

# Identity-sensitive Points-to Analysis for the Dynamic Behavior of JavaScript Objects

Shiyi Wei and Barbara G. Ryder

Department of Computer Science,  
Virginia Tech, Blacksburg, VA, USA.  
`{wei,ryder}@cs.vt.edu`

**Abstract.** JavaScript object behavior is dynamic and adheres to prototype-based inheritance. The behavior of a JavaScript object can be changed by adding and removing properties at runtime. Points-to analysis calculates the set of values a reference property or variable may have during execution. We present a novel, partially flow-sensitive, context-sensitive points-to algorithm that accurately models dynamic changes in object behavior. The algorithm represents objects by their creation sites and local property names; it tracks property updates via a new control-flow graph representation. The calling context is comprised of the receiver object, its local properties and prototype chain. We compare the new points-to algorithm with an existing JavaScript points-to algorithm in terms of their respective performance and accuracy on a client application. The experimental results on real JavaScript websites show that the new points-to analysis significantly improves precision, uniquely resolving on average 11% more property lookup statements.

**Keywords:** JavaScript, program analysis, points-to analysis

## 1 Introduction

Dynamic programming languages, including JavaScript, Ruby and PHP, are widely used in developing sophisticated software systems, especially web applications. These languages share several dynamic features, including dynamic code generation and dynamic typing, used in real-world programs [19]. For example, JavaScript code can be generated at runtime using *eval* and JavaScript functions can be variadic (i.e., functions can be called with different numbers of arguments). Despite the popularity of these dynamic languages, there is insufficient tool support for developing and testing programs because their dynamic features render many traditional analyses, and tools which depend on them, ineffective.

In addition, instead of class-based inheritance JavaScript supports prototype-based inheritance [15, 26] that results in a JavaScript object inheriting properties from a chain of (at least one) prototype objects. The model also allows the properties of a JavaScript object to be added, updated, or deleted at runtime. This means that JavaScript objects can exhibit different behaviors at different

times during execution. Moreover, object constructors may be polymorphic so that objects created by the same constructor may have distinct properties. These aspects of JavaScript further complicate tool building.

Some tools have been developed to support JavaScript software development (e.g., [18, 21]). Points-to analysis is the enabling analysis for such tools. Researchers have proposed several points-to analyses that handle different features of JavaScript (e.g., [8, 16, 25]). Nevertheless, there are opportunities to significantly improve the precision of points-to analysis through better modeling of object property set changes.

In this paper, we present a novel points-to algorithm that can accurately model JavaScript objects. Changes to object properties are tracked more accurately to reflect object run-time behavior at different program points. A new graph decomposition for control-flow graphs is used to better track object property changes. Prototype-based inheritance is more accurately modeled to locate delegated properties. The analysis identifies objects by their creation site as well as their local property names upon construction, more accurately than the per-creation-site representation. To distinguish polymorphic constructors, this analysis can incorporate dynamic information collected at runtime (see Section 3.6). Technically, the analysis is *partially flow-sensitive* (on our new control-flow graph structure) and *context-sensitive*, using a new form of object sensitivity [17].<sup>1</sup> Rather than using the receiver object creation site as a calling context in the analysis, we use an approximation of the receiver object and its properties at the call site (i.e., *object identity*).

In order to compare this algorithm with previous techniques, we instantiated our new points-to analysis as the static component of the JavaScript Blended Analysis Framework (*JSBAF*) [27]. Blended analysis collects run-time information through instrumentation to define the calling structure used by a subsequent static analysis, and to capture dynamically generated code. It has been demonstrated that blended analysis is practical and effective on JavaScript programs [27]. We measured the performance and accuracy of our new analysis on a statement-level points-to client (*REF* analysis) that calculates how many objects are returned by a property lookup (e.g., a read of  $x.p$ ).

The major contributions of this work are:

- We have designed a novel *identity-sensitive points-to analysis that accurately and safely handles dynamic changes in the behavior of JavaScript objects*. This algorithm presents a new program representation that enables partially flow-sensitive analysis, a more accurate object representation, and an expanded points-to graph that facilitates strong updates on the statements changing object properties.

---

<sup>1</sup> Informally, a flow-sensitive analysis follows the execution order of statements in a program; a flow-sensitive analysis can perform strong updates, but a flow-insensitive one cannot. A context-sensitive analysis distinguishes between different calling contexts of a method, producing different analysis results for each context [20, 23]. A context-insensitive analysis calculates one solution per method.

- Experimental results from our new analysis compared to a recent points-to analysis [25], both implemented in *JSBAF*, showed that *identity-sensitive analysis significantly improved precision*. On average over all the benchmarks (i.e., 12 popular websites), 48% of the property lookup statements were resolved to a single object by our new analysis, while the existing analysis [25] uniquely resolved only 37% of these statements. Although our analysis incurred a 135% time overhead on average to achieve the increased precision, it was able to analyze each of the programs in the benchmarks in under 5 minutes, attesting to its scalability in practice.

**Overview.** Section 2 defines the notion of object identity and then uses an example to illustrate the sources of imprecision in JavaScript points-to analysis. Section 3 describes our new points-to analysis algorithm and implementation. Section 4 presents the *REF* analysis and the experimental results. Section 5 discusses related work, and Section 6 offers conclusions and future work.

## 2 Background

In this section we define key concepts using an example to illustrate the sources of imprecision in current points-to analyses for JavaScript.

### 2.1 JavaScript Object Identity

JavaScript is a dynamically typed programming language whose object behavior can change as object properties are added or deleted at runtime. In strongly typed programming languages, the notion of *type* is used to abstract the possible behavior of an object (e.g., the class of an object in Java) [22]; however, in dynamically typed languages, the type of an object can change during execution. In order to avoid confusion, we call the type of a JavaScript object its *object identity*.<sup>2</sup>

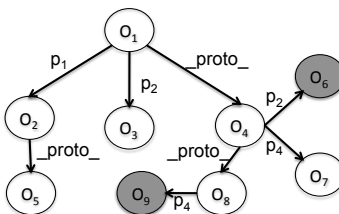
**Definition 1.** The *identity* of an object at a program point denotes all of its accessible properties and their non-primitive values.

The accessible properties of an object conform to the property lookup mechanism implemented in JavaScript. Every JavaScript object includes an internal reference to its prototype object from which it inherits non-local properties. A JavaScript object may have a sequence of prototype objects (i.e., a prototype chain) whose properties it can inherit. When reading a property  $p$  of an object  $o$ , the JavaScript runtime checks the local properties of  $o$  to see if  $o$  has a property named  $p$ . If not, the JavaScript runtime checks to see if the prototype object of  $o$  has a property named  $p$ , continuing to check along the prototype chain from object to object until the property is found (or not) [7].

<sup>2</sup> This general notion can be used for other dynamic languages and is related to structured typing for strongly typed languages [22].

**Definition 2.** *Identity-update statements* are: (1) property write statement (i.e.,  $x.p = y$  or  $x[p'] = y$ ), (2) property delete statement (i.e., *delete*  $x.p$  or *delete*  $x[p']$ ), and (3) an invocation that directly or indirectly results in execution of (1) and/or (2).

The *identity-update statements* are the set of statements in JavaScript that may affect the object identities. In Figure 1, we illustrate the identity of an object with an example. Figure 1 shows the objects connected to  $O_1$  at a program point. The local properties of object  $O_1$  are named  $p_1$  and  $p_2$  and  $O_4$  is its prototype object.  $O_7$  is visible from  $O_1$  by accessing  $O_1.p_4$  while  $O_6$  is not visible from  $O_1$  by accessing  $O_1.p_2$  because a local property named  $p_2$  exists for  $O_1$ . To sum up, the shaded nodes (i.e.,  $O_6$  and  $O_9$ ) are not accessible from  $O_1$  and the unshaded nodes constitute  $O_1$ 's identity.



**Fig. 1:** Object identity for  $O_1$ . (Unshaded nodes only)

## 2.2 Imprecision of Points-to Analysis

A flow-insensitive analysis may produce imprecise results when object identity changes, because it must safely approximate points-to relations and it cannot do strong updates. Existing context-sensitive analyses may produce imprecise results because they lack the power to distinguish between different object identities for the same JavaScript object. In Figure 2, we present a JavaScript example to illustrate the sources of imprecision of a flow-insensitive, context-insensitive points-to analysis resulting from several dynamic features of JavaScript. We also demonstrate that an existing context-sensitive analysis using the same object representation as [17] is ineffective at distinguishing the function calls in the example.

Lines 2-6 show a constructor function  $X()$ . Objects created by  $X()$  may or may not have the local property named  $p$  or  $q$  (lines 4 and 5) depending on the value of its argument. The statement in line 12 updates the value of local property  $p$  of an object pointed to by  $x$  if  $p$  exists; otherwise, the statement adds the local property named  $p$  to the object. Figures 3(a) and 3(b) show the points-to graphs that reflect the run-time behavior of this code. We use the line number to represent the object created (e.g., the object created at line 7 (*new X*) is  $O_7$ ). We focus on two program points in the execution, lines 10 and 15. The nodes  $O_7$ ,  $O_4$ , and  $O_3$  and  $O_9$  constitute the identity of  $O_7$  at line 10 and the nodes  $O_7$ ,  $O_{12}$ ,  $O_3$  and  $O_{14}$  constitute the identity of  $O_7$  at line 15. Note that

```

1 function P(){ this.p = new Y1(); }
2 function X(b){
3     this.__proto__ = new P();
4     if(b) { this.p = new Y2(); }
5     else this.q = new Y3();
6 }
7 var x = new X(true);
8 x.bar = function(v, z){ v.f = z; }
9 var z1 = new Z();
10 x.bar(x.p, z1);
11 ...
12 x.p = new A();
13 ...
14 var z2 = new Z();
15 x.bar(x.p,z2);

```

**Fig. 2:** JavaScript example.

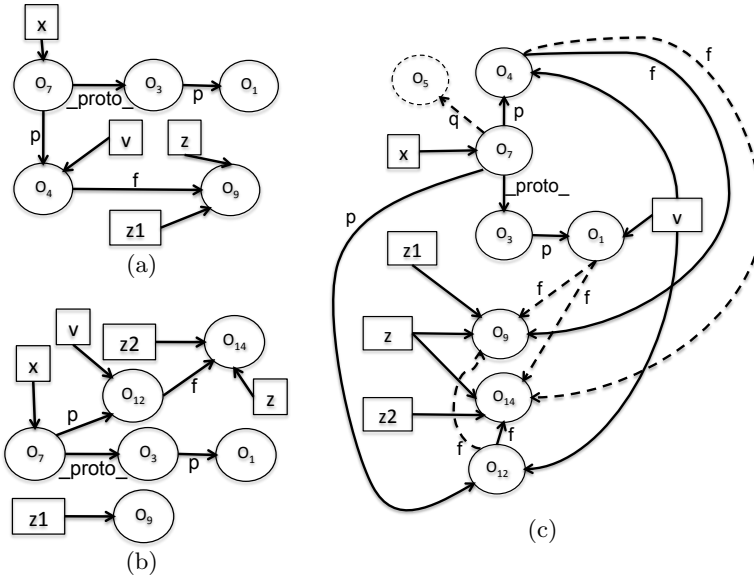
$O_1$  is not visible from  $O_7$  at lines 10 or 15 because of the existence of the local property named  $p$ . The identity of object  $O_7$  is different at these two program points.

Constructor polymorphism (lines 2-6), object property change (line 12) and function invocations (lines 10 and 15) in the example make precise static points-to analysis hard to achieve with current techniques. Figure 3(c) shows a points-to graph for the example built by a flow- and context-insensitive points-to analysis. Dashed nodes and edges are imprecise points-to relations that cannot exist at runtime.

There are several sources of imprecision. Line 7 creates an object pointed to by variable  $x$  by invoking the polymorphic constructor  $X()$ . Not knowing the value of  $b$ , static analysis conservatively builds all the points-to relations possible from execution of  $X()$ .<sup>3</sup> When reading the property  $p$  of  $x$  (line 10), static analysis returns objects  $O_4$  and  $O_1$  because a conservative analysis cannot distinguish whether or not  $O_4$  actually exists. Furthermore, because of the imprecise result of the read of  $x.p$ , invoking the  $bar()$  function results in imprecise property reference from  $O_1$  to  $O_9$ . Flow-insensitive points-to analysis simply adds  $O_{12}$  to  $O_7.p$  (line 12) because it cannot perform strong updates. Because the analysis does not distinguish which objects  $v$  and  $z$  point to on different calls of  $bar()$ , line 15 results in additional imprecision with respect to  $O_4.f$  and  $O_{12}.f$ .

Flow-sensitive analysis is not sufficient to resolve the imprecision in the example without an appropriate context for call sites. First, indirect assignment statements cannot be strongly updated in general. Second, assuming  $x.p$  is strongly updated to point to  $O_{12}$  at line 12, a context-insensitive analysis does not remove the imprecise edges  $\langle O_4, f \rangle, O_{14}$  and  $\langle O_{12}, f \rangle, O_9$  because calls to  $bar()$  (lines 10 and 15) are not distinguished by calling contexts. Object sensitivity [17]

<sup>3</sup> In this short example, constant propagation of parameters would help static analysis precision but clearly this is not always possible.



**Fig. 3:** Imprecision of static points-to analysis. (a) Run-time points-to graph at line 10. (b) Run-time points-to graph at line 15. (c) Flow- and context-insensitive points-to graph.

has been shown to perform better than call-string context sensitivity [23] for the idioms used in object-oriented languages [14]. However, object-sensitive analysis is not able to differentiate these two call sites because they have the same receiver object  $O_7$ , which has two different identities at these call sites. Our new points-to analysis is designed to handle these constructs more accurately and to address the challenges raised by object identity updating and prototype-based inheritance.

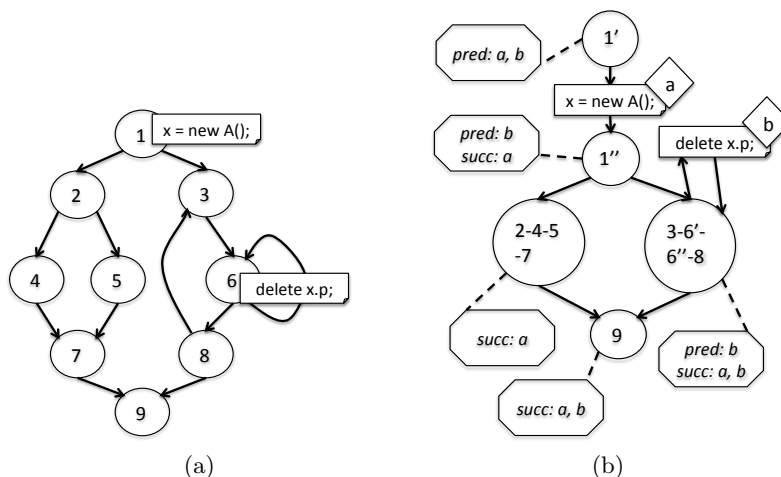
### 3 Identity-sensitive Points-to Analysis

In this section we will present our *identity-sensitive* points-to analysis for JavaScript. We will explain key ideas used in the analysis, including the intra-procedural program representation (i.e., the block-sensitive decomposition of control-flow graphs), the solution space (i.e., the annotated points-to graph with access path edges and in-construction nodes), the transfer functions of the identity-update statements as well as the identity-preserving statements, identity sensitivity (i.e., a form of context sensitivity based on object sensitivity that captures changes in object behavior during execution) and block sensitivity (i.e., a partial flow sensitivity performed on the transformed CFG). Finally, we will discuss the implementation details of our algorithm.

### 3.1 Identity-Preserving Block Graph

A flow-insensitive analysis ignores the control flow of a program while a flow-sensitive analysis typically uses an intra-procedural control-flow graph (CFG). Our analysis aims to provide a better model of a JavaScript object whose identity exhibits flow-sensitive characteristics (e.g., allowing addition and deletion of object properties at any program point). Cognizant of the possible overhead introduced by a fully flow-sensitive analysis, we designed a partially flow-sensitive analysis that only performs strong updates when possible on identity-update statements using a transformed CFG, called the *Identity-Preserving Block Graph (IPBG)*. Recall that the identity-update statements, including the property write (i.e., add or update a property) and delete (i.e., remove a property), directly change the object identity in JavaScript; all other statements (e.g., property read) are identity-preserving statements.

Figure 4 shows an example *IPBG* (Figure 4(b)) compared to its original CFG (Figure 4(a)). An *IPBG* is a transformed control-flow graph whose basic blocks are aggregated into region nodes according to whether or not they contain an object identity-update statement. The *IPBG* also contains identity-update statements as special singleton statement nodes (i.e., *identity-update* nodes). An example of a region node (i.e., *identity-preserving* node) is 2-4-5-7 in Figure 4(b) whereas node  $x = \text{new } A()$  is an example of a *identity-update* node. Note that in creating singleton nodes the algorithm breaks apart former basic blocks (e.g.,  $1 \rightarrow \{1', x = \text{new } A(), 1''\}$ ).



**Fig. 4:** *IPBG* generation. (a) CFG. (b) *IPBG*.

We first split any basic blocks in the CFG that contain at least an identity-update statement (see Definition 2 in Section 2), obtaining a *split-CFG*. Identity-update statements (1) and (2) can be detected syntactically and invocations that may result in an identity change (i.e., category (3)) are found by a linear

call graph traversal.<sup>4</sup> We then use a variant of the standard CFG construction algorithm [1] to build the split-CFG. The header nodes used include the standard headers [1] plus (i) any identity-update statement is a region header of an identity-update node containing only that statement, and (ii) any identity-preserving statement that immediately follows an identity-update statement is a region header of an identity-preserving node.

In an *IPBG*, identity-preserving region nodes are formed based on grouping nodes in the split-CFG that share the same control-flow relations with respect to identity-update nodes. The possible control-flow relations of node  $n_1$  and  $n_2$  in a split-CFG include: (1)  $n_1$  is a successor of  $n_2$ , (2)  $n_1$  is a predecessor of  $n_2$ , (3)  $n_1$  is both a successor and a predecessor of  $n_2$  (i.e.,  $n_1$  and  $n_2$  exist in a loop) and (4)  $n_1$  and  $n_2$  have no control-flow relation (e.g.,  $n_1$  and  $n_2$  are present in different branches). We label each node in a split-CFG with its relations to each identity-update node via depth-first searches. The set of labels form a signature for that node. If nodes share the same signature it means that they have the same control-flow relationship(s) to a (set of) identity-update statement(s) so that they can be collapsed to an identity-preserving node in the *IPBG*. Figure 4(b) shows the signatures of the identity-preserving regions in the generated *IPBG*;  $a$  and  $b$  represent the identity-update statements  $x=new A()$  and  $delete x.p$ , respectively. Basic blocks 2, 4, 5 and 7 are aggregated because they only appear as successors of  $x=new A()$  and have no control-flow relation to  $delete x.p$ . The region node 2-4-5-7 is not further aggregated with basic block 9 because 9 is a successor of  $delete x.p$  but 2-4-5-7 is not.

### 3.2 Points-to Graph Representation





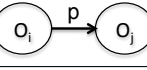

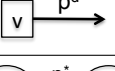
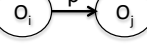
Our points-to graph representation includes constructs that facilitate the handling of strong updates by our analysis. Our algorithm design allows strong updates when possible for identity-update statements. In contrast, most flow-sensitive Java analysis algorithms cannot perform strong updates for indirect assignment statements (e.g.,  $x.p = y$ ) and few analyses consider property delete statements, which are uncommon in object-oriented languages. Two existing techniques help to enable strong updates for such statements in JavaScript: recency abstraction and access path maps.

Recency abstraction [2, 11] associates two memory-regions with each allocation site. The most-recently-allocated block, a concrete memory-region, allows strong updates and the not-most-recently-allocated block is a summary memory-region. We adapt the idea of recency abstraction to enable strong updates during analysis of constructor functions.

De *et al.* [6] performed strong updates at indirect assignments by computing the map from access paths (i.e., a variable followed a sequence of property accesses) to sets of abstract objects. This work demonstrated the validity of using access path maps to perform strong updates for indirect write statements in

---

<sup>4</sup> Our analysis requires a pre-computed call graph as input. See Section 3.6 for details.

points-to graph $G$	node $N$	variable $v$	
		abstract object $o$	
		in-construction object $@o$	
	edge $E$	variable reference $(v, \phi o)$	
		property reference $(\langle \phi o_i, p \rangle, \phi o_j)$	
		access path $(\langle v, p \rangle, \phi o)$	
	annotation $A$	$d$ annotation $p^d$	
		* annotation $p^*$	

**Table 1:** Expanded points-to graph with annotations.

Java. We adapt this approach to points-to analysis for JavaScript by expanding the points-to graph representation instead of using separate maps.

Table 1 lists the nodes, edges and annotations in our points-to graph. In addition to variable nodes  $v$  and abstract object nodes  $o$ , our points-to graph contains *in-construction object nodes*  $@o$ .<sup>5</sup> Details of the in-construction objects will be discussed in Section 3.3. For sake of simplicity, we use  $\phi o$  to represent either kind of object node (i.e.,  $o$  or  $@o$ ).

There are three kinds of edges. Variable reference and property reference edges exist in a traditional points-to graph. An *access path edge*,  $(\langle v, p \rangle, \phi o)$ , denotes that the property  $p$  of variable  $v$  refers to object  $\phi o$ .  $\langle v, p \rangle$  represents an access path with length of 2 (i.e., a variable followed by one field access  $v.p$ ).<sup>6</sup>

Our analysis calculates *may* pointer information, meaning that a points-to edge in the graph may or may not exist at runtime. To better approximate the identities of JavaScript objects, we introduce annotations on property reference edges as well as access path edges. The annotations help to calculate *must exist* information for object property names. In our analysis, the  $d$  annotation on a property name  $p$  (i.e.,  $p^d$ ) denotes that the local property named  $p$  must not exist. This annotation only applies to access path edges in our points-to graph. The other annotation,  $*$ , applies to both property reference edges and access path edges.  $p^*$  denotes that the local property named  $p$  may not exist. Property reference edges without annotation or access path edges without annotation

<sup>5</sup> Similar to the recency abstraction, an in-construction object always describes exactly one concrete object. In our analysis, it exists only during analysis of a constructor.

<sup>6</sup> The length of an access path is one more than the number of field accesses [6].

represent *must exist* information for the property names. We use  $p^\phi$  to represent any kind of  $p^d$ ,  $p^*$  or  $p$  edge. These annotated edges help us perform a more accurate property lookup (see Section 3.3).

$Pt(x)$  denotes the points-to set of  $x$  and  $Pt(< \phi o, p >)$  denotes the points-to set of the property  $p$  of  $\phi o$ .  $Pt(< v, p >)$  denotes the points-to set of access path  $v.p$ . We also define the operation  $Alias(v)$  which returns the set of variables  $W$  such that  $v$  and  $w \in W$  point to the same object.  $apset(v)$  denotes the set of all access path edges of  $v$  (i.e.,  $apset(v) = \forall q : \{(< v, q^\phi >, \phi o)\}$ ).

In addition to the points-to graph, we use a mapping data structure to store intermediate information in the analysis. The map  $M$  is used to record the list of property names when an object is constructed. An abstract object (e.g.,  $o$ ) is the key in  $M$  whose value is the set of local property names that exist when the constructor function of the abstract object returns (e.g.,  $\{p1, p2, p3\}$ ).

### 3.3 Points-to Analysis Transfer Functions

In this section we describe the data-flow transfer functions for the statements shown in Table 2.

**Object creation** ( $x = new X(a_1, a_2, \dots, a_n)$ ). In our analysis, an object creation statement (i.e., *new* statement) is modeled in three steps.  $x = new X$  creates an in-construction object  $@o_i$ . Then the invocation of the constructor  $new X((a_1, a_2, \dots, a_n))$  is modeled as a function call on  $@o_i$ . Upon the return of the constructor (i.e.,  $ret_X$ ), the analysis removes the in-construction object from the points-to graph and redirects all points-to relations from  $@o_i$  to an abstract object ( $remove(G, @o_i)$ ). If the local property set of the in-construction object matches that of an existing abstract object with the same allocation site, the in-construction object is merged into the abstract object; otherwise, a new abstract object is created to replace the in-construction object. There is at most one in-construction object for each creation site.<sup>7</sup>

The transfer function of the object creation statement ensures that abstract objects are based on their allocation site as well as their constructed local properties (i.e., an approximation of actual object identity); in other words, the objects created at the same allocation site that contain the same set of local property names share the same abstract object in our analysis. This object representation is more precise than using one abstract object per creation site.

**Property write** ( $x.p = y$ ). In general, strong updates cannot be performed on the property write statement because an abstract object may summarize multiple run-time objects; however, use of in-construction objects and access path edges enable strong updates in our analysis. In the points-to graph  $G$ , if  $x$  only refers to one object and the object is an in-construction object, we know

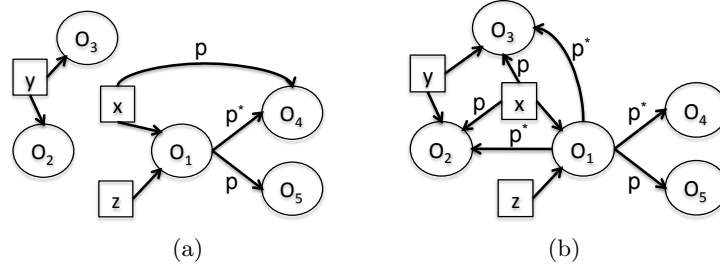
<sup>7</sup> Recursive constructor calls involve the creation of an in-construction object when the in-construction object for the same allocation site already exists (before it resolves into an abstract object). In our analysis, the existing in-construction object is resolved into a *special abstract object* whose set of properties upon construction is *unknown*. A fixed point calculation is done using the special abstract object.

Statement	Transfer function
$s_i : x = \text{new } X(a_1, a_2, \dots, a_n)$	<p>(1) <math>x = \text{new } X : (G - \text{apset}(x)) \cup (x, @o_i)</math></p> <p>(2) <math>\text{new } X((a_1, a_2, \dots, a_n)) : G \cup \{\text{invoke}(G, X, @o_i, a_1, a_2, \dots, a_n)\}</math></p> <p>(3) <math>\text{ret}_X : \text{remove}(G, @o_i)</math></p>
$x.p = y$	<p>(1) <i>if</i> <math> Pt(x)  = 1</math> and <math>\{\phi_{o_i} \in @O \mid \phi_{o_i} \in Pt(x)\} :</math>  <math>G - \{(\langle @o_i, p^\phi \rangle, \phi_{o_j}) \mid @o_i \in Pt(x) \wedge \phi_{o_j} \in Pt(\langle @o_i, p \rangle)\} \cup \{(\langle @o_i, p \rangle, \phi_{o_j}) \mid @o_i \in Pt(x) \wedge \phi_{o_j} \in Pt(y)\}</math></p> <p>(2) <i>otherwise</i> :  (2.1) <math>G - \{(\langle x, p^\phi \rangle, \phi_{o_i}) \mid \phi_{o_i} \in Pt(\langle x, p^\phi \rangle)\} \cup \{(\langle x, p \rangle, \phi_{o_j}) \mid \phi_{o_j} \in Pt(y)\}</math>  (2.2) <math>G \cup \{(\langle \phi_{o_i}, p^* \rangle, \phi_{o_j}) \mid \phi_{o_i} \in Pt(x) \wedge \phi_{o_j} \in Pt(y)\}</math>  (2.3) <math>G \cup \{(\langle \langle z, p^* \rangle, \phi_{o_i}) \mid z \in Alias(x) \wedge Pt(z, p^\phi) \neq \emptyset \wedge \phi_{o_i} \in Pt(y)\}</math></p>
<i>delete</i> $x.p$	<p>(1) <i>if</i> <math> Pt(x)  = 1</math> and <math>\{\phi_{o_i} \in @O \mid \phi_{o_i} \in Pt(x)\} :</math>  <math>G - \{(\langle @o_i, p^\phi \rangle, \phi_{o_j}) \mid @o_i \in Pt(x) \wedge \phi_{o_j} \in Pt(\langle @o_i, p^\phi \rangle)\}</math></p> <p>(2) <i>otherwise</i> :  (2.1) <math>G - \{(\langle x, p^\phi \rangle, \phi_{o_i}) \mid \phi_{o_i} \in Pt(\langle x, p^\phi \rangle)\} \cup \{(\langle x, p^d \rangle, \text{null})\}</math>  (2.2) <math>G \cup \{(\langle \phi_{o_i}, p^* \rangle, \phi_{o_j}) \mid \phi_{o_i} \in Pt(x) \wedge \phi_{o_j} \in Pt(\langle \phi_{o_i}, p \rangle)\} - \{(\langle \phi_{o_i}, p \rangle, \phi_{o_j}) \mid \phi_{o_i} \in Pt(x) \wedge \phi_{o_j} \in Pt(\langle \phi_{o_i}, p \rangle)\}</math>  (2.3) <math>G \cup \{(\langle \langle z, p^* \rangle, \phi_{o_i}) \mid z \in Alias(x) \wedge Pt(z, p) \neq \emptyset \wedge \phi_{o_i} \in Pt(\langle z, p \rangle)\} - \{(\langle z, p \rangle, \phi_{o_i}) \mid z \in Alias(x) \wedge Pt(z, p) \neq \emptyset \wedge \phi_{o_i} \in Pt(\langle z, p \rangle)\}</math></p> <p><math>(G - \text{apset}(x)) \cup \{(x, \phi_{o_i}) \mid \phi_{o_i} \in Pt(y)\}</math></p>
$x = y$	$(G - \text{apset}(x)) \cup \{(x, \phi_{o_i}) \mid \phi_{o_i} \in Pt(y)\}$
$x = y.p$	$(G - \text{apset}(x)) \cup \{(\langle x, \phi_{o_i} \rangle \mid \phi_{o_i} \in \text{lookup}(y, p))\}$
$x = y.m(a_1, a_2, \dots, a_n)$	$(G - \text{apset}(x)) \cup \{\text{invoke}(G, M, \phi_{o_i}, a_1, a_2, \dots, a_n) \mid \phi_{o_i} \in Pt(y) \wedge M \in \text{lookup}(y, m)\}$

**Table 2:** Transfer functions of program statements.

that  $x$  refers to a specific concrete object. The analysis then performs strong updates on the property reference edges by removing the points-to edges in  $G$  denoting  $@o_i.p$  (if they exist) and adding the new edges implied by  $Pt(y)$ . In other cases (i.e., the cardinality of  $Pt(x)$  is more than 1 or  $x$  refers to an abstract object), we use access path edges to enable strong updates on property write statements. First, the access path of  $x.p$  can be strongly updated by removing the access path edges in  $G$  denoting  $x.p$  (if they exist) and adding the new edges (e.g.,  $\langle x.p, o_j \rangle$  where  $o_j$  is referred to by  $y$ ). Second, the object(s)  $x$  points to are weakly updated (e.g., the edge  $\langle o_i.p^*, o_j \rangle$  is inserted if  $x$  points to  $o_i$  and  $y$  points to  $o_j$ ). The property reference edges are inserted with the  $*$  annotation because the property write statement may not affect all variables pointing to the updated object. Last, the access path edges of the variables that have a may alias relation to  $x$  need to be weakly updated. For example,  $\langle z.p^*, o_i \rangle$  is inserted to  $G$  if  $z$  may be an alias of  $x$ , and there exists at least an edge denoting  $x.p$  (with or without annotation).

In Figure 5, we show an example of the effects of property write statement on the points-to graph. Figure 5(a) illustrates the input points-to graph for the property write statement  $x.p = y$ . In Figure 5(b), our analysis performs a strong update on the access path  $x.p$  (i.e., delete  $\langle x.p, O_4 \rangle$  and add  $\langle x.p, O_2 \rangle$ ,  $\langle x.p, O_3 \rangle$ ) and inserts the edges  $\langle O_1.p^*, O_2 \rangle$ ,  $\langle O_1.p^*, O_3 \rangle$  (i.e., weak updates).



**Fig. 5:** Property write example. (a) Input points-to graph. (b) Updated points-to graph.

**Property delete** (*delete  $x.p$* ). The transfer function of the delete statement is similar to the property write statement. Our analysis strongly updates the access path edges by removing the existing edges and adding a new edge (i.e.,  $\langle x.p^d, null \rangle$ ) that denotes  $x$  must not have a local access path  $x.p$ . When performing weak updates on the property reference edges of an object  $o_i$  that is referred to by  $x$ , all existing edges denoting  $o_i.p$  should be annotated by  $*$  because the property named  $p$  may not exist locally for  $o_i$ . The same rule applies when updating the access path edges of the aliases of  $x$ .

**Direct write** ( $x = y$ ). The effects of direct variable assignment on the points-to graph are relatively straightforward.  $x = y$  creates points-to edges from  $x$  to all objects pointed to by  $y$ . Note that we perform weak updates on direct

assignments. Although the analysis removes all the access path edges of  $x$  from the points-to graph (i.e.,  $G - apset(x)$ ), soundness is ensured because lookups through the abstract objects reflect *less precise, yet safe approximations*. Also, the access path edges of  $y$  cannot be copied to  $x$  because access path edges can only be added via strong updates.

**Property read** ( $x = y.p$ ). JavaScript enforces an asymmetry between reading and writing property values. When writing the value of a property or deleting a property, JavaScript always uses the local property, ignoring the prototype object. When reading a property of a variable (e.g.,  $x = y.p$ ), recall that JavaScript supports prototype-based inheritance. In some existing points-to analyses for JavaScript, when reading property  $p$  of an object, the property lookup mechanism is modeled by reporting *all* properties named  $p$  in the prototype chain of the object to ensure analysis safety.

---

**Procedure 1** Optimized object property lookup:  $lookup(v, p)$

---

**Output:** accessible objects  $v.p$ :  $P$

```

1: if  $Pt(< v, p >) \neq \emptyset$  then
2:    $P \cup Pt(< v, p >) \cup Pt(< v, p^* >)$ 
3:   return
4: else if  $Pt(< v, p^* >) \neq \emptyset$  or  $Pt(< v, p^d >) \neq \emptyset$  then
5:    $P \cup Pt(< v, p^* >)$ 
6:   for each object  $\phi o$  in  $lookup(v, \_proto\_)$  do
7:      $S.push(\phi o)$ 
8:   end for
9: else
10:  for each object  $\phi o$  in  $Pt(v)$  do
11:     $S.push(\phi o)$ 
12:  end for
13: end if
14: while  $S$  is not empty do
15:   $\phi o_i \leftarrow S.pop()$ 
16:   $P \cup Pt(< \phi o_i, p >) \cup Pt(< \phi o_i, p^* >)$ 
17:  if  $|Pt(< \phi o_i, p >)| = 0$  and  $(Pt(< \phi o_i, \_proto\_ >) \neq null$  or  $Pt(< \phi o_i, \_proto\_^* >) \neq null)$  then
18:    for each object  $\phi o_j$  in  $Pt(< \phi o_i, \_proto\_ >) \cup Pt(< \phi o_i, \_proto\_^* >)$  do
19:       $S.push(\phi o_j)$ 
20:    end for
21:  end if
22: end while

```

---

Procedure 1 illustrates our potentially more precise property lookup procedure enabled by our edge types and their annotations. This worklist algorithm iterates through all the accessible objects in the points-to graph when property  $p$  of variable  $v$  is read. Intuitively, it favors the use of access path edges in property lookup because they reflect the results of strong updates, before examining property reference edges. Lines 1 to 12 initialize the algorithm upon three conditions.

(1) If there exist access path edges for  $v.p$  without annotation (i.e., property  $p$  must exist locally), the objects in the  $Pt(< v, p >)$  and  $Pt(< v, p^* >)$  are considered to be accessible properties (line 2) and the algorithm returns (line 3). (2) If there exist access path edges for  $v.p$  with either annotation, the algorithm needs to lookup objects in the prototype chain. In this case, the objects in the  $Pt(< v, p^* >)$  (if  $v.p^*$  exists) are considered to be accessible properties (line 5) and the algorithm pushes all the immediate prototype objects of  $v$  onto the worklist (lines 6 to 8). (3) Otherwise (i.e., no access path edge for  $v.p$  exists), only the abstract objects are used for looking up so that all the objects in the  $Pt(v)$  are pushed onto the worklist (lines 10 to 12). Lines 14 to 22 iterate the worklist. All the objects in  $Pt(< \phi o, p >)$  and  $Pt(< \phi o, p^* >)$  are considered to be accessible properties by our analysis (Line 16). Since an edge annotated with  $*$  means that the property may not exist locally, the algorithm will continue looking up the prototype chain, until it reaches at least one points-to edge named  $p$  without annotation or the end of the prototype chain (Line 17 to 21). Thus, instead of finding *all* the properties named  $p$  in the prototype chain (i.e.,  $lookup\_all(v, p)$ ), our algorithm can stop when it finds an *existing* property  $p$  (i.e., a property named  $p$  without annotation).

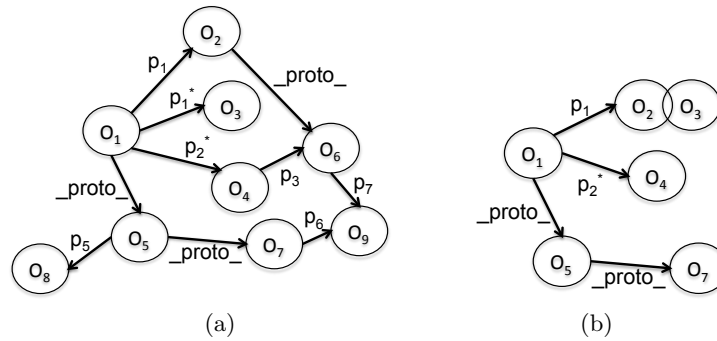
This new property lookup algorithm  $lookup(v, p)$  mimics the run-time property lookup mechanism of JavaScript while still assuring the safety of our analysis. For the example in Figure 5(b),  $lookup(z, p)$  results in  $O_2$  and  $O_3$  through the access path while  $lookup(x, p)$  results in  $O_2, O_3, O_4$  and  $O_5$  through the abstract object  $O_1$ . In Table 2, the transfer function of the property read statements refers to this optimized object property lookup algorithm. Because we perform weak updates on the property read statements, similar to direct writes, the analysis removes all the access path edges of  $x$  from the points-to graph to ensure safety.

**Method invocation** ( $x = y.m(a_1, a_2, \dots, a_n)$ ). The method invocation (e.g.,  $x = y.m(a_1, a_2, \dots, a_n)$ ) resolves for every receiver object pointed to by  $y$ . The invoked methods are determined by reading the property  $y.m$  through our optimized lookup algorithm. Upon the return of method invocation,  $x$  is weakly updated by removing all its access path edges from  $G$ .

### 3.4 Identity Sensitivity

Identity sensitivity for JavaScript is a new form of context sensitivity derived from the notion of object sensitivity for languages such as Java [17]. In object sensitivity, each method is analyzed separately for each object on which it may be invoked. For strongly typed languages like Java, often object sensitivity identifies objects in the analysis by their creation sites. Calls of a method using two different receiver objects (i.e., created at different sites) will result in two separate analyses of the method, even if the calls originated from the same call site. However, this is insufficient for JavaScript analysis, because object behavior may change dynamically at any program point during execution.

Ideally, identity sensitivity would analyze each method separately for each object identity on which it may be invoked. However, the graph representation



**Fig. 6:** Approximate object identity as a context. (a) Object identity of  $O_1$ . (b) Approximate object identity of  $O_1$ .

of object identity may contain many edges and nodes both locally and along prototype chains (e.g., object identity of  $O_1$  in Figure 6(a)), which would be prohibitively expensive to use as a context. Therefore, we use an *approximation* of the object identity of the receiver object to differentiate calls that will be analyzed separately. Our approximation consists of the object, its local properties and their object values plus its chain of prototype objects. In Figure 6(b) we show the approximation corresponding to the object identity of object  $O_1$  in Figure 6(a). Note that the edges with the same local property name (annotated and not annotated) in the points-to graph are merged in the approximate object identity (e.g.,  $\langle O_1, p_1 \rangle, O_2$ ) and  $\langle O_1, p_1^* \rangle, O_3$  in Figure 6(a)). An object-sensitive analysis groups the calls using a receiver object created at the same allocation site and our identity-sensitive analysis more accurately groups the calls where receiver objects have the same approximate object identity. We intend to study the effects of using different object identity approximations as calling contexts in future work.

### 3.5 Block-sensitive Analysis

Our new points-to analysis algorithm is a fixed point calculation on the call graph, initialized with an empty points-to graph on entry to the JavaScript program, in which every constituent *IPBG* is traversed in a flow-sensitive manner. Essentially, we have designed the points-to algorithm to emphasize precision for the object identity information in the points-to graph and the *IPBG* to hide control flow not relevant to identity updates.

More specifically, our analysis solves for the points-to graph on exit of each *IPBG* node. The transfer function for a node in the *IPBG* is one of two kinds: (1) for a identity-update node perform strong update of the changed property, if possible (as in Table 2), or (2) for a identity-preserving node perform a flow-insensitive analysis of the statements in that node, using an initial points-to graph (*IN*) and storing the fixed point reached in points-to graph *OUT*.

(a)					(b)			
$\cup$	$\emptyset$	$v.p^d$	$v.p^*$	$v.p$	$\cup$	$\emptyset$	$o.p^*$	$o.p$
$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$o.p^*$	$o.p^*$
$v.p^d$	$\emptyset$	$v.p^d$	$v.p^*$	$v.p^*$	$o.p^*$	$o.p^*$	$o.p^*$	$o.p^*$
$v.p^*$	$\emptyset$	$v.p^*$	$v.p^*$	$v.p^*$	$o.p$	$o.p^*$	$o.p^*$	$o.p$
$v.p$	$\emptyset$	$v.p^*$	$v.p^*$	$v.p$				

**Table 3:** Union rules. (a) Access path edges union rules. (b) Property reference edges union rules.

Normally in a points-to analysis, we would form *IN* as a union of the *OUT* points-to graphs of predecessors of a node. In our algorithm, we need to maintain the invariant of our annotated property edges, namely that a property name without an annotation means that property exists and a property name with the *d* annotation means that property must not exist.

Table 3 shows the union rules for the access path edges and property reference edges when two points-to graphs are unioned. For the access path edges: (1) if access path  $v.p^\phi$  does not exist in at least one predecessor, then  $v.p^\phi$  does not exist after union; (2) if  $v.p^d$  or  $v.p$  exists in both predecessors, then  $v.p^d$  or  $v.p$  respectively exists after union; (3) otherwise,  $v.p^*$  exists after union. For the property reference edges: (1) if  $o.p$  exists in both predecessors, then  $o.p$  exists after union; (2) otherwise, if  $o.p$  or  $o.p^*$  exists in at least one predecessor, then  $o.p^*$  exists after union. These rules ensure analysis safety when property lookup is performed.

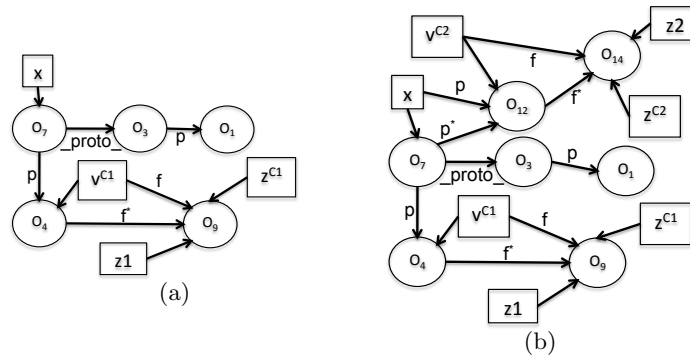
### 3.6 Implementation of Identity-sensitive Analysis in *JSBAF*

Our new points-to analysis was implemented with a client as the static component of the *JavaScript Blended Analysis Framework (JSBAF)*, a general-purpose analysis framework for JavaScript [27]. This framework was designed to strongly couple dynamic and static analyses to account for the effects of the dynamic features of JavaScript. *JSBAF* can be applied to analyze a JavaScript program (i.e., JavaScript code on a webpage) automatically in the presence of a good test suite. The dynamic phase gathers run-time information by executing tests. A trace of each test is collected, including call statements, object creations, variadic function calls with parameters, and dynamically created code. The implementation separates each trace into its constituent page traces. Each subtrace on a page is analyzed separately in the static phase. Data-flow solutions from different page subtraces are combined into a entire solution for that page.

Blended analysis uses only the observed calling structure as a basis to model the JavaScript program. Knowledge of unexecuted calls or object creations can be used to prune other unexecuted code sharing the same control dependence. For example, knowing  $Y3()$  is not called at line 5 in Figure 2, blended analysis prunes this unexecuted statement so that the imprecise node  $O_5$  and its connected edges will not be created. Thus, blended analysis is unsafe because not all executions are explored, but sound on the observed executions.

Our points-to algorithm was implemented on the *IBM T.J. Watson Libraries for Analysis (WALA)* open-source static analysis framework<sup>8</sup> which contains several existing static points-to analysis algorithms. *WALA* has been extended to enable blended analysis by providing dynamic information (i.e., a run-time collected call graph, dynamically generated code, object creation sites) [27].

We chose this implementation platform because blended analysis has been demonstrated to be more efficient and effective in analyzing real JavaScript programs.<sup>9</sup> Our algorithm takes as inputs the run-time collected calling structure (i.e., call graph<sup>10</sup>) and source code including dynamically generated code. Code pruning was performed on function bodies so that the code in polymorphic constructors and variadic functions was specialized. Hence, constructor polymorphism was handled by our improved object representation combined with dynamic information (i.e., objects created at the same allocation site with different sets of property names are represented as separate abstract objects).



**Fig. 7:** Blended identity-sensitive points-to analysis. (a) points-to graph at line 10. (b) points-to graph at line 15.

**Example.** In comparison to the inaccurate points-to solution of a flow- and context-insensitive analysis for the JavaScript code in Figure 2, we now demonstrate the results of our identity-sensitive points-to analysis in the context of blended analysis. Figures 7(a) and 7(b) show the points-to graphs obtained at lines 10 and 15, respectively. Because blended analysis executes the program and does not observe an object created by the constructor  $Y3$ , the code at line 5 is pruned so that our analysis does not generate the inaccurate node  $O_5$  nor the edge  $(\langle O_7, q \rangle, O_5)$ . For the call statement at line 10, our points-to analysis calculates the identity approximation of  $O_7$ , namely  $C1: \{O_7, p:O_4, \_proto_: O_3\}$ . Also, when looking up  $x.p$  at line 10, our algorithm returns  $O_4$  because there is

<sup>8</sup> <http://wala.sourceforge.net/>

<sup>9</sup> A static analysis was not able to finish analyzing most webpages in [27].

<sup>10</sup> Each node in the call graph is associated with the object creations observed during its execution.

no annotation on the property reference edge so that further lookup through the prototype chain is not necessary. Note that the points-to graph in Figure 7(a) is as precise as the run-time points-to graph (Figure 3(a)).

At line 12,  $x.p$  is strongly updated via the access path edge  $\langle x, p \rangle, O_{12}$ . For the call statement at line 15, our points-to analysis calculates the identity approximation of  $O_7$ ,  $C2: \{O_7, p:[O_4, O_{12}], \_proto.: O_3\}$ . Our points-to algorithm distinguishes this call site from line 10 because  $O_7$  has a different identity here. The lookup of  $x.p$  at line 15 follows the access path edge so that the node  $O_{12}$  is returned. Thus, in this example our analysis results in none of the inaccurate edges in the flow- and context-insensitive analysis (Figure 3(c)) and reflects the run-time behavior of JavaScript objects (Figure 3(b)).

## 4 Evaluation

In this section, we present experiments using *JSBAF* with our identity-sensitive points-to analysis compared to an existing points-to analysis [25], evaluating both with a *REF* client.

### 4.1 Experimental Design

**REF Analysis.** To evaluate the precision and performance of our points-to analysis, we implemented a JavaScript reference analysis (*REF*). The *REF* client calculates the set of objects returned by property lookup at a property read statement (i.e.,  $x = y.p$ ) or call statement (i.e.,  $x = y.p(\dots)$ ).<sup>11</sup> For each of these statements  $s$  in a function being analyzed in calling context  $c$ , we compute  $REF(s, c)$ , the set of objects returned by a property lookup for each  $o.p$  where  $o$  is pointed to by  $y$ . The cardinality of the *REF* set depends on the precision of the points-to graph and the property lookup operation; the smaller the set returned, the more useful for program understanding, for example.

In Figure 2, assume we add the function property

$$x.foo = function()\{var a = this.p; return a;\}$$

Effectively,  $foo()$  returns the property lookup result for  $this.p$ . If  $x.foo()$  is called at line 11 before the property update statement  $x.p=new A()$ , it will return  $O_4$ . If  $x.foo()$  is called at line 13 after  $x.p=new A()$ , it will return  $O_{12}$ . For an analysis that is flow-insensitive or that cannot distinguish these call sites by calling context, the return value of each of these function calls will contain at least two objects (i.e.,  $O_4$  and  $O_{12}$ ).

**Comparison with Points-to Analysis in [25].** We use the term *Corr* to refer to a blended version of correlation-tracking points-to analysis [25] (see Section 5 for more details) and its *REF* client. To demonstrate the additional

<sup>11</sup> All source code instances of property lookups (e.g.,  $return y.p$ ) occur as one of these two statements in the *WALA* intermediate code.

Website	<i>Corr</i>			<i>CorrBSIS</i>		
	1	2-4	$\geq 5$	1	2-4	$\geq 5$
facebook.com	38%	52%	10%	50%	47%	3%
google.com	32%	51%	17%	53%	42%	5%
youtube.com	41%	47%	12%	54%	41%	5%
yahoo.com	48%	46%	6%	52%	45%	3%
wikipedia.org	29%	45%	26%	43%	39%	18%
amazon.com	45%	52%	3%	46%	51%	3%
twitter.com	32%	53%	15%	39%	49%	12%
blogspot.com	35%	34%	31%	53%	36%	11%
linkedin.com	34%	49%	17%	44%	50%	6%
msn.com	40%	36%	24%	48%	37%	15%
ebay.com	30%	40%	30%	46%	40%	14%
bing.com	41%	34%	25%	54%	37%	9%
<b>Average</b>	<b>37%</b>	<b>45%</b>	<b>18%</b>	<b>48%</b>	<b>43%</b>	<b>9%</b>

**Table 4:** *REF* analysis precision.

precision of our analysis over *Corr*, we applied the correlation extraction transformation to our JavaScript benchmarks before performing our points-to analysis. We use the term *CorrBSIS* to refer to a blended version of this augmented new points-to analysis and its *REF* client. For each algorithm, an object property lookup returns a *REF* set whose cardinality  $|REF(s, c)|$  is calculated. For *Corr*, the *lookup\_all()* approximate algorithm described in Section 3.3 is used. For *CorrBSIS*, we use our optimized lookup algorithm *lookup()* in Procedure 1.

**Benchmarks.** We conducted the experiments with the benchmarks collected from 12 websites among the top 25 most popular sites on *alexa*, reusing website traces originally used in [27]. The results in [27] showed that the collected traces covered a large portion of the executable JavaScript code in those websites, including dynamically generated code. Although the benchmarks we used cover the most popular websites, it will require further investigation to determine how representative these benchmarks are of other websites. The experimental results were obtained on a 2.53GHz Intel Core 2 Duo MacBook Pro with 4GB memory running the Mac OS X 10.5 operating system.

## 4.2 Experimental Results

**Improved *REF* Precision.** Table 4 shows the *REF* client results for the 12 websites. Columns 2-4 present the results for *Corr* and columns 5-7 present the results for *CorrBSIS*. For each website, columns 2 & 5, 3 & 6, and 4 & 7 in Table 4 correspond to the percentage of property lookup statements that return 1 object, 2-4 objects, and more than 4 objects, respectively. The result shown for each website is averaged over the corresponding percentage numbers for all the webpages in that domain; for example, the 38% entry for *facebook.com* in column 2 is the average for *Corr* over the 27 webpages analyzed of the percentage of property lookup statements returning only 1 object.

Website	<i>Corr</i>	<i>CorrBSIS</i>	overhead
facebook	17.4	45.9	163%
google	13.0	30.4	134%
youtube	31.2	75.3	141%
yahoo	28.5	54.1	90%
wiki	16.0	40.1	151%
amazon	15.1	24.2	61%
twitter	38.1	94.5	148%
blog	15.9	42.4	137%
linkedin	27.8	62.0	167%
msn	34.4	57.9	68%
ebay	8.3	27.2	227%
bing	22.1	50.4	128%
<b>Average</b>	<b>22.3</b>	<b>50.3</b>	<b>135%</b>

**Table 5:** *REF* analysis cost (in seconds) on average per webpage.

Comparing columns 2-4 with 5-7 in Table 4 for each website, we see the relative precision improvement of *CorrBSIS* over *Corr*. For *REF* analysis, the best result is that the lookup returns only one object and the property lookup is more precise if the number of objects returned is smaller. On average over all the websites, *Corr* reported **37%** of the property lookup statements were resolved to a single object, while *CorrBSIS* improved this metric to **48%**, a *significant improvement*. In addition, *REF* analysis results may become too approximate to be useful if too many objects are returned. Although **18%** of the statements returned more than 4 objects for *Corr*, *CorrBSIS* reduced that number to **9%**. These improved precision results indicate the potential for greater practical use of identity-sensitive points-to information by client analyses.

We also investigated the average number of objects returned by a property lookup statement. For each website, we calculated the number of objects per statement on average over all its webpages. Over all the benchmarks, *Corr* produced on average **2.8** objects and *CorrBSIS* only reported on average **2.3** objects. Intuitively, this means that on average fewer objects at each property lookup statement must be examined to gain better understanding of the code.

***REF* Performance.** An analysis approach is practical if it scales to real-world programs, such as JavaScript code from actual websites. Because *CorrBSIS* is partially flow-sensitive and context-sensitive, it is important to demonstrate that this analysis is scalable. Table 5 shows the time performance of *Corr* versus *CorrBSIS*.<sup>12</sup> Columns 2 and 3 present the average webpage analysis time for each website, averaging over all of its webpages. Both *Corr* and *CorrBSIS* completely analyzed all the benchmark programs. On average over all the websites, *Corr* completely analyzed a webpage in **22.3** seconds, while *CorrBSIS* did so in **50.3**

<sup>12</sup> The time cost in Table 5 reflects the performance of the static phase of blended analysis. In the experiments, the dynamic phase of *Corr* and *CorrBSIS* is the same for both analyses. The work in [27] has demonstrated that the static phase dominates the blended analysis cost.

Website	No. of functions	% of functions w/ update(s)	% of identity-update stmt	No. of contexts
facebook	2123	9%	8%	4.0
google	1002	17%	6%	6.7
youtube	1329	7%	6%	3.9
yahoo	3810	4%	4%	2.4
wiki	270	10%	19%	4.8
amazon	729	6%	6%	1.9
twitter	618	15%	5%	3.4
blog	583	14%	14%	6.1
linkedin	920	8%	11%	3.6
msn	1537	8%	8%	2.8
ebay	581	18%	13%	7.5
bing	1131	7%	11%	4.9
<b>Average</b>	<b>1219</b>	<b>10%</b>	<b>9%</b>	<b>4.3</b>

**Table 6:** Benchmark and context statistics. (Total contexts per website is approximately column 2 times column 5.)

seconds, incurring an 135% average time overhead per webpage, acceptable for a research prototype implementation which has not been optimized.

**Discussion.** We collected data characterizing benchmark program structure and complexity to relate these characteristics to observed analysis precision and performance. The entries in Table 6 all represent averages per webpage that are averaged over an entire website. Column 2 shows the average number of functions in a JavaScript program. Column 3 shows the percentage of functions containing at least one identity-update statement. Column 4 shows the percentage of statements that are identity-update statements. Column 5 shows the number of contexts produced by *CorrBSIS* as a multiplier for column 2. On average over all the websites, **10%** of the functions contained local identity-update statement(s); these averages ranged from **4%** for *yahoo.com* to **18%** for *msn.com*. This suggests that *the identity-update statements are localized in a relatively small portion of the JavaScript program (e.g., in constructor functions)*. Manual inspection of several websites (i.e., *facebook*, *google* and *youtube*) revealed there were significant object behavior changes in JavaScript code outside of constructors. On average over all the websites, **9%** of the statements were identified as identity-update statements. The relatively small number of identity-update statements means that our *IPBG* contained many fewer nodes than the corresponding CFGs; therefore the flow-sensitive analysis was more practical in cost on the *IPBGs*.

Now we compare the analysis precision observed in Table 4 with the number of contexts generated on average per function per page (column 5 in Table 6) to observe the effect of identity sensitivity. *google*, *blog*, and *ebay* were the websites for which *CorrBSIS* improved precision the most, whereas *amazon*, *yahoos*, *twitter*, and *msn* were the websites for which *CorrBSIS* produced similar results to *Corr*. For the former websites, *CorrBSIS* generated the greatest number of contexts per function per webpage. For the latter websites, *CorrBSIS* generated the

fewest. We observe strong correlation between the precision gain and the number of contexts generated by *CorrBSIS*, demonstrating that *identity sensitivity significantly increased analysis precision on these benchmarks*, and suggesting that identity sensitivity will be an effective form of context sensitivity for JavaScript analysis.

As shown in Table 5, the *CorrBSIS* time overhead differed significantly for different websites, from 61% (*amazon.com*) to 227%(*ebay.com*). We investigate several program characteristics to reason about such differences. First, the *IPBGs* created by *CorrBSIS* determine the efficiency of the flow-sensitive analysis. On average over all the websites, an *IPBG* was comprised of about **6** nodes, explaining why *CorrBSIS* scaled on real websites. Functions with large numbers of nodes in their *IPBG* usually contained multiple identity-update statements and complex control flow. The largest number of nodes for an *IPBG* was **23** in *linkedin*. Second, the websites with the least performance overhead from *CorrBSIS* were *amazon*, *msn* and *yahoo*. These websites contained a relatively small percentage of update statements (i.e., all below average) and *CorrBSIS* generated the lowest number of contexts for them. The website that incurred the most overhead (i.e., *ebay*) contained **13%** update statements, (i.e., the third highest percentage in our benchmarks), and the greatest number of contexts per function (i.e., **7.5**) generated by *CorrBSIS*. These results support the reasoning that more complex block structure and more context comparisons contribute to the higher overhead for *CorrBSIS*.

## 5 Related Work

In this section due to space limitations, we present only the work most closely related to our identity-sensitive points-to algorithm.

**Related analyses of JavaScript.** Several approaches were proposed to analyze JavaScript programs. Sridharan *et al.* presented a points-to analysis for JavaScript that focused on handling correlated dynamic property accesses [25]. Correlated property accesses were identified and then extracted into a function. Using the property name as the calling context, points-to analysis tracking correlation was shown to be more precise and efficient than a field-sensitive Andersen’s points-to analysis. In our experiments, *CorrBSIS* was augmented by correlation analysis (i.e., *Corr*) demonstrating a significant improvement in the analysis precision.

Jensen *et al.* presented a static analysis that can precisely model prototype chains [12]. In their analysis, the *absent* set indicated potentially missing properties. The property edges annotated with \* play a similar role in our analysis. Jensen’s analysis is context-sensitive similar to 1-object-sensitivity used in Java [17]. The static flow-sensitive analysis presented in [12] was not scalable on large JavaScript programs, whereas our experiments showed the *CorrBSIS* was practical for blended analysis of real-world websites.

Several points-to analyses were proposed to handle other important challenges introduced by JavaScript. Guarnieri and Livshits designed a points-to

analysis to detect security and reliability issues in JavaScript widgets [8]. They used a subset of JavaScript language, *JavaScript<sub>SAFE</sub>*, that can be statically approximated. Guarnieri *et al.* presented a static taint analysis based on a points-to analysis finding security vulnerabilities in real-world websites [10]. The points-to algorithm focused on addressing features of JavaScript including object creations and accesses through constructed property names. In these analyses, prototyping was modeled as *lookup\_all* rather than our more accurate property lookup algorithm *lookup*. The points-to analysis in [10] and our algorithm are both implemented in *WALA*.

A hybrid analysis (i.e., a combination of static and dynamic analyses) is attractive when analyzing JavaScript programs. Chugh *et al.* presented a staged information flow analysis for JavaScript [4]. The approach analyzed the static code and incrementally analyzed the dynamically generated code. A similar approach was proposed by Guarnieri and Livshits [9]; their experiments studied the performance of the incremental analysis. In addition to collecting the dynamically generated/loaded code, blended analysis uses run-time information to make the analysis more precise (e.g., polymorphic constructors are distinguished via object initialization and some unexecuted code is pruned).

Several type-based points-to techniques have been proposed for JavaScript that support dynamic features such as prototype-based inheritance (e.g., [3, 5, 13]). It is difficult to compare our analysis with them as to practicality, because no empirical evidence on large JavaScript programs was presented.

Software tools supporting large JavaScript software including libraries are desirable. Schafer *et al.* provided an IDE support for JavaScript programming [21]. Points-to analysis was used to calculate code completion suggestions. Points-to analysis precision is crucial to determine the effectiveness of the tool. Madsen *et al.* presented a static analysis of JavaScript focusing on frameworks and libraries [16]. A novel use analysis was proposed to analyze libraries precisely and a points-to analysis was used to find aliases in the program. Our work is complementary to these techniques, in that more precise points-to results would make them more practical.

**Context-Sensitive Analysis.** Our identity-sensitive analysis is inspired by object sensitivity. Milanova *et al.* first introduced object sensitivity and implemented a object-sensitive points-to analysis for Java using a receiver object represented by its creation site as the calling context [17]. The experiments in [14] showed object sensitivity is the better choice as a calling context when analyzing an object-oriented language. Changes to object properties in JavaScript render object creation sites insufficient to represent object behavior, whereas identity sensitivity captures object behavior changes better.

Smaragdakis *et al.* formalized object sensitivity, summarizing its variations [24]. They introduced type sensitivity where object type was used as the calling context. For dynamically-typed languages like JavaScript, type is a run-time notion, encapsulated in the idea of object identity used as a calling context.

## 6 Conclusion

JavaScript object behavior is difficult to analyze well because of prototype-based inheritance and allowed changes to object properties during execution. In this paper, we introduced an identity-sensitive points-to analysis that models object behavior changes accurately by using a hierarchical program representation emphasizing identity-update statements, by defining identity sensitivity, a better context sensitivity mechanism for a dynamic language, and by enhancing the points-to graph representation for improving object property lookups. We implemented our new points-to algorithm as the static phase of a blended analysis in *JSBAF*. Experimental results on a *REF* client showed our analysis, *CorrB-SIS*, significantly improved the precision of a previous good JavaScript points-to analysis *Corr* [25]. For example, 48% of property lookups were resolved to a single object by our analysis versus 37% by *Corr*. Although our research prototype implementation incurred on average 135% overhead versus *Corr* on the popular website benchmarks used, further optimization will improve the performance, which is in the practical range.

In future work, we intend to investigate variations of identity sensitivity, to study the effects of different object identity approximations on analyzing JavaScript programs. We are interested in exploring the capability of identity-sensitive analysis to support program understanding. We also plan to generalize the proposed techniques to other dynamic programming languages.

## References

1. A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques and Tools*. Addison Wesley, 1986.
2. G. Balakrishnan and T. Reps. Recency-abstraction for heap-allocated storage. In *Proceedings of the 13th international conference on Static Analysis, SAS'06*, pages 221–239, 2006.
3. R. Chugh, D. Herman, and R. Jhala. Dependent types for JavaScript. In *Proceedings of the ACM international conference on Object oriented programming systems languages and applications, OOPSLA '12*, pages 587–606, 2012.
4. R. Chugh, J. A. Meister, R. Jhala, and S. Lerner. Staged information flow for JavaScript. In *Proceedings of the 2009 ACM SIGPLAN conference on Programming language design and implementation, PLDI '09*, pages 50–62, 2009.
5. R. Chugh, P. M. Rondon, and R. Jhala. Nested refinements: a logic for duck typing. In *Proceedings of the 39th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages, POPL '12*, pages 231–244, 2012.
6. A. De and D. D'Souza. Scalable flow-sensitive pointer analysis for Java with strong updates. In *Proceedings of the 26th European conference on Object-Oriented Programming, ECOOP'12*, pages 665–687, 2012.
7. D. Flanagan. *JavaScript: The Definitive Guide*. O'Reilly Media, Inc., 2006.
8. S. Guarnieri and B. Livshits. Gatekeeper: mostly static enforcement of security and reliability policies for JavaScript code. In *Proceedings of the 18th conference on USENIX security symposium, SSYM'09*, pages 151–168, 2009.

9. S. Guarnieri and B. Livshits. Gulfstream: staged static analysis for streaming JavaScript applications. In *Proceedings of the 2010 USENIX conference on Web application development*, WebApps'10, pages 6–6, 2010.
10. S. Guarnieri, M. Pistoia, O. Tripp, J. Dolby, S. Teilhet, and R. Berg. Saving the world wide web from vulnerable JavaScript. In *Proceedings of the 2011 International Symposium on Software Testing and Analysis*, ISSTA '11, pages 177–187, 2011.
11. P. Heidegger and P. Thiemann. Recency types for analyzing scripting languages. In *Proceedings of the 24th European conference on Object-oriented programming*, ECOOP'10, pages 200–224, 2010.
12. S. H. Jensen, A. Møller, and P. Thiemann. Type analysis for JavaScript. In *Proceedings of the 16th International Symposium on Static Analysis*, SAS '09, pages 238–255, 2009.
13. B. S. Lerner, P. Joe Gibbs, A. Guha, and K. Shriram. TeJaS: Retrofitting type systems for JavaScript. In *Proceedings of the 9th symposium on Dynamic languages*, DLS '13, 2013.
14. O. Lhoták and L. Hendren. Context-sensitive points-to analysis: is it worth it? In *Proceedings of the 15th international conference on Compiler Construction*, CC'06, pages 47–64, 2006.
15. H. Lieberman. Using prototypical objects to implement shared behavior in object-oriented systems. In *Conference proceedings on Object-oriented programming systems, languages and applications*, OOPSLA '86, pages 214–223, 1986.
16. M. Madsen, B. Livshits, and M. Fanning. Practical static analysis of JavaScript applications in the presence of frameworks and libraries. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, ESEC/FSE 2013, pages 499–509, 2013.
17. A. Milanova, A. Rountev, and B. G. Ryder. Parameterized object sensitivity for points-to analysis for Java. *ACM Trans. Softw. Eng. Methodol.*, 14(1):1–41, Jan. 2005.
18. Orion. <http://www.eclipse.org/orion/>.
19. G. Richards, S. Lebresne, B. Burg, and J. Vitek. An analysis of the dynamic behavior of JavaScript programs. In *Proceedings of the 2010 ACM SIGPLAN conference on Programming language design and implementation*, PLDI '10, pages 1–12, 2010.
20. B. G. Ryder. Dimensions of precision in reference analysis of object-oriented programming languages. In *Proceedings of the 12th international conference on Compiler construction*, CC'03, pages 126–137, 2003.
21. M. Schafer, M. Sridharan, J. Dolby, and F. Tip. Effective smart completion for JavaScript. Technical Report RC25359, IBM, 2013.
22. R. Sethi. *Programming Languages, Concepts & Constructs, 2nd edition*. Addison Wesley, 1996.
23. M. Sharir and A. Pnueli. Two approaches to interprocedural data flow analysis. *Program Flow Analysis: Theory and Applications*, pages 189–234, 1981.
24. Y. Smaragdakis, M. Bravenboer, and O. Lhoták. Pick your contexts well: understanding object-sensitivity. In *Proceedings of the 38th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, POPL '11, pages 17–30, 2011.
25. M. Sridharan, J. Dolby, S. Chandra, M. Schäfer, and F. Tip. Correlation tracking for points-to analysis of JavaScript. In *Proceedings of the 26th European conference on Object-Oriented Programming*, ECOOP'12, pages 435–458, 2012.

26. P. Wegner. Dimensions of object-based language design. In *Conference proceedings on Object-oriented programming systems, languages and applications*, OOPSLA '87, pages 168–182, 1987.
27. S. Wei and B. G. Ryder. Practical blended taint analysis for JavaScript. In *Proceedings of the 2013 International Symposium on Software Testing and Analysis*, ISSTA 2013, pages 336–346, 2013.