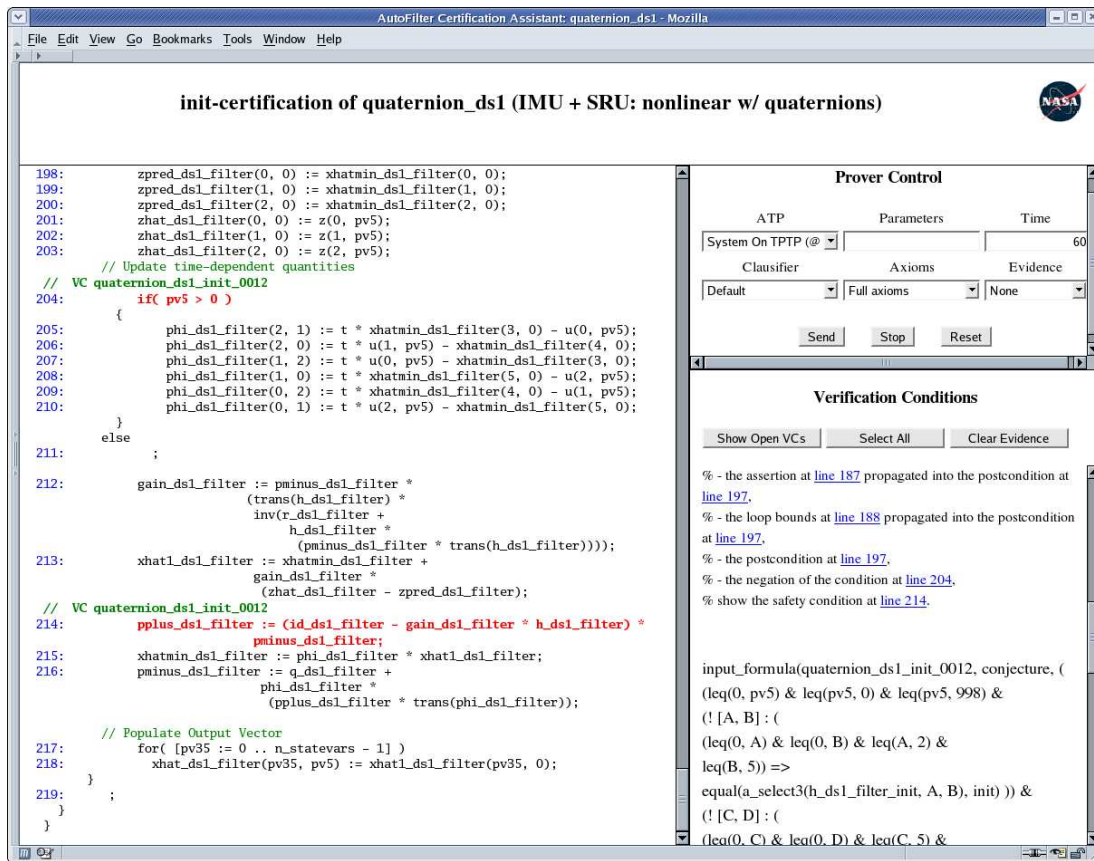


Fig. 6. Certification assistant: linking from VC

can be resolved. This traceability is also mandated by relevant standards such as DO-178B.

In practice, safety checks are often carried out during code reviews [33], where reviewers look in detail at each line of the code and check the individual safety properties statement by statement. To support this, linking works in both directions: clicking on a statement or annotation displays all VCs to which it contributes (i.e., which are labeled with its line number). Figure 7 shows the result of clicking on the label for line 220; the unproven verification condition indicates that this line of code has not been completely cleared yet.

## VI. CONCLUSIONS

We believe that there is a natural synergy between code generation and evidence-based certification. To gain trust in a black-box generator, it is necessary to have evidence that the generated code satisfies some desirable criteria. So long as the evidence is in a form that can be independently scrutinized, the generator can provide that evidence itself without loss of assurance. Since certification is ultimately a human process it is important to support both machine and human checking of evidence.

We have implemented a safety-proof based extension to two code generators that integrates the generation of safety proofs, safety explanations, and a browser-based assistant that allows tracing between the various generated artifacts. In more recent work [34], we have modularized the certification system in

such a way that it can be customized for third-party code generators.

Our long-term vision is that the system will support the construction of a safety case for the generated code, incorporating information about the generator itself, the code derivation, diverse forms of evidence, and customizable documentation.

## REFERENCES

- [1] I. Stürmer, D. Weinberg, and M. Conrad, "Overview of existing safeguarding techniques for automatically generated code," *SIGSOFT Software Engineering Notes*, 30(4):1–6, Jul. 2005.
- [2] V. L. Winter and J. M. Boyle, "Proving refinement transformations for deriving high-assurance software," in *Proc. High-Assurance Systems Engineering Workshop*. IEEE Comp. Soc. Press, 1996, pp. 68–77.
- [3] D. Basin, "The next 700 synthesis calculi," 2002, FME '02 Invited Talk.
- [4] Y. Bertot and P. Castéran, *Interactive Theorem Proving and Program Development - Coq'Art: The Calculus of Inductive Constructions*, EATCS Texts in Theoretical Computer Science. Springer-Verlag, 2004.
- [5] RTCA Special Committee 167, "Software considerations in airborne systems and equipment certification," RTCA, Inc., Tech. Rep., Dec. 1992.
- [6] M. Stickel, R. Waldinger, M. Lowry, T. Pressburger, and I. Underwood, "Deductive composition of astronomical software from subroutine libraries," in *Proc. 12th CADE*, LNAI 814. Springer, 1994, pp. 341–355.
- [7] G. C. Necula, "Proof-carrying code," in *Proc. 24th POPL*. ACM Press, 1997, pp. 106–119.
- [8] A. W. Appel, "Foundational proof-carrying code," in *Proc. 16th Annual IEEE Symp. Logic in Computer Science*. IEEE Comp. Soc. Press, 2001, pp. 247–258.
- [9] N. A. Hamid, Z. Shao, V. Trifonov, S. Monnier, and Z. Ni, "A syntactic approach to foundational proof-carrying code," in *Proc. 17th Annual IEEE Symp. Logic in Computer Science*. IEEE Comp. Soc. Press, 2002, pp. 89–100.

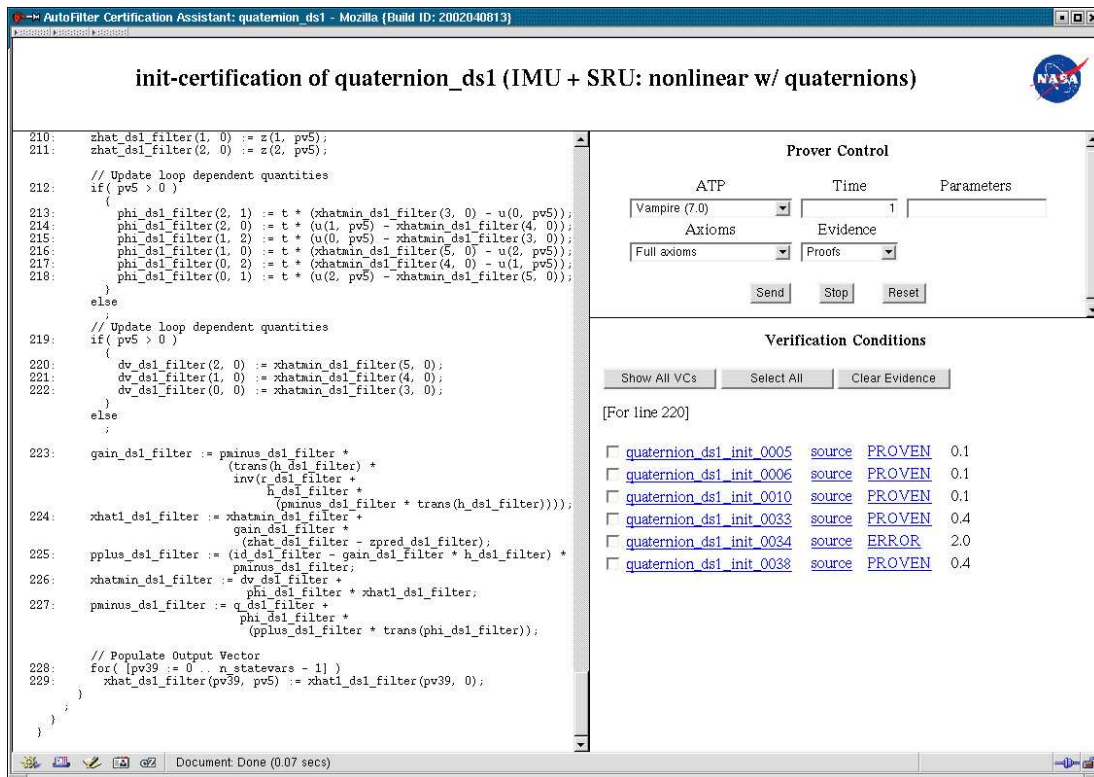


Fig. 7. Certification assistant: linking from code

- [10] E. Denney and B. Fischer, “Correctness of source-level safety policies,” in *Proc. FM 2003: Formal Methods*, LNCS 2805. Springer, 2003, pp. 894–913.
- [11] —, “Certifiable program generation,” in *Proc. Conf. Generative Programming and Component Engineering*, LNCS 3676. Springer, 2005, pp. 17–28. Invited Talk.
- [12] M. Whalen, J. Schumann, and B. Fischer, “Synthesizing certified code,” in *Proc. FME 2002: Formal Methods—Getting IT Right*, LNCS 2391. Springer, 2002, pp. 431–450.
- [13] B. Fischer and J. Schumann, “AutoBayes: A system for generating data analysis programs from statistical models,” *J. Functional Programming*, 13(3):483–508, May 2003.
- [14] J. Whittle and J. Schumann, “Automating the implementation of Kalman filter algorithms,” *ACM Trans. Mathematical Software*, 30(4):434–453, Dec. 2004.
- [15] M. Lowry, T. Pressburger, and G. Roşu, “Certifying domain-specific policies,” in *Proc. 16th ASE*. IEEE Comp. Soc. Press, 2001, pp. 118–125.
- [16] G. Sutcliffe and C. Suttner, “The TPTP problem library,” [www.tptp.org](http://www.tptp.org).
- [17] E. Denney, B. Fischer, and J. Schumann, “An empirical evaluation of automated theorem provers in software certification,” *International J. AI Tools*, 15(1):81–107, Feb. 2006.
- [18] B. Beckert and J. Posegga, “lean<sup>TP</sup>: Lean tableau-based deduction,” *J. Automated Reasoning*, 15(3):339–358, 1995.
- [19] W. McCune and O. Shumsky-Matlin, “Ivy: A Preprocessor and Proof Checker for First-Order Logic,” in *Computer-Aided Reasoning: ACL2 Case Studies*, Advances in Formal Methods 4, M. Kaufmann, P. Manolios, and J. Strother Moore, Eds. Kluwer Academic Publishers, 2000, pp. 265–282.
- [20] W. McCune, “Otter 3.3 Reference Manual,” Argonne National Laboratory, Argonne, USA, Tech. Rep. ANL/MS-C-263, 2003.
- [21] F. J. Pelletier, G. Sutcliffe, and C. B. Suttner, “The Development of CASC,” *AI Communications*, 15(2-3):79–90, 2002.
- [22] G. Sutcliffe and C. Suttner, “The CADE-15 ATP System Competition,” *J. Automated Reasoning*, 23(1):1–23, 1999.
- [23] G. Sutcliffe, “The CADE-16 ATP System Competition,” *J. Automated Reasoning*, 24(3):371–396, 2000.
- [24] —, “The IJCAR-2004 Automated Theorem Proving Competition,” *AI Communications*, 18(1), 2005.
- [25] S. Berghofer and T. Nipkow, “Proof terms for simply typed higher order logic,” in *Proc. 13th Intl. Conf. Theorem Proving in Higher Order Logics*, LNCS 1869, 2000, pp. 38–52.
- [26] W. Wong, “Validation of HOL proofs by proof checking,” *Formal Methods in System Design: An International Journal*, 14(2):193–212, Mar. 1999.
- [27] G. Sutcliffe and D. Belfiore, “Semantic Derivation Verification,” in *Proc. 18th Florida Artificial Intelligence Research Symposium*. AAAI Press, 2005, pp. 641–646.
- [28] C. Weidenbach, U. Brahm, T. Hillenbrand, E. Keen, C. Theobald, and D. Topic, “SPASS Version 2.0,” in *Proc. 18th CADE*, LNAI 2392. Springer, 2002, pp. 275–279.
- [29] S. Schulz, “E: A Brainiac Theorem Prover,” *AI Communications*, 15(2-3):111–126, 2002.
- [30] G. Sutcliffe, E. Denney, and B. Fischer, “Practical proof checking for program certification,” in *Proc. CADE-20 Workshop on Empirically Successful Classical Automated Reasoning (ESCAR’05)*, Jul. 2005.
- [31] E. Denney and R. P. Venkatesan, “A generic software safety document generator,” in *Proc. 10th Intl. Conf. Algebraic Methodology and Software Technology*, LNCS 3097. Springer, 2004, pp. 102–116.
- [32] R. Fraer, “Tracing the origins of verification conditions,” in *Proc. 5th Intl. Conf. Algebraic Methodology and Software Technology*, LNCS 1101. Springer, 1996, pp. 241–255.
- [33] S. Nelson and J. Schumann, “What makes a code review trustworthy?” in *Proc. Thirty-Seventh Annual Hawaii Intl. Conf. System Sciences (HICSS-37)*. IEEE, 2004.
- [34] E. Denney and B. Fischer, “A generic annotation inference algorithm for the safety certification of automatically generated code,” in *Proc. Conf. Generative Programming and Component Engineering*. ACM Press, 2006.