# Analyzing JavaScript Web Applications

## Anders Møller

**Aarhus University**

Joint work with Simon Holm Jensen, Peter Thiemann, and Magnus Madsen

# JavaScript:  the lingua franca of Web 2.0

# JavaScript is a *dynamic language*

- Object-based
- Prototype-based inheritance
- First-class functions, closures
- Runtime types
- …

## NO STATIC TYPE CHECKING

# How JavaScript was designed

"I hacked the JS prototype
in ~1 week. And it showed!
Mistakes were frozen early."

– Brendan Eich, inventor of JavaScript

# JavaScript

# Static type analysis to the rescue!

```
var x = {
    a: 42,
    b: function() {return a;}
}
print(x.b())
```

an *abstract state* for each program point,
analyze it further to detect likely programming errors

# **<u>Potential</u> programming errors**

1. invoking a non-function value (e.g. undefined) as a function
2. reading an absent variable
3. accessing a property of `null` or `undefined`
4. reading an absent property of an object
5. writing to variables or object properties that are never read
6. calling a function object both as a function and as a constructor, or passing function parameters with varying types
7. calling a built-in function with an invalid number of parameters, or with a parameter of an unexpected type

etc…

# Flow of control and data can be subtle

```
function Person(n) {
    this.setName(n);
    Person.prototype.count++;
}
Person.prototype.count = 0;
Person.prototype.setName = function(n) { this.name = n; }
function Student(n,s) {
    this.b = Person;
    this.b(n);
    delete this.b;
    this.studentid = s.toString();
}
Student.prototype = new Person;

var t = 100026.0;
var x = new Student("Joe Average", t++);
var y = new Student("John Doe", t);
y.setName("John Q. Doe");
```

declares a "class" named Person

declares a "static field" named count

declares a shared method named setName

declares a "sub-class" named Student

creates two Student objects…

does y have a `setName` method at this program point?

# An abstract state
# (as produced by our analysis)

# General and widely used approaches

- *The monotone framework*
  [Kam & Ullman '77]

- *The functional approach*
  [Sharir & Pnueli '81]

- *IFDS*
  [Reps, Horwitz, and Sagiv '95]

- ...

# Our approach

- Abstract interpretation using the monotone framework

- The recipe:
  1. construct a **control flow graph** for the program to be analyzed

  2. define an appropriate **dataflow lattice** (abstraction of data)

  3. define **transfer functions** (abstraction of operations)

# Control flow graphs

- declare-variable[$x$]

- read-variable[$x,v$]

- write-variable[$v,x$]

- constant[$c,v$]

- read-property[$v_{obj},v_{property},v_{result}$]

- write-property[$v_{obj},v_{property},v_{value}$]

- delete-property[$v_{obj},v_{property},v_{result}$]

- if[$v$]

- entry[$f,x_1,...,x_n$], exit, exit-exc

- call[$w,v_0,...,v_n$], construct[$w,v_0,...,v_n$], after-call[$v$]

- return[$v$]

- throw[$v$],  catch[$x$]

- <op>[$v_1,v_2$], <op>[$v_1,v_2,v_3$]

- ...

- Convenient representation of JavaScript programs

- *Nodes* describe primitive instructions, *edges* describe control-flow

- Each *x* is a program variable

- Each *v* is a temporary variable (i.e. a register)

12

# Analysis lattice

the analysis lattice

abstract states

abstract objects

abstract values

# Abstract values

object labels
(allocation sites)

$$\text{Value} = \text{Undef} \times \text{Null} \times \text{Bool} \times \text{Num} \times \text{String} \times \mathcal{P}(L)$$

$$\text{Undef} = \begin{array}{c} \text{undef} \\ | \\ \bot \end{array} \qquad \text{Null} = \begin{array}{c} \text{null} \\ | \\ \bot \end{array} \qquad \text{Bool} = \begin{array}{c} \text{bool} \\ \diagup \diagdown \\ \text{true} \quad \text{false} \\ \diagdown \diagup \\ \bot \end{array}$$

$$\text{Num} = \begin{array}{c} \text{Num} \\ \text{INF} \quad \text{UInt} \quad \text{NotUInt} \\ -\text{Inf} \ +\text{Inf} \ \text{NaN} \ 0 \ \cdots \ 4294967295 \quad \cdots -42 \ -1.87 \ 1.2 \ \cdots \\ \bot \end{array}$$

$$\text{String} = \begin{array}{c} \text{string} \\ \text{UIntString} \quad \text{NotUIntString} \\ \text{"0"}\cdots\text{"4294967295"} \ \text{"foo"}\cdots\text{"bar"} \\ \bot \end{array}$$

Example:  $(\bot, \text{null}, \text{true}, 42.0, \bot, \{\ell_7, \ell_9\})$

# Abstract objects

property names including [[Prototype]]

$$Obj = (\mathcal{P} \hookrightarrow Value \times Absent \times Attributes \times Modified) \times \mathcal{P}(ScopeChain)$$

(explained later)

describes the

[[Scope]] property

$$Absent = \begin{array}{c} absent \\ | \\ \bot \end{array} \qquad Modified = \begin{array}{c} modified \\ | \\ \bot \end{array}$$

$$Attributes = ReadOnly \times DontDelete \times DontEnum$$

$$ReadOnly = \begin{array}{c} \top \\ RO \quad notRO \\ \bot \end{array} \qquad DontDelete = \begin{array}{c} \top \\ DD \quad notDD \\ \bot \end{array} \qquad DontEnum = \begin{array}{c} \top \\ DE \quad notDE \\ \bot \end{array}$$

# Abstract states

heap

$$\text{State} = (L \hookrightarrow \text{Obj}) \times \text{Stack} \times \mathcal{P}(L) \times \mathcal{P}(L)$$

(explained later)

temporary variables

the current
activation record

stack-reachable
objects

$$\text{Stack} = (T \rightarrow \text{Value}) \times \mathcal{P}(\text{ExecutionContext}) \times \mathcal{P}(L)$$
$$\text{ExecutionContext} = \text{ScopeChain} \times L \times L$$
$$\text{ScopeChain} = L^*$$

the `this` object

for resolving variables

the variable object,
for variable declarations

16

# The analysis lattice

$$AnalysisLattice = C \times N \rightarrow State$$

contexts

(for context sensitivity)

the flow graph nodes

# Transfer functions

Example:   **read-property**$[v_{obj}, v_{property}, v_{result}]$

1. Coerce $v_{obj}$ to objects

2. Coerce $v_{property}$ to strings

3. Descend the object prototype chains
   (using the [[Prototype]] property)
   to find the relevant properties

4. Join the property values

5. Assign the result to $v_{result}$

# Weak vs. strong updates

- For a **write-property**$[v_{obj}, v_{property}, v_{value}]$ node, $v_{obj}$ refers to one or more abstract objects (identified by their allocation sites)

- Each abstract object generally describes *multiple* concrete objects

- So **write-property** must conservatively be modeled by *joining* $v_{value}$ into the existing value of $v_{property}$ at $v_{obj}$ (i.e. a *weak update*)

- This is bad for precision!

- *Strong update* (overwriting instead of joining) is possible whenever the abstract object is known to represent a single concrete object

# Recency abstraction

- For each allocation site $\ell$
  maintain **two** abstract objects:
  - $\ell$ @ corresponds to the *most recently*
      allocated object originating from $\ell$
  - $\ell$ * older objects from $\ell$
- $\ell$ @ always describes at most one concrete
  object and hence permits strong updating!
- To make this work, we just need some extra
  bookkeeping in the transfer functions

# Interprocedural analysis with *maybe-modified*



At function exits, **restore unmodified parts of the heap (and the stack)** from the call node

# Observing redundancy

```
...
TaskControlBlock.prototype.markAsRunnable = function () {
  this.state = this.state | STATE_RUNNABLE;
};
...
```

- Why is this function (from `richards.js`, V8) visited **18** times by the analyzer???

- Mostly, new dataflow that arrives at the function entry (and triggers re-analysis) is irrelevant to the function body!

# Lazy propagation

- Defer propagation of field values that are not known to be relevant to the current function

- Use a placeholder value:  *unknown*

- When analyzing a function, assume initially that no fields are referenced

- When a field is referenced, recover its proper value

$\Rightarrow$ irrelevant dataflow isn't propagated

$\Rightarrow$ *unknown* implies *unmodified*

# An example



- Each •
  represents an
  abstract state

- For simplicity,
  no context
  sensitivity here

# Formalization of lazy propagation

*How do we express the idea
more concisely and formally?*

(necessary for reasoning about its properties
and for obtaining a good implementation)

1) Start with a basic analysis framework where
   transfer functions are expressed via an
   **abstract data type** (ADT)

2) Introduce lazy propagation by a systematic
   modification of the ADT **(without touching
   the transfer functions!)**

# The basic lattice (simplified)

object labels (allocation sites)     functions

$$\text{Value} = \mathcal{P}(L) \times \mathcal{P}(F) \times \text{Base}$$

$$\text{Obj} = P \rightarrow \text{Value}$$

property names (fields)

$$\text{State} = L \rightarrow \text{Obj}$$

$$\text{CallGraph} = \mathcal{P}(C \times N \times C \times F)$$

set of call edges

$$\text{AnalysisLattice} = (C \times N \rightarrow \text{State}) \times \text{CallGraph}$$

contexts     nodes (primitive statements)

# AnalysisLattice **as an abstract data type (ADT)**

reads a field value

$$getfield : C \times N \times L \times P \rightarrow \mathsf{Value}$$

$$getcallgraph : () \rightarrow \mathsf{CallGraph}$$ reads the call graph

reads an abstract state

$$getstate : C \times N \rightarrow \mathsf{State}$$

$$propagate : C \times N \times \mathsf{State} \rightarrow ()$$ intra-procedural flow

$$funentry : C \times N \times C \times F \times \mathsf{State} \rightarrow ()$$

$$funexit : C \times N \times C \times F \times \mathsf{State} \rightarrow ()$$

inter-procedural flow

- The transfer functions can only access the AnalysisLattice element through these operations
- (We'll skip their definitions here…)

# Introducing lazy propagation

## – a systematic modification of the lattice and the ADT operations

property values can now be "unknown"!

$$Obj = P \to (Value\downarrow_{unknown})$$

$$CallGraph = C \times N \times C \times F \to (State\downarrow_{none})$$

$$AnalysisLattice = (C \times N \to (State\downarrow_{none})) \times CallGraph$$

now distinguishing between the *unreachable* state and the all-*unknown* state

each call edge is now labelled with an abstract state

$$recover : C \times N \times L \times P \to Value$$

used for recovering "unknown" values

28

# *getfield'* (read a field)

$a.\mathit{getfield'}(c \in C,\, n \in N,\, \ell \in L,\, p \in P):$

$\vdots$

    $v := a.\mathit{getfield}(c, n, \ell, p)$    ←   *getfield* from the basic framework

    **if** $v =$ unknown **then**

        *// the field value has been reduced to* unknown*, so recover the real value*

        $v := a.\mathit{recover}(c, n, \ell, p)$

    **end if**    ←   call *recover* if the value is "unknown"

    **return** $v$

$\vdots$

# *funentry'* (flow at function entry)

$a.funentry'(c_1 \in C,\ n_1 \in N,\ c_2 \in C,\ f_2 \in F,\ s \in \textsf{State})$:
  **let** $(m, g) = a$ **and** $u = m(c_2, entry(f_2))$

    $\vdots$

  // *introduce* unknown *field values*
  $s' := \bot_{\textsf{State}}$
  **if** $u \neq$ none **then**
    **for all** $\ell \in L, p \in P$ **do**
      **if** $u(\ell)(p) \neq$ unknown **then**
        // *the field has been referenced*
        $s'(\ell)(p) := s(\ell)(p)$
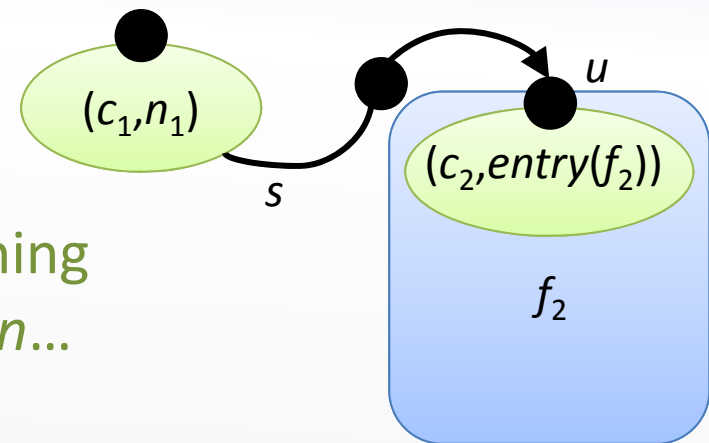      **end if**
    **end for**
  **end if**
  // *propagate the resulting state into the function entry*
  $a.propagate'(c_2, entry(f_2), s')$
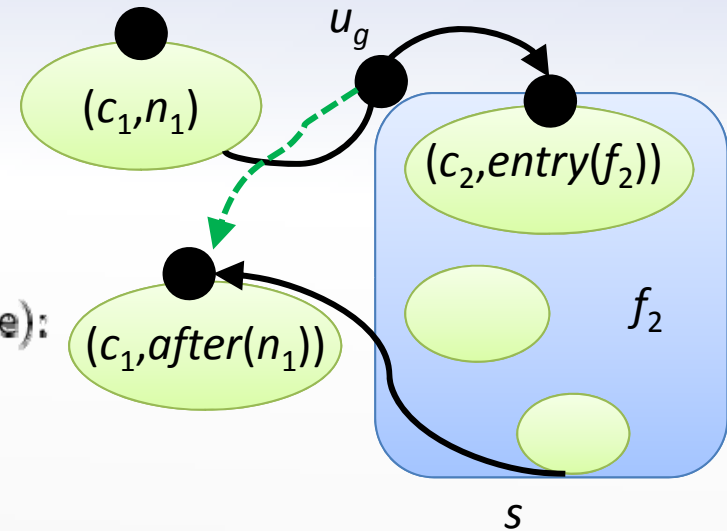
    $\vdots$

set everything
to *unknown…*

… except fields that are
known to be referenced

join *s'* into the function entry state

$(c_1, n_1)$

$s$

$(c_2, entry(f_2))$

$u$

$f_2$

# *funexit'* (flow at function exit)



$a.funexit'(c_1 \in C, n_1 \in N, c_2 \in C, f_2 \in F, s \in \mathbf{State})$:
  **let** $(\_, g) = a$ **and** $u_g = g(c_1, n_1, c_2, f_2)$
  $s' := \bot_{\mathbf{State}}$
  **for all** $\ell \in L, p \in P$ **do**
    **if** $s(\ell)(p) = $ unknown **then**
      *// the field has not been accessed, so restore its value from the call edge state*
      $s'(\ell)(p) := u_g(\ell)(p)$
    **else**
      $s'(\ell)(p) := s(\ell)(p)$
    **end if**
  **end for**
  $a.propagate'(c_1, after(n_1), s')$

*unknown*
implies
not modified within the function

join *s'* into the node after the call

31

# Theoretical properties of lazy prop.

- **Precision** is at least as good as before

- **Soundness** (wrt. language semantics)
  is preserved

- **Recovery** does not affect amortized complexity

- **Number of fixpoint iterations** increases in some situations and decreases in other

# Experiments

- >200 **small test cases**, to get into the obscure corner cases of JavaScript

- A few larger benchmarks: **Google's V8 benchmark suite** (500-1800 lines of code)

- Also tested on the **SunSpider benchmarks**

# Experiments

Some results for `richards.js` from V8:

- the analysis guarantees for **95%** of the call/construct instructions that they always succeed

- **1** location where an absent variable is read, with **0** spurious warnings

- **93%** of all read/write/delete-property operations will never attempt to coerce null or undefined into an object

- **6** functions dead (guaranteed unreachable)

# Experimental results

| | LOC | Blocks | Iterations | | | Time (seconds) | | | Memory (MB) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | *lazy* | *basic+* | *basic* | *lazy* | *basic+* | *basic* | *lazy* | *basic+* | *basic* |
| richards.js | 529 | 478 | 1399 | 2782 | 2663 | 3.8 | 4.6 | 5.6 | 3.7 | 6.4 | 11.05 |
| benchpress.js | 463 | 710 | 5097 | 12581 | 18060 | 5.4 | 13.4 | 33.2 | 7.8 | 24.0 | 42.02 |
| delta-blue.js | 853 | 1054 | 63611 | $\infty$ | $\infty$ | 136.7 | $\infty$ | $\infty$ | 140.5 | $\infty$ | $\infty$ |
| cryptobench.js | 1736 | 2857 | 17213 | 43848 | $\infty$ | 22.1 | 99.4 | $\infty$ | 42.8 | 127.9 | $\infty$ |
| 3d-cube.js | 342 | 545 | 2009 | 4147 | 7116 | 4.0 | 5.3 | 14.1 | 6.2 | 10.6 | 18.4 |
| 3d-raytrace.js | 446 | 575 | 6749 | 30323 | $\infty$ | 8.2 | 24.8 | $\infty$ | 10.1 | 16.7 | $\infty$ |
| crypto-md5.js | 296 | 392 | 646 | 1004 | 5358 | 1.8 | 2.0 | 4.5 | 2.7 | 3.6 | 6.1 |
| access-nbody.js | 179 | 149 | 317 | 523 | 551 | 1.0 | 1.3 | 1.8 | 0.9 | 1.7 | 3.2 |

$\infty$ means >512MB

*basic:* naive monotone framework

*basic+:* basic extended with *maybe-modified* (and *copy-on-write*)

*lazy:* basic extended with *lazy propagation*

# Summary

★ ***Static analysis*** is a useful tool for reasoning about programs written in a scripting language such as JavaScript

★ ***Lazy propagation*** ensures that only relevant information is propagated from one function to another

    – reduces the amount of data being propagated

    – may improve precision:  non-referenced fields respect interprocedurally realizable paths