

# Analyzing Test Completeness for Dynamic Languages

Christoffer Quist  
Adamsen  
Aarhus University, Denmark  
cqqa@cs.au.dk

Gianluca Mezzetti  
Aarhus University, Denmark  
mezzetti@cs.au.dk

Anders Møller  
Aarhus University, Denmark  
amoeller@cs.au.dk

## ABSTRACT

In dynamically typed programming languages, type errors can occur at runtime. Executing the test suites that often accompany programs may provide some confidence about absence of such errors, but generally without giving any guarantees. We present a program analysis that can check whether a test suite has sufficient coverage to prove a given type-related property, which is particularly challenging for program code with overloading and value dependent types. The analysis achieves a synergy between scalable static analysis and dynamic analysis that goes beyond what can be accomplished with the static analysis alone. Additionally, the analysis provides a new coverage adequacy metric as a measure of the completeness of a test suite for a family of type-related properties.

Based on an implementation for Dart, we demonstrate how such a hybrid static/dynamic program analysis can be used for measuring the quality of a test suite with respect to showing absence of type errors and inferring sound call graph information, specifically for program code that is difficult to handle for traditional static analysis techniques.

## 1. INTRODUCTION

A well-known quote from Dijkstra is: “Program testing can be used to show the presence of bugs, but never to show their absence” [13]. This is of course true in general, but, as also observed decades ago [19], there are situations where test suites are *complete* in the sense that they allow verification of correctness properties that hold in any execution. As a trivial example, for a straight-line Java program that takes a string as input, a single test execution suffices to detect any null pointer bugs that may exist. More generally, test suites may be complete with respect to local properties at specific program points. In this paper, we explore how such a notion of test completeness can be exploited for soundly type checking challenging programming patterns in dynamically typed languages.

One of the advantages of programming languages with dynamic or optional typing, such as, Dart [14], TypeScript [33], Typed Racket [43], and Reticulated Python [45], is flexibility: they allow dynamic programming patterns, such as, value-based overloading [44, 1] and other forms of value dependent types, that do not fit into traditional static type systems. The cost of this flexibility is that type-related errors are not detected until runtime and that type information is not avail-

able statically for guiding auto-completion and navigation in IDEs or optimizers in compilers. Our goal is to provide program analysis techniques that can automatically infer sound and precise type information for such code. More specifically, we focus on Dart and aim to ensure absence of runtime type errors and to infer precise call graph information. Two kinds of runtime type errors are particularly interesting in Dart programs: subtype-violation errors at implicit downcasts and message-not-understood errors at field and method lookup operations. Programmers can benefit from knowing to what extent their tests guarantee absence of such runtime errors. Likewise, call graphs are useful for enabling dead code elimination and other optimizations, as well as for various kinds of IDE support [16].

*Example* Figure 1 shows some Dart code from a sample application shipped with the *box2d* library<sup>1</sup> that uses the *vector\_math* library,<sup>2</sup> which illustrates a style of programming that is common with dynamically typed languages and that is challenging for static type analysis. The `cross` function is overloaded in the sense that its behavior and return type depend on the runtime types of its parameters: the branch in line 4 returns type `vec3` or the runtime type of the parameter `out`, line 7 returns type `double`, and line 16 returns type `vec2` or the type of `out`. Moreover, if the types of the parameters at a call do not match any of the cases, then a failure occurs in line 20 if assertions are enabled and otherwise `null` is returned, which will likely trigger a runtime error later in the program. The `cross` function is called, for example, in lines 32 and 33. In Dart’s checked mode execution, the assignments to `dp2perp` and `dp1perp` will fail if the values being returned are not of type `vec3`. In production mode execution, the assignments will succeed, but the applications of `*` in line 34 will fail if `dp2perp` or `dp1perp` is of type `double` and the other argument is of type `vec2` or `vec3`, and `+` will fail if the left-hand argument and the right-hand argument do not have the same type. (Lines 26–27 show the implementation of `+` for the case where the left-hand argument is of type `vec2`.) How can the programmer be certain that none of these potential type-related errors can occur in the application? Running a few tests may give some confidence, but it is difficult to know when enough testing has been done.

Traditional static analysis is not suitable for reasoning about such code since very high precision (e.g. context sensitivity and path sensitivity) would be needed, which is difficult to achieve together with sufficient scalability. For the exam-

<sup>1</sup><https://github.com/financeCoding/dartbox2d>

<sup>2</sup>[https://github.com/google/vector\\_math.dart](https://github.com/google/vector_math.dart)

```

1 dynamic cross(dynamic x, dynamic y,
2               [dynamic out=null]) {
3   if (x is vec3 && y is vec3) {
4     return x.cross(y, out);
5   } else if (x is vec2 && y is vec2) {
6     assert(out == null);
7     return x.cross(y);
8   } else if (x is num && y is vec2) {
9     x = x.toDouble();
10    if (out == null) {
11      out = new vec2.zero();
12    }
13    var tempy = x * y.x;
14    out.x = -x * y.y;
15    out.y = tempy;
16    return out;
17  } else if (x is vec2 && y is num) {
18    ...
19  } else {
20    assert(false);
21  }
22  return null;
23 }

24 class vec2 {
25   ...
26   vec2 operator+(vec2 other) =>
27     new vec2.raw(x + other.x, y + other.y);
28   ...
29 }

30 ...
31 // solve the linear system
32 vec3 dp2perp = cross(dp2, normal);
33 vec3 dp1perp = cross(normal, dp1);
34 vec3 tangent = dp2perp * duv1.x + dp1perp * duv2.x;
35 ...

```

**Figure 1:** Motivating example, code from the *vec\_tor\_math* library (lines 1–29) and the *box2d* library (lines 31–34).

ple described above, a context-insensitive analysis may be able to show that `cross` can return values of type `vec3`, `vec2`, or `double` (assuming that `out` is always also of one of those types), but that information is not precise enough to rule out the various type-related errors. Two-phase typing [44] has been proposed to address similar challenges but is not yet practical for realistic programs. Conversely, traditional dynamic analysis is also not suitable as it does not give any soundness guarantees. One notable exception, which has inspired our work, is the constraint-based dynamic type inference technique by An et al. [2]. It can infer sound type information for Ruby programs using dynamic analysis, however, it requires full path coverage of every function in the program being analyzed, which is rarely realistic.

**This work** We propose a hybrid of dynamic and static analysis: one component is a lightweight static analysis that provides information about where another component, a dynamic analysis, has adequate coverage to provide sound conclusions for a given test suite. Our key insight is that such a combination of static and dynamic techniques can help to determine when a test suite has sufficient coverage to guarantee type-related correctness properties as in the example.

The static analysis has two parts: a dependence analysis and a type analysis (technically, a points-to analysis). It is context- and path-insensitive and thereby scales to large

programs, and it is relatively easy to implement; notably it requires simpler modeling of native functions than what would be required by a fully static analysis approach.

In summary, the contributions of this paper are:

- We define a notion of test completeness as a sufficient condition for a test suite to have enough coverage to guarantee a given type-related correctness property. Using a lightweight static analysis to approximate test completeness, we then demonstrate how a hybrid static/dynamic analysis can produce sound type information for challenging dynamic programming patterns that resist traditional static analysis.
- Based on an implementation, GOODENOUGH,<sup>3</sup> for the Dart programming language, we evaluate the ability to ensure absence of runtime type errors and to produce precise call graphs, compared to a more traditional technique that uses the static analysis information only. Across a variety of Dart programs, we find numerous examples where precision is improved. Specifically, the analysis is able to guarantee for 81% of the expressions that all types that can possibly appear at runtime are in fact observed by execution of the test suite. The experiments additionally show that the limiting factor of the precision in some cases is the test coverage and in other cases the precision of the dependence analysis, which suggests opportunities for future improvements of the technique.

## 2. THE DART PROGRAMMING LANGUAGE

In this paper we use Dart as subject for presenting and evaluating our technique. The Dart programming language was introduced by Google in 2013 and has later been standardized by Ecma [14]. The language is now widely used by Google and elsewhere, primarily as an alternative to JavaScript but also for server-side high performance computation and, more recently, embedded devices. Dart supports both object oriented and functional programming using a Java-like syntax.

From our perspective, the most interesting aspect of the language is its type system. Type annotations are optional (the default is called `dynamic`). The static type checker can warn about likely type errors in the parts of the program where type annotations have been provided, while ignoring the parts without type annotations. Warnings produced by the type checker do not preclude executing the program. Moreover, even for fully annotated programs, the type checker is deliberately unsound. These design choices give Dart the flexibility of dynamically typed programming languages, but they also mean that type errors can occur at runtime.

Dart programs run in either *production mode* where type annotations are entirely ignored, or *checked mode* where implicit type casts are performed at assignments to non-`dynamic` variables or fields, checking that the runtime type is a subtype of the annotated type. Two kinds of runtime type errors can occur: a *subtype-violation* error occurs if a type cast fails, and a *message-not-understood* error occurs if attempting to access a field or method that does not exist. We here use the terminology from Ernst et al. [15] who recently presented a formalization of a core of Dart with a focus on its type system and the causes of unsoundness.

<sup>3</sup>Named in honor of J. B. Goodenough [19].

The example code in Figure 1 shows how programmers use optional types in practice. Some variables and methods have type annotations (e.g. `dp2persp` in line 32 and `operator+` in line 26), whereas in other parts (e.g. `cross`) the programmer chose to use `dynamic` (that annotation could in fact be omitted without changing the meaning of the code).

### 3. OVERVIEW

Given a Dart program with a test suite, we wish to know for any expression  $e$  whether the test suite has sufficient coverage to explore all possible types  $e$  may evaluate to. In our setting, a test suite is simply a finite set of program inputs. To simplify the discussion we assume that program execution is deterministically determined by the input.

**Definition 1** (Test completeness). A test suite  $T$  is *complete* with respect to the type of an expression  $e$ , written  $\text{COMPLETE}_T(e)$ , if execution of  $T$  covers all possible types  $e$  may have at runtime.

The analysis we present conservatively approximates completeness (which is clearly undecidable): if the analysis reports that the completeness condition is satisfied, then this is indeed the case, but the analysis may sometimes be unable to prove completeness even though the condition semantically holds.

Our approach involves four components: 1) a dynamic execution of the test suite, 2) a static dependence analysis, 3) a static type analysis, and 4) a test completeness analysis. In this section we present an overview of the approach and explain the purpose of each component. The subsequent sections give more details about the components.

#### 3.1 Combining over- and under-approximation

Our starting point is a well-known fact about dynamic execution and static analysis: executing a test suite constitutes an under-approximation of the program behavior, whereas a static analysis (in the abstract interpretation style) can provide an over-approximation. If the two agree, for example regarding the possible types of an expression  $e$ , then the completeness condition is trivially satisfied. Obviously, no further tests can then possibly reveal new types of  $e$ .

For the program code shown in lines 36–39, a simple static analysis is able to show that the only possible value of  $x$  on line 38 is the `A` object created on line 37. That information tells us not only that line 38 cannot fail with a message-not-understood runtime error, but also that the only possible callee is the `m` method in class `A`.

For scalability reasons we wish to use only context-insensitive and path-insensitive static analysis. (An analysis is *context-insensitive* if it does not distinguish between different calling contexts inter-procedurally [39], and it is *path-insensitive* if it does not distinguish between the different paths that may lead to a given program point intra-procedurally [6].) The precision of such an analysis may suffice for large parts of a program, but programs written in dynamically typed languages often contain code that require more heavyweight program analysis techniques, such as, re-

finement typing [44] or context sensitivity [4]. The following example shows a typical case of *value-dependent types*.

```

40 class A {
41   m() { ... }
42 }
43 class B {}
44
45 f() {
46   var t = 42;
47   A x = g(t);
48   x.m();
49 }
50 g(a) {
51   var r;
52   if (a > 0) {
53     r = new A();
54   } else {
55     r = new B();
56   }
57   return r;
58 }

```

Here, `A` has a `m` method and `B` does not. The function `g` is overloaded as its argument determines whether it returns an object of type `A` or `B` (assume that some other code not shown here contains a call to `g` with argument `0`). The call in line 47 clearly always returns an `A` object, but this fact cannot be obtained by a context-insensitive static analysis alone (it would infer that the type is either `A` or `B`). Nor is it obvious by executing a test suite covering `f` that `A` is the only possibility. If the call to `g` instead returned a `B` object, then the program would fail at runtime, in checked mode with a subtype-violation error at line 47 and in production mode with a message-not-understood error at line 48.

#### 3.2 Exploiting tests and dependencies

Our key insight is that it is possible through a *combination* of lightweight static analysis and execution of a test suite to obtain completeness guarantees for the kind of code shown in the previous example.

**The dependence analysis component** One component in our system is a context-insensitive and path-insensitive *dependence analysis* that over-approximates the dependencies of each expression in the given program. Unlike traditional dependence analyses, this one considers both value and type dependencies. (This dependence analysis is described in more detail in Section 4.) For example, it infers that the *type* of `r` in line 57 depends (only) on the *value* of the parameter `a`. It also tells us that the parameter passed to `g` on line 47 has no dependencies (it is a constant). By combining these pieces of information, we see that *a single concrete execution of line 47 suffices* for learning all the possible types of the return value at that call. Thus, we run the test suite, and if it covers line 47 we learn that the only possible type is `A`—in other words, the test suite is complete with respect to the type of the return value of this particular call. Notice that the static analysis alone does not have this information; we need the concrete execution also.

**The type analysis component** The `bar` function shown in lines 59–71 is a Dart version of a Ruby example by An et al. [2]. Assume that elsewhere in the program, there are calls to `bar` with arguments `true` and `false`. A purely static analysis would require path sensitivity to be able to prove that `y` is always a number in line 67 (so that the `+` operation is guaranteed to succeed) and a string in line 69 (so that it has a `length` field).

We now show how we can obtain the necessary precision without path sensitive static analysis. The dependence analysis gives us that the type

```

59 bar(p) {
60   var y;
61   if (p) {
62     y = 3;
63   } else {
64     y = "hello";
65   }
66   if (p) {
67     y + 6;
68   } else {
69     y.length;
70   }
71 }

```

of  $y$  in line 67 and the type of  $y$  in line 69 can only depend on the value of the parameter  $p$ . As the actual value of  $p$  is unknown to the dependence analysis, we need more information to prove type safety of lines 67 and 69. For this reason, we include another component in our system: a context-insensitive and path-insensitive *type analysis* that provides an over-approximation of the possible types of all expressions. For the `bar` example, the type analysis tells us that  $p$  is of type `bool`, i.e. its value can only be `true` or `false`. Now, notice that by combining this information with the dependence information we see that if executing the test suite includes call to `bar` with both these two values then the test suite is complete with respect to the type of  $y$  in lines 67 and 69. We thereby know that runtime type errors cannot occur in those lines.

The Rubydust technique by An et al. [2] is able to infer sound types for the `bar` function if all possible paths inside the function are executed by the tests. For this particular example, Rubydust can therefore, like our technique, infer sound and precise types using only two executions of the function. However, our technique differs in several important aspects: (1) Rubydust infers sound types *if* all possible paths are executed, but it does not know whether or not that condition is satisfied (in this example, the control-flow of the function suggests that there are four paths, but only two are possible because of the branch correlation); in contrast, our technique is able to tell that the two executions suffice. (2) In this example, the fact that the two branches are correlated is quite obvious and could be inferred by a very simple static analysis, and that information could easily be incorporated into Rubydust. However, our technique is capable of reaching the conclusion about the type of  $y$  without reasoning explicitly about branch correlations. (3) If `bar` contained additional branches with code not affecting  $y$ , then the Rubydust soundness criterion would require more than two executions, whereas our technique would still only require two, to show type safety for the operations on  $y$ .

As these examples suggest, we use two lightweight static analyses: a dependence analysis and a type analysis. A central part of our approach for inferring test completeness facts is combining the information provided by these analyses with the information obtained from executing the test suite. In Section 6 we explain how this can be achieved, a key step being a mechanism for substituting type properties according to dependencies.

### 3.3 Using the inferred completeness facts

**Program correctness** For Dart programs we are particularly interested in test completeness at assignments, calls, and field lookup operations. As mentioned in Section 2, running Dart programs may encounter message-not-understood errors at property lookup operations and subtype-violation errors at implicit casts, even in programs that are accepted without warnings by the static type checker. For a programmer, knowing that his test suite is complete for such properties gives confidence about the correctness of the program.

**Test adequacy** The notion of test completeness with respect to type properties directly gives rise to a new metric for test adequacy alongside statement coverage, branch coverage, path coverage, etc. [49]:

**Definition 2** (Type coverage). For a given set of expressions  $X$  in a program and a test suite  $T$ , the *type coverage* of  $T$ , denoted  $\mathcal{C}_T(X)$ , is computed as

$$\mathcal{C}_T(X) = \frac{|\{x \in X \mid \text{COMPLETE}_T(x)\}|}{|X|}$$

As  $X$ , one may select, for example, the set of expressions that are subjected to implicit casts in a unit of code of interest, e.g. a method, class, or library. A type coverage of 100% then means that the test suite is adequate for revealing all implicit cast runtime errors that are possible in that unit of code.

Traditional coverage metrics are not suitable for giving such guarantees. For example, full statement coverage or full branch coverage is not always sufficient, and full path coverage is impossible to achieve whenever loops are involved [49]. Other techniques that focus on branches and paths, such as, basis path testing using cyclomatic complexity [32], have similar limitations.

By selecting  $X$  as the set of receiver expressions of method calls (e.g.  $x$  in line 48) and the function expressions at function calls (e.g.  $g$  in line 47) in a unit of code, then 100% type coverage implies that the call graph has been fully exercised in that code. A programmer may use such information to guide further testing. For example, if a test suite has full statement coverage for two classes  $C$  and  $D$ , maybe only 30% of the part of the call graph that involves  $C$  has been covered, while the number is 95% for  $D$ , in which case the programmer should perhaps prioritize adding tests targeting  $C$  rather than  $D$ . We leave it to future work to evaluate the practical usefulness of reporting such type coverage numbers to developers and how the type coverage metric correlates with other metrics and with errors; for the rest of this paper we focus on the use of test completeness in checking type safety and inferring call graphs.

**Type filtering** A well-known trick in points-to analysis and dataflow analysis is to use type tests (e.g., casts) that appear in the program as *type filters* [40, 4, 20]. The type-related completeness facts inferred by our analysis can of course be exploited for such type filtering: if we have inferred that, for example,  $y$  is definitely a string in line 69 then a subsequent analysis may use that fact to filter away spurious values that are not strings. As part of our evaluation in Section 7, we demonstrate that performing type filtering based on the completeness facts inferred by our analysis can have a substantial impact on the precision of type analysis and call graph analysis.

The lightweight type analysis we use as one of our components can directly use this mechanism itself, and increasing precision of this component may lead to more completeness facts being inferred. It therefore makes sense to repeat the entire analysis, boosting precision using type filtering based on type information inferred in the previous iteration. For now, this remains an intriguing observation, however; although it is possible to construct examples where the completeness analysis becomes more precise by such a feedback mechanism, we have not yet encountered much need in practice.

**Optimizations** The inferred completeness facts may also be used for optimizations, for example, removal of runtime type checks, dead-code elimination, replacement of indirect calls by direct calls, inlining, and type specialization [9, 30, 20, 23, 26]; we leave such opportunities for future work.

## 4. DEPENDENCE ANALYSIS

As motivated in the preceding section, a key component of our technique is a static dependence analysis. Unlike traditional dependence analyses as used in, for example, optimization [17], information flow analysis [36] and program slicing [42] we are interested in both value and type dependencies. We therefore introduce a general *dependence relation* denoted  $\triangleleft$ , which is a binary relation over abstractions of runtime states at different program points (inspired by the notion of abstract dependence by Mastroeni and Zanardini [31]). For example, we have seen that the type of  $y$  at the program point after line 67 depends on the value of  $p$  at line 59, which we write as  $\text{TYPE}_y[67] \triangleleft \text{VAL}_p[59]$ . The dependence relation is computed using a whole-program analysis, however, all dependencies are intra-procedural in the sense that they relate variables and parameters within the same function.

More generally, the dependence information we need can be expressed via different abstractions of runtime states:

**Definition 3** (Type, value, and top abstraction). The *type* abstraction  $\text{TYPE}_x$  for a variable  $x$  maps a state  $\sigma$  to the runtime type of  $x$  in  $\sigma$ . The *value* abstraction  $\text{VAL}_x$  instead maps  $\sigma$  to the value of  $x$  in  $\sigma$ . The *top* abstraction  $\top$  is the identity function on program states (we use this abstraction later to express dependencies that are unknown due to analysis approximations).

The dependence relation can now be expressed as a relation of the form  $\pi[c] \triangleleft \pi'[c']$  such that  $\pi, \pi' \in \Pi$  and  $c \in \mathbb{C}$  where  $\Pi$  is a family of state abstractions and  $\mathbb{C}$  is a set of program points. (The program point associated with a given line is the one immediately after that line.) We want the dependence relation being computed to conservatively approximate all dependencies in the following sense.

**Property 1** (Valid dependence relation). Given any two executions of the program that both reach a program point  $c$  inside a function  $f$  with entry program point  $c_0$ , let  $\sigma$  and  $\sigma'$  be the states at  $c$  for the two executions, respectively, and similarly let  $\sigma_0$  and  $\sigma'_0$  be the states at  $c_0$  when  $f$  was entered. If there exists some state abstraction  $\pi \in \Pi$  where  $\pi(\sigma) \neq \pi(\sigma')$  then there must exist some  $\pi' \in \Pi$  where  $\pi[c] \triangleleft \pi'[c_0]$  and  $\pi'(\sigma_0) \neq \pi'(\sigma'_0)$ .

Intuitively, if two executions disagree on  $\pi$  at program point  $c$ , then they must also disagree on  $\pi'$  at the entry of the function containing  $c$ , where  $\pi'$  is some state abstraction that  $\pi$  depends on. For example, if two executions disagree on the type of  $r$  in line 50 then they must also disagree on the value of  $a$  in line 57, so for our choice of state abstractions, the dependence relation must include the fact that the type of  $r$  depends on the value of  $a$ .

**Program representation** To concisely explain how the dependence analysis works we represent a program as a control flow graph  $(\mathbb{C}, \hookrightarrow)$  where  $\mathbb{C}$  is now a set of nodes corresponding to primitive instructions of the different kinds as shown in Figure 2 and  $\hookrightarrow \subseteq \mathbb{C} \times \mathbb{C}$  is the intra-procedural control flow relation between nodes. We let  $\hookrightarrow^+$  be the transitive closure of  $\hookrightarrow$ .

We assume nested expressions have been flattened, so all primitive instructions operate directly on local variables (which include function parameters) and nodes are uniquely identified by the line number. Every call is represented by

<code>entry[f, p<sub>1</sub>, ..., p<sub>n</sub>]</code>	the entry of function $f$ with parameters $p_1, \dots, p_n$
<code>call[f, x<sub>1</sub>, ..., x<sub>n</sub>]</code>	calls function $f$ with arguments $x_1, \dots, x_n$
<code>aftercall[x]</code>	writes the returned value to $x$
<code>const[x, p]</code>	writes constant value $p$ to $x$
<code>new[x, D]</code>	writes new $D$ instance to $x$
<code>assign[x, y]</code>	writes value of $y$ to $x$
<code>load[x, y, f]</code>	writes value of $y.f$ to $x$
<code>store[x, f, y]</code>	writes value of $y$ to $x.f$
<code>is[x, D, y]</code>	writes <i>true</i> to $x$ if value of $y$ is a subtype of $D$ , <i>false</i> otherwise
<code>binop[x, y, <math>\oplus</math>, z]</code>	writes the result of the operation $y \oplus z$ to $x$
<code>if[x]</code>	a branch on condition $x$
<code>phi[x, x<sub>1</sub>, x<sub>2</sub>, B]</code>	writes $\phi(x_1, x_2)$ to $x$ after branch
<code>return[x]</code>	returns the value stored in $x$

Figure 2: Primitive instructions.

two nodes: a `call` node and an associated `aftercall` node. The `call` node is labeled by the name of the function being called (we explain in Section 5 how indirect function/method calls are resolved), and for simplicity we assume every call has a single callee. The program representation is in SSA form [12], and every  $\phi$ -node `phi[x, x1, x2, B]` is associated with the set  $B$  of all `if`-nodes that dominate  $x_1$  or  $x_2$ . Due to the use of SSA, each variable name uniquely corresponds to a node where the variable is assigned its value (where parameters are considered to be assigned at the function entry node), and we sometimes use the variable name to denote the program point after that node. Furthermore, we let  $\text{entry}_c$  denote the entry program point of the function containing node  $c$ . The meaning of the other kinds of nodes is straightforward.

**The dependence analysis** We compute the dependence relation  $\triangleleft$  as the smallest fixpoint of the rules shown in Figure 3 (in addition,  $\triangleleft$  is reflexive and transitive).

The `ASSIGN` rule for an assignment `assign[x, y]` establishes that the type and value of  $x$  depend on the type and value of  $y$ , respectively. More precisely, for  $\pi = \text{TYPE}_x$  we see that the type of  $x$  at the program point where  $x$  has been defined, i.e.  $\text{TYPE}_x[x]$ , depends on the type of  $y$  at the program point where  $y$  has been defined, i.e.  $\text{TYPE}_y[y]$ , and similarly for the value. The `BINOP` rule shows that the value of the result depends on the values of the two operands; however, the *type* of the result is always fixed by the operator so it does not depend on the operands. The `IS` rule for the instruction `is[x, D, y]` shows that the *value* of  $x$  depends on the *type* of  $y$ . (The *type* of  $x$  is always boolean.) A rule for `const[x, p]` can be entirely omitted, since neither the type nor the value of constants depend on anything.

To keep the analysis lightweight we choose to entirely avoid tracking dependencies on the heap. This is modeled using the  $\top$  abstraction. In the `LOAD` rule for `load[x, y, f]`, the type and value of  $x$  conservatively depend on the entire state at the entry of the function, i.e.  $\top[\text{entry}_c]$ . With this coarse model of the heap, it is safe to omit a rule for `store[x, f, y]`. The `NEW` rule captures the fact that the value being produced by `new[x, D]` is unknown to the analysis (whereas the type is trivially  $D$ , without any dependencies).

The `PHI` rule models dependencies that are due to the data *and* control flow at branches, which are represented by `if[y]` and `phi[x, x1, x2, B]` nodes. First, the type and value

$$\begin{array}{c}
\text{(ASSIGN)} \\
\frac{\text{assign}[x, y] \in \mathbb{C} \quad \pi \in \{\text{TYPE}, \text{VAL}\}}{\pi_x[x] \triangleleft \pi_y[y]} \\
\\
\text{(LOAD)} \\
\frac{c = \text{load}[x, y, f] \in \mathbb{C} \quad \pi \in \{\text{TYPE}, \text{VAL}\}}{\pi_x[x] \triangleleft \top[\text{entry}_c]} \\
\\
\text{(PHI)} \\
\frac{\text{phi}[x, x_1, x_2, B] \in \mathbb{C} \quad \pi \in \{\text{TYPE}, \text{VAL}\} \quad \text{if}[y] \in B \quad l \in \{x_1, x_2, y\}}{\pi_x[x] \triangleleft d} \quad d = \begin{cases} \pi_l[l] & \text{if } l \in \{x_1, x_2\} \\ \text{VAL}_l[l] & \text{if } l = y \end{cases} \\
\\
\text{(CALL)} \\
\frac{c = \text{call}[f, x_1, \dots, x_n] \in \mathbb{C} \quad c_e = \text{entry}[f, p_1, \dots, p_n] \in \mathbb{C} \quad c \hookrightarrow c_a \quad c_e \hookrightarrow^+ c_r \quad d = \begin{cases} \text{VAL}_{x_i}[x_i] & \text{if } \pi_y[y] \triangleleft \text{VAL}_{p_i}[c_e] \\ \text{TYPE}_{x_i}[x_i] & \text{if } \pi_y[y] \triangleleft \text{TYPE}_{p_i}[c_e] \\ \top[\text{entry}_c] & \text{if } \pi_y[y] \triangleleft \top[c_e] \end{cases}}{\pi_x[x] \triangleleft d} \\
c_a = \text{aftercall}[x] \in \mathbb{C} \quad c_r = \text{return}[y] \in \mathbb{C}
\end{array}$$

Figure 3: Rules for type and value dependence.

of  $x$  depend on the type and value, respectively, of both  $x_1$  and  $x_2$ . Second, the type and value of  $x$  also depend on the value of the branch conditions in  $B$ . It is standard to compute control flow dependence using post-dominance information [12]. A node  $c \in \mathbb{C}$  is control flow dependent on  $c' \in \mathbb{C}$  if  $c'$  determines whether  $c$  is executed, i.e. (i) there is a path from  $c$  to  $c'$  where every node in the path is post-dominated by  $c'$ , and (ii)  $c$  is not post-dominated by  $c'$ . Since we assume the control flow graph to be on SSA form, control flow dependence for variables is already explicit in the  $\phi$ -nodes.

The CALL rule exploits the dependencies computed between the return value and the parameters of the callee, to obtain dependencies between the resulting value and the arguments at the call site. In this rule,  $c$  is a call node and  $c_a$  is the associated `aftercall` node at the call site, while  $c_e$  is the callee `entry` node and  $c_r$  is the `return` node. The dependencies for the result  $x$  are found simply by substituting the dependencies of the return value  $y$  according to the parameter binding.

Finally, the DEF rule has the consequence that the type or value of a variable  $x$  at any program point  $c$  can always be determined from its type or value at the definition of  $x$ , and the top abstraction at  $c$  depends on the top abstraction at the entry of the function. This rule is strictly not necessary, but it simplifies the use of  $\triangleleft$  in Section 6.

Notice that even though type abstractions are conservative approximations of value abstractions, all combinations of dependencies between types and values can occur; in particular, it is possible that the *value* of a variable  $x$  depends on the *type* of another variable  $y$  (due to the IS rule), and that the *type* of  $x$  depends on the *value* of  $y$  (due to the PHI rule).

**Proposition 1** (Soundness of computed dependencies). For any program, the relation  $\triangleleft$  defined by the least fixpoint of the dependence rules satisfies Property 1.

*Example* For the example in lines 40–58, the dependence analysis infers, in particular, that  $\text{TYPE}_r[57] \triangleleft \text{VAL}_a[50]$  but *not*  $\text{TYPE}_r[57] \triangleleft \top[50]$ , and that there is no  $\pi$  such that

$\text{TYPE}_t[47] \triangleleft \pi[45]$ . Therefore the CALL rule gives that there is also no  $\pi'$  such that  $\text{TYPE}_x[x] \triangleleft \pi'[45]$ , i.e., the type of  $x$  does not depend on the input to  $f$ .

## 5. TYPE ANALYSIS

The static type analysis component serves two purposes: it resolves indirect calls for the dependence analysis in Section 4, and it computes an over-approximation of the possible types of every expression, which we use in the test completeness analysis in Section 6. The type analysis is simply a context-insensitive and path-insensitive Andersen-style points-to analysis [3] that simultaneously tracks functions/methods and types. This is a well-known analysis technique (see e.g. the survey by Sridharan et al. [40]), so due to the limited space we omit a detailed description of how this component works. We choose a lightweight analysis for the reasons given in Sections 1 and 3.

In Dart, all values are objects, including primitive values and functions. The analysis abstracts primitive values by their type (e.g. `bool` or `String`) and treats each function and method as having its own type. As output, for every program point  $c$  the analysis returns an abstract state  $\hat{\sigma}_c$  that over-approximates all concrete states that may appear at  $c$  when the program is executed. We define  $\hat{\sigma}_c(e)$  to be the set of types the expression  $e$  may evaluate to according to  $\hat{\sigma}_c$ .

*Example* Consider a call to the `cross` function from Figure 1

```
72 x = cross(y, z);
```

in a context where the type analysis finds that the type of  $y$  is either `vec3` or `vec2` and the type of  $z$  is `vec3`. That is,  $\hat{\sigma}_c(y) = \{\text{vec3}, \text{vec2}\}$  and  $\hat{\sigma}_c(z) = \{\text{vec3}\}$  where  $c$  is the program point at the call. (This example uses a direct call, so the connection between the call and the callee is trivial.) Assuming there are other calls to `cross` elsewhere in the program, the context-insensitive type analysis has only imprecise information about the possible return types for this function. However, the dependence analysis has inferred that the return type only depends on the types of the parameters.

$$\begin{array}{c}
\text{(BASE)} \\
\frac{\forall (a_1, \dots, a_n) \in \pi_1(\hat{\sigma}_c) \times \dots \times \pi_n(\hat{\sigma}_c) : \\
\exists t \in T, \sigma \in \llbracket t \rrbracket_c : \forall i = 1, \dots, n : \pi_i(\sigma) = a_i}{T \vdash \pi_1, \dots, \pi_n[c]}
\end{array}
\quad
\begin{array}{c}
\text{(INDUCTIVE-TOENTRY)} \\
\frac{\{\pi'_1, \dots, \pi'_m\} = \{\pi' \mid \pi_i[c] \triangleleft \pi'[entry_c]\} \quad T \vdash \pi'_1, \dots, \pi'_m[entry_c]}{T \vdash \pi_1, \dots, \pi_n[c]}
\end{array}$$

$$\begin{array}{c}
\text{(INDUCTIVE-TOCALL)} \\
\frac{
\begin{array}{c}
c = \text{call}[f, \mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{C} \quad c_e = \text{entry}[f, \mathbf{p}_1, \dots, \mathbf{p}_n] \in \mathbb{C} \quad c \hookrightarrow c_a \quad \{\pi'_1, \dots, \pi'_m\} = \\
c_a = \text{aftercall}[\mathbf{x}] \in \mathbb{C} \quad c_r = \text{return}[\mathbf{y}] \in \mathbb{C} \quad c_e \hookrightarrow^+ c_r \quad \{ \text{VAL}_{\mathbf{x}_i} \mid \pi_{\mathbf{y}}[\mathbf{y}] \triangleleft \text{VAL}_{\mathbf{p}_i}[c_e] \} \cup \\
\pi \in \{ \text{TYPE}, \text{VAL} \} \quad T \vdash \pi'_1, \dots, \pi'_m[c] \quad \{ \text{TYPE}_{\mathbf{x}_i} \mid \pi_{\mathbf{y}}[\mathbf{y}] \triangleleft \text{TYPE}_{\mathbf{p}_i}[c_e] \} \cup \\
\{ \top \mid \pi_{\mathbf{y}}[\mathbf{y}] \triangleleft \top[c_e] \}
\end{array}
}{T \vdash \pi_{\mathbf{x}}[c_a]}
\end{array}$$

Figure 4: Rules for test completeness.

This allows the test completeness analysis, presented in the following section, to conclude that two executions suffice to cover all possible types of  $\mathbf{x}$ : one where  $\mathbf{y}$  has type `vec3` and one where it has type `vec2`. This example demonstrates the power of combining dependence analysis and type analysis.

Now consider a slight modification of the example, where the type of both  $\mathbf{y}$  and  $\mathbf{z}$  is either `vec3` or `vec2`. Our technique then requires four executions of the call to ensure that all relevant combinations are covered. This suggests that *relational* dependence and type information may be valuable: since we only obtain test completeness guarantees at a call when all combinations of dependencies have been exercised, it may be beneficial to know that only some combinations are relevant. This is an opportunity to be explored in future work.

## 6. TEST COMPLETENESS

In this section we present a proof system for test completeness, i.e., for proving  $\text{COMPLETE}_T(\mathbf{x})$  for a given variable  $\mathbf{x}$  in some program with a test suite  $T$  (see Definition 1). (We use the program representation from Section 4, so we assume without loss of generality that the expression of interest is simply a variable  $\mathbf{x}$ .)

The proof rules, shown in Figure 4, rely on three ingredients: the execution of  $T$ , the dependence relation  $\triangleleft$  from Section 4, and for each program point  $c$  the abstract state  $\hat{\sigma}_c$  produced by the type analysis from Section 5. Each test input  $t \in T$  gives rise to a program execution, which can be seen as a sequence of concrete states at different program points. We write  $\llbracket t \rrbracket_c$  for the set of states encountered at program point  $c$  when  $t$  is executed.

A simple judgment  $T \vdash \pi[c]$  intuitively means that the test suite  $T$  is complete with respect to the state abstraction  $\pi$  (see Definition 3) at program point  $c$ . In other words, the judgment holds if for every abstract value in  $\pi(\sigma_c)$  where  $\sigma_c$  is a runtime state at  $c$  in *some* execution of the program, there exists a test in  $T$  that also encounters that abstract value at  $c$ . In particular, we are interested in the type abstraction  $\text{TYPE}_{\mathbf{x}}$  and the program point where  $\mathbf{x}$  is defined. We therefore aim for the following connection between proof judgments and test completeness for type properties.

**Proposition 2** (Soundness of test completeness analysis).

$$T \vdash \text{TYPE}_{\mathbf{x}}[\mathbf{x}] \text{ implies } \text{COMPLETE}_T(\mathbf{x})$$

To show  $\text{COMPLETE}_T(\mathbf{x})$  for some variable  $\mathbf{x}$ , we thus attempt to derive  $T \vdash \text{TYPE}_{\mathbf{x}}[\mathbf{x}]$ .

More generally, judgments may involve *multiple* state abstractions:  $T \vdash \pi_1, \dots, \pi_n[c]$ . The intuitive interpretation of such a judgment is that  $T$  is complete with respect to the product of the state abstractions  $\pi_1, \dots, \pi_n$  at program point  $c$ .

We now briefly describe the rules from Figure 4.

The BASE rule corresponds to the observation in Section 3.1 that completeness sometimes can be shown using the information from the type analysis. To understand the rule in its full generality, consider first this special case:

$$\frac{\forall a \in \hat{\sigma}_c(\mathbf{x}) : \exists t \in T, \sigma \in \llbracket t \rrbracket_c : \text{TYPE}_{\mathbf{x}}(\sigma) = a}{T \vdash \text{TYPE}_{\mathbf{x}}[c]}$$

This rule states that a test  $T$  is complete for the type of  $\mathbf{x}$  if for all the types  $a \in \hat{\sigma}_c(\mathbf{x})$  that can be observed according to the type analysis (Section 5), there exists an execution that reaches program point  $c$  with a concrete state  $\sigma$  where  $\mathbf{x}$  has type  $a$ .

A first step from this special case toward the general BASE rule is to generalize it from using the type abstraction  $\text{TYPE}$  to use any state abstraction  $\pi \in \Pi$ . For this, we introduce notation for lifting a state abstraction  $\pi$  to operate on abstract states produced by the type analysis: we define  $\pi(\hat{\sigma}_c) = \{\pi(\sigma_c) \mid \hat{\sigma}_c \text{ is the type abstraction of } \sigma_c\}$ . Note that  $\pi(\hat{\sigma}_c)$  can be infinite, for example if  $\pi = \top$ . This corresponds to a completeness condition that requires an infinite set of executions, so in that situation we can simply give up. Another interesting case is when  $\pi = \text{VAL}_{\mathbf{x}}$  for some variable  $\mathbf{x}$ . This occurs when our analysis is required to prove  $T \vdash \text{VAL}_{\mathbf{p}}[c]$  in the example lines 59–71 in Section 3 (where  $c$  is the program point at the entry of `bar`). Because  $\mathbf{p}$  is a boolean according to the type analysis,  $\text{VAL}_{\mathbf{p}}(\hat{\sigma}_c)$  contains only the two values `true` and `false`, so two executions suffice to prove  $T \vdash \text{VAL}_{\mathbf{p}}[c]$ .

The last step to the general BASE rule is to account for judgments with multiple state abstractions  $\pi_1, \dots, \pi_n$ . For this case, we simply require observations of all combinations of abstract values in  $\pi_1(\hat{\sigma}_c) \times \dots \times \pi_n(\hat{\sigma}_c)$ .

The INDUCTIVE-TOENTRY rule simply uses the dependence relation  $\triangleleft$  to substitute a completeness property at a program point  $c$  with a completeness property at the function entry  $entry_c$ .

The INDUCTIVE-TOCALL rule mimics the CALL rule from the dependence analysis (Figure 3), substituting completeness properties at calls according to the dependencies of

the callee. This corresponds to the reasoning used for the example in Section 3.2 for the call in line 47.

Completeness proofs can be obtained by proving ground facts using the BASE rule, and then deriving all the others using the inductive rules whenever their premises are satisfied. It is easy to see that this procedure terminates, as the induction proceeds in program order.

**Using test completeness for type filtering** As suggested in Section 3.3, the completeness facts being inferred can be plugged in as type filters in a type analysis. For example, we can run the type analysis described in Section 5 a second time, now using type filtering.

Let  $X$  be the set of observed runtime types for a variable  $x$  where  $\text{COMPLETE}_T(x)$ . By Proposition 2 this is the set of all the possible types  $x$  may have in any execution. During the type analysis, for every abstract state, we can filter out those types of  $x$  that are not in  $X$ . Removing those spurious types may improve precision throughout the program.

## 7. EXPERIMENTAL EVALUATION

Our experimental evaluation addresses the following four research questions.

- Q1** To what extent is the technique capable of showing *test completeness* for realistic Dart programs and test suites? More specifically, what is the *type coverage* being computed for such programs and test suites?
- Q2** Does the hybrid static/dynamic approach result in better precision for soundly checking absence of runtime type errors (message-not-understood and subtype-violation errors) and producing precise call graph information (outdegree of calls, reachable functions), compared to an analysis that does not exploit the test suites?
- Q3** How important is the dependence analysis, which is a central component in our technique, for the precision of the computed type coverage?
- Q4** In situations where the analysis cannot show test completeness, is the cause likely to be insufficient tests or lack of precision in the static analysis components?

**Implementation and benchmarks** Our implementation, GOODENOUGH, consists of the four components listed in Section 3, including the type filtering technique that uses completeness facts to improve the type analysis. For logging the runtime types of expressions we use a small runtime instrumentation framework inspired by Jalangi [38]. We keep the runtime overhead low by only observing types at variable reads and calls in the application code, and we do not instrument libraries.

The evaluation is based on 27 benchmarks, ranging from small command-line and web applications to sample applications that demonstrate the use of libraries. The benchmarks are listed in the first column of Table 1. The second column shows the size of each benchmark in number of lines of code (excluding unreachable library code). We use the existing test suites available for the command-line applications; for the web applications we obtain test suites by manually exercising the applications via a browser for one minute each. To obtain some degree of confidence that the static analysis components are sound, we check that all runtime observations agree with the static analysis results.

Table 1: Results.

Program	LOC	$C_T(X)$						PI
			>0%	1%	5%	30%	100%	
a_star	1 025	83.6%	-	-	0.31%	7.53%	77.27%	
archive	1 682	87.8%	11.11%	-	-	0.17%	3.45%	
bad.aliens	1 393	77.6%	4.26%	18.18%	-	0.16%	22.22%	
box2d	7 675	70.9%	4.13%	10.10%	3.87%	3.32%	22.52%	
csslib	5 911	82.4%	6.45%	1.54%	2.02%	2.69%	7.85%	
Dark	12 853	78.4%	25.88%	1.41%	-	0.55%	10.63%	
dart_regexp_tester	2 208	84.8%	-	-	-	-	-	
Dartrix	1 323	80.2%	-	-	0.50%	1.49%	66.67%	
dartrocket	9 641	74.7%	-	-	-	-	-	
frappe	4 723	77.8%	-	-	-	-	-	
graphviz	2 648	80.7%	-	6.25%	-	-	-	
markdown	1 350	89.9%	-	-	1.98%	8.57%	63.72%	
qr	4 906	87.2%	-	-	-	0.88%	16.00%	
petitparser	2 981	69.0%	-	-	0.31%	3.68%	33.62%	
pop-pop-win	12 558	71.8%	-	-	0.08%	0.86%	17.01%	
rgb_cube	831	86.1%	-	-	-	-	-	
solar	439	85.1%	-	-	-	-	-	
solar3d	7 181	86.5%	71.74%	6.90%	2.88%	4.65%	17.65%	
spaceinvaders	1 026	75.9%	-	-	-	-	-	
speedometer	267	91.0%	-	-	-	-	-	
spirodrow	923	86.5%	4.76%	5.56%	-	4.54%	-	
stats	1 162	84.2%	35.71%	9.09%	-	-	-	
sunflower	359	86.4%	75.00%	33.33%	-	-	-	
three	6 590	83.3%	42.05%	5.29%	20.93%	7.91%	8.00%	
todomvc_vanilladart	2 786	78.4%	2.08%	5.00%	0.15%	0.93%	7.58%	
tutorial.kanban	578	70.3%	-	-	-	-	-	
vote	6 377	76.7%	17.29%	7.69%	15.16%	18.08%	31.10%	

The implementation, benchmarks, test suites, and details of the experimental results are available online.<sup>4</sup>

### Q1: Type coverage

To answer the first question we run GOODENOUGH on each benchmark with its test suite  $T$  and measure the type coverage  $C_T(X)$ , where  $X$  is chosen as the set of all the expressions in the program. (Due to Proposition 2, the type coverage numbers we report are safe in the sense that they underestimate the semantic definition of type coverage given by Definition 2.) The results are shown in the third column of Table 1.

Type coverage is generally high, with an average of 81%. In other words, the analysis is able to guarantee for 81% of the expressions that all types that can possibly appear at runtime are in fact observed by execution of the test suite.

### Q2: Type errors and call graphs

To investigate whether our hybrid static/dynamic approach to test completeness can be useful for showing absence of type errors and inferring sound call graph information, we compare GOODENOUGH with a variant that does not exploit the runtime observations from the test suite. We select five metrics for measuring the precision:

- MNU** (*message-not-understood*): reported number of warning for possibly failing access to a field or a method;
- UIC** (*unsafe implicit cast*): reported number of warning for subtype violation that may arise from an implicit downcast (only meaningful for checked mode execution);
- CGS** (*call-graph size*): number of edges in the call graph;
- RF** (*reached functions*): number of possibly reached functions;
- PI** (*polymorphic invocation*): number of call sites with multiple potential callees.

<sup>4</sup>URL omitted for double-blind reviewing.

The columns MNU, UIC, RF, CGS, and PI in Table 1 show the percentage change of the resulting numbers for the hybrid analysis compared to the fully-static analysis.

We see that the hybrid analysis is able to greatly improve the precision for some benchmarks, while it gives marginal or no improvements in others. This is not surprising, considering the different nature of the benchmarks and their programming style. For a few benchmarks, the naive fully-static analysis already obtains optimal precision for MNU or UIC, leaving no room for improvements with the hybrid technique. Interestingly, 19 out of 27 benchmarks improve for at least one metric, meaning that the hybrid approach shows advantages across a variety of programs. This confirms that type filters based on inferred completeness facts can have a substantial impact on the precision of the type analysis, as discussed in Section 3.

An important design decision has been that the static analysis components are context- and path-insensitive to keep the system lightweight. By manually studying some of the examples where our hybrid approach obtains high precision, we see that a fully-static alternative would require a high degree of context sensitivity to reach the same conclusions, and it is well known that such analyses do not scale well.

For example, a context-insensitive static analysis is insufficient to reason precisely about the return type of functions similar to `cross` from the `vector_math` library discussed in Section 1. In contrast, `GOODENOUGH` finds 19 calls to `cross` where the test suite is complete, which enables the static type analysis to filter the inferred types for the return of the calls and thereby avoid several spurious warnings.

Another example is found the `Dark` benchmark:

```

73 var canvas;
74 GL.RenderingContext gl;
75 void startup() {
76   canvas = querySelector("#game");
77   gl = canvas.getContext("webgl");
78   if (gl==null)
79     gl = canvas.getContext("experimental-webgl");
80   if (gl==null) {
81     crash("No webgl", "Go to [...]");
82     return;
83   }
84   ...
85 }
86 void start(Level _level) {
87   ...
88   GL.Texture colorLookupTexture = gl.createTexture();
89   ...
90 }
```

According to the type analysis, `canvas` can be *any* subtype of the `HTML Element` class, and the calls to `getContext` return objects of type `CanvasRenderingContext2D` or `RenderingContext`. With this information, one would conclude that in checked mode execution subtype-violation errors may occur at the two assignments to `gl`, and in production mode the `createTexture` invocation may result in a message-not-understood error since `createTexture` is not declared by `CanvasRenderingContext2D`. Although `canvas` can be an arbitrary `HTML` element, the call is monomorphic. The results of the calls to `getContext` only depend on the value of the arguments, which are constant. Executing the two branches dynamically is enough to prove type completeness for both calls, which reveals that both calls return objects of type `RenderingContext`, not `CanvasRenderingContext2D`. Type filtering uses this fact to remove spurious dataflow. This example also shows the im-

portance of test coverage: a single execution does *not* suffice to cover all types at both calls.

While a lightweight type analysis (e.g. the one from Section 5) cannot reach the conclusion that only `RenderingContext` objects may be returned by `getContext`, a more precise analysis could. In this case, 2-CFA context sensitivity would be needed if we did not exploit the dynamic executions, since the argument to `getContext` is passed to the native function `getContext_Callback_2_`. In other cases, we see that 5-CFA would be necessary to reach the same precision as our hybrid approach.

To correctly handle the example above with a fully static analysis, it would additionally be necessary to precisely model the implementation of `getContext_Callback_2_`: it returns a `RenderingContext` object when the argument is `"webgl"` or `"experimental-webgl"`, and a `CanvasRenderingContext2D` object when the argument is `"2d"`. The Dart SDK contains more than 10 000 external functions, many of which are highly overloaded in such ways. As an example, the method `getParameter` of `RenderingContext` from the WebGL API returns 15 (unrelated) types depending on which of 87 constants the argument matches.<sup>5</sup> Another example is the library function `Element.tag`, which calls a native function that has more than 200 cases involving different types, which a fully static analysis would need detailed models for. Thus, our hybrid approach avoids a massive task of modeling such details, as it only needs the dependence information.

We also observe that test completeness is useful for reasoning about the types of numeric expressions without the need for context sensitivity, constant folding and constant propagation. Indeed, the type of many numerical expressions only depends on the types of the operands. The following simple example from `Dark` is one of many where cast failures may occur in checked mode execution according to a naive type analysis:

```

91 void move(double iX, double iY, double passedTime) {
92   ...
93   double frictionXZ = pow(0.0005, passedTime);
94   double frictionY = pow(0.2, passedTime);
```

The declared return type of `pow` is `num`, which has subtypes `int` and `double`. Covering the two calls to `bar` with a single execution is enough to ensure completeness.

### Q3: Dependence analysis

The dependence analysis finds that 73 % of all function return expressions have no type dependencies on the function entry state. This means that a single execution of those expressions is enough to cover all their possible types. At other 16 % of the return expressions, the return type depends on the types or values of one or more parameters but not on the heap, and for the remaining 11 % the dependence analysis finds dependencies on the heap.

To investigate the importance of the dependence analysis, we have repeated the experiments from Q1 and Q2 using a weaker dependence analysis that does not exploit the call graph produced by the type analysis. Instead, the type and value dependencies of the return value at the call site of all indirect calls are conservatively set to `top`. This single restriction causes a reduction in the number of completeness facts being proven and a significant degradation in the precision

<sup>5</sup><https://developer.mozilla.org/en-US/docs/Web/API/WebGLRenderingContext/getParameter>

metrics: type coverage drops from 81% to 77%, and the precision of type error detection and call graph construction drops to almost the same as the naive fully-static analysis (leaving only a 0.7% improvement of the UIC metric for the Dark benchmark). These results demonstrate that the dependence analysis is indeed a critical component.

#### Q4: Reasons for inability to show test completeness

To answer Q4, we have investigated the behavior of our analysis in selected places where GOODENOUGH is not able to prove test completeness. It is naturally difficult to quantify the reasons, yet some patterns are clear. Our preliminary investigation has identified two main reasons: missing coverage by the test suites, and too coarse heap modeling in the dependence analysis. On the other hand, imprecision in the type analysis does not appear to be significant.

As we have already seen in the *Dark* example in the answer to Q2, coverage is indeed important to prove test completeness. In that case two executions on specific browser versions were needed. In many other cases, we see that simply improving the statement coverage of the test suite would likely improve the type coverage significantly.

For this first design of our technique, we have chosen a simplistic approach to dependence analysis, in particular, by using the top abstraction at all field load operations. Consider the following example, which uses the `parse` function from a Dart HTML5 parser<sup>6</sup>.

```

95 f() {
96   DivElement div = parse("<div></div>").firstChild;
97 }
98 Document parse(input, ...) {
99   var p = new HtmlParser(input, ...);
100   return p.parse();
101 }

```

The `parse` function always returns a tree of HTML elements whose structure is determined solely by the input to `parse`. A dependence analysis that is able to track dependencies on the heap could in principle determine that the lookup of `firstChild` always has the same type, or equivalently, that the expression has no type dependencies. We see similar cases in many of the libraries being used in the benchmarks. This observation suggests that extending the dependence analysis to also track dependencies involving the heap may be a promising direction for future work.

## 8. RELATED WORK

**Types in dynamic languages** Many techniques have been proposed to statically reconstruct precise type information for dynamic languages, for example, for optimization purposes [26, 23]. The type inference approach by An et al. [2] is discussed in Section 3.2. Dependent types and flow-based type checking can deal with common dynamic patterns, such as, value-based overloading [10, 29, 35, 44, 41, 21]. These techniques require programmers to provide detailed type annotations. Advanced static analysis has been used to precisely infer and check types in JavaScript [4], however, this has not yet been proven to scale to realistic programs.

**Test adequacy and coverage metrics** Numerous criteria for deciding whether a test suite is adequate and metrics for measuring coverage have been proposed (the key concepts

being described by Zhu et al. [49]), so due to the limited space we can only mention the high-level relations to our work. The focus in the literature is typically on using coverage metrics to guide the effort of testing programs. Common to most of those techniques, including the seminal work by Goodenough and Gerhart [19], is that they do not come with tool support for checking whether a test suite has adequate coverage to guarantee properties like absence of runtime type errors. The general idea of “test completeness” has a long history [19, 25, 24, 8], but until now without the connection to types in dynamic languages.

**Hybrid analysis** Many hybrids of static and dynamic analysis have been developed, typically aiming to combine the strengths of both parts [48, 11, 7, 22, 27], or to use dynamic executions to guide static analysis [46, 37, 34]. Notably, the method by Yorsh et al. [48] uses automated theorem proving to check that a generalized set of concrete states obtained by dynamic execution covers all possible executions. This involves a notion of abstraction-based adequacy, which is reminiscent of our notion of type coverage. Predicate-complete testing is another related idea [5]. Most hybrid analysis create tests on-the-fly and do not exploit preexisting test suites, but for the dynamic programming patterns that are our primary target, automatically creating useful tests is by itself a considerable challenge. One approach that does use test suites and is designed for a dynamic programming language is blended analysis [46], however, it is unsound in contrast to our technique.

**Dependence analysis** Many variations of the concept of dependence described in Section 4 appear in the literature. From a theoretical point of view, our definitions fit nicely into the framework proposed by Mastroeni and Zanardini [31]. In the field of information flow [18], dependence plays an important role in reasoning about non-interference. In program slicing [47] and compiler optimization (e.g. [17]), program dependence graphs [28] model the dependencies between values used in different statements. The novelty of the dependence analysis in Section 4 is to capture dependencies not only between values but also between types.

## 9. CONCLUSION

We have presented the notions of test completeness and type coverage together with a hybrid program analysis for reasoning about adequacy of test suites and proving type-related properties. Moreover, we have demonstrated using our implementation GOODENOUGH how this analysis technique is suitable for showing absence of type errors and inferring sound call graph information, in Dart code that is challenging for fully-static analysis.

Many interesting opportunities for future work exist. We plan to explore how better heap modeling in the dependence analysis and relational dependence and type analysis can improve precision further. In another direction, we intend to perform an empirical study of how the type coverage metric correlates with other metrics and with programming errors, and to use type coverage to guide automated test generation. It may also be interesting to use the type analysis results for program optimizations and to apply our approach to other dynamic languages.

<sup>6</sup><https://github.com/dart-lang/html>

## References

- [1] B. Åkerblom and T. Wrigstad. Measuring polymorphism in Python programs. In *Proc. 11th Symposium on Dynamic Languages (DLS)*, 2015.
- [2] J. D. An, A. Chaudhuri, J. S. Foster, and M. Hicks. Dynamic inference of static types for Ruby. In *Proc. 38th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, 2011.
- [3] L. O. Andersen. *Program analysis and specialization for the C programming language*. PhD thesis, University of Copenhagen, 1994.
- [4] E. Andreassen and A. Møller. Determinacy in static analysis for jQuery. In *Proc. ACM International Conference on Object Oriented Programming Systems Languages & Applications (OOPSLA)*, 2014.
- [5] T. Ball. A theory of predicate-complete test coverage and generation. In *Proc. Formal Methods for Components and Objects, 3rd International Symposium (FMCO)*, 2004.
- [6] T. Ball and S. K. Rajamani. Bebop: a path-sensitive interprocedural dataflow engine. In *Proc. ACM SIGPLAN-SIGSOFT Workshop on Program Analysis For Software Tools and Engineering (PASTE)*, 2001.
- [7] N. E. Beckman, A. V. Nori, S. K. Rajamani, R. J. Simmons, S. Tetali, and A. V. Thakur. Proofs from tests. *IEEE Transactions on Software Engineering*, 36(4):495–508, 2010.
- [8] T. A. Budd and D. Angluin. Two notions of correctness and their relation to testing. *Acta Informatica*, 18(1):31–45, 1982.
- [9] C. Chambers and D. Ungar. Customization: Optimizing compiler technology for SELF, a dynamically-typed object-oriented programming language. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, 1989.
- [10] R. Chugh, D. Herman, and R. Jhala. Dependent types for JavaScript. In *Proc. ACM International Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA)*, 2012.
- [11] C. Csallner and Y. Smaragdakis. Check 'n' crash: combining static checking and testing. In *Proc. 27th International Conference on Software Engineering (ICSE)*, 2005.
- [12] R. Cytron, J. Ferrante, B. K. Rosen, M. N. Wegman, and F. K. Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Transactions on Programming Languages and Systems*, 13(4):451–490, 1991.
- [13] E. W. Dijkstra. Notes on structured programming. Technical Report EWD249, Technological University Eindhoven, 1970.
- [14] Ecma International. *Dart Programming Language Specification, ECMA-408, 3rd Edition*, June 2015.
- [15] E. Ernst, A. Møller, M. Schwarz, and F. Strocchio. Message safety in Dart. In *Proc. 11th Symposium on Dynamic Languages (DLS)*, 2015.
- [16] A. Feldthaus, M. Schäfer, M. Sridharan, J. Dolby, and F. Tip. Efficient construction of approximate call graphs for JavaScript IDE services. In *Proc. 35th International Conference on Software Engineering (ICSE)*, 2013.
- [17] J. Ferrante, K. J. Ottenstein, and J. D. Warren. The program dependence graph and its use in optimization. *ACM Transactions on Programming Languages and Systems*, 9(3):319–349, 1987.
- [18] J. A. Goguen and J. Meseguer. Security policies and security models. In *IEEE Symposium on Security and Privacy*, 1982.
- [19] J. B. Goodenough and S. L. Gerhart. Toward a theory of test data selection. *ACM SIGPLAN Notices – International Conference on Reliable Software*, 10(6):493–510, 1975.
- [20] M. Gorbovitski, Y. A. Liu, S. D. Stoller, T. Rothamel, and K. T. Tekle. Alias analysis for optimization of dynamic languages. In *Proc. 6th Symposium on Dynamic Languages (DLS)*, 2010.
- [21] A. Guha, C. Saftoiu, and S. Krishnamurthi. Typing local control and state using flow analysis. In *Proc. Programming Languages and Systems – 20th European Symposium on Programming (ESOP)*, 2011.
- [22] A. Gupta, R. Majumdar, and A. Rybalchenko. From tests to proofs. *International Journal on Software Tools for Technology Transfer*, 15(4):291–303, 2013.
- [23] B. Hackett and S. Guo. Fast and precise hybrid type inference for JavaScript. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, 2012.
- [24] W. Howden. Reliability of the path analysis testing strategy. *IEEE Transactions on Software Engineering*, SE-2(3):208–215, Sept 1976.
- [25] W. E. Howden. Weak mutation testing and completeness of test sets. *IEEE Transactions on Software Engineering*, (4):371–379, 1982.
- [26] M. N. Kedlaya, J. Roesch, B. Robotmili, M. Reshadi, and B. Hardekopf. Improved type specialization for dynamic scripting languages. In *Proc. 9th Symposium on Dynamic Languages (DLS)*, 2013.
- [27] K. Knowles and C. Flanagan. Hybrid type checking. *ACM Trans. Program. Lang. Syst.*, 32(2), 2010.
- [28] B. Korel. The program dependence graph in static program testing. *Information Processing Letters*, 24(2):103–108, 1987.
- [29] B. S. Lerner, J. G. Politz, A. Guha, and S. Krishnamurthi. TeJaS: retrofitting type systems for JavaScript. In *Proc. 9th Symposium on Dynamic Languages (DLS)*, 2013.

- [30] F. Logozzo and H. Venter. RATA: rapid atomic type analysis by abstract interpretation - application to JavaScript optimization. In *Proc. Compiler Construction, 19th International Conference (CC)*, 2010.
- [31] I. Mastroeni and D. Zanardini. Data dependencies and program slicing: from syntax to abstract semantics. In *Proc. ACM SIGPLAN Symposium on Partial Evaluation and Semantics-based Program Manipulation (PEPM)*, 2008.
- [32] T. J. McCabe. A complexity measure. *IEEE Transactions on Software Engineering*, 2(4):308–320, 1976.
- [33] Microsoft. TypeScript language specification, February 2015. <http://www.typescriptlang.org/Content/TypeScript%20Language%20Specification.pdf>.
- [34] M. Naik, H. Yang, G. Castelnovo, and M. Sagiv. Abstractions from tests. In *Proc. 39th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, pages 373–386, 2012.
- [35] A. Rastogi, N. Swamy, C. Fournet, G. M. Bierman, and P. Vekris. Safe & efficient gradual typing for TypeScript. In *Proc. 42nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, 2015.
- [36] A. Sabelfeld and A. C. Myers. Language-based information-flow security. *IEEE Journal on Selected Areas in Communications*, 21(1):5–19, 2003.
- [37] M. Schäfer, M. Sridharan, J. Dolby, and F. Tip. Dynamic determinacy analysis. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, 2013.
- [38] K. Sen, S. Kalasapur, T. Brutch, and S. o. Gibbs. Jalangi: A selective record-replay and dynamic analysis framework for JavaScript. In *Proc. Symposium on the Foundations of Software Engineering (FSE)*, 2013.
- [39] M. Sharir and A. Pnueli. Two approaches to interprocedural dataflow analysis. In *Program Flow Analysis: Theory and Applications*, pages 189–233. Prentice-Hall, 1981.
- [40] M. Sridharan, S. Chandra, J. Dolby, S. J. Fink, and E. Yahav. Alias analysis for object-oriented programs. In *Aliasing in Object-Oriented Programming. Types, Analysis and Verification*, volume 7850 of *LNCS*, pages 196–232. Springer, 2013.
- [41] A. Takikawa, D. Feltey, E. Dean, M. Flatt, R. B. Findler, S. Tobin-Hochstadt, and M. Felleisen. Towards practical gradual typing. In *Proc. 29th European Conference on Object-Oriented Programming (ECOOP)*, 2015.
- [42] F. Tip. A survey of program slicing techniques. *Journal of Programming Languages*, 3(3), 1995.
- [43] S. Tobin-Hochstadt and M. Felleisen. The design and implementation of typed Scheme. In *Proc. 35th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, 2008.
- [44] P. Vekris, B. Cosman, and R. Jhala. Trust, but verify: Two-phase typing for dynamic languages. In *Proc. 29th European Conference on Object-Oriented Programming (ECOOP)*, 2015.
- [45] M. M. Vitousek, A. M. Kent, J. G. Siek, and J. Baker. Design and evaluation of gradual typing for Python. In *Proc. 10th ACM Symposium on Dynamic Languages (DLS)*, 2014.
- [46] S. Wei and B. G. Ryder. Practical blended taint analysis for JavaScript. In *Proc. International Symposium on Software Testing and Analysis (ISSTA)*, 2013.
- [47] M. Weiser. Program slicing. In *Proc. 5th International Conference on Software Engineering (ICSE)*, 1981.
- [48] G. Yorsh, T. Ball, and M. Sagiv. Testing, abstraction, theorem proving: better together! In *Proc. ACM/SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*, 2006.
- [49] H. Zhu, P. A. V. Hall, and J. H. R. May. Software unit test coverage and adequacy. *ACM Computing Surveys*, 29(4):366–427, 1997.