

1 Correct Compilation of Proofs About Embedded 2 Programs

3 **Audrey Seo***

4 University of Washington, USA

5 **Chris Lam***

6 University of Illinois Urbana-Champaign, USA

7 **Dan Grossman**

8 University of Washington, USA

9 **Talia Ringer**

10 University of Illinois Urbana-Champaign, USA

11 — Abstract —

12 Embedded programming languages and program logics make it possible to prove properties about
13 programs written in a variety of languages, all within a proof assistant. These embeddings consist
14 of implementations of the language and logic inside of the proof assistant language, sometimes
15 combined with frameworks for effective reasoning. By chaining these embeddings with verified
16 program compilers, proof engineers can reason about the source program, and get strong guarantees
17 about the target program by composition. But something is lost in that composite workflow: there
18 is no machine-checkable, source-independent proof object about the target program.

19 In this paper, we explore a different workflow that produces proof objects at the target level. The
20 workflow involves compiling embedded programs, specifications, and proof objects all at once, from
21 an embedded source language and logic to an embedded target language and logic. We implement
22 such a proof compiler in the Coq proof assistant, moving from an imperative language with variables
23 and functions, to a language where the only memory is a single call stack. Our proof compiler is
24 formally verified—it is constructive and correct by construction, with its type specifying the relation
25 between the source and target proofs. We believe this is the first proof compiler across program
26 logics that is formally proven correct.

27 **2012 ACM Subject Classification** Software and its engineering → Compilers; Software and its
28 engineering → Formal software verification

29 **Keywords and phrases** proof transformations, certified compilation, program logics, proof engineering

30 **Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

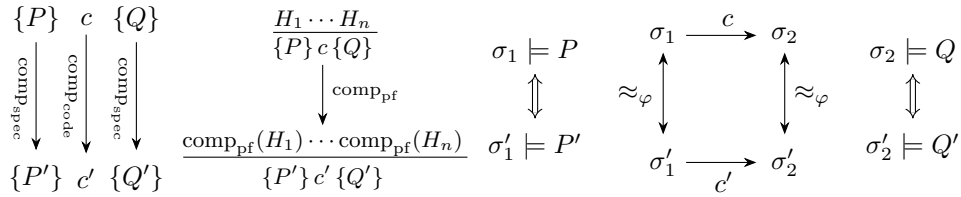
31 **Acknowledgements** This research was developed with funding from the Defense Advanced Research
32 Projects Agency. The views, opinions, and/or findings expressed are those of the author(s) and
33 should not be interpreted as representing the official views or policies of the Department of Defense
34 or the U.S. Government.

35 **1 Introduction**

36 Proof engineers sometimes need to escape the confines of the proof assistant language. They
37 may wish to reason instead about programs written in languages with features that the proof
38 assistant lacks, like effects or concurrency. In this position, they may turn to a framework
39 that embeds program logics for these languages directly within the proof assistant [41], like
40 Iris [21, 23], VST [9], CHL [10], or SEPREF [26]. By chaining these program logics with a
41 verified compiler, they can get strong guarantees about compiled programs [9].

* Co-first authors





■ **Figure 1** The idea behind proof compiler correctness. As our proof compiler’s actual correctness proof is via dependent types, its representation is merely the transformation of proof trees. Here, we feature code and spec compiler correctness. Section 4 describes the assumptions each arrow needs to satisfy in this diagram.

42 In this paper, we explore a different approach to getting strong guarantees about compiled
 43 programs. Our approach is to compile embedded programs, specifications, and proofs all at
 44 once—in the style of certificate translation [25] and proof-transforming compilation [33, 37, 17].
 45 The key benefit to this approach is that it transforms proofs about source programs directly
 46 to machine-checkable proofs about their compiled counterparts. Furthermore, it does so in a
 47 way that preserves the original specification, up to changes in the abstraction level.

48 To that end, this paper presents a formally verified proof compiler in Coq, which we
 49 call POTPIE (Proof Object Transformation Preserving Imp Embeddings). POTPIE compiles
 50 embedded programs, specifications, and proofs all at once, moving from an imperative
 51 language with variables and functions, to a language where the only memory is a single call
 52 stack. The proofs about compiled programs that POTPIE produces are machine-checkable,
 53 and reduce away references to the source code whenever possible.

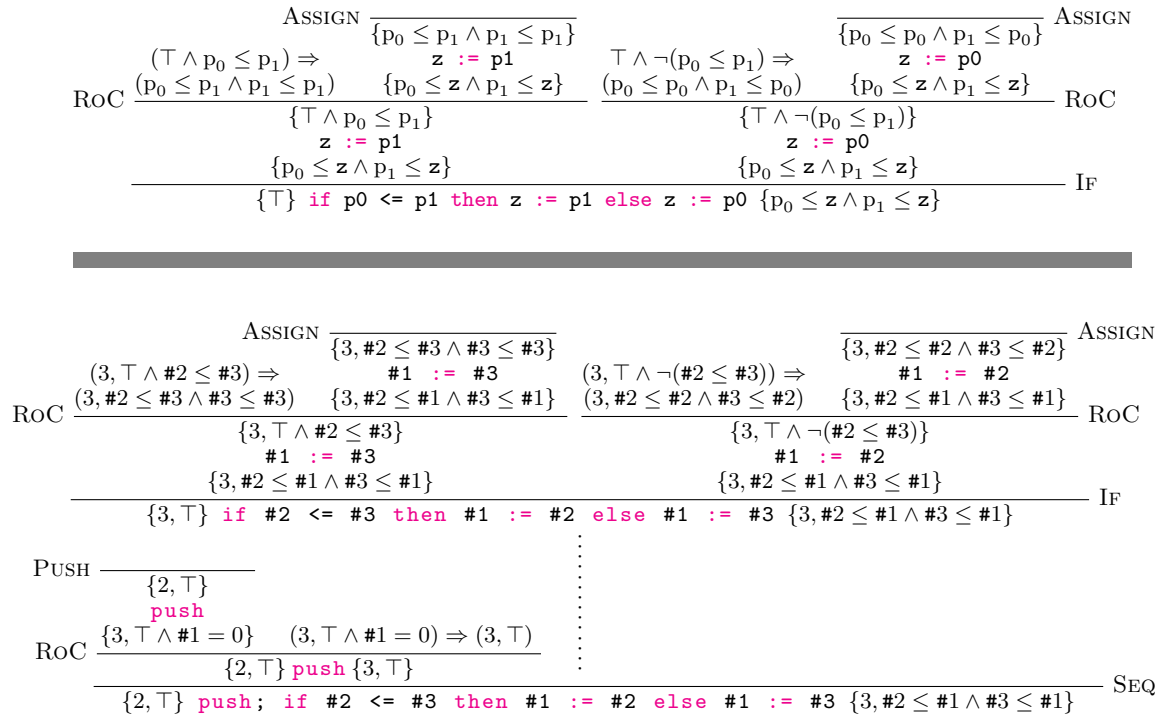
54 POTPIE is constructive and correct by construction. In particular, it is defined as a
 55 dependently typed function that compiles embedded source programs, specifications, and
 56 proofs in a way that is guaranteed by virtue of its type signature to respect the correspondence
 57 between the source and target language and logic. Accomplishing this was a significant proof
 58 engineering undertaking, which we discuss in detail. Our main contributions are twofold:

- 59 1. We introduce POTPIE, a formally verified, correct-by-construction proof compiler in
 60 Coq that compiles embedded programs, specifications, and proofs (Sections 2, 3, and
 61 4). POTPIE is, to the best of our knowledge, the first formally verified proof compiler
 62 across program logics (Section 6 discusses other projects, including those that did proof
 63 compilation but with only an informal proof of proof compiler correctness).
- 64 2. We describe the proof engineering considerations key to building a proof compiler like
 65 POTPIE (Section 5). This includes a discussion of how to address open challenges related
 66 to reducing away remaining references to the source language, compiling proofs efficiently,
 67 and removing the sole remaining dependency on total program compiler correctness.

68 Throughout the paper, we annotate claims supported by formalizations with a circled number,
 69 like (42). These icons are links to the POTPIE code, and are detailed further in `GUIDE.md`. Our
 70 hope moving forward is to reopen inquiry into work on producing source-independent proofs
 71 about compiled programs, with an eye on eventually linking proof objects about components
 72 of complex systems in which components may not share a common source language, and on
 73 eventually removing the need to prove the program compiler itself correct (Section 7).

74 2 Overview

75 POTPIE compiles source programs, specifications, and proofs to target programs, specifica-
 76 tions, and proofs. It does so soundly, meaning that the compiled proof proves the compiled



■ **Figure 2** The Hoare tree for the program `max` (above) and its compiled version (below). For brevity, we use p_i and p_i as shorthand for the i^{th} parameter in the specification and code, respectively.

77 specification about the compiled program, and that the compiled specification is the same
78 as the source specification up to the change in abstraction from the source language to the
79 target language. This is illustrated in Figure 1.

80 Our source and target languages are IMP and STACK. While both are imperative, IMP
81 uses variables that are local to a function and a distinct construct for accessing function
82 parameters. STACK, on the other hand, maintains only a single stack for storing values,
83 so the compiler must map variables and parameters to stack positions as well as manage
84 the stack size (via push and pop commands) to compile function calls. Due to this change
85 in abstraction, their specification languages (which encode the Hoare triple pre and post
86 conditions) are also different. For more details about the languages and logics, see Section 3.

87 To demonstrate end-to-end proof compilation, we consider a program that computes the
88 max of two numbers, along with a specification and proof that the specification is met:

```

89 {⊤}
90 1
91 2 if (param 0 ≤ param 1) then
92 3   z := param 1
93 4 else
94 5   z := param 0
95 6 {param 0 ≤ z ∧ param 1 ≤ z}

```

97 This program takes in two parameters (`param 0` and `param 1`), and sets the variable `z` to
98 the maximum of its arguments. The program's precondition, \top , indicates that there are
99 no requirements on its input. The program's postcondition, $z \geq \text{param } 0 \wedge z \geq \text{param } 1$,
100 indicates `z` is greater than both of the inputs.

101 We represent the proof of this specification explicitly as a proof object that corresponds
102 to a Hoare tree (Figure 2, top). Our goal is to compile this program, specification, and

23:4 Correct Compilation of Proofs About Embedded Programs

$$\begin{array}{ll}
a ::= \mathbb{N} \mid x \mid \text{param } k \mid a_1 + a_2 \mid a_1 - a_2 & a ::= \mathbb{N} \mid \#k \mid a_1 + a_2 \mid a_1 - a_2 \\
\quad \mid f(a_1, a_2, a_3, \dots, a_n) & \quad \mid f(a_1, a_2, a_3, \dots, a_n) \\
b ::= T \mid F \mid \neg b \mid a_1 \leq a_2 \mid b_1 \wedge b_2 \mid b_1 \vee b_2 & b ::= T \mid F \mid \neg b \mid a_1 \leq a_2 \mid b_1 \wedge b_2 \mid b_1 \vee b_2 \\
i ::= \text{skip} \mid x := a \mid i_1; i_2 & i ::= \text{skip} \mid \text{push} \mid \text{pop} \mid \#k := a \mid i_1; i_2 \\
\quad \mid \text{if } b \text{ then } i_1 \text{ else } i_2 \mid \text{while } b \text{ do } i & \quad \mid \text{if } b \text{ then } i_1 \text{ else } i_2 \mid \text{while } b \text{ do } i \\
\lambda ::= (f, n_{\text{args}}, x_{\text{ret}}, i_{\text{body}}) & \lambda ::= (f, n_{\text{args}}, i_{\text{body}}, \text{return } a_{\text{ret}} \ n_{\text{pop}}) \\
p ::= (\{\lambda_1, \dots, \lambda_n\}, i_{\text{main}}) & p ::= (\{\lambda_1, \dots, \lambda_n\}, i_{\text{main}})
\end{array}$$

■ **Figure 3** IMP (left) and STACK (right) language syntax, where a describes arithmetic expressions, b boolean expressions, i imperative statements, λ function definitions, and p whole programs, which consist of a set of functions and a “main” body. The evaluation of the main body yields the result of program. In the syntax for STACK functions, the construct $\text{return } a_{\text{ret}} \ n_{\text{pop}}$ means that the function returns the result of evaluating a_{ret} and then pops off n_{pop} variables from the stack.

$$\frac{\Lambda(f) = (f, k, y, i) \quad \langle \sigma, \Lambda, \Delta, a \rangle \Downarrow_a c_j \quad \langle \emptyset, \Lambda, [c_1, \dots, c_k], i \rangle \Downarrow_i \sigma' \quad \sigma'(y) = c_y}{\langle \sigma, \Lambda, \Delta, f(a_1, a_2, a_3, \dots, a_k) \rangle \Downarrow_a c_y} \text{FUN} \quad \frac{\sigma(x) = c}{\langle \sigma, \Lambda, \Delta, x \rangle \Downarrow_a c} \text{VAR} \quad \frac{|\Delta| \geq i \quad \Delta[i] = c}{\langle \sigma, \Lambda, \Delta, \text{param } i \rangle \Downarrow_a c} \text{PARAM}$$

■ **Figure 4** Big step semantics for the relevant cases of the IMP language.

103 proof object from the source language and logic to the target language and logic. The code
104 compiler $\text{comp}_{\text{code}}$ is just a conventional compiler that makes memory-layout decisions (i.e.,
105 where parameters and variables go on the stack). The specification compiler $\text{comp}_{\text{spec}}$ needs
106 to know the behavior of the code compiler so that the target specifications use the correct
107 stack positions and thus the right Hoare triples can be produced. The proof compiler comp_{pf}
108 is then just a recursive traversal over the source proof, using the code and specification
109 compilers in the natural way. This results in Figure 2, bottom. Note that the target proof
110 object is independent of the source language and logic.

111 3 Programs, Specifications, and Proofs

112 This section presents our six languages, with Section 3.1 describing programming languages,
113 Section 3.2 describing specification languages, and Section 3.3 describing proof languages.

114 3.1 Programs

115 We designed IMP and STACK to be similar enough to facilitate building a formally verified
116 proof compiler in a single pass, yet still have interesting properties and changes in abstraction
117 that are reminiscent of real-world source and target languages. IMP and STACK have similar
118 syntax (Figure 3), except for variables vs. stack slots ($\#k$) and the addition of push and
119 pop in STACK. Furthermore, parameters ($\text{param } i$) in IMP cannot be assigned to.

120 The basic infrastructure for both our IMP and STACK language semantics is based on the
121 materials for Xavier Leroy’s course on mechanized semantics [29]. Taking inspiration from
122 his approach, we define a big step semantics for our languages. Our semantics judgments for
123 IMP are of the form $\langle \sigma, \Lambda, \Delta, e \rangle \Downarrow r$, where σ is the variable environment, Λ is the function
124 environment, and Δ contains the values of parameters. The result, r , depends on the type of
125 e : it is another variable environment if e is a statement, a boolean value if e is a boolean

$$\begin{array}{c}
\boxed{\langle \sigma, \Lambda, a \rangle \Downarrow_a (\sigma', c)} \\
\frac{1 \leq k \leq |\sigma| \quad \sigma[k] = c}{\langle \sigma, \Lambda, \#k \rangle \Downarrow_a (\sigma, c)} \text{STKVAR} \quad \frac{\langle \sigma, \Lambda, a_1 \rangle \Downarrow_a (\sigma', c_1) \quad \langle \sigma', \Lambda, a_2 \rangle \Downarrow_a (\sigma'', c_2)}{\langle \sigma, \Lambda, a_1 + a_2 \rangle \Downarrow_a (\sigma'', c_1 + c_2)} \text{STKADD} \\
\frac{\Lambda(f) = (f, k, i, \text{return } a_{\text{ret}} \ n) \quad \forall j \leq k, \langle \sigma_{j-1}, \Lambda, a_j \rangle \Downarrow_a (\sigma_j, c_j) \quad \sigma' = [c_1, \dots, c_k] : \sigma_k \quad \langle \sigma', \Lambda, i \rangle \Downarrow_i \sigma'' \quad \langle \sigma'', \Lambda, a_{\text{ret}} \rangle \Downarrow_a (\sigma''', c_{\text{ret}}) \quad \sigma_f = \text{pop}(\sigma''', k + n)}{\langle \sigma_0, \Lambda, f(a_1, a_2, a_3, \dots, a_k) \rangle \Downarrow_a (\sigma_f, c_{\text{ret}})} \text{STKFUN} \\
\frac{}{\langle \sigma, \Lambda, \text{push} \rangle \Downarrow_i [0] : \sigma} \text{STKPUSH} \quad \frac{}{\langle [v] : \sigma, \Lambda, \text{pop} \rangle \Downarrow_i \sigma} \text{STKPOP} \quad \boxed{\langle \sigma, \Lambda, i \rangle \Downarrow_i \sigma'}
\end{array}$$

■ **Figure 5** Big step semantics for the relevant cases of the STACK language. Here, σ is a stack and Λ is the function environment. We denote concatenation by $\sigma_1 : \sigma_2$, and stack length by $|\sigma|$.

$$\begin{array}{c}
M ::= T \mid F \mid p_n \ e_{\text{list}} \\
\mid M \wedge M \mid M \vee M \\
\frac{}{\sigma \vDash T} \text{TRUE} \quad \boxed{\sigma \vDash M} \\
\frac{\text{map_eval}_\sigma \ a_{\text{list}} \ v_{\text{list}} \ p_{\text{list}} \ v_{\text{list}}}{\sigma \vDash p_{\text{list}} \ a_{\text{list}}} \text{N-ARY}
\end{array}$$

■ **Figure 6** Syntax (left) and semantics (right) for base assertions for both IMP and STACK. map_eval_σ is a relation from lists of expressions to lists of values. The semantic interpretation is parametric over the types of v , σ , and map_eval_σ . Interpretations for \wedge and \vee are standard.

126 expression, and a natural number if e is an arithmetic expression. Figure 4 contains the
127 inference rules for IMP that are either non-standard or are different from STACK's semantics.
128 Note we allow for function calls inside of arithmetic expressions. Functions are call-by-value,
129 and the function body (a command statement) is then evaluated on an empty variable
130 environment with the evaluated arguments as its parameter environment. Lastly, when
131 a call completes, it returns the value stored in the function's return variable in the local
132 environment. These semantics are formalized in Coq (13).

133 For STACK, our semantics judgments have the form $\langle \sigma, \Lambda, e \rangle \Downarrow r$, where Λ is again the
134 function environment, but σ is a stack. As in IMP, the type of r depends on the type of e ,
135 and r is a stack if e is a statement. But if e is an arithmetic expression, r is a *tuple* of a stack
136 and a natural number, since arithmetic expressions in STACK can have side effects (similarly
137 for booleans). Figure 5 contains the inference rules for STACK that are either non-standard,
138 or that are different from IMP's semantics. STACK adds several complexities—instead of
139 having an abstract environment where any variable can be queried and used, variables are
140 stored on a stack that can be directly manipulated via pushing, popping, and assigning.
141 Further, instead of each function evaluating on its own private environment, function calls
142 involve pushing variables and arguments onto the stack and later popping them. While in
143 this paper, we usually deal with programs that do not modify memory outside of their own
144 stack frame, there is no semantic restriction that prevents this. Later, we will utilize that our
145 compiled programs and functions *do* behave nicely in this respect and thus do not interfere
146 with the execution of other code. The semantics are formalized in Coq (14).

$$\begin{array}{c}
\frac{}{\{P\} \text{ skip } \{P\}} \text{IMPSKIP} \quad \frac{}{\{P[x \rightarrow a]\} x := a \{P\}} \text{IMPASSIGN} \quad \frac{\{P\} i_1 \{R\} \quad \boxed{\{P\} c \{Q\}}}{\{R\} i_2 \{Q\}} \text{IMPSEQ} \\
\frac{\{P \wedge b\} i_1 \{Q\} \quad \{P \wedge \neg b\} i_2 \{Q\}}{\{P\} \text{ if } b \text{ then } i_1 \text{ else } i_2 \{Q\}} \text{IMPIF} \quad \frac{\{P \wedge b\} i \{P\}}{\{P\} \text{ while } b \text{ do } i \{P \wedge \neg b\}} \text{IMPWHILE} \quad \frac{P_1 \Rightarrow P_2 \quad \{P_2\} i \{Q_1\} \quad Q_1 \Rightarrow Q_2}{\{P_1\} i \{Q_2\}} \text{RoC}
\end{array}$$

■ **Figure 7** Hoare deduction rules for the IMP language, where RoC stands for “rule of consequence.”

147 3.2 Specifications

148 Our program logics have the same base for assertions over the respective languages (Figure 6,
149 left). Here, p_n is an n -ary predicate over expressions, which we denote e . We offer support
150 for predicates of any length by allowing predicates that take lists of expressions. This specific
151 structure enables two important aspects of our system:

- 152 1. It makes the logic compilation very easy. In order to compile these assertions, all that
153 has to be done is to compile the expressions.
- 154 2. It permits the user of POTPIE to be as expressive as they wish, without adding new
155 relations or operators to the languages. For example, we do not provide a multiply
156 operator, but anyone writing proofs could use it in their specifications by defining a 2-ary
157 predicate. At the same time, POTPIE does not have to know extra mathematics.

158 The evaluation rules for assigning a truth value to a formula (Figure 6, right) define predicates
159 as parameterized over the value types. For example, if we have the proposition $p_1 a$ where a
160 is a language expression that evaluates to v , p_1 then evaluates on v as a Coq proposition.

161 We can instantiate the framework described in Figure 6 using the arithmetic and boolean
162 semantics relations as the evaluators. For IMP we refer to this instantiation as SM_e , and for
163 STACK as TM_e . We then use this to construct the following specification grammars:

$$\begin{array}{l}
164 \quad SM ::= SM_e \mid SM \wedge SM \mid SM \vee SM \quad (1) \\
165 \quad TM ::= (n, TM_e) \mid TM \wedge TM \mid TM \vee TM \\
166
\end{array}$$

Because the minimum stack size required by the compilation might not be captured by
formula itself, we also want to specify a minimum stack size in STACK specifications. This is
represented by the following judgement:

$$\frac{|\sigma| \geq n \quad \sigma \vDash TM_e}{\sigma \vDash (n, TM_e)} TM_e$$

167 For implementation details, see Section 5.1.

168 3.3 Proofs

169 Both IMP and STACK use a Hoare logic as their program logic, and we have proven them
170 sound (16) (17). In our Coq implementation, these are defined as inductive types that project
171 the logics into **Type**, because we need to access, match over, and transform the contents of
172 the proof in order to compile it.

173 Our IMP logic is standard (Figure 7; see the original paper on Hoare logic for more [18]).
174 While later on we need to introduce some termination conditions (see Section 5.2) in order
175 to properly compile the proofs, our IMP logic as a whole imposes no such condition.

$$\begin{array}{c}
\frac{}{\{A\} \text{ skip } \{A\}} \text{ STKSKIP} \quad \frac{\text{preservesStack}(b, n) \quad \{n, P \wedge b\} i \{n, P\}}{\{n, P\} \text{ while } b \text{ do } i \{n, P \wedge \neg b\}} \text{ STKWHILE} \quad \boxed{\{n, P\} c \{m, Q\}} \\
\frac{}{\{n, P\} \text{ push } \{n + 1, \text{inc}(P) \wedge \#1 = 0\}} \text{ STKPUSH} \quad \frac{}{\{(n + 1, \text{inc}(P) \wedge Q)\} \text{ pop } \{(n, P)\}} \text{ STKPOP} \\
\frac{\text{preservesStack}(a, n)}{\{n, P[\#k \rightarrow a]\} \#k := a \{n, P\}} \text{ STKASSIGN} \quad \frac{\{n_1, P\} i_1 \{n_2, Q\} \quad \{n_2, Q\} i_2 \{n_3, R\}}{\{n_1, P\} i_1; i_2 \{n_3, R\}} \text{ STKSEQ} \\
\frac{\text{preservesStack}(b, n_1) \quad \{n_1, P \wedge b\} i_1 \{n_2, Q\} \quad \{n_1, P \wedge \neg b\} i_2 \{n_2, Q\}}{\{n_1, P\} \text{ if } b \text{ then } i_1 \text{ else } i_2 \{n_2, Q\}} \text{ STKIF}
\end{array}$$

■ **Figure 8** Hoare rules for the STACK language. RoC is omitted, as it is identical to the IMP rule.

176 The STACK proof rules (Figure 8) are similar to the IMP ones, except for obvious language
177 differences (stack index accesses instead of variables, the addition of push and pop). There
178 is one other major way that the STACK rules differ from their IMP counterparts: the
179 preservesStack premise in the assign, if, and while cases, which specifies that the output
180 stack of the expression must be the same as the input stack. This is a simplifying assumption
181 that we add to avoid reasoning about how arbitrary functions can change the stack, since in
182 the STACK semantics for function calls (Figure 5), there are no guarantees about side effects.
183 Function calls could affect the entire stack, so even the addition rule is not pure. This can
184 present a welter of issues for how to track all the changes introduced by expressions, as well
185 as how to alter the Hoare logic rules to still keep them sound in such a setting.

186 For any expression e , preservesStack e says that all function call subexpressions of e
187 do not modify the rest of the stack frame. Because preserving stacks is dependent on the
188 function environment, the preservesStack predicate is implicitly parametric on a function
189 environment. Since function calls are the only source of effects in expressions, this means
190 that we have the following lemma:

191 ► **Lemma 1** (Stack Preservation ①). *If for STACK expression e we have preservesStack(e),*
192 *then for all stacks σ, σ' , if $\langle \sigma, \Lambda, e \rangle \Downarrow (\sigma', v)$ for some value v , we must have $\sigma = \sigma'$.*

193 While this assumption is fine for programs with function calls that do not change the stack
194 (such as the codomain of our compiler), if we want a more general program logic to reason
195 about the target language we need a more robust system to reason about both memory usage
196 (à la separation logic) and behavior within assertions that may change the state.

197 4 Proof Compilation

198 Now that we have our programming languages, specification languages, and program logics
199 defined, we describe in this section how we compile them in POTPIE.

200 In order to define our three compilers, we need to define an equivalence between the
201 environments of IMP (i.e., the variable store σ and parameter array Δ) and the stacks of
202 STACK (which we will represent by σ_s). This equivalence is entirely dependent on our choice
203 of a mapping between variables and stack slots. Intuitively, since parameters are always at
204 the top of the stack at the beginning of a function call, and are then pushed down as space

205 for local variables is allocated, the parameters should be “after” (i.e., appended to the end
206 of) the local variables. In other words:

207 ► **Definition 1.** *Let V be a finite set of variable names, and let $\varphi : V \rightarrow \{1, \dots, |V|\}$ be
208 bijective with inverse φ^{-1} . Then for all variable stores σ , parameter stores Δ , and stacks σ_s ,
209 we say that σ and Δ are φ -equivalent to σ_s , written $(\sigma, \Delta) \approx_\varphi \sigma_s$, if*

- 210 ■ for $1 \leq i \leq |V|$, we have $\sigma_s[i] = \sigma(\varphi^{-1}(i))$, and
- 211 ■ for $|V| + 1 \leq i \leq |V| + |\Delta|$, we have $\sigma_s[i] = \Delta[i - |V|]$.

212 Note that this definition implies that $|V| + |\Delta| \leq |\sigma_s|$ while also saying nothing about
213 the stack indices beyond $|V| + |\Delta|$. This will be important for ensuring that our compiler
214 correctness theorems also work for evaluating function bodies.

215 Building on this mapping from environments to stacks, Sections 4.1, 4.2, and 4.3 describe
216 our program, specification, and proof compilers, respectively.

217 4.1 Compiling Programs

218 Once we have a way to relate environments and stacks, our program compiler is fairly standard.
219 Given IMP program $(\{f_1, \dots, f_n\}, i_{\text{main}})$, where V contains the free variables of i_{main} , we
220 can construct bijective $\varphi : V \rightarrow \{1, \dots, |V|\}$. Once we have φ , syntactic compilation of code
221 is simple: replace a variable x with $\#(\varphi(x))$ and a parameter, param k , with $\#(|V| + k)$,
222 producing c' , which is STACK code. Finally, add $|V|$ pushes to the start of c' to make enough
223 room for local variables. While this compiler is simple, the transformation does allow us
224 to examine how program logics compile across abstraction gaps. Compiling functions and
225 function environments follows directly from program compilation.

226 We have proven the following correctness theorem for our simple compiler (for clarity we
227 have omitted assumptions about the well-formedness of compiler inputs—e.g., functions are
228 called with the right number of arguments—please see our formalization):

229 ► **Theorem 1 (Compiler Correctness (15)).** *For IMP arithmetic expressions e_a , boolean expres-
230 sions e_b , and commands c , function environment Λ , variable environment σ , and compiled
231 function environment Λ' , we have the following:*

- 232 ■ $\langle \sigma, \Lambda, \Delta, e_a \rangle \Downarrow_a v \wedge (\sigma, \Delta) \approx_\varphi \sigma_s \implies \langle \sigma_s, \Lambda', \text{comp}_{\text{expr}}(e_a) \rangle \Downarrow_a (\sigma_s, v)$
- 233 ■ $\langle \sigma, \Lambda, \Delta, e_b \rangle \Downarrow_a v \wedge (\sigma, \Delta) \approx_\varphi \sigma_s \implies \langle \sigma_s, \Lambda', \text{comp}_{\text{expr}}(e_b) \rangle \Downarrow_a (\sigma_s, v)$
- 234 ■ $\langle \sigma, \Lambda, \Delta, c \rangle \Downarrow_i \sigma' \wedge (\sigma, \Delta) \approx_\varphi \sigma_s \wedge (\sigma', \Delta) \approx_\varphi \sigma'_s \implies \langle \sigma_s, \Lambda', \text{comp}_{\text{code}}(c) \rangle \Downarrow_i \sigma'_s$

235 Because the relation \approx_φ ignores the stack beyond index $|V| + |\Delta|$, we can be assured that
236 function execution is correct as well since our compiler correctness theorems are robust to
237 the rest of the call stack. Further, if we assume termination of source programs, we get the
238 backwards direction of compiler correctness for free via determinism: if a compiled program
239 evaluates to a value, then the source program will evaluate to the same value ((11), (12)).

240 Our proof of correctness differs slightly from more traditional compiler correctness
241 theorems, as our semantics is big-step as opposed to small-step. Because of this, we don't
242 use the traditional simulation diagram to model the usual observable behavior mimicry.

243 4.2 Compiling Specifications

244 Once we have a program compiler and $\varphi : V \rightarrow \{1, \dots, |V|\}$, compiling specifications is also
245 relatively simple: recurse over the source logic formula and compile the leaves, i.e., IMP
246 expressions. If k is then number of function arguments, give each assertion a minimal stack
247 size, $|V| + k$, to ensure well-formedness of the IMP expressions within the specification, which
248 is given as the maximum value of φ plus k , where k is the number of arguments.

■ **Table 1** The proof engineering effort that went into stating and formalizing POTPIE, as well as creating the infrastructure to support the code and specification languages, the logics, and all of the compilers. (Here, “specs” means the number of **Definitions**, **Fixpoints**, and **Inductives**.)

Category	IMP			STACK			Base Assertions	Compiler			Total
	Lang	Logic	WF	Lang	Logic	Frame		Code	Spec	Proof	
LOC	776	1,832	3,465	3,142	1,064	5,344	921	3,858	2,680	6,256	29,338
Theorems	10	56	98	82	16	204	37	109	54	165	831
Specs	43	34	51	55	49	49	30	34	32	117	494

249 For any given specification compilation function $\text{comp}_{\text{spec}}^{\varphi,k}$ for given φ and k , we want
 250 this $\text{comp}_{\text{spec}}^{\varphi,k}$ to be sound. In this context, sound means that if a variable and parameter
 251 environment satisfy a formula P , then we want the translation of those environments to also
 252 satisfy $\text{comp}_{\text{spec}}^{\varphi,k}(P)$. More formally:

253 ► **Theorem 2.** *A specification compilation function $\text{comp}_{\text{spec}}^{\varphi,k}$ is sound if for all $P, \sigma, \Delta, \sigma_s$
 254 such that $(\sigma, \Delta) \approx_{\varphi} \sigma_s$, we have $\sigma, \Delta \models P \Leftrightarrow \sigma_s \models \text{comp}_{\text{spec}}^{\varphi,k}(P)$.*

255 We have formalized that our translation of specifications is sound for surely terminating
 256 specifications in Coq as ⑧ and ⑨. Because of the way that our specifications are designed,
 257 this proof followed relatively straightforwardly from the proof of compiler correctness.

258 4.3 Compiling Proofs

259 Because the control flow of the compiled program mirrors that of the source, the structure of
 260 the Hoare tree can largely stay the same. Therefore, in order to compile a proof, we need to
 261 prove two more results: (1) the constructors for each Hoare rule in our Hoare proof tree types
 262 commute across compilation, and (2) the implications inside of the rule of consequence hold
 263 across compilation. These two conditions, in addition to the correctness of the specification
 264 compiler, are sufficient to prove our correct-by-construction proof compiler. More formally:

265 ► **Theorem 3.** *Given an IMP Hoare proof of $\{P\}c\{Q\}$ (denoted as x), code compiler
 266 $\text{comp}_{\text{code}}$, and sound specification compiler $\text{comp}_{\text{spec}}$, a proof compiler comp_{pf} is sound if
 267 $\text{comp}_{\text{pf}}(x)$ is a sound STACK Hoare proof of the triple $\{\text{comp}_{\text{spec}}(P)\} \text{comp}_{\text{code}}(c) \{\text{comp}_{\text{spec}}(Q)\}$.*

268 For the precise formulation of the statement of the proof compiler as well as a discussion of
 269 its implementation details, see Section 5.

270 5 Proof Engineering

271 This was a substantial proof engineering effort; Table 1 summarizes the number of lines of
 272 code, proofs, and specifications for different components of POTPIE. Here we detail important
 273 design decisions (Section 5.1), assumptions (Section 5.2), and challenges (Section 5.3).

274 5.1 Design

275 We discuss five design decisions from engineering POTPIE: (1) making POTPIE correct
 276 by construction, (2) building modular assertion types, (3) separating language and logical
 277 reasoning, (4) using custom induction principles, and (5) building higher-order relations.

23:10 Correct Compilation of Proofs About Embedded Programs

278 **Correct-by-Construction Proof Compilation.** We made the deliberate decision to make
 279 our proof compiler correct by construction. That is, POTPIE is a function in Coq whose type
 280 specifies its correctness [⑦](#):

```

281
282 1   hl_compile : ∀ (P Q : log) (i : imp) (fenv : fun_env) (HL : hl P i Q fenv),
283 2     ∀ (map : list ident) (args : nat) (P' Q' : AbsState) (i' : stk) (fenv' : fun_env_stk),
284 3     (* ... well-formedness and termination conditions ... *) →
285 4     logic_transrelation args map P P' →
286 5     logic_transrelation args map Q Q' →
287 6     i' = compile_code i → fenv' = compile_fenv fenv →
288 7     hl_stk P' i' Q' fenv'.
  
```

290 where the type above is simplified for clarity, and where the omitted conditions are described
 291 in Section 5.2. In other words, this function takes an IMP Hoare logic proof $\mathbb{H}\mathbb{L}$ that shows
 292 $\{P\}i\{Q\}$ in function environment fenv , and compiles it to a STACK Hoare logic proof
 293 that shows $\{P'\}i'\{Q'\}$ in function environment fenv' , where each of P' , i' , Q' , and fenv'
 294 are correctly related to their IMP counterparts. We chose this design for its elegance and
 295 simplicity. Notably, it echoes three approaches to encoding proofs in proof assistants:

- 296 1. The **LCF Approach**: The hl and $\mathit{hl_stk}$ types that encode IMP and STACK Hoare logic
 297 proofs are each a datatype whose constructors are precisely the rules of the corresponding
 298 logic, hence all terms of that type must be valid proofs that are correct by construction—as
 299 with terms that have the ML thm type in the *LCF approach* [15].
 - 300 2. The **de Bruijn Criterion**: These correct-by-construction proofs are themselves explicit
 301 proof terms in Coq’s proof term language that represent Hoare trees, and that can be
 302 checked by a small proof checking kernel—thereby satisfying the *de Bruijn criterion* [3, 4].
 - 303 3. The **Curry-Howard Correspondence**: Thanks to Coq’s rich type system, we can
 304 encode the specification that each Hoare logic proof term proves as its type, by making the
 305 hl and $\mathit{hl_stk}$ types dependent—taking advantage of (generalized) *Curry-Howard* [14, 19].
- 306 From these combined approaches, we get correct-by-construction, machine-checkable, explicit
 307 proof objects coupled with simple specifications that correspond to types in Coq. But this
 308 design is not without cost. Most notably, it tightly couples proofs and computation, which
 309 poses challenges for efficiency and reduction (see Section 5.3). We are currently weighing
 310 whether to decouple proofs and computation, considering what we may lose in the process.

311 **Modular Assertions.** Our implementation of assertions in Coq are modular, in that the
 312 datatype $\mathit{LogicProp}$ that defines the syntax for assertions M is parameterized over two types:
 313 an expression type A and a value type V . As a result, we can provide any one of our
 314 expression types, its corresponding value type, and the semantics relation for that type, and
 315 end up with a model of this logic. Because of this modularity, the atoms SM_e and TM_e
 316 referenced in Section 3.2 are actually each two instantiations of the $\mathit{LogicProp}$ framework,
 317 one each for arithmetic and boolean in both IMP and STACK.

318 **Separating Language Reasoning and Logical Reasoning.** Our system separates *language*
 319 *semantic* reasoning and *logical* reasoning. The language semantics reasoning occurs entirely
 320 inside of the expressions e that are arguments to the logical predicates p , whereas the logical
 321 reasoning occurs only once the expressions have been evaluated to values. This is needed
 322 because, when compilation changes the semantics of programs and programmatic functions
 323 are included in arithmetic statements, program semantics can no longer be abstracted away
 324 in specifications, as different languages may have different semantics. With this design, the
 325 proof engineer can separate their logical reasoning from the language. For example, if they

326 need to use the fact that 5 is less than 7, then they can do so in the logical fragment. If
 327 they need to use the fact that the language expression $5 + 2$ evaluates to 7, then that is also
 328 expressible, and distinct from what it means purely mathematically.

329 **Custom Induction Principles.** One tried-and-true [44, 13] proof engineering design principle
 330 we employed was the use of custom induction principles. For example, for arithmetic
 331 expressions in IMP, we represented function arguments as lists of arbitrary lengths. We then
 332 realized that Coq’s automatically generated induction principles were too weak. With the
 333 help of Clément Blaudeau [7], we defined a better induction principle ⑤:

```
334
335 1  ∀ (P : aexp → Prop),
336 2    (* ... other cases ... *) →
337 3    (∀ (f : ident) (aexps : list aexp), Forall P aexps → P (f · aexps)) →
338 4    ∀ (a : aexp), P a.
```

340 where `Forall P aexps` is inhabited whenever the inductive motive `P` holds on all elements of
 341 the list `aexps`. This produces a stronger inductive hypothesis in the function application
 342 case. We used a similar trick to define an induction principle for IMP that gives informative
 343 inductive hypotheses about arithmetic and boolean subexpressions. We hope, going forward,
 344 that third-party automation [43, 30] for producing better induction principles for containers
 345 will make it into Coq’s standard library, sparing proof engineers this kind of effort.

346 **Higher-Order Relations.** We used higher-order relations (relations that take in other
 347 relations) to facilitate engineering and proving. This allowed us to lift relations from
 348 arithmetic and boolean expressions to the assertions those expressions were embedded within.
 349 From an engineering perspective, this reduced the number of names we had to remember,
 350 since we could simply lift expression-level relations without defining new names. From a
 351 proof perspective, it gave us some adequacy theorems “for free,” by lifting proofs about
 352 expressions to the relation level. For example, this is how we proved Theorem 4, which
 353 allowed us to extend adequacy up to types that were parameterized on that type:

354 ► **Theorem 4 (Higher-Order Adequacy For Free ②).** *If V, A are types and $l, l' : \text{LogicProp } V A$,*
 355 *for any $\varphi : A \rightarrow A$ and $\psi : A \rightarrow A \rightarrow \text{Prop}$, if for all $a_1, a_2 \in A$ we have $a_2 = \varphi a_1$ if and*
 356 *only if $\psi a_1 a_2$, then $l' = \text{transform_prop_exprs } l \varphi \Leftrightarrow \text{transformed_prop_exprs } \psi l l'$.*

357 5.2 Assumptions

358 POTPIE produces proof objects that are both correct by construction and checkable against
 359 the compiled specification. Its *trusted computing base* (TCB) consists of just the Coq kernel,
 360 the mechanized semantics, and two localized axioms for reasoning about dependent types.
 361 Its compiler correctness proof provides some confidence that the semantics are implemented
 362 correctly, at least with respect to one another. To keep the TCB small, it takes as input
 363 additional proof obligations, which we plan to relax in the near future.

364 **Localized Axioms.** Both axioms that we assume are localized instances of Uniqueness of
 365 Identity Proofs (UIP) over particular types, without broadly implying UIP. UIP, which is
 366 consistent with Coq, states that all equality proofs are equal to each other for *all* types—we
 367 instead assume that all equality proofs are equal to each other for *two particular types*. We
 368 place this restriction on UIP because assuming UIP in general would be computationally
 369 undesirable, since UIP does not have a constructive interpretation in Coq. We assume UIP
 370 for `AbsEnvs` (the Coq implementation of SM in Equation 1) and function environments.

371 ► **Axiom 1** (UIP for `AbsEnv`s ③). *If $a, a' : \text{AbsEnv}$, then if $p_1, p_2 : a = a'$, we have $p_1 = p_2$.*

372 ► **Axiom 2** (UIP for Function Environments ④). *If Λ, Λ' are function environments, and we*
 373 *have proofs $p_1, p_2 : \Lambda = \Lambda'$, then $p_1 = p_2$.*

374 Axiom 1 is reasonable to assume since we do not ever compute on the contents of
 375 `AbsEnv`s; they really are computationally irrelevant. Standing in the way of provable UIP is
 376 the presence of functions (representing predicates in our base assertions from Section 3.1) in
 377 the constructors of `LogicProp`, one of the types that makes up `AbsEnv`. But we never compute
 378 on the contents of these functions. Furthermore, even though Axiom 1 implies UIP for
 379 `LogicProp`, by requiring that `LogicProp`'s parameters V and A have decidable equality ⑩,
 380 UIP provably holds for V and A , so Axiom 1 does not imply UIP more generally. Likewise,
 381 Axiom 2 is a reasonable assumption because the function environment in both languages
 382 is considered to be global—there can only ever be one function environment. Furthermore,
 383 as function environments are functions whose domain and codomain both have decidable
 384 equality and thus UIP, Axiom 2 also does not imply UIP more generally.

385 We may be able to avoid these axioms by adapting the definitions of `AbsEnv` and function
 386 environments, or by rephrasing some inductive relations as `Fixpoints` that project into `Prop`.
 387 We have done the latter for one of our later theorems, but have found adapting earlier
 388 definitions and relations to be invasive. Assuming UIP over these two types made it possible
 389 to write proofs by dependent induction with minimal changes to our proof development.

390 **Proof Obligations.** To minimize our TCB, `POTP`IE requires some proof obligations from
 391 the caller, beyond just the source proof that is being compiled. One such proof obligation is
 392 to prove that Hoare proofs must have only surely terminating assertions. This simplifies the
 393 backwards program compiler correctness proof, which follows from forward correctness and
 394 determinism. It also simplifies the semantics of our Hoare logic. However, this comes with a
 395 tradeoff. Beyond the additional user proof obligation imposed, this assumption also prohibits
 396 calling functions that may have non-terminating behavior, ruling out entire classes of useful
 397 proofs. We hope to sort out the semantics of our Hoare logic in the face of potentially
 398 nonterminating preconditions and postconditions so that we can relax this in the future.

399 We also take well-formedness of the source code and function environments as proof
 400 obligations, requiring that, for example, functions are called on the correct number of
 401 arguments. We plan to relax this assumption as well. One potential way to relax this
 402 assumption is to replace our current program compiler with one that rejects ill-formed source
 403 programs. Another potential way is to prove an adequacy condition that lets us loosen the
 404 proof obligation to proving an equality with the result of calling our current compiler. These
 405 assumptions are documented in the assumptions required to call the proof compiler ⑦.

406 5.3 Challenges

407 In developing `POTP`IE, we encountered several interesting challenges, both specific to our
 408 project and also more broadly related to Coq. Here we describe some of these challenges and
 409 how we addressed them: obtaining source-independent proofs, efficiently compiling proofs,
 410 expanding our source logic, and handling mostly-correct compilers.

411 **Getting Source-Independent Proofs.** One major benefit of compiling proof objects from
 412 a source logic to a target logic is that the compiled proofs can be free of references to
 413 the source language and logic. We are well on our way to source-independence—but fully
 414 achieving this introduces its own challenges. For example, due to the tight coupling of our

415 correctness proof and computation in POTPIE, source code well-formedness principles may
 416 show up in compiled proofs. In addition, opaque proofs in Coq may prohibit reduction that
 417 would otherwise remove source program and proof references. We see three paths toward
 418 full source-independence: (1) relaxing proof obligations for the user, (2) decoupling the
 419 correctness proof from the computation, and (3) setting opaque proofs to transparent. Each
 420 path comes with its own challenges—the first to engineering, the second to elegance and
 421 simple specifications, and the third to efficiency. We are actively exploring these paths.

422 **Efficient Proof Compilation.** One of the factors that prevented us from computing source-
 423 independent proofs was the inefficiency of Coq’s term reduction, especially when combined
 424 with proof script automation. Proof script automation through specialized tactics can be
 425 effective for reducing engineering effort. But such automation rarely creates simple proof
 426 terms, which is important when computation matters, as it does for POTPIE. For example,
 427 automation often introduces highly nested proof terms that may take minutes to reduce. In
 428 one case, reducing such a term even crashed CoqIDE!

429 One way to get more efficient computation is to use opaque proofs, which is the default
 430 in Coq when one uses `Qed` instead of `Defined` to end a proof. For us, this was often not an
 431 option, since opaque proofs block reduction and thus computation. We needed transparent
 432 proof terms since computation was relevant, so we often used `Defined`. The only way to
 433 ensure our terms computed while staying transparent was to *manually* optimize proof terms.
 434 We used this process to do so:

- 435 1. Prove a given theorem using automation.
- 436 2. Print the term automatically generated by the proof script for that theorem. Locate
 437 inefficiencies in that term, and match those inefficiencies to tactics in the proof script.
- 438 3. Simplify the term, typically by working at the tactic level.
- 439 4. Repeat until the term efficiently reduces.

440 The details varied by the source of inefficiency. For example, we often had to simplify proof
 441 terms generated using the `inversion` tactic, which performs (possibly dependent) pattern
 442 matching to pull premises into the proof state. Consider the `WHILE` case of the mutually
 443 inductive well-formedness relation for function applications ⑤:

```
444
445 1 | fun_app_while : ∀ fenv wf_funcs b i,
446 2   fun_app_bexp_well_formed fenv wf_funcs b → (* If the guard b is well formed, *)
447 3   fun_app_imp_well_formed fenv wf_funcs i → (* and so is the expression i, *)
448 4   fun_app_imp_well_formed fenv wf_funcs (WHILE b i). (* then so is the loop. *)
```

450 In proofs with hypotheses applying `fun_app_imp_well_formed` to some loop `WHILE b i`, it was
 451 helpful to invert those hypotheses to derive the well-formedness of the guard `b` and the `IMP`
 452 expression `i`. This pattern was so common that, for large relations like this one, we had
 453 defined custom inversion tactics. These custom tactics very quickly led to case explosion. To
 454 optimize case-exploding proofs, we refactored inversion into opaque lemmas, like ⑥:

```
455
456 1 Lemma inv_fun_app_imp_while : ∀ fenv funcs b i,
457 2   fun_app_imp_well_formed fenv funcs (WHILE b i) →
458 3   fun_app_bexp_well_formed fenv funcs b ∧ fun_app_imp_well_formed fenv funcs i.
```

460 We then left the rest of those proofs transparent. This workflow brought down the reduction
 461 time of the term that had previously crashed CoqIDE to the order of seconds.

462 Some of our reduction woes may be preventable by replacing some of our inductive
 463 relations with `Fixpoints` that project into `Prop`. We would also like a notion of semi-opaque
 464 subterms, on which Coq’s full reduction (normalization) process could act like the `simpl`

23:14 Correct Compilation of Proofs About Embedded Programs

465 tactic (by *refolding* [8] constants) without ever fully reducing. Going forward, we would love
466 more work on making Coq’s reduction faster, easier to control, and simpler to debug. Finally,
467 we would love tools for automatic proof term optimization.

468 **Expanding Our Stack Logic** For convenience of implementation, the STACK program logic
469 semantics we introduced in Section 3.2 is limited in the programs it can describe. While it is
470 sufficient to handle proofs compiled from IMP, we would eventually love to be able to link
471 those proofs with other STACK programs’ proofs. To do this, we plan to derive a program
472 logic semantics that fully describes what happens with side effects for STACK. For example,
473 if we know that a function may access up to the fifth index outside of its stack frame, we
474 could remove all knowledge in our assertions having to do with the affected stack indices. We
475 could even define our current STACK logic semantics as a special case of the richer semantics,
476 so proofs produced by POTPIE could still interface with proofs that use the full STACK logic.

477 **Handling Sometimes-Incorrect Compilation** POTPIE currently assumes that the program
478 compiler itself is verified. This is a reasonable start, but we would love to avoid the need
479 for a verified program compiler. In fact, we suspect that one of the key benefits of proof
480 compilation ought to be handling compilation that is sometimes *incorrect*—for example,
481 unsafe optimizations, so long as they do not impact the truth of the compiled specification.

482 We attempted to build POTPIE without this assumption initially, but we found this
483 challenging. In order to translate the rule of consequence, POTPIE as implemented must
484 ensure that arbitrary logical implications hold through the logic translation, which requires a
485 proof of compiler correctness; this is the only place in which that proof is needed.

486 This challenge arises from including functions calls in arithmetic expressions. Unlike in
487 traditional uses of Hoare logic, the semantics of arithmetic do not entirely line up with our
488 semantics—one must also reason about the control flow of functions. This means that the
489 implications inside an application of consequence cannot be proven just using arithmetic,
490 but rather must reason about the semantics of the program. When compiling programs,
491 the semantics and level of abstraction change, and so soundly compiling proofs relies on
492 those semantics changing soundly. We have three options: (1) add a proof obligation that all
493 source and compiled implications hold, (2) rely on a proof of compiler correctness, or (3)
494 rework the rule of consequence so that it no longer requires a proof of compiler correctness.

495 For now, we choose option (2), since it is easy to generalize and, unlike option (1), avoids
496 exposing confusing proof obligations about target proofs to users. One path toward option
497 (3) that we are considering is specifying that all implications must be of the form of a Hoare
498 triple for a function. Such a change would allow us to view all expressions inside of our
499 specifications exclusively as deeply embedded language constructs and eject all Coq-level
500 reasoning to the relation in our predicate construct. The end result would still require some
501 target-level proof obligations from the user, but those proof obligations would be simpler than
502 those from option (1). We will continue to investigate option (3), alongside ways to relax the
503 proof obligations in option (1), to later smoothly handle sometimes-incorrect compilation.

504 **6 Related Work**

505 **Proof-Transforming Compilation.** Early work compiling proofs positions itself as an exten-
506 sion of *proof-carrying code* [35], in which the compiler produces proofs of important properties
507 about compiled code. Work in 2006 [5] stated a theorem relating source and target program
508 logics reasoning about languages similar to those we use in our work. There was not yet

509 transformation of proof objects themselves. Transformation of proof objects followed shortly
510 after, going by the names *proof-transforming compilation* and *certificate translation*.

511 Early proof-transforming compilation work [33] transformed proofs about Java-like pro-
512 grams to proofs about bytecode. Proofs were written in Hoare-style program logics defined
513 over the respective source and target languages, and implemented by way of a prototype tar-
514 getting proofs written in XML. Later work [37] implemented proof-transforming compilation
515 from Eiffel to bytecode, formalizing the specification compiler in Isabelle/HOL, and writing
516 a hand-written proof of correctness of the proof compiler. A subsequent thesis [17] showed
517 how to embed the compiled proofs about bytecode into Isabelle/HOL. Our work is the first
518 we know of to formally verify the proof of correctness of the proof compiler.

519 Certificate translation [25] compiles proofs in a program logic as well, but focuses on
520 compiler optimization passes, like dead code elimination or constant propagation. While it
521 does not come with an accompanying formal proof development, revisiting work on certificate
522 translation may prove fruitful should we move to support optimization passes as well.

523 **Proof Term Transformations.** There is some work on proof term transformations that hap-
524 pen directly over proof assistant proofs, not inside of any program logic. Proof transformations
525 for higher-order logic were introduced in 1987 to bridge automation and usability [38]. They
526 have since been used for many different purposes, like automatic generalization [16, 22, 20],
527 reuse [31], and repair [40] of proofs. Our work implements a proof term transformation
528 directly within the Coq proof assistant for an embedded program logic, which empowers us
529 to fully prove the transformation correct within the proof assistant itself.

530 **Correct Compilation.** Our work is part of the broader space of correct compilation. The
531 two main approaches to correct compilation in a proof assistant are *certified* and *certifying*
532 compilation. *Certified compilation* refers to formally proving compilers correct. One example
533 is the CompCert verified C compiler [28, 27], which was proven correct in Coq, and has been
534 shown to lack bugs present in other compilers [45]. The CakeML [24] verified implementation
535 of ML in HOL4 includes a verified compiler in HOL4. Oeuf [32] and CertiCoq [2] are
536 both certified compilers for programs written in Gallina, the language underpinning Coq.
537 *Certifying compilation*, in contrast, produces proofs that a given compilation is correct
538 alongside the compiled code. The COGENT language, for example, has a certifying compiler
539 that, for a given compiled program from COGENT to C, proves that the compiled C code
540 correctly implements a generated high-level semantics embedded in Isabelle/HOL [1, 42]. A
541 comprehensive survey of work in certified and certifying compilation is in QED at Large [41].

542 Our work is both certifying (the proof compiler produces a target proof object, showing
543 that the target program proves the target specification) and certified (the proof compiler is
544 proven correct, so that the source and target programs, specifications, and proofs must be
545 analogous). It provides a clear specification of correctness that must always be satisfied, but
546 also allows for proofs to be checked at the target level. *Rupicola* [39], which phrases correct
547 compilation from proof assistants to low-level code as proof search, also has both certified and
548 certifying components. Our addition to this rich space is a technique that correctly compiles
549 proofs from source to target embedded program logic, producing target proof objects.

550 Proof assistants are not the only way to certify correct compilation. *Translation validation*
551 is a method of certifying that a compiler preserves the behavior of the original source program
552 for a particular compilation run. Through the use of simulation relations, one can check
553 whether two programs are semantically equivalent [36]. Our work instead focuses on producing
554 proof objects that can be checked in a target program logic.

555 **Program Logics.** Our work correctly compiles programs and proofs from a source embedded
 556 language and program logic to a target embedded language and program logic. Our source
 557 language and logic were developed using Xavier Leroy’s course materials [29] as a base.
 558 Frameworks based on embedded program logics have seen a lot of success in helping proof
 559 engineers write proofs in a proof assistant about code with features that the proof assistant
 560 lacks. Some example frameworks for reasoning with embedded program logics are Iris [21, 23],
 561 VST-Floyd [9], Bedrock [11, 12], YNot [34], CHL [10], and SEPREF [26].

562 The VST project that VST-Floyd is a part of is especially related to our work, as it aims
 563 to create machine-checked proofs about systems that hold all the way down to machine code.
 564 VST provides three main tools to produce verified software written in C: a program logic
 565 (Verifiable C), a Coq tactic and lemma library for proof automation (VST-Floyd), and a
 566 semantic model of programs (which interacts heavily with CompCert) [6]. By composition
 567 with CompCert, source C programs verified in the VST program logic are guaranteed to
 568 preserve their specifications once compiled, all the way down to assembly code.

569 Our primary contribution to this space is producing a proof object at the target level—
 570 one that is still guaranteed to relate correctly to the source proof of the source program.
 571 In contrast, chaining program logics with certified compilers preserves guarantees about
 572 compiled code, but does not produce a target proof object. We hope that producing target
 573 proof objects will help overcome barriers for understanding and checking target proofs
 574 independently, forming composite guarantees about compiled components written in different
 575 source languages, optimizing and repairing compiled proofs, and keeping proofs about
 576 compiled code immune to changes in the source language or logic.

577 For now, relative to the program logics found in practical frameworks like Iris and VST,
 578 our embedded program logics are fairly simple. And unlike VST chained with CompCert,
 579 our source and target languages are also fairly similar. We hope to extend our work to
 580 handle more practical program logics and lower-level target languages in the future, while
 581 still producing correct proof objects at the target level.

582 **7 Conclusions & Future Work**

583 We defined and formally verified a correct-by-construction proof compiler, POTPIE, in Coq.
 584 POTPIE produces proofs about the target STACK language, given proofs about the source
 585 IMP language. Both proofs about IMP and their compiled STACK counterparts are explicit,
 586 machine-checkable proof terms in Coq, guaranteed to prove their respective specifications,
 587 and to be correctly related. This is a major step toward allowing proof engineers to reason
 588 directly about source programs, and still obtain source-independent proofs about compiled
 589 programs. We believe this opens up two particularly promising directions for future work:

- 590 **1. Linking proofs:** What if proof engineers could write proofs about components of systems
 591 implemented in different languages, and easily get machine-checkable proofs at the target
 592 level about the entire system? This is the equivalent of linking for proofs, and proof
 593 compilation offers a path toward it. The main barriers left are removing the remaining
 594 source-dependent fragments of compiled proofs, implementing proof compilers for multiple
 595 languages, and reasoning about the interactions between compiled components.
- 596 **2. Correctness in the face of unsafe compilation:** What if proof engineers could write
 597 proofs about source programs, and from those proofs derive strong guarantees about
 598 compiled programs, even when the compiler itself is sometimes unsafe? With proof
 599 compilation, there is a chance for freedom from relying on certified and certifying program
 600 compilers, without sacrificing satisfaction of important specifications. The one thing

601 standing in the way is the dependency on full program compiler correctness that our proof
 602 compiler imposes to support the rule of consequence; this is at least in part surmountable.
 603 We are excited to explore these directions.

604 ——— References ———

- 605 1 Sidney Amani, Alex Hixon, Zilin Chen, Christine Rizkallah, Peter Chubb, Liam O'Connor,
 606 Joel Beeren, Yutaka Nagashima, Japheth Lim, Thomas Sewell, Joseph Tuong, Gabriele Keller,
 607 Toby Murray, Gerwin Klein, and Gernot Heiser. Cogent: Verifying high-assurance file system
 608 implementations. In *International Conference on Architectural Support for Programming
 609 Languages and Operating Systems*, pages 175–188, Atlanta, GA, USA, April 2016. doi:
 610 10.1145/2872362.2872404.
- 611 2 Abhishek Anand, Andrew W. Appel, Greg Morrisett, Matthew Weaver, Matthieu Sozeau,
 612 Olivier Savary Belanger, Randy Pollack, and Zoe Paraskevopoulou. CertiCoq: A verified
 613 compiler for Coq. In *CoqPL*, 2017. URL: [http://www.cs.princeton.edu/~appel/papers/
 614 certicoq-coqpl.pdf](http://www.cs.princeton.edu/~appel/papers/certicoq-coqpl.pdf).
- 615 3 Henk Barendregt and Erik Barendsen. Autarkic computations in formal proofs. *Journal of
 616 Automated Reasoning*, 28(3):321–336, 2002. doi:10.1023/A:1015761529444.
- 617 4 Henk Barendregt and Freek Wiedijk. The challenge of computer mathematics. *Philosophical
 618 Transactions of the Royal Society of London A: Mathematical, Physical and Engineering
 619 Sciences*, 363(1835):2351–2375, 2005. doi:10.1098/rsta.2005.1650.
- 620 5 Gilles Barthe, Tamara Rezk, and Ando Saabas. Proof obligations preserving compilation. In
 621 *Formal Aspects in Security and Trust: Thrid International Workshop, FAST 2005, Newcastle
 622 upon Tyne, UK, July 18-19, 2005, Revised Selected Papers*, pages 112–126. Springer, 2006.
- 623 6 Lennart Beringer and Andrew W. Appel. Abstraction and subsumption in modular verification
 624 of C programs. *Formal Methods in System Design*, 58(1):322–345, October 2021. doi:
 625 10.1007/s10703-020-00353-1.
- 626 7 Clément Blaudeau. Untitled, 2022. URL: <https://pastebin.com/BAvg3Jdh>.
- 627 8 Pierre Boutillier. *De nouveaux outils pour calculer avec des inductifs en Coq*. Theses, Université
 628 Paris-Diderot - Paris VII, February 2014. URL: <https://theses.hal.science/tel-01054723>.
- 629 9 Qinxiang Cao, Lennart Beringer, Samuel Gruetter, Josiah Dodds, and Andrew W. Appel. VST-
 630 Floyd: A Separation Logic Tool to Verify Correctness of C Programs. *Journal of Automated
 631 Reasoning*, 61(1):367–422, June 2018. doi:10.1007/s10817-018-9457-5.
- 632 10 Haogang Chen, Daniel Ziegler, Tej Chajed, Adam Chlipala, M. Frans Kaashoek, and Nikolai
 633 Zeldovich. Using crash hoare logic for certifying the fscq file system. In *Proceedings of the 25th
 634 Symposium on Operating Systems Principles, SOSP '15*, page 18–37, New York, NY, USA,
 635 2015. Association for Computing Machinery. doi:10.1145/2815400.2815402.
- 636 11 Adam Chlipala. Mostly-automated verification of low-level programs in computational separa-
 637 tion logic. *SIGPLAN Not.*, 46(6):234–245, jun 2011. doi:10.1145/1993316.1993526.
- 638 12 Adam Chlipala. The bedrock structured programming system: Combining generative metapro-
 639 gramming and hoare logic in an extensible program verifier. *SIGPLAN Not.*, 48(9):391–402,
 640 sep 2013. doi:10.1145/2544174.2500592.
- 641 13 Adam Chlipala. Formal reasoning about programs, 2017. URL: [http://adam.chlipala.net/
 642 frap/](http://adam.chlipala.net/frap/).
- 643 14 H. B. Curry. Functionality in combinatory logic. *Proceedings of the National Academy of
 644 Sciences of the United States of America*, 20(11):584–590, 1934. URL: [http://www.jstor.
 645 org/stable/86796](http://www.jstor.org/stable/86796).
- 646 15 Michael J. C. Gordon, Robin Milner, and Christopher P. Wadsworth. Edinburgh LCF: A
 647 mechanised logic of computation. In *Lecture Notes in Computer Science*, volume 78. 1979.
 648 doi:10.1007/3-540-09724-4.
- 649 16 Robert W Hasker and Uday S Reddy. Generalization at higher types. In *Proceedings of the
 650 Workshop on the λ Prolog Programming Language*, pages 257–271, 1992.

- 651 17 Bruno Hauser. Embedding proof-carrying components into Isabelle. Master's thesis, ETH,
652 Swiss Federal Institute of Technology Zurich, Institute of Theoretical . . . , 2009.
- 653 18 C. A. R. Hoare. An axiomatic basis for computer programming. *Commun. ACM*,
654 12(10):576–580, oct 1969. doi:10.1145/363235.363259.
- 655 19 W. A. Howard. *The formulae-as-types notion of constructions*. Academic Press, 1980.
- 656 20 Einar Broch Johnsen and Christoph Lüth. Theorem reuse by proof term transformation.
657 In *Theorem Proving in Higher Order Logics: 17th International Conference, TPHOLs 2004,*
658 *Park City, Utah, USA, September 14-17, 2004. Proceedings*, pages 152–167. Springer, Berlin,
659 Heidelberg, 2004. doi:10.1007/978-3-540-30142-4_12.
- 660 21 Ralf Jung, David Swasey, Filip Sieczkowski, Kasper Svendsen, Aaron Turon, Lars Birkedal,
661 and Derek Dreyer. Iris: Monoids and invariants as an orthogonal basis for concurrent reasoning.
662 *SIGPLAN Not.*, 50(1):637–650, jan 2015. doi:10.1145/2775051.2676980.
- 663 22 Thomas Kolbe and Christoph Walther. *Proof Analysis, Generalization and Reuse*, pages
664 189–219. Springer, Dordrecht, 1998. doi:10.1007/978-94-017-0435-9_8.
- 665 23 Robbert Krebbers, Amin Timany, and Lars Birkedal. Interactive proofs in higher-order
666 concurrent separation logic. *SIGPLAN Not.*, 52(1):205–217, jan 2017. doi:10.1145/3093333.
667 3009855.
- 668 24 Ramana Kumar, Magnus O. Myreen, Michael Norrish, and Scott Owens. CakeML: A verified
669 implementation of ML. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on*
670 *Principles of Programming Languages*, POPL '14, pages 179–191, New York, NY, USA, 2014.
671 ACM. doi:10.1145/2535838.2535841.
- 672 25 César Kunz. *Certificate Translation Alongside Program Transformations*. PhD thesis, ParisTech,
673 Paris, France, 2009.
- 674 26 Peter Lammich. Refinement to imperative HOL. *J. Autom. Reason.*, 62(4):481–503, apr 2019.
675 doi:10.1007/s10817-017-9437-1.
- 676 27 Xavier Leroy. Formal certification of a compiler back-end, or: programming a compiler with
677 a proof assistant. In *33rd ACM symposium on Principles of Programming Languages*, pages
678 42–54. ACM Press, 2006. URL: <http://xavierleroy.org/publi/compiler-certif.pdf>.
- 679 28 Xavier Leroy. Formal verification of a realistic compiler. *Communications of the ACM*,
680 52(7):107–115, July 2009. doi:10.1145/1538788.1538814.
- 681 29 Xavier Leroy. Coq development for the course “Mechanized semantics”, 2019-2021. URL:
682 <https://github.com/xavierleroy/cdf-mech-sem>.
- 683 30 Bohdan Liesnikov, Marcel Ullrich, and Yannick Forster. Generating induction principles and
684 subterm relations for inductive types using metacoq. *arXiv preprint arXiv:2006.15135*, 2020.
- 685 31 Nicolas Magaud. Changing data representation within the Coq system. In *International*
686 *Conference on Theorem Proving in Higher Order Logics*, pages 87–102. Springer, 2003.
- 687 32 Eric Mullen, Stuart Pernsteiner, James R. Wilcox, Zachary Tatlock, and Dan Grossman. Oeuf:
688 Minimizing the Coq extraction TCB. In *CPP*, pages 172–185, New York, NY, USA, 2018.
689 ACM. doi:10.1145/3167089.
- 690 33 Peter Müller and Martin Nordio. Proof-transforming compilation of programs with abrupt
691 termination. In *SAVCBS*, pages 39–46, 01 2007. doi:10.1145/1292316.1292321.
- 692 34 Aleksandar Nanevski, Greg Morrisett, Avraham Shinnar, Paul Govereau, and Lars Birkedal.
693 Ynot: Dependent types for imperative programs. *SIGPLAN Not.*, 43(9):229–240, sep 2008.
694 doi:10.1145/1411203.1411237.
- 695 35 George C. Necula. Proof-carrying code. In *Proceedings of the 24th ACM SIGPLAN-SIGACT*
696 *Symposium on Principles of Programming Languages*, POPL '97, pages 106–119, New York,
697 NY, USA, January 1997. Association for Computing Machinery. doi:10.1145/263699.263712.
- 698 36 George C. Necula. Translation validation for an optimizing compiler. In *Proceedings of the*
699 *ACM SIGPLAN 2000 Conference on Programming Language Design and Implementation*,
700 PLDI '00, page 83–94, New York, NY, USA, 2000. Association for Computing Machinery.
701 doi:10.1145/349299.349314.

- 702 37 Martin Nordio, Peter Müller, and Bertrand Meyer. Formalizing proof-transforming compilation
703 of eiffel programs. Technical report, ETH Zurich, 2008.
- 704 38 Frank Pfenning. *Proof Transformations in Higher-Order Logic*. PhD thesis, Carnegie Mellon
705 University, 1987.
- 706 39 Clément Pit-Claudel, Jade Philipoom, Dustin Jamner, Andres Erbsen, and Adam Chlipala.
707 Relational compilation for performance-critical applications: Extensible proof-producing
708 translation of functional models into low-level code. In *Proceedings of the 43rd ACM SIGPLAN
709 International Conference on Programming Language Design and Implementation, PLDI 2022*,
710 page 918–933, New York, NY, USA, 2022. Association for Computing Machinery. doi:
711 10.1145/3519939.3523706.
- 712 40 Talia Ringer. *Proof Repair*. PhD thesis, University of Washington, 2021.
- 713 41 Talia Ringer, Karl Palmskog, Ilya Sergey, Milos Gligoric, and Zachary Tatlock. Qed at large: A
714 survey of engineering of formally verified software. *Foundations and Trends® in Programming
715 Languages*, 5(2-3):102–281, 2019. URL: <http://dx.doi.org/10.1561/25000000045>, doi:10.
716 1561/2500000045.
- 717 42 Christine Rizkallah, Japheth Lim, Yutaka Nagashima, Thomas Sewell, Zilin Chen, Liam
718 O’Connor, Toby Murray, Gabriele Keller, and Gerwin Klein. A framework for the automatic
719 formal verification of refinement from Cogent to C. In *Interactive Theorem Proving*, pages
720 323–340, Cham, 2016. Springer. doi:10.1007/978-3-319-43144-4_20.
- 721 43 Enrico Tassi. Deriving Proved Equality Tests in Coq-Elpi: Stronger Induction Principles for
722 Containers in Coq. In John Harrison, John O’Leary, and Andrew Tolmach, editors, *10th
723 International Conference on Interactive Theorem Proving (ITP 2019)*, volume 141 of *Leibniz
724 International Proceedings in Informatics (LIPIcs)*, pages 29:1–29:18, Dagstuhl, Germany, 2019.
725 Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. URL: [http://drops.dagstuhl.de/opus/
726 volltexte/2019/11084](http://drops.dagstuhl.de/opus/volltexte/2019/11084), doi:10.4230/LIPIcs.ITP.2019.29.
- 727 44 Doug Woos, James R. Wilcox, Steve Anton, Zachary Tatlock, Michael D. Ernst, and Thomas
728 Anderson. Planning for change in a formal verification of the Raft consensus protocol. In
729 *Proceedings of the 5th ACM SIGPLAN Conference on Certified Programs and Proofs, CPP
730 2016*, pages 154–165, New York, NY, USA, 2016. ACM. doi:10.1145/2854065.2854081.
- 731 45 Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. Finding and understanding bugs in c
732 compilers. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language
733 Design and Implementation, PLDI ’11*, page 283–294, New York, NY, USA, 2011. Association
734 for Computing Machinery. doi:10.1145/1993498.1993532.