ZHE ZHOU, Purdue University, USA ROBERT DICKERSON, Purdue University, USA BENJAMIN DELAWARE, Purdue University, USA SURESH JAGANNATHAN, Purdue University, USA

Programmers often leverage data structure libraries that provide useful and reusable abstractions. Modular verification of programs that make use of these libraries naturally rely on specifications that capture important properties about how the library expects these data structures to be accessed and manipulated. However, these specifications are often missing or incomplete, making it hard for clients to be confident they are using the library safely. When library source code is also unavailable, as is often the case, the challenge to infer meaningful specifications is further exacerbated. In this paper, we present a novel data-driven abductive inference mechanism that infers specifications for library methods sufficient to enable verification of the library's clients. Our technique combines a data-driven learning-based framework to postulate candidate specifications, along with SMT-provided counterexamples to refine these candidates, taking special care to prevent generating specifications that overfit to sampled tests. The resulting specifications form a minimal set of requirements on the behavior of library implementations that ensures safety of a particular client program. Our solution thus provides a new multi-abduction procedure for precise specification inference of data structure libraries guided by client-side verification tasks. Experimental results on a wide range of realistic OCaml data structure programs demonstrate the effectiveness of the approach.

 $\label{eq:CCS} \textit{Concepts:} \bullet \textit{Computing methodologies} \rightarrow \textit{Classification and regression trees}; \bullet \textit{Software and its engineering} \rightarrow \textit{Software verification and validation}; Automated static analysis.}$ 

Additional Key Words and Phrases: Automated Verification, Data-Driven Specification Inference, Data Structures, Decision Tree Learning, Counterexample Guided Refinement

## **ACM Reference Format:**

Zhe Zhou, Robert Dickerson, Benjamin Delaware, and Suresh Jagannathan. 2021. Data-Driven Abductive Inference of Library Specifications. *Proc. ACM Program. Lang.* 5, OOPSLA, Article 116 (October 2021), 29 pages. https://doi.org/10.1145/3485493

## **1 INTRODUCTION**

Using a specification of a library's methods in the verification of its clients is a hallmark of modular reasoning. Because these specifications encapsulate the interface between the client and the library, each may be independently verified without access to the other's implementation. This modularity is particularly beneficial when the library function is complex or its source code is unavailable. All too often, though, such specifications are either missing or incomplete, preventing the verification of clients without making (often unwarranted) assumptions about the behavior of the library. This problem is further exacerbated when libraries expose rich datatype functionality, which often leads to specifications that rely on inductive invariants [Dillig et al. 2013; Itzhaky et al. 2014] and complex

Authors' addresses: Zhe Zhou, Purdue University, USA; Robert Dickerson, Purdue University, USA; Benjamin Delaware, Purdue University, USA; Suresh Jagannathan, Purdue University, USA.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2021 Copyright held by the owner/author(s). 2475-1421/2021/10-ART116 https://doi.org/10.1145/3485493 structural relations. One solution to this problem is to automatically *infer* missing specifications. Unfortunately, while significant progress has been made in specification inference over the past several years [Albarghouthi et al. 2016; Miltner et al. 2020; Padhi et al. 2016; Zhu et al. 2016], existing techniques have not considered inference in the frequently occurring case of client programs that make use of data structure libraries with unavailable implementations.

To highlight the challenge, consider the following simple program, which concatenates two stacks together using four operations provided by a Stack library: push, top, is\_empty and tail.

```
1 let rec concat s1 s2 =
2 if Stack.is_empty s1 then s2
3 else Stack.push (Stack.top s1) (concat (Stack.tail s1) s2)
```

To ensure the correctness of this client function, its author may wish to verify that (a) the top element of the output stack is always the top element of one of the input stacks; and, (b) every element of the output stack is also an element of one of the input stacks and vice-versa. In order to express this behavior in a form amenable for automatic verification, we need some mechanism to encode the semantics of stacks in a decidable logic. To do so, we rely on a pair of *method predicates*, "*a* is the head of stack s", hd(s, a), and "*a* is a member of stack s", mem(s, a), to write our postcondition:

$$\forall u, (hd(v, u) \implies hd(s_1, u) \lor hd(s_2, u)) \land (mem(v, u) \iff mem(s_1, u) \lor mem(s_2, u))$$

$$(\phi_{\text{concat}})$$

We assume these method predicates are associated with (possibly blackbox) implementations that we can use to check the specifications in which they appear. For example, *hd* may be defined in terms of the stack operations top and is\_empty, while the implementation of *mem* might additionally use tail. The variable v in  $\phi_{concat}$  is used to represent the output stack of concat. The above assertion claims that the head of the output stack must be the head of either s<sub>1</sub> or s<sub>2</sub> and that any element found in the output must be a member of either s<sub>1</sub> or s<sub>2</sub>.

By treating method predicates (*hd* and *mem*) and library functions (push, top, is\_empty, and tail) as uninterpreted function symbols, it is straightforward to generate verification conditions (VCs), e.g. using weakest precondition inference, which can be handed off to an off-the-shelf SMT solver like Z3 to check. However, the counter-examples returned by the theorem prover may be spurious, generated by incorrect assumptions about library method behavior in the absence of any constraints on these behaviors outside the client VCs. For example, the prover might assume the formula  $\neg hd(s_1, top(s_1))$  is valid, i.e. that the result of  $top(s_1)$  is not the head of  $s_1$ . This claim is obviously inconsistent with the client's expectation of top's semantics, but it is not disallowed by any constraints in the SMT query. Using this assumption, Z3 may return the following counterexample:  $\exists s u.hd(s, u) \land \neg mem(s, u)$ . This counterexample, which is obviously incongruous with the intended semantics of *hd* and *mem*, occurs because the expected relationship between *hd* and *mem* is lost when the predicates are embedded as uninterpreted functions in the SMT query.

To overcome this problem, we need stronger specifications for the library methods, defined in terms of these method predicates, that are sufficient to imply the desired client postcondition. In particular, these specifications should rule out spurious unsafe executions such as the counter-example given above. In the (quite likely) scenario that such library specifications are not already available, a reasonable fallback is to infer *some* specifications for these functions that are strong enough to ensure the safety of the client. Traditional approaches to specification inference usually adopt a closed-world assumption in which specifications of library methods are discovered in isolation, independent of the client context in which they are being used. Such assumptions are not applicable here since (a) we do not have access to the library's method implementations and (b) the



Fig. 1. Elrond pipeline.

nature of the specifications we need to infer are impacted by the verification demands of the client. In this setting, some form of data-driven inference [Miltner et al. 2020; Padhi et al. 2016; Zhu et al. 2016] can be beneficial. Such an approach may be tailored to the client context in which the library methods are used, postulating candidate specifications for library methods based on observations of their input-output behavior. Unfortunately, completely blackbox data-driven approaches are susceptible to overfitting on the set of observations used to train them, and can thus discount reasonable and safe behaviors of the underlying library functions.

To address the problem of overfitting, we might instead consider attacking this problem from a purely logical standpoint, treating specification inference as an instance of a multi-abductive inference problem [Albarghouthi et al. 2016] that tries to find formulae  $R_{push}$ ,  $R_{top}$ ,  $R_{tail}$  and  $R_{is\_empty}$ such that  $\bigwedge R_i \not\models \bot$  and yet which are sufficient to prove the desired verification condition. While such problems have been previously solved over linear integer arithmetic constraints [Albarghouthi et al. 2016] using quantifier elimination, these prior techniques cannot be directly applied to formulae with uninterpreted function symbols like the method predicates (e.g., *hd* and *mem*) used to encode library method specifications in our setting.

In this work, we combine aspects of these data-driven and abductive approaches in a way that addresses the limitations each approach has when considered independently. Our technique uses SMT-provided counterexamples to generate infeasible interpretations of these predicates (similar to other abductive inference methods) while using concrete test data to generate feasible interpretations (similar to data-driven inference techniques). This combination yields a novel CEGIS-style inference methodology that allows us to postulate specifications built from method predicates sufficient to prove the postcondition in a purely blackbox setting. The specifications learned by this procedure are guaranteed to be both *consistent* with the observed input-output behavior of the blackbox library implementations and *safe* with respect to the postcondition of the client program. As there may be many such specifications, we also endeavor to find a *maximal* one that is at least as weak as every other safe and consistent specification, in order to avoid overfitting to observed library behaviors. Our algorithm applies another data-driven weakening procedure to find these maximal specifications.

To demonstrate the effectiveness of our approach, we have implemented a fully automated abductive specification inference pipeline in OCaml called Elrond (see Figure 1). This pipeline takes as input (a) an OCaml client program that may call blackbox library code defined over algebraic datatypes like lists, trees, heaps, etc.; (b) assertions about the behavior of this client program; and, (c) a set of method predicates (e.g., *hd* or *mem*), along with their (possibly blackbox) implementations, that are used to synthesize library method specifications. It combines tests and counterexample-guided refinement techniques to either generate a set of maximal specifications for the library

methods used by the client program, or a counterexample that demonstrates a violation of the postcondition. The notion of "weakest" used in our definition of maximal is bounded by the "shape" of specifications (e.g., the number of quantified variables, the set of method predicates, etc.) and a time bound. Our results over a range of sophisticated data-structure manipulating programs, including those drawn from e.g., Okasaki [1999], show that Elrond is able to discover maximally-weak specifications (as determined by an oracle executing without any time constraints) for the vast majority of applications in our benchmark suite within one hour.

Our key contribution is thus a new abductive inference framework that is a fusion of automated data-driven methods and counterexample-guided refinement techniques, tailored to specification inference for libraries that make use of rich algebraic datatypes. Specifically, we:

- **1.** Frame client-side verification as a multi-abduction inference problem that searches for library method specifications that are both consistent with the method's implementation and sufficient to verify client assertions.
- **2.** Devise a novel specification weakening procedure that yields the weakest specification among the collection of all safe and consistent ones with respect to a given set of quantified variables and method predicates.
- **3.** Evaluate our approach in a tool, Elrond, which we use to analyze a comprehensive set of realistic and challenging functional (OCaml) data structure programs. An artifact containing this tool and our benchmark suite is publicly available [Zhou et al. 2021b].

The remainder of the paper is structured as follows. The next section presents an overview of our approach using a detailed example to motivate its key ideas . A formal characterization of the problem is given in Section 3. Section 4 defines how a data-driven learning strategy can be used to perform inference. A detailed presentation of the algorithm used to manifest these ideas in a practical implementation is given in Section 5. Details of our implementation and evaluation results are explained in Section 6. Related work and conclusions are given in Section 7 and Section 8.

# 2 OVERVIEW AND MOTIVATION

We divide the inference of maximal library specifications into two stages, which are represented as the "Specification Inference" and "Weakening" components in Figure 1. Both stages leverage datadriven learning to overcome the lack of a purely logical abduction procedure for our specification language. The initial inference stage learns a set of safe and consistent specifications from a combination of concrete tests and verifier-provided counterexamples. The next stage then weakens these specifications by iteratively augmenting this data set with additional safe behaviors until a set of maximal specifications are found.

Figure 2 provides a more detailed depiction of the initial specification inference stage in Elrond. Starting from an initial set of maximally permissive specifications, this stage iteratively refines the set of candidate specifications until either a set of safe and consistent solutions or a counterexample witnessing an unsafe execution is found.

Each iteration first uses a property-based sampler, e.g. QuickCheck [Claessen and Hughes 2000], to look for executions of the blackbox library implementations that are inconsistent with the current set of inferred specifications. The reliance on a generator to provide high-quality tests provides yet more motivation for the subsequent weakening phase, in order to ensure that the final abduced specifications are not overfitted to or otherwise biased by the tests provided by the generator. At the same time, we also observe that the sorts of shape properties (e.g., membership and ordering) used in our specifications and assertions are relatively under-constrained and are thus amenable to property-based random sampling. We do not ask, for example, for QuickCheck to generate inputs satisfying non-structural properties like "*a list whose 116<sup>th</sup> element is equal to 5.*"



Fig. 2. The details of Elrond's Specification Inference component from Figure 1. The blackbox library implementation and method predicates are used by multiple components; their corresponding arrows within the diagram are omitted for clarity.

Any tests that are disallowed by the current solution are passed to a learner which uses them to generalize the current specification. If no inconsistencies are detected, Elrond attempts to verify the client against the candidate specifications using a theorem prover. If the inferred specifications are sufficient to prove client safety, the loop exits, returning the discovered solution. If not, the verifier has identified a model that represents a *potential* safety violation. The model is then analyzed in an attempt to extract test inputs that trigger a safety violation. If we are unable to find such a counterexample, the model is most likely incongruous with the semantics of the method predicates and is thus spurious. In this case, the model is passed to the learner so that it can be used to strengthen candidate specifications, preventing this and similar spurious counterexamples from manifesting in subsequent iterations.



Fig. 3. The details of Elrond's Weakening component from Figure 1. Arrows for method predicates are again omitted for clarity.

While the previous loop is guaranteed to return safe and consistent solutions, it may find specifications that are nonetheless *too strong* with respect to the underlying library implementation. This occurs when the property-based sampler fails to find a test that identifies an inconsistent specification, which may happen when the input space of a library function is very large. To combat overfitting specifications to test data, candidate solutions are iteratively weakened using the data-driven counterexample-guided refinement loop depicted in Figure 3. The data in this phase is supplied by the underlying theorem prover rather than a concrete test generator. Each iteration of the refinement loop first attempts to find a safe execution of the client program that is disallowed by the current set of specifications. If no such execution can be found, the specifications are maximal

and the loop terminates. Otherwise, the identified execution is passed to a learner, which uses it to generalize the candidate solution so that the execution is permitted before continuing the refinement loop. The learner always generalizes candidate specifications, maintaining the invariant that the current solution is consistent with all previously observed library behaviors.

## 2.1 Elrond in Action

To illustrate our approach in more detail, we apply it to the stack concatenation example from the introduction. Given the postcondition  $\phi_{concat}$  and the implementation of concat from Section 1, Elrond generates a formula that can be simplified to the following implication:

$$\begin{aligned} \forall \mathbf{s}_{1}, \mathbf{s}_{2}, v, v_{\text{top}}, v_{\text{tail}}, v_{\text{concat}}, v_{\text{is}\_empty}, \\ (\mathbf{R}_{\text{is}\_empty}(\mathbf{s}_{1}, v_{\text{is}\_empty}) \land \mathbf{R}_{\text{top}}(\mathbf{s}_{1}, v_{\text{top}}) \land \mathbf{R}_{\text{tail}}(\mathbf{s}_{1}, v_{\text{tail}}) \land \mathbf{R}_{\text{push}}(v_{\text{top}}, v_{\text{concat}}, v)) \\ \implies \\ ((v_{\text{is}\_empty} = \bot \land \forall u, (hd(v_{\text{concat}}, u) \implies hd(v_{\text{tail}}, u) \lor hd(\mathbf{s}_{2}, u)) \land \\ (mem(v_{\text{concat}}, u) \iff mem(v_{\text{tail}}, u) \lor mem(\mathbf{s}_{2}, u))) \implies \\ \forall u, (hd(v, u) \implies hd(\mathbf{s}_{1}, u) \lor hd(\mathbf{s}_{2}, u)) \land (mem(v, u) \iff mem(\mathbf{s}_{1}, u) \lor mem(\mathbf{s}_{2}, u))) \end{aligned}$$

The four predicates in the premise of this formula correspond to the four library functions (push, top, is\_empty and tail) invoked in a recursive call to concat. The specification of a blackbox library function f in our assertion logic is represented as a *placeholder predicate*: an uninterpreted predicate that relates the parameters of f to its return value. For a library function f, we adopt a naming convention of  $R_f$  and  $v_f$  for its placeholder predicates and return values, respectively. The predicate  $R_{top}(s_1, v_{top})$  in the above formula, for example, says the variable  $v_{top}$  holds the return value of the call to Stack.top s1 in concat. The result of the recursive call to concat is similarly denoted as  $v_{concat}$ . The conclusion of the formula encodes the expected verification condition for a recursive call to concat (Stack.tail s1) s2 satisfies  $\phi_{concat}$ , then the result of concat must also satisfy  $\phi_{concat}$ . The remainder of this section refers to the premise and conclusion of this implication as  $\Sigma_{concat}$  and  $\Phi_{concat}$ , respectively.

Method Predicates. From a logical standpoint, the method predicates used in  $\Phi_{concat}$  are simply uninterpreted function symbols which have no intrinsic semantics. This representation allows our specifications to use predicates whose semantics may be difficult to encode directly in the logic. Embedding recursively defined predicates like mem, for example, requires particular care [Zhu et al. 2016]. In order to ensure that the specifications inferred by Elrond are tethered to reality, users must also supply Elrond with implementations (possibly blackbox) of these predicates. One possible implementation for *hd* and *mem* is:

```
1 let hd(s, u) =
2 if Stack.is_empty s then false else (Stack.top s) == u
3
4 let rec mem(s, u) =
5 if Stack.is_empty s then false else (Stack.top s = u) ||
6 (mem(Stack.tail s, u))
```

where Stack.is\_empty, Stack.top and Stack.tail refer to blackbox implementations of stack library methods.

While it is possible to naïvely include a method predicate for each library method, such an approach may not be useful for verification. Some functions may be irrelevant to client assertions,

Proc. ACM Program. Lang., Vol. 5, No. OOPSLA, Article 116. Publication date: October 2021.

unnecessarily increasing the set of possible specifications that must be considered. Conversely, the library may not include functions for desirable predicates; e.g., the Stack library does not provide a Stack.mem function, although it is quite relevant for verifying our running example.

$$\begin{split} \mathsf{R}_{\mathsf{is\_empty}}(\mathsf{s},v) &\mapsto \forall u, (v \implies \neg mem(\mathsf{s},u)) \land (\neg v \land hd(\mathsf{s},u) \implies mem(\mathsf{s},u)) \\ \mathsf{R}_{\mathsf{top}}(\mathsf{s},v) &\mapsto \forall u, mem(\mathsf{s},v) \land (v = u \iff hd(\mathsf{s},u)) \\ \mathsf{R}_{\mathsf{tail}}(\mathsf{s},v) &\mapsto \forall u, (mem(\mathsf{s},u) \implies (mem(v,u) \lor hd(\mathsf{s},u))) \land \\ & ((mem(v,u) \lor hd(v,u)) \implies mem(\mathsf{s},u)) \\ \mathsf{R}_{\mathsf{push}}(\mathsf{x},\mathsf{s},v) &\mapsto \forall u, (mem(v,u) \land mem(\mathsf{s},u) \implies \neg(\mathsf{x} = u \lor hd(v,u))) \land \\ & (mem(v,u) \land \neg mem(\mathsf{s},u) \implies (\mathsf{x} = u \land hd(v,u))) \land \\ & ((\mathsf{x} = u \lor hd(v,u) \lor hd(\mathsf{s},u) \lor mem(\mathsf{s},u)) \implies mem(v,u)) \end{split}$$

Fig. 4. Candidate verification interface for the concat example.

Solution Space. Our ultimate goal is to find a mapping from each placeholder predicate in  $\sum_{concat}$  to an interpretation that entails  $\Phi_{concat}$ . We refer to such a mapping as a *verification interface*; Figure 4 presents a potential verification interface for our running example. Not every mapping that ensures the safety of concat is reasonable, however. At one extreme, interpreting every predicate as  $\perp$  ensures the safety of the client, but does not capture the behavior of any sensible stack implementation. Our goal, then, is to find interpretations that are general enough to cover a range of possible implementations. From a purely logical perspective, this an instance of a multi-abductive inference problem that tries to find the weakest interpretations of  $R_{push}$ ,  $R_{top}$ ,  $R_{tail}$  and  $R_{is\_empty}$  in terms of predicates mem and hd such that the interpretation condition. While solutions to the multi-abduction problem have been developed for domains that admit quantifier elimination, e.g. linear integer arithmetic constraints Albarghouthi et al. [2016], there is no purely logical solution for formulae involving equalities over uninterpreted functions. An additional challenge in our setting is that we seek to infer specifications consistent with the library's implementation, a requirement that is absent in [Albarghouthi et al. 2016].

## 2.2 Data-Driven Abduction

We overcome these challenges by adopting a data-driven approach to abducing maximal library specifications, framing the problem as one of training a Boolean classifier on a set of *example behaviors* for each function f. Under this interpretation, a classifier represents a specification of the "acceptable" behaviors of f. Thus, the goal of the specification inference stage of our algorithm is to learn a set of classifiers that recognize the behaviors of each f that a) are consistent with f's underlying implementation and b) preserve the safety of the client program. As discussed at the start of this section, this algorithm uses both an SMT-based verifier and a property-based test generator as sources of training data for our learner. The former identifies example behaviors consistent with  $\neg(\Sigma_{concat} \implies \Phi_{concat})$ ; these are labelled as "negative" examples, so that the learned specifications can help the solver rule out behaviors that are inconsistent with the semantics of the method predicates or that produce unsafe executions (i.e., interpretations that would violate the postcondition). Example behaviors drawn from tests are labelled as "positive", so that our learner is biased towards explanations that are consistent with the (unknown) implementations of library functions. Notably, our algorithm generalizes this data-driven abduction procedure for

individual functions to the multi-abduction case, ensuring that discovered interpretations are globally consistent over all library methods.



Fig. 5. Elements in the circle labelled  $\Sigma \Rightarrow \Phi$  represent positive examples, while the orange  $\neg(\Sigma \Rightarrow \Phi)$  circle contains negative examples. Our algorithm generates a mapping  $\Delta$ , which maps placeholder predicates to specifications, separating the two.

Figure 5 depicts the space of example behaviors for our learner, as well four potential verification interfaces. Each of these represents a potential solution in the *hypothesis space* for this learner, which is tasked with building a classifier that separates negative (-) and positive examples (+) for each library function. The dashed purple line, labelled  $\Delta^{unsafe}$ , represents an unsafe verification interface that allows a client program to violate the desired postcondition. The remaining red lines represent the range of safe verification interfaces. The two dashed red lines represent the verification interfaces that are sufficient to verify the client, but which are suboptimal.  $\Delta^{safe}$  is safe but inconsistent with the observed behaviors of the library implementation, and is thus overly restrictive.  $\Delta_{cons}^{safe}$  is safe and consistent, but not maximal, as there exists a weaker verification interface ( $\Delta^{max}$ ) in the hypothesis space that is still safe and consistent. Intuitively, the goal of our first phase is to identify  $\Delta_{cons}^{safe}$ , which is then weakened by the second phase to produce  $\Delta^{max}$ .

*Hypothesis space.* Our learner limits the shape of solutions it considers so that inferred specifications are both amenable to automatic verification and strong enough to verify specified postconditions. To enable automated verification of client programs, potential specifications are required to be prenex universally-quantified propositional formulae over datatype values and variables representing arguments to the predicates under consideration. Some possible specifications of the library function push in our running example include:

$$push(x, l) = v: \begin{cases} \forall u, mem(l, u) \implies mem(v, u), \quad \forall u, \neg mem(v, u), \\ \forall u, u = x, \qquad \qquad \forall u, u = x \iff hd(v, u), \\ \forall u, mem(l, u) \land hd(l, u), \qquad \dots \end{cases}$$

which contains, among other candidates, the desired specification.

All atomic literals in generated formulae are applications of uninterpreted method predicates and equalities over quantified variables, parameters, and return values of functions. The literals in the above formulae are simply applications of *hd* and *mem* to *l*, *v* and *u*, and the equality u = x. We automatically discard equalities between terms of different types, e.g. l = x. The *feature set* for the predicate R<sub>push</sub>, i.e. the set of atomic elements used to construct its specification, is thus: {*hd*(*l*, *u*), *mem*(*l*, *u*), *hd*(*v'*, *u*), *mem*(*v*, *u*), x = u}.

Proc. ACM Program. Lang., Vol. 5, No. OOPSLA, Article 116. Publication date: October 2021.

$R_{top}(s_1, v_{top})$	∀u	$hd(s_1, u)$	$\mathit{mem}(s_1, u)$	$v_{top} = u$		
	u = a	false	false	true		
	u≠a	false	false	false		
$R_{tail}(s_1, v_{tail})$	∀u	$hd(s_1, u)$	$mem(s_1, u)$	$hd(v_{tail}, u)$	$mem(v_{tail}, u)$	
	u = a	false	false	false	true	
	u≠a	false	false	false	false	
$R_{push}(v_{top}, v_{concat}, v) \mid \forall u$		$hd(v_{concat}, u)$	$\mathit{mem}(v_{concat}, u)$	hd(v, u)	mem(v, u)	$v_{top} = u$
	u = a	false	true	false	true	true
	u≠a	false	false	false	false	false
$R_{is\_empty}(s_1, \nu_{is\_empty})$	∀u	$hd(s_1, u)$	$mem(s_1, u)$	$v_{is\_empty}$		
	u = a	false	false	false		
	u≠a	false	false	false		

Table 1. Potential negative feature vectors extracted from CEX.

*Training Data.* We now consider how to represent library behaviors in a form that is useful for learning a solution in our hypothesis space. To illustrate our chosen representation, consider the counterexample produced by an off-the-shelf theorem prover when asked to verify the formula  $\Sigma_{concat} \implies \Phi_{concat}$  from our running stack example, where the set of candidate specifications are initialized to true (i.e.,  $\forall f. R_f \mapsto \top$ ):

$$\forall l \ u, \neg hd(l, u) \land v_{top} = a \land (mem(l, u) \iff ((l = v \lor l = v_{concat} \lor l = v_{tail}) \land u = a)) \quad (Cex)$$

Intuitively, this counterexample asserts that the stacks v,  $v_{tail}$ , and  $v_{concat}$  contain exactly one element, the constant a, and the other two stacks,  $s_1$  and  $s_2$ , are empty. This assertion indeed violates the second conjunct of the postcondition,  $mem(v, u) \iff mem(s_1, u) \lor mem(s_2, u)$ , but it is inconsistent with the expected semantics of the library functions, and can thus be safely ignored. The verifier generates this counterexample because the interpretations of the placeholder predicates in  $\Sigma_{concat}$  are too permissive. In order for the verifier to rule out this counterexample, the placeholder predicates need to be strengthened to rule out this inconsistent behavior.

Ignoring how we identify this counterexample as spurious for now, note that there are many ways to strengthen these specifications. One approach is to focus on one particular function at a time. For example, we could choose to refine the specification of  $R_{top}$  so that it guarantees that  $v_{top}$  is a member of  $s_1$ . Alternatively, we could focus on  $R_{tail}$ , ensuring the members of  $v_{tail}$  are also contained by  $s_1$ . In general, however, it may be necessary to strengthen multiple specifications at once. Therefore, instead of focusing on one specification at a time, we learn refined specifications simultaneously.

The first step to refining our placeholder specifications is to extract data from CEX in a form that can be used to train a classifier. We do this by using the assignments to the arguments of the placeholder predicates in a counterexample to build *feature vectors* that describe the valuations of method predicates and equalities in the unsafe execution. Table 1 presents the feature vectors extracted from CEX. The first column of this table indicates the particular placeholder predicates that can be strengthened to rule out this counterexample. The second column gives the feature vectors for a particular instantiation of the quantified variables (*e.g.*, u) of the placeholders. The subsequent columns list applications of method predicates, with the rows underneath listing the valuation of these predicates in the offending run. The second row corresponds to the assertion

$R_{top}(s_1, v_{top})$	∀u	$hd(s_1, u)$	$\mathit{mem}(s_1, u)$	$v_{top} = u$		
	u = a	true	true	true		
	u = b	false	false	false		
$R_{tail}(s_1, v_{tail})$	∀u	$hd(s_1, u)$	$mem(s_1, u)$	$hd(v_{tail}, u)$	$mem(v_{tail}, u)$	
u = a		true	true	false	false	
	u = b		false	false	false	
$R_{push}(v_{top}, v_{concat}, v)$	∀u	$hd(v_{concat}, \mathbf{u})$	$mem(v_{concat}, u)$	<i>hd</i> ( <i>v</i> , u)	mem(v, u)	$v_{top} = u$
	u = a	false	false	true	true	true
	u = b	true	true	false	true	false
$R_{is\_empty}(s_1, v_{is\_empty})$	∀u	$hd(s_1, u)$	$mem(s_1, u)$	Vis_empty		
	u = a	true	true	false		
	u = b	false	false	false		

Table 2. Positive feature vectors extracted from executing concat([a], [b]).

Table 3. Disjoint positive and negative feature vectors of R<sub>top</sub>.

$R_{top}(s_1, v_{top})$	$hd(s_1, u)$	$\textit{mem}(s_1,u)$	$v_{top} = u$
-	false	false	true
+	true	true	true
+	false	false	false

that  $\neg hd(s_1, a) \land \neg mem(s_1, a) \land v_{top} = a$ , for example. A strengthening of the specifications that disallows any one of these interpretations will also rule out the corresponding unsafe run of the program. Put another way, each row corresponds to a potential *negative feature vector*, and a classifier (i.e., specification) for the corresponding placeholder that disallows this feature will disallow the counterexample.

The designation of these features as *potentially* negative is deliberate, as we only want to disallow features that are inconsistent with the implementation of the library functions. As an example, the first feature vector for  $R_{top}$  (the second row of Table 1) states that the result of top s1 is not the head element of the input stack (since  $v_{top} = u$  is true and  $hd(s_1, u)$  is false), and thus is inconsistent with *any* reasonable implementation of top. In contrast, the next feature vector is compatible with an execution of top where, e.g. top([1]) = 1 and u = 2. The second feature vector represents a behavior that is consistent with the underlying library implementation and that should be allowed by the learned specification. The consistency checker generates this positive training data via random testing of the client program. To see how we extract positive feature vectors from training data, consider the execution of concat with the inputs  $s_1 = [a]$  and  $s_2 = [b]$ , which produces the following assignment to program variables:

$$v_{\text{is empty}} = \bot, v_{\text{top}} = a, v_{\text{tail}} = [], v_{\text{concat}} = [b], v = [a; b]$$

Similar to how we built negative feature vectors from CEX, we can construct feature vectors for each function specification from these assignments. Table 2 illustrates the feature vectors corresponding to this assignment. Consider the second row where u is instantiated with a: under this assignment,  $hd(s_1, u) \equiv hd([a], a)$  is true, as are  $mem(s_1, u) \equiv mem([a], a)$  and  $v_{top} = u \equiv a = a$ .

*Classification.* In order to train a classifier, we need to label the extracted feature vectors as either positive or negative. In other words, we need to identify behaviors that should (and should not) be allowed by the inferred specification. Assigning labels is not as straightforward as labelling the feature vectors extracted from counterexamples as negative and those extracted from testing as positive, as the two sets can overlap. We observe this in the vectors of  $R_{top}$  from Table 1 and Table 2: the interpretation  $hd(s_1, u) \mapsto false; mem(s_1, u) \mapsto false; v_{top} = u \mapsto false; occurs in both tables. Intuitively, we do not want to strengthen the specification of <math>R_{top}$  to rule out this interpretation, as the positive sample is a witness that this execution is consistent with the implementation of  $R_{top}$ . Ultimately, therefore, the specification must be relaxed to allow the execution. In cases where a negative feature vector conflicts with a positive feature vector, we identify the potential negative feature vector as positive and remove it from the learner's negative feature vector set. This strategy is similar to the one used by Miltner et al. [2020] to deal with inductiveness counterexamples.

Thankfully, as long as the counterexample set contains at least one feature vector not known to be consistent with the underlying library implementation, we can strengthen the specifications to disallow it. The first feature vector for  $R_{top}$  in Table 1 represents one such infeasible execution. This vector encodes the case where *u* is the output of top but is not a member or head of the input stack. Clearly, no reasonable implementation would support such an interpretation. We use this observation to label as "negative" those feature vectors that are extracted from a counterexample but do not appear in the set drawn from a concrete execution. Table 3 shows the partition of positive and negative feature vectors for  $R_{top}$ .

Given this labelled set of positive and negative feature vectors, our data-driven learner builds a separator over the training data. One such classifier (formula) for the data in Table 3 is  $v_{top} = u \implies hd(s_1, u)$ . Substituting similarly learned specifications for the other library functions in  $\Sigma_{concat}$  equips the SMT solver with enough constraints to rule out CEX while maintaining the invariant that the learned specifications are also consistent with the underlying library implementations. Additional iterations of this counterexample-guided refinement loop gather additional positive and negative features, eventually producing the specifications for the library functions presented in Figure 4.

*Identifying Spurious Counterexamples.* Thus far, we have only considered spurious counterexamples generated by the safety checker. Counterexamples, however, can also result from an incorrect client assertion. For example, suppose the client (unsoundly) asserts :

$$\forall u, mem(v, u) \implies mem(s_2, u)$$

This assertion is wrong, assuming reasonable implementations of top and push, as the elements of the result stack  $\nu$  can also come from s<sub>1</sub>. We distinguish counterexamples corresponding to actual safety violations by first checking if all the feature vectors extracted from the counterexample are included in the set of known positive feature vectors. For example, given this unsound assertion, the verifier may produce the following counterexample:

$$\forall l \ u, (hd(l, u) \iff ((l = v \lor l = s_1) \land u = a)) \land v_{top} = a \land$$

$$(mem(l, u) \iff ((l = v \lor l = s_1) \land u = a))$$
(Cex')

Table 4 shows all of the feature vectors for  $R_{top}$  that are extracted from this counterexample . Since these are a subset of the positive feature vectors from Table 2, there are no feature vectors that can be labeled as negative, and there are thus no new bad behaviors that the learner can use to generate a refined specification mapping that rejects the counterexample. In this scenario, our algorithm tries to discover concrete values of  $s_1$  and  $s_2$  consistent with Cex'; that is,  $s_1$  only

$R_{top}(s_1, \nu_{top})$	∀u	$hd(s_1, u)$	$\textit{mem}(s_1,u)$	$v_{top} = u$
	u=a	true	true	true
	u≠a	false	false	false

Table 4. Feature vectors for Rtop extracted from CEX'.

Table 5. A weakening feature vector of R<sub>push</sub> in Figure 4

$R_{push}(v_{top}, v_{concat}, v) \mid \forall u \mid hd(v_{cor})$	$_{cat}$ , u) $mem(v_{concat}, u)$	) $hd(v, u)$	mem(v, u)	$v_{top} = u$
u = 1   tru	e true	true	true	true

contains a and  $s_2$  is empty. When called with these parameters, concat([a], []) will return  $v \equiv [a]$ , which Elrond returns as a witness of an unsafe execution.

Note that this situation may also occur when the feature set is not large enough, as the specifications in the corresponding hypothesis space are not expressive enough to identify a spurious counterexample. Thus, if we are not able to find inputs that trigger a safety violation, we grow the feature set by increasing the number of quantified variables (*e.g.* from u to u, v) so that Elrond can explore a richer space of specifications.

*Weakening.* While the above strategy is guaranteed to find a safe and consistent verification interface when one exists, the solutions it produces may still be suboptimal, as illustrated in Figure 5. For example, the first conjunct of the specification for push in Figure 4 states that any existing member of both the input and output stacks should not be the same as the element being added to the stack; that is, push always produces a stack with no duplicates. This specification is too restrictive, however, as it disallows reasonable behaviors such as push([1;2], 1) = [1;1;2]. However, if our sampler never generates an observation corresponding to this behavior, e.g.  $x \equiv 1$ ,  $s \equiv [1; 2]$  and  $v_{push} \equiv [1; 1; 2]$ , the candidate specification for top produced by Elrond's first phase will incorrectly disallow it.

In other words, our reliance on testing to identify and label negative feature vectors may result in initial specifications that are overfitted to the examples enumerated by the test generator. There are two potential reasons such a positive example might be missed: (1) the input space of the program might be too large for a test generator to effectively explore, and (2) the provided implementation may simply not exhibit this behavior (e.g., it may be the case that the implementation of push that we are trying to verify against does indeed remove duplicates). While exhaustive or more effective enumeration can address the first cause, it cannot remedy the second. Elrond's weakening phase helps ameliorate both issues.

Our weakening algorithm iteratively weakens candidate specifications, focusing on one library function at a time. To weaken the specification of push, for example, we fix the specifications of the other library functions to their assignments in the current verification interface, and then try to find a maximal weakening of  $R_{push}$  that admits a larger set of implementations of push. To do so, Elrond attempts to discover additional *weakening feature vectors* for  $R_{push}$ , or feature vectors corresponding to behaviors disallowed by the current specification but which would not lead to a violation of client safety. One possible weakening feature vector for our current example is shown in Table 5. Here, the head of both  $s_1$  and the result of the recursive call is 1; this scenario is mistakenly disallowed by the specification of push in Figure 4.

Elrond repeatedly queries the verifier to identify weakening feature vectors for push, indicating that the current specification is maximal when none can be found. It then moves on to the next function specification, iteratively weakening each until a fixpoint is reached. Figure 6 shows the maximal verification interface Elrond builds by weakening the candidate specifications in Figure 4. Compared with Figure 4, the specification of push now permits duplicate elements in stacks. In addition, the specification of is\_empty has been simplified by removing the redundant conjunction  $\neg v \land hd(s, u) \implies mem(s, u)$ , as  $hd(s, u) \implies mem(s, u)$  can never be violated by a concrete stack value thanks to the observed semantics of hd and mem.

$$\begin{split} & \mathsf{R}_{\mathsf{is}\_\mathsf{empty}}(\mathsf{s},v) \mapsto \forall u, (v \implies \neg \mathit{mem}(\mathsf{s},u)) \\ & \mathsf{R}_{\mathsf{top}}(\mathsf{s},v) \mapsto \forall u, v = u \iff (\mathit{hd}(\mathsf{s},u) \land \mathit{mem}(\mathsf{s},u)) \\ & \mathsf{R}_{\mathsf{tail}}(\mathsf{s},v) \mapsto \forall u, (\mathit{mem}(\mathsf{s},u) \implies (\mathit{mem}(v,u) \lor \mathit{hd}(\mathsf{s},u))) \land \\ & ((\mathit{mem}(v,u) \lor (\mathit{hd}(v,u) \land \mathit{hd}(\mathsf{s},u))) \implies \mathit{mem}(\mathsf{s},u)) \\ & \mathsf{R}_{\mathsf{push}}(\mathsf{x},\mathsf{s},v) \mapsto \forall u, (\mathit{mem}(v,u) \land \neg \mathsf{x} = u \implies \mathit{mem}(\mathsf{s},u) \land \neg \mathit{hd}(v,u)) \land \\ & (\mathsf{x} = u \lor \mathit{mem}(\mathsf{s},u) \implies \mathit{mem}(v,u)) \land (\neg \mathit{mem}(v,u) \land \mathit{hd}(\mathsf{s},u)) \end{split}$$

Fig. 6. Maximal specifications for Stack operations. The special variable  $\nu$  represents the result of the method.

### **3 PROBLEM FORMULATION**

Having completed a high-level tour of Elrond in action, we now present a precise description of the specification synthesis problem and our data-driven inference procedure. We consider functional programs that use data structure libraries providing functions to access and construct instances of inductively-defined algebraic datatypes (e.g., list, stacks, trees, heaps, tries, etc.).

In the remainder of the paper, we use  $(\Sigma, \Phi)$  to refer to the *verification query* whose validity we are attempting to establish. These structures serve the same role as verification conditions in a typical verification framework. The first component of this query,  $\Sigma$ , is a conjunction of applications of specification placeholders  $(R_f)$  to arguments; these represent the library method calls made by the client program. The second component,  $\Phi$ , represents the client program's pre- and post-conditions, encoded as sentences built from logical connectives  $(\wedge, \vee, \Longrightarrow)$  over prenex universally-quantified propositional formulae. Each verification query corresponds to a control flow path in the client program; the full algorithm considers the conjunction of all these verification queries at once. To keep the formalization and the description of our algorithms concise, our description considers a single verification query in isolation. The extension to sets of verification queries is provided in the full version of the paper [Zhou et al. 2021a].

*Definition 3.1 (Problem Definition).* A given verification query  $(\Sigma, \Phi)$  with unknown library functions *F* has the form  $\Sigma \implies \Phi$ , where:

$$\Sigma \equiv (\bigwedge_{f \in F, i=0}^{n} R_{f}(\vec{x}_{i}) \land \bigwedge_{x, x' \in \bigcup_{0 \le j \le i} \vec{x}_{j} \cup C} x = x')$$

Here, the equality constraints are either between program variables (x, x') or between variables and constants *C* of some base type (e.g., Booleans and integers). Each  $R_f(\vec{x}_i)$  in  $\Sigma$  is an application of a placeholder predicate to some arguments; the conjunction of these placeholder applications and equality constraints represents a sequence of library method invocations in one control-flow path of the client. To model the input and output behaviors of the blackbox implementations of library functions and method predicates, our formalization relies on a pair of partial functions with the same signature as the implementations. We use partial functions to reflect the fact that we can only observe a subset of the full behaviors of these implementations when searching for specifications.

Definition 3.2 (Specification Configuration). Let *P* be a set of method predicates and *F* be the set of functions in a library used by the client. Let  $\Gamma_f$  be a partial function from the domain of  $f \in F$  to its codomain, and  $\Gamma_p$  be a partial function with the same signature as  $p \in P$ . Let  $\Gamma_P = \bigcup_{p \in P} \Gamma_p$  and  $\Gamma_F = \bigcup_{f \in F} \Gamma_f$ . A specification configuration is a 5-tuple ( $(\Sigma, \Phi), P, F, \Gamma_P, \Gamma_F$ ), where ( $\Sigma, \Phi$ ) is the verification query extracted from the client.

**Example.** The specification configuration of our running example consists of a verification query ( $\Sigma_{concat}$ ,  $\Phi_{concat}$ ), a method predicate set {*memhd*}, and library functions {push, is\_empty, top, tail}. The partial functions in  $\Gamma_F$  and  $\Gamma_P$  abstract over observations on their corresponding blackbox implementations; for an execution that produces the feature vectors shown in Table 2, they are:

$$\begin{split} \Gamma_{mem} &\equiv \lambda (l \in \{s_1, s_2, v_{tail}, v_{concat}, v\}), u. \\ & (u = a \land (l = s_1 \lor l = v)) \lor (u = b \land (l = s_2 \lor l = v_{concat} \lor l = v)) \\ \Gamma_{hd} &\equiv \lambda (l \in \{s_1, s_2, v_{tail}, v_{concat}, v\}), u. \\ & (u = a \land (l = s_1 \lor l = v)) \lor (u = b \land (l = s_2 \lor l = v_{concat})) \\ \Gamma_{push} &\equiv \lambda x, (l \in \{v_{concat})\}).v \\ \Gamma_{is\_empty} &\equiv \lambda (l \in \{s_1\}), false \\ \Gamma_{top} &\equiv \lambda (l \in \{s_1\}), a \\ \Gamma_{tail} &\equiv \lambda (l \in \{s_1\}), v_{tail} \end{split}$$

where stack arguments are limited to values in this particular execution, e.g.  $s_1$  or  $v_{tail}$ .

Given a specification configuration as input, the output of our verification pipeline is a *verification interface* ( $\Delta$ ), a logical interpretation of the method predicates that maps each placeholder predicate for a function  $f \in F$  to a universally-quantified propositional formula over the parameters and result of f. We impose two requirements on  $\Delta$ . The first is *safety*: an underlying theorem prover (e.g., a SMT solver) must be able to prove  $\Sigma[\Delta] \implies \Phi$ , where  $\Sigma[\Delta]$  denotes the formula constructed by replacing all occurrences of specification placeholders with their interpretation in  $\Delta$ , and  $\Sigma \implies \Phi$ is the verification query built from a client program:

Definition 3.3 (Safe Verification Interface). For a given verification query  $(\Sigma, \Phi)$ , a verification interface  $\Delta$  is safe when:

- (1) it makes the VC valid:  $\Sigma[\Delta] \models \Phi$ , and
- (2) is not trivial:  $\Sigma[\Delta] \not\models \bot$

In addition to safety, we also desire that any proposed mapping  $\Delta$  be *consistent* with the provided implementations of method predicates and library functions, *i.e.* that  $\Delta$  must accurately represent their observed behavior. Formally:

Definition 3.4 (Interface Consistency). A verification interface  $\Delta$  is consistent with  $\Gamma_P$  and  $\Gamma_F$  when all specifications in  $\Delta$  are consistent with the inputs on which  $\Gamma_P$  and  $\Gamma_F$  are defined. Formally,

$$\forall f \in F, \Gamma_f \in \Gamma_F, \Gamma_f(\vec{\alpha}) = \nu \Longrightarrow \Delta(R_f)(\vec{\alpha}, \nu)[\Gamma_P]$$

The expression  $\Delta(R_f)(\vec{\alpha}, \nu)$  denotes the instantiation of the formula bound to  $R_f$  in  $\Delta$  with the input arguments  $\vec{\alpha}$  and observed output  $\nu$ . The expression  $\Delta(R_f)(\vec{\alpha}, \nu)[\Gamma_P]$  replaces all free occurrences of p in  $\Delta(R_f)(\vec{\alpha}, \nu)$  with  $\Gamma_p$  where  $\Gamma_p \in \Gamma_P$ .

This definition thus relates the observed behavior of a library method on test data, encoded by  $\Gamma_P$  and  $\Gamma_F$  with its logical characterization provided by  $\Delta$ . Note that there may be many possible verification interfaces for a given specification configuration. In order to identify the best such interface, we use an ordering based on a natural logical inclusion property:

*Definition 3.5 (Interface Order* >). The verification interface  $\Delta'$  is weaker (>) than  $\Delta$  when,

- (1) The two interfaces contain the same functions:  $dom(\Delta) = dom(\Delta')$
- (2) They are not equal:  $\exists R_f \in dom(\Delta), \Delta(R_f) \iff \Delta'(R_f)$
- (3) The specifications in  $\Delta'$  are logically weaker that those in  $\Delta: \forall R_f \in dom(\Delta), \Delta(R_f) \implies \Delta'(R_f)$

Intuitively, weaker verification interfaces are preferable because they place fewer restrictions on the behavior of the underlying implementation. Given an ordering over verification interfaces, we seek to find the weakest safe and consistent interface, *i.e.* one that imposes the fewest constraints while still enabling verification of the client program.

Definition 3.6 (Maximal Verification Interface). For a specification configuration  $((\Sigma, \Phi), P, F, \Gamma_P, \Gamma_F), \Delta$  is a maximal verification interface when:

- (1)  $\Delta$  is safe for the verification query  $(\Sigma, \Phi)$ .
- (2)  $\Delta$  is consistent with  $\Gamma_P$  and  $\Gamma_F$ .
- (3) For a given bound on the number of quantified variables k used by the specifications in Δ, there is no safe and consistent interface Δ' whose specifications use at most k quantified variables such that Δ' > Δ.

We now refine our expectation for the output of our verification pipeline to be not just any safe and consistent verification interface, but also a *maximal* one. Notice that our notion of maximality is parameterized by the number of quantified variables used in the interpretation. As this bound increases, we can always find a weaker specification mapping. Thus we frame our definition of maximality to be relative to the number of quantified variables in the specification.

## 4 LEARNING LIBRARY SPECIFICATIONS

As Section 2 outlined, Elrond frames the search for a safe verification interface as a data-driven learning problem. At a high level, the goal of learning is to build a *classifier* (a function from unlabeled data to a label) from a set of labeled data. More precisely, our goal is to learn classifiers for each of the library functions in a specification configuration that can correctly identify any input and output behavior that could induce an unsafe execution in the client.

Our first challenge is to find an encoding of program executions that is amenable to a data-driven learning framework. To begin, we need to identify the salient *features* used by a classifier to make its decisions.

Definition 4.1 (Feature). A feature of a function for a set of variables  $\vec{x}$  is a method predicate applied to elements of  $\vec{x}$  or equalities between variables in  $\vec{x}$ .

A feature is similar to a literal in first-order logic, but does not allow for method predicates as arguments (e.g. hd(l, mem(l, u))) or constant arguments (e.g. hd(l, 3)).

Definition 4.2 (Feature Set). The feature set of a function f with method predicates P and quantified variables  $\vec{u}$ , denoted as  $S \equiv FSet(P, f(\vec{a}_f) = v_f, \vec{u})$ , is a list of all well-typed features in P for the set of variables  $\vec{a}_f \cup \{v_f\} \cup \vec{u}$  which is minimally linearly independent:

$$\forall \eta' \notin \mathcal{S}, \exists \vec{\eta} \subseteq \mathcal{S}, \eta' \iff \bigwedge \vec{\eta}$$

**Example.** The feature set for the function top :  $list_a \rightarrow a$  from the Stack library, where *a* is some base type, for predicate set  $P \equiv \{hd, mem\}$ , equality operation  $=_a$  and quantified variables  $\vec{u} \equiv \{u : a\}$  is  $FSet(P, top(l) = v, \vec{u}) \equiv [hd(l, u), mem(l, u), v=_a u]$ . Note that the features mem(v, u) and  $l =_a u$  are not included in this set because they are not well-typed. The feature hd(l, v), on the other hand, is omitted because it can be represented by  $hd(l, u) \land v=_a u$  and is thus not linearly independent with respect to the other features in the set. We use *feature vectors* to encode the features of observed tests:

*Definition 4.3 (Feature Vector).* A feature vector is a vector of Booleans that represents the value of each feature in the feature set for some test.

We also need to define the *hypothesis space* of possible solutions considered by our learning system. To easily integrate learned classifiers into the underlying theorem prover, we choose to represent such solutions as Boolean combinations over terms consisting of applications of interpreted base relations and uninterpreted functions. In order to preserve decidability, we limit this space to a subset of effectively propositional sentences. This limitation was expressive enough for all of our benchmarks.<sup>1</sup>

Definition 4.4 (Hypothesis Space). The hypothesis space of specifications for a library function f, method predicate set P, and quantified variables  $\vec{u}$  is the set of formulas in prenex normal form with the quantifier prefix  $\forall \vec{u}$ , and whose bodies are built from  $FSet(P, f(\vec{a}) = v, \vec{u})$ , the logical connectives { $\land, \lor, \neg, \Longrightarrow$ }, and Boolean constants  $\top$  (TRUE) and  $\perp$  (FALSE). The hypothesis space of f over P and  $\vec{u}$  is denoted  $Hyp(P, f(\vec{a}) = v, \vec{u})$ .

In order to classify feature vectors, we ascribe them a semantics in logic:

Definition 4.5 (Unitary classifier). For a given feature vector fv in a feature set  $S \equiv FSet(P, f(\alpha) = v, \vec{u})$ , the logical embedding of fv is a formula encoding the assignment to its features:

$$\llbracket fv \rrbracket \equiv \forall \vec{u}, \bigwedge_{i=0}^{|S|} S[i] \iff fv[i]$$

We say that a classifier  $\phi$  labels a feature vector fv as positive when  $\llbracket fv \rrbracket \implies \phi$ , and negative otherwise.

**Example.** Given the classifier  $\phi \equiv \forall u, hd(s_1, u) \implies v_{top} = u$  for the top function from Section 2, the first row in Table 3 corresponds to the feature vector  $fv^- \equiv \{hd(s_1, u) \mapsto \text{false}; mem(s_1, u) \mapsto \text{false}; v_{top} = u \mapsto \text{true}\}$ . The unitary classifier for fv is  $[\![fv^-]\!] \equiv \forall u, \neg hd(s_1, u) \land \neg mem(s_1, u) \land v_{top} = u$ . fv is labelled negative by  $\phi$ , as  $[\![fv^-]\!] \implies \phi$ . The other two feature vectors in Table 3 are labeled as positive by  $\phi$ .

Proc. ACM Program. Lang., Vol. 5, No. OOPSLA, Article 116. Publication date: October 2021.

116:16

<sup>&</sup>lt;sup>1</sup>The main sorts of properties that we do not support as a consequence of this choice are those which use quantifier alternation, e.g., for every element in a stream, there exists another larger element that appears after it:  $(\forall u, \exists v, mem(l, u) \implies (mem(l, v) \land ord(l, u, v) \land u \le v))$ .

Definition 4.6 (Classification). For a given classifier  $\phi$  and feature set S, it is straightforward to partition the feature vectors of S into positive ( $\phi^+$ ) and negative ( $\phi^-$ ) sets:

$$\begin{split} \phi^+ &\equiv \quad \{ fv \in 2^S \mid \llbracket fv \rrbracket \implies \phi \} \\ \phi^- &\equiv \quad \{ fv \in 2^S \mid \llbracket fv \rrbracket \implies \phi \} \end{split}$$

Notice that these two sets are trivially disjoint:  $\phi^+ \cap \phi^- = \emptyset$ .

For a particular configuration, we can straightforwardly lift this partitioning to verification interfaces:

$$\Delta(R_f)(\vec{\alpha}, v)^+ \equiv \{ fv \in 2^{FSet(P, f(\vec{\alpha}) = v, \vec{u})} \mid \llbracket fv \rrbracket \implies \Delta(R_f)(\vec{\alpha}, v) \}$$
  
$$\Delta(R_f)(\vec{\alpha}, v)^- \equiv \{ fv \in 2^{FSet(P, f(\vec{\alpha}) = v, \vec{u})} \mid \llbracket fv \rrbracket \implies \Delta(R_f)(\vec{\alpha}, v) \}$$

## 4.1 Learning Safe and Consistent Verification Interfaces

We now confront the challenge of how to generate training data from a specification configuration in a way that guarantees the safety of the learned formulas (classifiers). To do so, we extract feature vectors from a set of *logical samples*:

Definition 4.7 (Sample). A sample *s* of a formula  $\phi$  is an instantiation of its quantified variables and a Boolean-valued interpretation for each application of a method predicate to those variables in  $\phi$ , which we denote as  $s \models \phi$ . The positive and negative samples of a verification query ( $\Sigma$ ,  $\Phi$ ) are samples of  $\Phi$  and  $\neg \Phi$ , respectively.

Intuitively, the positive samples of a verification query correspond to safe executions of a client program, while negative samples represent potential violations that safe verification interfaces need to prevent. For example, CEx from Section 2 corresponds to the following negative sample of  $\Phi_{concat}^2$ :

$$\{ s_1 \mapsto l_0; s_2 \mapsto l_1; \nu \mapsto l_2; \nu_{top} \mapsto a; \nu_{tail} \mapsto l_3; \nu_{concat} \mapsto l_4; \nu_{is\_empty} \mapsto \bot;$$

$$hd(l, u) \equiv \{ \}; mem(l, u) \equiv \{ (l_2, a); (l_3, a); (l_4, a) \} \}$$

$$(s^-)$$

and the following sample, extracted from a concrete input and client execution result, is positive:

$$\{ s_1 \mapsto l_0; s_2 \mapsto l_1; v \mapsto l_2; v_{top} \mapsto a; v_{tail} \mapsto l_3; v_{concat} \mapsto l_4; v_{is\_empty} \mapsto \bot;$$
(s<sup>+</sup>)  
 
$$hd(l, u) \equiv \{ (l_0, a); (l_1, b); (l_4, b); (l_2, a) \}; mem(l, u) \equiv \{ (l_0, a); (l_1, b); (l_4, b); (l_2, a); (l_2, b) \} \}$$

Although they come from different sources, both samples provide the values of variables (e.g., the value of  $v_{is\_empty}$  in both  $s^-$  and  $s^+$  is  $\perp$ ) and the values of predicate applications (e.g.,  $hd(v_{concat}, v_{top})$  is true in  $s^-$  sample and false in  $s^+$ ).

Using  $[\cdot]$ , we can *extract* a collection of feature vectors under a feature set S from a sample s:

$$\chi_{\mathcal{S}}(s) \equiv \{ fv \in 2^{\mathcal{S}} \mid s \models \llbracket fv \rrbracket \}$$

For example, the feature vectors extracted from  $s^-$  and  $s^+$  for the feature set of top are shown in the second and third rows of Table 1 and Table 2, respectively.

Definition 4.8 (Classifier Consistency). For a verification query  $(\Sigma, \Phi)$ , we say that a verification interface  $\Delta$  is consistent with a negative sample  $s^-$  if at least one of the library specifications in  $\Delta$  classifies one or more features extracted from that sample as negative:

$$\exists R_f(\vec{\alpha}, \nu) \in \Sigma, \ \exists f\nu \in \chi_{FSet(P, f(\vec{\alpha}) = \nu, \vec{u})}(s^-), f\nu \in \Delta(R_f)(\vec{a}, \nu)^-$$

<sup>&</sup>lt;sup>2</sup>The interpretation of method predicates are represented as binary relations.

Similarly,  $\Delta$  is consistent with a positive sample  $s^+$  if *all* specifications in  $\Delta$  positively identify *every* feature vector extracted from  $s^+$ :

$$\forall R_f(\vec{\alpha}, \nu) \in \Sigma, \ \forall f\nu \in \chi_{FSet(P, f(\vec{\alpha}) = \nu, \vec{u})}(s^+), f\nu \in \Delta(R_f)(\vec{a}, \nu)^+$$

**Example.** The verification interface  $\Delta$  from Figure 4 is consistent with ( $s^-$ ), as the specification of the top function labels as negative the following feature vector of  $s^-$ ,  $fv^- \equiv \{hd(s_1, u) \mapsto false; mem(s_1, u) \mapsto false; v_{top} = u \mapsto true\}$  to negative. Furthermore,  $\Delta$  is also consistent with all the feature vectors extracted from  $s^+$ .

THEOREM 4.9. For a given specification configuration  $((\Sigma, \Phi), P, F, \Gamma_P, \Gamma_F)$  and verification interface  $\Delta, \Sigma[\Delta] \implies \Phi$  is valid iff  $\Delta$  is consistent with all negative samples  $s^-$ ;  $\Delta$  is a consistent interface iff  $\Delta$  is consistent with all positive samples  $s^+$  entailed by  $\Gamma_F$  and  $\Gamma_P$ .<sup>34</sup>

# 4.2 Learning Maximal Verification Interfaces

While Theorem 4.9 identifies the conditions under which a verification interface is safe and consistent, it does not ensure that it is maximal. We frame the search for a maximal solution as a learning problem for a single function specification assuming all other specifications are fixed .

Definition 4.10 (Weakest safe specification). For a given verification query  $(\Sigma, \Phi)$  and safe and consistent verification interface  $\Delta$ ,  $\phi$  is the weakest safe specification of f iff

- (1)  $\Delta[R_f \mapsto \phi]$  is safe
- (2) For a given bound on the number of quantified variables k allowed in the specification of f, there is no other specification φ' with k quantified variables that makes Δ[R<sub>f</sub> → φ'] safe such that φ ⇒ φ'.

Definition 4.11 (Sample with respect to library function). For a verification query  $(\Sigma, \Phi)$ , and safe verification interface  $\Delta$ , a sample *s* is positive (resp. negative) with respect to library function *f* when *s* is positive (resp. negative) and consistent with the specifications of all other library functions in the domain of  $\Delta$ :

$$s \models \Sigma[\Delta[R_f \mapsto \top]] \land \Phi$$
$$s \models \Sigma[\Delta[R_f \mapsto \top]] \land \neg \Phi \quad (\text{resp.})$$

Ideally, the weakest safe specification of f would be able to positively classify every such positive sample. Because this is not possible in general due to the intrinsic granularity of the hypothesis space, we must limit ourselves to covering some subset of this space instead. The weakening relation between classifiers can be viewed as the difference between the sets of positive and negative feature vectors induced by classifiers:

Definition 4.12 (Weakening feature vector and samples). For a verification query  $(\Sigma, \Phi)$  and safe verification interface  $\Delta$ , a weakening feature vector fv distinguishes between a weaker safe specification and  $\Delta(R_f)$ :

$$fv \in \Delta(R_f)^-$$
 and  $\Sigma[\Delta[R_f \mapsto \llbracket fv \rrbracket \lor \Delta(R_f)]] \implies \Phi$ 

A sample *s* is a weakening sample when *s* is positive with respect to f and includes some weakening feature vector. Intuitively, such weakening samples can be used to safely generalize f. If we cannot find any weakening sample, then the specification must have converged to a maximal one.

Proc. ACM Program. Lang., Vol. 5, No. OOPSLA, Article 116. Publication date: October 2021.

116:18

<sup>&</sup>lt;sup>3</sup>All positive samples  $s^+$  entailed by  $\Gamma_F$  and  $\Gamma_P$  means every positive sample consistent with the observations encoded by  $\Gamma_F$  and  $\Gamma_P$ .

<sup>&</sup>lt;sup>4</sup>Proofs for all theorems are available in the full version of the paper [Zhou et al. 2021a] .

Algorithm 1: Multi-Abductive Inference Algorithm.

**Inputs** :Specification configuration  $((\Sigma, \Phi), F, P, \Gamma_F, \Gamma_P)$ **Output** : Maximal verification interface  $\Delta$  or counterexample CEX 1  $\vec{u} \leftarrow 0$ : <sup>2</sup> while true do match SpecInfer( $(\Sigma, \Phi), F, P, \Gamma_F, \Gamma_P, \vec{u}$ ): 3 case Fail s<sup>-</sup> do 4 **match** *ExtractCex*(*s*<sup>-</sup>): 5 case CEX do return CEX; 6 **case** None **do**  $\vec{u} \leftarrow \vec{u} \cup FreshVariable()$ ; 7 **case** Fail None **do**  $\vec{u} \leftarrow \vec{u} \cup FreshVariable()$ ; 8 case  $\Delta$  do repeat 10  $\Delta_0 \leftarrow \Delta;$ 11 **for**  $f \in F$  **do**  $\Delta(R_f) \leftarrow Weaken((\Sigma, \Phi), P, \Delta, f, \vec{u})$ ; 12 until  $\Delta = \Delta_0$ ; 13 return  $\Delta$ ; 14

Example. The sample

$$\{ s_1 \mapsto l_0; s_2 \mapsto l_1; v \mapsto l_2; v_{top} \mapsto a; v_{tail} \mapsto l_3; v_{concat} \mapsto l_4; v_{is\_empty} \mapsto \bot; \\ hd(l, u) \equiv \{ (l_0, a); (l_1, a); (l_2, a); (l_4, a) \}; mem(l, u) \equiv \{ ((l_0, a); (l_1, a); (l_2, a); (l_4, a) \} \}$$

is a positive sample with respect to push since it makes both input stacks and the output stack contain a and also have the head element a, which entails  $\Phi_{concat}$ . It is also consistent with all the specifications in the verification interface from Figure 4 outside of the one for push. The feature vector fv shown in Table 5 is extracted from this sample, and it is not included in  $\Delta(R_{push})^+$ . Moreover,  $\Sigma[\Delta[R_{push} \mapsto [\![fv]\!] \vee \Delta(R_{push})]] \implies \Phi$ , thus fv is a weakening feature vector for push, and the above sample is a weakening sample for push.

THEOREM 4.13. For a given specification configuration  $((\Sigma, \Phi), P, F, \Gamma_P, \Gamma_F)$  and safe verification interface  $\Delta$ ,  $\Delta(R_f)$  is weakest safe specification of f if and only if there are no weakening samples for f.

## 5 ALGORITHM

Algorithm 1 presents the abductive inference procedure depicted in Figure 1. The algorithm maintains a list of variables  $\vec{u}$  that defines the hypothesis space it is currently exploring. The main loop of the algorithm searches for a maximal verification interface in this hypothesis space, starting from an initially empty set of variables (line 1). We first try to use the *SpecInfer* subalgorithm (Algorithm 2) to infer a verification interface that is safe and consistent with the specification configuration. If inference fails with a sample, we search for corresponding inputs that trigger a safety violation using random testing via the *ExtractCex* subroutine (line 8). If we are unable to find such unsafe inputs, we extend  $\vec{u}$  with a fresh variable (lines 7-8) and restart the loop under this enhanced hypothesis space. Otherwise, *SpecInfer* has found a safe and consistent verification interface,  $\Delta$ , which we then refine to a maximal solution via a weakening loop. This loop iteratively weakens the specification of each library function in  $\Delta$  using the *Weaken* subalgorithm (Algorithm 3), and returns  $\Delta$  once it has reached a fixed-point.

### 5.1 Specification Inference

Algorithm 2: Safe and Consistent Specification Inference (SpecInfer)

**Inputs** :Specification configuration  $(\Sigma, \Phi), F, P, \Gamma_F, \Gamma_P$ , variables in hypothesis space  $\vec{u}$ Output : Consistent safe verification interface or reports failure 1 for  $f \in F$  do  $\pi_f, \omega_f \leftarrow \emptyset, \emptyset$ ; <sup>2</sup>  $\Delta \leftarrow \{R_f \mapsto \forall \vec{u}, \top\};$ while true do 3 while true do 4 match Sample( $\Delta(R_f), \Gamma_F, \Gamma_P$ ): 5 case None do break; 6 case Some s<sup>+</sup> do 7 for  $f \in F$  do 8  $\pi_f \leftarrow FvecFromSample(FSet(P, f(\vec{\alpha}_f) = v_f, \vec{u}), s^+) \cup \pi_f;$ 0  $\omega_f \leftarrow \omega_f \setminus \pi_f$ 10 **for**  $f \in F$  **do**  $\Delta(R_f) \leftarrow Learner(\pi_f, \omega_f)$ ; 11 match  $Verify((\Sigma[\Delta] \implies \Phi))$ : 12 case OK do 13 match *Verify*( $\neg \Sigma[\Delta]$ ): 14 case OK do return Fail None; 15 case *Sat* \_ do return  $\Delta$ ; 16 case Sat s<sup>-</sup> do 17 **for**  $f \in F$  **do**  $\omega'_f \leftarrow FVecFromSample(FSet(P, f(\vec{\alpha}_f) = v_f, \vec{u}), s^-) \setminus \pi_f$ ; 18  $\mbox{if } \bigwedge_{f \in F} \omega'_f = \emptyset \ \mbox{then return } Fail\, s^- \ ; \\$ 19 else for  $f \in F$  do  $\omega_f \leftarrow \omega_f \cup \omega'_f$ ; 20 for  $f \in F$  do  $\Delta(R_f) \leftarrow Learner(\pi_f, \omega_f)$ ; 21

Our data-driven multi-abduction specification inference algorithm (*SpecInfer*) given by Algorithm 2 assumes three key components: a property-based random sampler *Sample*, a satisfiability checker *Verify*, and a classifier learner *Learner*. *Sample* takes a first-order formula as input and attempts to find a counterexample via random sampling, returning data instances which violate this specification based on these counterexamples. *Verify* returns OK when given a valid formula, and produces an assignment of variables and member predicates (i.e. a sample) that invalidates its input formula otherwise. *Learner* takes two disjoint sets of feature vectors,  $\pi$  and  $\omega$ , as inputs, and returns a formula that classifies elements of  $\pi$  and  $\omega$  as positive and negative, respectively. Elrond uses a decision tree algorithm [Ruggieri 2002] as the underlying learning framework.

In more detail, for each library function f, SpecInfer maintains a pair of positive and negative sets of feature vectors  $\pi_f$  and  $\omega_f$ ; these are initially empty (line 1). The algorithm holds the current candidate verification interface in the variable  $\Delta$ , which initially stores trivial specifications for each library function. SpecInfer implements a two-step refinement loop which first calls Sample to test the current candidate verification interface  $\Delta$  against  $\Gamma_F$  and  $\Gamma_P$  in order to identify any library function behaviors that are inconsistent with their current specifications (line 5). If so, Sample returns a positive sample s<sup>+</sup> that we then use to augment the set of positive feature vectors (line 7). We subsequently remove any feature vectors appearing in  $\pi_f$  from  $\omega_f$  in order to ensure the

Algorithm 3: Weakening (Weaken)

```
Inputs : (\Sigma, \Phi), P, \Delta, f, \vec{u}
    Output: Weakest safe specification of f
1 W_f \leftarrow \forall \vec{u}, \perp; \pi, \omega \leftarrow \emptyset, \emptyset;
    while true do
2
          match Verify(\Sigma[\Delta[R_f \mapsto \forall \vec{u}, \top]] \land \Phi \implies \Sigma[\Delta[R_f \mapsto \Delta(R_f) \lor W_f \lor [\![\omega]\!]]):
 3
                 case OK do
 4
                       if W_f = \forall \vec{u}, \perp then return \Delta(R_f) else return W_f \vee \Delta(R_f);
 5
                 case Sat s do
 6
                        W_f, \pi, \omega \leftarrow Update(s, W_f, \pi, \omega);
 7
                        W_f, \pi, \omega \leftarrow SafetyLoop (W_f, \pi, \omega);
 8
    Procedure Update(s, W_f, \pi, \omega)
9
          for fv \in FVecFromModel(Fset(P, f(\alpha_f) = v_f, \vec{u}), s) do
10
                 match Verify((\Sigma[\Delta[R_f \mapsto W_f \lor \Delta(R_f) \lor \llbracket fv \rrbracket)] \implies \Phi)):
11
                       case OK do \pi \leftarrow \pi \cup \{fv\};
12
                    case Sat _ do \omega \leftarrow \omega \cup \{fv\};
13
                 W_f \leftarrow Learner(\pi, \omega)
14
          return W_f, \pi, \omega
15
    Procedure SafetyLoop(W_f, \pi, \omega)
16
           match Verify((\Sigma[\Delta[R_f \mapsto W_f \lor \Delta(R_f)]] \implies \Phi)):
17
                 case OK do return W_f, \pi, \omega;
18
                 case Sat s do
19
                        W_f, \pi, \omega \leftarrow Update(s, W_f, \pi, \omega);
20
                        return SafetyLoop(W_f, \pi, \omega)
21
```

two sets are disjoint before invoking *Learner* to ascertain whether  $\Delta$  is consistent with the newly observed behaviors (lines 9-11).

SpecInfer then queries Verify to determine if the current verification interface is sufficient to ensure the desired safety property (line 12). We then verify that the specifications are not contradictory (line 14) before returning them, in order to ensure that they are not vacuously safe. Alternatively, if the verifier returns a sample  $s^-$  witnessing an unsafe execution that  $\Delta$  was not strong enough to disallow, we extract sets of feature vectors for each library function f from  $s^-$ , filtering out any that also occur in  $\pi_f$  (line 18) to maintain the invariant that  $\pi_f$  and  $\omega_f$  are disjoint. If no new negative feature vectors are found, *SpecInfer* cannot ensure that  $s^-$  is spurious in the current hypothesis space, and it thus fails, returning Fail  $s^-$  to the main algorithm. Otherwise, the algorithm updates  $\omega_f$  with the new negative feature vectors and strengthens  $\Delta$  to rule out  $s^-$  by invoking *Learner* and the refinement loop restarts (line 21).

THEOREM 5.1 (SPECINFER IS TOTAL AND SOUND). Algorithm 2 always halts and returns a verification interface that is safe and consistent.

## 5.2 Weakening

Algorithm 3 shows the *Weaken* procedure that is used to generalize the specification of a library function f from the initial solution returned by *SpecInfer*. This procedure takes a verification query, method predicate set, verification interface, library function, and a list of quantified variables as

input, and infers a weakest safe specification for f. Notably, the arguments of *Weaken* do not include  $\Gamma_P$  or  $\Gamma_P$ , reflecting that our approach to weakening is purely logical.

Weaken maintains a formula  $W_f$  which is used to weaken the specification of f; the start of the function initializes this formula to false (line 1). Like SpecInfer, Weaken maintains two sets of feature vectors:  $\pi$ , which holds weakening feature vectors, and  $\omega$ , which holds feature vectors representing library behaviors that cause client safety violations. Both sets are initially empty (line 1). The body of Weaken uses a counterexample-guided refinement loop to iteratively construct  $W_{f}$ . The loop first queries the verifier to ensure that every safe execution consistent with the candidate specifications besides  $\Delta(R_f)$  is either a) currently allowed by  $W_f \vee \Delta(R_f)$  or b) involves a behavior of f that can trigger a safety violation in other contexts (line 3). The clauses of the formula representing each of these pieces are highlighted with the corresponding color. When this query is valid, there is no weakening sample for f; thus, by Theorem 4.13, the current specification for f is the weakest safe one, and Weaken returns  $W_f \vee \Delta(R_f)$  (or  $\Delta(R_f)$  when  $W_f \equiv \bot$ ) (line 5). Alternatively, Verify will provide a sample s witnessing a safe execution; this sample is passed to the Update subroutine to further weaken  $W_f$ . Update first uses Verify to identify feature vectors extracted from s that trigger safety violations in other contexts. After inserting safe and unsafe feature vectors into  $\pi$  and  $\omega$ , respectively (lines 12-13), Update uses Learner to refine  $W_f$ . The weakened specification  $W_f \lor \Delta(R_f)$  may violate client program safety, however. We use a learning based procedure for its intrinsic generalization ability: when we add a new feature vector to  $\pi$ , the learned classifier will not only classify *that* new feature vector to positive, but also potentially others. This ability results in faster convergence, since we do not need to explicitly enumerate each weakening feature vector. However, to avoid the learned classifier from being weakened to the point that it mistakenly classifies unsafe feature vectors as safe, Weaken then invokes SafetyLoop to filter any unsafe weakenings (line 8). SafetyLoop implements a similar strategy as Algorithm 2 to progressively refine  $W_f$  until it guarantees safety. After this step, the loop restarts in order to identify additional weakening possibilities .

THEOREM 5.2 (WEAKEN IS SOUND). For given  $(\Sigma, \Phi)$ ,  $P, \Delta, f, \vec{u}$ , if  $\Delta$  is safe, Algorithm 3 will halt with the weakest safe specification for f.

THEOREM 5.3 (RELATIVE COMPLETENESS). Algorithm 1 either returns either a maximal verification interface or a concrete counterexample.

#### 6 IMPLEMENTATION AND EVALUATION

We have implemented a verification pipeline based on the above approach, called Elrond[Zhou et al. 2021b], that targets OCaml programs which rely on libraries to manipulate algebraic data types. Elrond consists of 7267 lines of OCaml and uses Z3 [de Moura and Bjørner 2008] as its backend solver. Elrond's frontend generates verification queries from client programs via weakest liberal precondition predicate transformers. Elrond does not automatically infer inductive invariants for recursive client functions, and it expects client programs to provide such invariants, in addition to any pre- and post-conditions.

Our experimental evaluation of Elrond addresses five key questions:

- **Q1**: Is Elrond able to find specifications sufficient to verify a range of properties and client programs in a reasonable amount of time?
- Q2: Can Elrond identify unsafe client programs?
- Q3: Can Elrond efficiently find maximal solutions?
- Q4: Is Elrond able to find useful intermediate generalizations of initial specifications?
- Q5: Does weakening improve the quality of the inferred specifications?

Table 6. Experimental results. The columns in the table can be divided into three groups. The first group presents the number of distinct library functions (|F|), the number of library function applications (|R|), and the size of the method predicate set (|P|) for each benchmark. The second column describes the number of quantified variables ( $|\vec{u}|$ ), the number of counter-examples generated (|cex|), and the time in seconds (time<sub>c</sub>) needed for Elrond to find a consistent and safe verification interface. Times indicate how long it took for counter-example generation, sampling, feature vector extraction, and labelling and DT learning. They do not include the time taken to generate the initial verification query. Red entries in the time column indicate how long it took to identify safety violations in the five unsafe benchmarks. The last group lists the number of gathered (#Gather) and total positive ( $|\phi^+|$ ) feature vectors in the space of weakenings, the time needed by the weakening phase (time<sub>w</sub>), and the time needed for the SMT solver to find a sample allowed by a weakened solution but not the initial one (time<sub>d</sub>). Blue entries in the  $|\phi^+|$  column indicate a lower bound. "Limit" entries in the time<sub>w</sub> column indicate that the one hour time bound was reached. "Max" entries in the time<sub>d</sub> column indicate that the initial solution was already maximally weak.

	( F ,  R )	P	$ \vec{u} $	cex	$time_c(s)$	#Gather / $ \phi^+ $	$time_w(s)$	$time_d(ms)$
	(6,10)	2	2	26	1.85	366 / 1828	47.9	123.5
0110110	(4, 5)	2	1	10	0.28	235 / 1536	9.2	18.2
queue	(4, 5)	3	2	23	0.87	2076 / 30054	318.9	42.7
	(6, 10)	1	1	8	0.22	53 / 114	1.9	Max
	(4,5)	2	1	11	0.61	29 / 13	0.5	18.5
stack	(4, 5)	3	2	39	2.94	3220 / 22716	Limit	169.3
	(4, 5)	2	1	4	0.08			
	(3,8)	2	2	71	11.21	1790 / 26588	Limit	895.8
	(2, 21)	1	1	10	0.37	172 / 432	21.5	21.3
heap	(2, 21)	3	2	141	107.18	702 / 177388	Limit	328.3
	(2, 21)	3	2	192	152.99	512 / 190215	Limit	149.7
	(3,8)	1	1	10	0.37	464 / 3308	46.0	Max
	(4, 10)	1	1	4	0.15	50 / 110	1.3	21.1
stream	(4, 10)	2	2	23	1.57	640 / 1376	1038.0	70.6
stream	(4, 10)	3	2	18	1.69	1755 / 31276	741.8	24.3
	(4, 10)	3	2	12	0.37			
	(1, 21)	4	2	91	77.75	1369 / 59797035	Limit	1268.7
	(2,8)	1	1	11	0.30	181 / 448	11.2	21.1
	(2,8)	2	1	21	0.81	3327 / 27609	1741.9	23.2
	(2,8)	2	2	63	8.28	2496 / 28495	Limit	314.1
set	(2,6)	3	1	15	0.39	3644 / 22526	1203.3	19.5
	(2,6)	3	1	25	0.69	3091 / 20090	968.7	19.2
	(1, 21)	1	1	10	0.27	228 / 772	34.0	Max
	(2,8)	2	1	8	0.14			
	(2,6)	3	3	29	4.69			
trio	(4, 18)	2	1	12	0.47	167 / 404	13.6	28.0
trie	(4, 18)	2	1	15	0.31			

All reported data was collected on a Linux server with an Intel(R) Core(TM) i7-8700 CPU @ 3.20GHz and 64GB of RAM.

To answer these questions, we have evaluated Elrond on a corpus<sup>5</sup> of client programs drawn from Okasaki [1999], the OCaml standard library [Leroy et al. 2014], Verified Functional Algorithms [Appel 2018] and Software Foundations [Pierce et al. 2010]. Our benchmarks cover a range

<sup>&</sup>lt;sup>5</sup>All benchmarks, post-conditions, and inferred library specifications from our evaluation are provided in the full version of the paper [Zhou et al. 2021a] .

of abstract data types manipulating a diverse set of algebraic data types, including queues (bankers queues and batched queues), list-based stacks, heaps (leftist heaps and splay heaps), streams, sets (backed by trees, red black trees, and lists), and tries. The underlying representation of each algebraic data type provides a set of method predicates related to ordering, membership, and uniqueness (no duplicate elements) that can be queried as part of a test. These sorts of shape properties are relatively under-constrained and are thus amenable to property-based random sampling. We used Elrond to verify several different properties for each data type, including membership, ordering, distinct elements, and sorting.

*Example 6.1.* One of our benchmarks uses the following specification for an insert program that inserts an element x into an unbalanced set s using a binary tree for its underlying representation:

$$(\forall u, (root(s, u) \implies (u < x \implies root(v_{insert}, u)) \land (u \ge x \implies root(v_{insert}, x))) \land (mem(v_{insert}, u) \iff mem(s, u) \lor x = u)$$

Given this specification, Elrond infers the following specification for the maket(x, l, r) = v function used by add that constructs a new tree from x and the left and right subtrees 1 and r:

$$\forall u, (mem(v, u) \iff mem(1, u) \lor mem(r, u) \lor (x = u)) \land (root(v, u) \iff x = u) \\ \land (root(1, u) \implies mem(1, u)) \land (root(r, u) \implies mem(r, u))$$

The first line of this specification captures the key semantic properties of the tree, while the second line encodes a key relationship between the *mem* and *hd* method predicates needed by the solver to verify the specification.

The detailed results of our experiments are shown in Table 6. The first group of columns in Table 6 describes the salient features of our benchmarks. Each client specification uses between 1 and 4 member predicates, and the client programs make between 5 and 21 calls to library functions. In order to evaluate Elrond's ability to identify faulty specifications, our experiments also included five unsafe client programs. Elrond can infer specifications for all the benchmarks with valid assertions, and returns concrete counter-examples for clients with unsafe post-conditions.

The second group of columns in Table 6 presents our evaluation of Elrond's ability to discover an initial safe and consistent verification interface (Q1 and Q2). These columns show that Elrond is relatively efficient at finding safe specifications, with none of the benchmarks taking longer than three minutes to learn an initial solution (Q1). As expected, more complicated benchmarks (i.e. those with more function calls or member predicates) required more SMT-provided counterexamples and took longer to complete, as did benchmarks requiring a larger hypothesis space (i.e. specifications with more quantified variables). Notably, Elrond was able to quickly generate concrete counterexamples witnessing safety violations (Q2) in the five benchmarks with invalid postconditions (these benchmarks have red entries in the time<sub>c</sub> column).

The final group of columns in Table 6 addresses the questions dealing with Elrond's weakening phase (Q3 - Q5). Unsurprisingly, weakening requires more time than the initial inference phase; while the latter needs to identify a *single* solution, the former needs to account for *all* possible solutions in order to select the best one. When evaluating this phase, we choose to use a one hour time bound for each experiment. If this time bound was reached, we had Elrond return the current weakened solution. Three-fourths (16/22) of our safe benchmarks were able to find maximal solutions from the initial solution within this limit (Q3). In general, the time taken to weaken a solution was correlated with the complexity of the benchmark.

To further investigate how effective Elrond was at exploring the space of candidate weakenings (Q4), we calculated the ratio of the total number of feature vectors gathered during weakening

(#Gather) against the total number of positive feature vectors admitted by the final maximal specification ( $|\phi^+|$ ). The latter number represents the set of vectors that a naïve exhaustive enumeration would need to consider in order to find our solution. In our experiments, Elrond only needed to consider at most 40% of the full search space.<sup>6</sup> To get  $|\phi^+|$  for the six benchmarks on which Elrond returned a partially weakened solution, we increased the time bound to 24 hours and ran Elrond until it converged on the maximal solution. Elrond was unable to converge under this longer bound for three of these benchmarks. For these three, we report the total number of feature vectors in the partial solution, which is indicated using a blue entry in the  $|\phi^+|$  column.

In an attempt to quantify how well Elrond was able to generalize an initial solution, we asked the SMT solver to identify samples permitted under the weakened solution but not the original (Q5). The more general a weakened solution, the larger this space should be, so the solver should be able to more readily identify one of its elements. The results of this experiment are presented under the time<sub>d</sub> column. The initial solution was already maximal for three of our benchmarks, which is indicated via a "Max" entry in the time<sub>d</sub> column. In general, the solver was able to quickly find such samples for the remaining nineteen benchmarks. This search took longer for the benchmarks that hit the time bound, suggesting those solutions are either closer to the initial specifications or are otherwise more complicated for the SMT solver to handle.



Fig. 7. Run times of the first 1000 weakening iterations in benchmarks.

When running our experiments without a time bound, we observed that individual iterations of the weakening loop took longer to complete over time. As an example, in the first hour of the second stack benchmark, Elrond performed 46 weakening steps per minute on average, but only averaged 33 after that. Based on this observation, we conjectured that it became harder for Elrond to find further generalizations as the weakening phase progresses. To test this hypothesis, we measured the time required by each of the first 1000 iterations of the weakening loop for the six benchmarks that hit the time bound. The results of this experiment are presented in Figure 7; the y-axis uses a logarithmic scale in order to account for times that range from 0.1s to 50s. The trajectory agrees with our hypothesis, and suggests diminishing returns for running the weakening procedure for long periods of time.

Finally, in order to show that Elrond produces useful specifications (**Q1** and **Q5**), we used the Coq proof assistant to verify that implementations of the library functions from our benchmarks satisfied their inferred specifications.<sup>7</sup> In total, Elrond inferred 68 specifications across all our experiments; we were able to verify 64 of these in Coq. Of these, three had initial specifications that were too strong; it was only after the weakening phase that the specifications could be used for verification. The four specifications that we were not able to verify came from the benchmarks

<sup>&</sup>lt;sup>6</sup>In the first stack benchmark, Elrond gathers more feature vectors than the total number because the algorithm may generate false feature vectors, as discussed in Section 5.

<sup>&</sup>lt;sup>7</sup>The corresponding Coq developments are provided in the full version of the paper [Zhou et al. 2021a] .

whose weakening phase timed out, i.e., these specifications were not maximal. Taken together, these points suggest that the weakening phase results in more general (and thus more useful) specifications, and that Elrond is highly effective in finding meaningful specifications.

## 7 RELATED WORK

Data-Driven Approaches. There have been several recent data-driven approaches that use learning to infer specifications of library functions. Zhu et al. [2016] automatically infer specifications that use a fixed set of features (analogous to our method predicates) to identify relationships between the input and outputs of a function. Padhi et al. [2016] use program synthesis to automatically learn features on demand when inferring preconditions for data-structure manipulating library functions. This approach is further extended by Miltner et al. [2020] to synthesize [Osera and Zdancewic 2015] representation invariants that are sufficient to verify specifications of the abstract datatype operations. Rather than inferring specifications of libraries in isolation, we consider the complementary problem of discovering those specifications in service of verifying library clients. Broadly speaking, this change in perspective differentiates our approach from these prior works insofar as we cannot appeal to any source of ground truth (e.g., SMT verification as in Padhi et al. [2016] or bounded model checking as in Miltner et al. [2020]). This leads to the need for a weakening procedure to facilitate generalization to ultimately aid client-side verification. However, techniques like Miltner et al. [2020] could be used alongside ours, for example, by leveraging abduced specifications to infer and verify representation invariants.

Automated Verification. Encoding verification conditions in a logic for which efficient solvers exist (e.g. SMT) is ubiquitous in the automatic program verification community. In this setting, the standard approach to reasoning about clients of user-defined functions is to rely on some manually written axiomatization of those functions, paying particular care in order to ensure that the underlying solver will terminate [Itzhaky et al. 2014, 2013]. In the case that a specification is incomplete, e.g., when defining new functions, manual intervention is required to extend the axiomatization. Informally, our approach can be thought of as filling in the missing parts of specifications by using the latent semantics of method predicates. More recently, Vazou et al. [2017] introduced refinement reflection in order to enable SMT-based reasoning about arbitrary user-defined functions. There, the semantics of a function are embedded directly into the logic as a set of equations, and users can *manually* construct equational proofs about their behavior using a library of proof combinators. The authors introduce a proof search algorithm to help automate the construction of these proofs. While this algorithm is complete when a proof exists for a bounded unfolding of function definitions, users are still required to provide instantiations of lemmas and induction hypotheses to completely automate program verification. Our approach uses data-driven methods and counterexample-guided search to generate specifications without the need to reflect the implementation, which in our setting is unavailable, into the solver's underlying logic.

*Abductive Inference.* As noted in Section 2, our logical formulation of specification inference is an instance of an abductive inference problem. This observation has been previously exploited to develop inference algorithms for loop invariants [Dillig et al. 2013] and specifications of functions in a client program [Albarghouthi et al. 2016]. For a given program, both algorithms rely on an abduction procedure to iteratively strengthen loop invariants (resp. function specifications) until they are strong enough to prove a user-provided post-condition. While completely automated, these approaches critically rely on an abduction procedure for the underlying specification logic, in particular the first-order theory of linear integer arithmetic in their experiments. To the best of our knowledge, no such abduction procedure exists for the theory of equalities with uninterpreted function symbols that is commonly used to specify recursive functions over algebraic datatypes.

Our approach provides an alternative solution that combines data-driven methods with SMT-based counterexample-guided refinement to discover library method specifications to aid client-side verification.

While Bastani et al. [2015] consider a similar abductive inference problem – discovering minimal sufficient assumptions in order to analyze client programs when some portion of the code is missing – there are several major differences between our two approaches: 1) their specification domain is limited to CFL reachability (alias and taint specifications), while our hypothesis space involves shape properties (membership, ordering, etc.) over algebraic datatypes; 2) we guarantee specification consistency with respect to our observations on the blackbox implementation of the libraries; and (3) our approach incorporates an explicit weakening procedure to generalize candidate solutions.

*Specification Inference.* Similar to our problem setup, Nguyen et al. [2014] and Su et al. [2018] infer specifications of library APIs from client information. However, both works try to infer specifications of correct library usage under the assumption that the client is itself correct. In contrast, we leverage abductive methods to provide guarantees that any library implementation must satisfy to ensure the client is safe. Our approach does not assume clients are always safe – e.g, Elrond was able to identify safety violations in 5 of our benchmarks. Another distinguishing feature of our work is our focus on generating maximally weak specifications to overcome overfitted specifications. The work of Pandita et al. [2012] learns specifications from comments in natural language but does not provide safety and consistency guarantees.

Qin et al. [2010] also infers specifications for unknown procedures in a rich domain involving shape properties of datatypes. However, their technique does not query library implementations and does not generalize inferred specifications, which can potentially lead to specifications that are overfitted to the client safety property. In contrast to our work, their method also requires users to understand the underlying representation of the datatype used by the library method, and present explicit interpretations of predicates used in specifications sufficient to verify the client. Finally, we apply an additional algorithmic weakening procedure to refine inferred specifications.

*CHC solving.* While it is possible to frame our problem as an instance of data-driven CHC solving à la Zhu et al. [2018], this would not address the important challenge of generating maximally weak specifications, a concern that is not considered by Horn-clause semantics [Albarghouthi et al. 2016].

## 8 CONCLUSIONS

This paper presents a novel data-driven approach to infer specifications using a minimal set of assumptions that are nonetheless consistent with provided blackbox library implementations and sufficient to verify client assertions. We demonstrate that our technique, manifested in a tool called Elrond, is highly effective in identifying sophisticated specifications that enable verification of challenging functional data structure programs.

## ACKNOWLEDGEMENTS

We thank Pedro Abreu and the anonymous reviewers for their detailed comments and suggestions. This material is based upon work supported by the NSF under Grants CCF-SHF 1717741, CCF-FMiTF 2019263, and CCF-1755880.

## REFERENCES

Aws Albarghouthi, Isil Dillig, and Arie Gurfinkel. 2016. Maximal Specification Synthesis. In *Proceedings of the 43rd* Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (St. Petersburg, FL, USA) (POPL '16). Association for Computing Machinery, New York, NY, USA, 789–801. https://doi.org/10.1145/2837614.2837628 Andrew Appel. 2018. Software Foundations Volume 3: Verified Functional Algorithms.

- Osbert Bastani, Saswat Anand, and Alex Aiken. 2015. Specification Inference Using Context-Free Language Reachability. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Mumbai, India) (*POPL '15*). Association for Computing Machinery, New York, NY, USA, 553–566. https://doi.org/10.1145/2676726. 2676977
- Koen Claessen and John Hughes. 2000. QuickCheck: A Lightweight Tool for Random Testing of Haskell Programs. In Proceedings of the Fifth ACM SIGPLAN International Conference on Functional Programming (ICFP '00). Association for Computing Machinery, New York, NY, USA, 268–279. https://doi.org/10.1145/351240.351266
- Leonardo de Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In *Tools and Algorithms for the Construction and Analysis of Systems*, C. R. Ramakrishnan and Jakob Rehof (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 337–340. https://doi.org/10.1007/978-3-540-78800-3\_24
- Isil Dillig, Thomas Dillig, Boyang Li, and Ken McMillan. 2013. Inductive Invariant Generation via Abductive Inference. In Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages and Applications (Indianapolis, Indiana, USA) (OOPSLA '13). Association for Computing Machinery, New York, NY, USA, 443–456. https://doi.org/10.1145/2509136.2509511
- Shachar Itzhaky, Anindya Banerjee, Neil Immerman, Ori Lahav, Aleksandar Nanevski, and Mooly Sagiv. 2014. Modular Reasoning about Heap Paths via Effectively Propositional Formulas. *SIGPLAN Not.* 49, 1 (Jan. 2014), 385–396. https: //doi.org/10.1145/2578855.2535854
- Shachar Itzhaky, Anindya Banerjee, Neil Immerman, Aleksandar Nanevski, and Mooly Sagiv. 2013. Effectively-Propositional Reasoning about Reachability in Linked Data Structures. In Proceedings of the 25th International Conference on Computer Aided Verification - Volume 8044 (Saint Petersburg, Russia) (CAV 2013). Springer-Verlag, Berlin, Heidelberg, 756–772. https://doi.org/10.5555/2958031.2958053
- Xavier Leroy, Damien Doligez, Alain Frisch, Jacques Garrigue, Didier Rémy, and Jérôme Vouillon. 2014. The OCaml system release 4.02. *Institut National de Recherche en Informatique et en Automatique* 54 (2014).
- Anders Miltner, Saswat Padhi, Todd Millstein, and David Walker. 2020. Data-Driven Inference of Representation Invariants. In Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation (London, UK) (PLDI 2020). Association for Computing Machinery, New York, NY, USA, 1–15. https://doi.org/10.1145/3385412.3385967
- Hoan Anh Nguyen, Robert Dyer, Tien N. Nguyen, and Hridesh Rajan. 2014. Mining Preconditions of APIs in Large-Scale Code Corpus. In Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (Hong Kong, China) (FSE 2014). Association for Computing Machinery, New York, NY, USA, 166–177. https://doi.org/10.1145/2635868.2635924
- Chris Okasaki. 1999. Purely Functional Data Structures. Cambridge University Press, USA.
- Peter-Michael Osera and Steve Zdancewic. 2015. Type-and-Example-Directed Program Synthesis. In Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (Portland, OR, USA) (PLDI '15). Association for Computing Machinery, New York, NY, USA, 619–630. https://doi.org/10.1145/2737924.2738007
- Saswat Padhi, Rahul Sharma, and Todd Millstein. 2016. Data-Driven Precondition Inference with Learned Features. In Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (Santa Barbara, CA, USA) (PLDI '16). Association for Computing Machinery, New York, NY, USA, 42–56. https://doi.org/10.1145/2908080. 2908099
- Rahul Pandita, Xusheng Xiao, Hao Zhong, Tao Xie, Stephen Oney, and Amit Paradkar. 2012. Inferring method specifications from natural language API descriptions. In 2012 34th International Conference on Software Engineering (ICSE). 815–825. https://doi.org/10.1109/ICSE.2012.6227137
- Benjamin C Pierce, Chris Casinghino, Marco Gaboardi, Michael Greenberg, Cătălin Hriţcu, Vilhelm Sjöberg, and Brent Yorgey. 2010. Software foundations. *Webpage: http://www.cis.upenn.edu/bcpierce/sf/current/index.html* (2010).
- Shengchao Qin, Chenguang Luo, Guanhua He, Florin Craciun, and Wei-Ngan Chin. 2010. Verifying Heap-Manipulating Programs with Unknown Procedure Calls. In *Formal Methods and Software Engineering*, Jin Song Dong and Huibiao Zhu (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 171–187. https://doi.org/10.1007/978-3-642-16901-4\_13
- S. Ruggieri. 2002. Efficient C4.5 [classification algorithm]. IEEE Transactions on Knowledge and Data Engineering 14, 2 (2002), 438–444. https://doi.org/10.1109/69.991727
- Jingyi Su, Mohd Arafat, and Robert Dyer. 2018. Using Consensus to Automatically Infer Post-Conditions. In Proceedings of the 40th International Conference on Software Engineering: Companion Proceedings (Gothenburg, Sweden) (ICSE '18). Association for Computing Machinery, New York, NY, USA, 202–203. https://doi.org/10.1145/3183440.3195096
- Niki Vazou, Anish Tondwalkar, Vikraman Choudhury, Ryan G. Scott, Ryan R. Newton, Philip Wadler, and Ranjit Jhala. 2017. Refinement Reflection: Complete Verification with SMT. Proc. ACM Program. Lang. 2, POPL, Article 53 (Dec. 2017), 31 pages. https://doi.org/10.1145/3158141
- Zhe Zhou, Robert Dickerson, Benjamin Delaware, and Suresh Jagannathan. 2021a. Data-Driven Abductive Inference of Library Specifications (Full Version). arXiv:2108.04783 [cs.PL]

- Zhe Zhou, Robert Dickerson, Benjamin Delaware, and Suresh Jagannathan. 2021b. OOPSLA2021 Artifact: Data-Driven Abductive Inference of Library Specifications. https://doi.org/10.5281/zenodo.5130646
- He Zhu, Stephen Magill, and Suresh Jagannathan. 2018. A Data-Driven CHC Solver. In Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation (Philadelphia, PA, USA) (PLDI 2018). Association for Computing Machinery, New York, NY, USA, 707–721. https://doi.org/10.1145/3192366.3192416
- He Zhu, Gustavo Petri, and Suresh Jagannathan. 2016. Automatically Learning Shape Specifications. In Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (Santa Barbara, CA, USA) (PLDI '16). Association for Computing Machinery, New York, NY, USA, 491–507. https://doi.org/10.1145/2908080.2908125