

The Paradox of Flow Analysis

Or: What We Talk About
When We Talk About Higher-Order Flow Analysis

David Van Horn



On Flow Analysis (*k*CFA):

- ▶ Background
- ▶ Complexity
- ▶ Paradox
- ▶ Resolution
- ▶ Exploitation

Part I

Background

Q: What is program analysis?

A: Prediction of the future.

underpins:

verification, optimization, transformation, understanding, ...

global register allocation, global constant subexpression elimination, loop invariant detection, redundant assignment detection, dead code elimination, constant propagation, range analysis, code hoisting, induction variable elimination, copy propagation, live variable analysis, loop unrolling, loop jamming, type inference

Q: Which are higher-order programming languages?

A: Most of them.

Object-oriented: Java, Smalltalk, C#, C++, JavaScript, Eiffel, Scala, Ruby, Perl, Clojure, Python, ...

Functional: Common Lisp, Erlang, F#, Haskell, JavaScript, Scheme, SML, OCaml, ...

Q: What is a higher-order programming language (HOPL)?

A: Computation as values.

```
var cos = deriv(sin)
```

```
[cos, sin, sqr].map(function (f) { return f(5) })
```

Q: What makes program analysis difficult in a HOPL?

A: Mutual recursion between data and control.

- ▶ To do data-flow analysis, we need a control-flow analysis.
- ▶ To do control-flow analysis, we need a data-flow analysis.

```
var j = [cos, sin, sqr] [i] (5)
```

Evaluation:

```
//  $(\mathbb{R} \rightarrow \mathbb{R}) \rightarrow (\mathbb{R} \rightarrow \mathbb{R})$   
// Produce the numerical derivative of f.  
function deriv(f) {  
  function (x) { (f(x+ $\epsilon$ ) + f(x- $\epsilon$ )) / 2* $\epsilon$ ; }  
};  
  
deriv(sqr);    //=> ⟨function (x) {..}, [f ↦ sqr]⟩  
deriv(sin)    //=> ⟨function (x) {..}, [f ↦ sin]⟩
```

Idealized JavaScript:

$$E ::= x \qquad x \in \{x, y, \dots\}$$
$$\begin{array}{l} | \text{function } (x) \{ E \} \\ | E; E \\ | E(E) \end{array}$$

Evaluation:

$ev : \text{Exp} \times \text{Env} \rightarrow \text{Val}$

$ev(x, \rho) = \rho(x)$

$ev(\text{function } (x) \{E\}, \rho) = \langle \text{function } (x) \{E\}, \rho \rangle$

$ev(E_0; E_1, \rho) = ev(E_0, \rho); ev(E_1, \rho)$

$ev(E_0(E_1), \rho) =$

let $ev(E_0, \rho) = \langle \text{function } (x) \{E_0\}, \rho_0 \rangle$ in

let $ev(E_1, \rho) = v$ in

$ev(E_0, \rho_0[x \mapsto v])$

Labeled Idealized JavaScript:

$$\begin{array}{l} E^l ::= x^l \\ \quad | \text{function}^l(x) \{ E^l \} \\ \quad | E^l; E^l \\ \quad | E^l(E^l)^l \end{array} \quad x \in \{x, y, \dots\}$$

Abstract Evaluation (OCFA):

Computes a cache C that maps locations ℓ and variables x to sets of “abstract” values, i.e. function terms.

$$\begin{aligned} C(\ell) &= \{ \text{function } (x) \{..\}, \\ &\quad \text{function } (y) \{..\} \} \\ C(x) &= \{ \text{function } (x) \{..\}, \\ &\quad \text{function } (y) \{..\} \} \end{aligned}$$

Soundness:

- ▶ if E^ℓ evaluates to $\langle f, \rho \rangle$, then $f \in C(\ell)$,
- ▶ if x is bound to $\langle f, \rho \rangle$, then $f \in C(x)$.

Abstract Evaluation (OCFA):

```
//  $\forall X . X \rightarrow X$   
// The identity function.  
function id(x) {  
    x  
};  
  
id(sqr);           //=> { sqr }  
id(sin)           //=> { sqr, sin }
```

Abstract Evaluation (OCFA):

```
//  $\forall X, Y. (X \rightarrow Y) \rightarrow (X \rightarrow Y)$ 
// The identity function (for functions).
function id(f) {
  function (x) { f(x) }
};

id(sqr);           //=> { function (x) {..} }
id(sin)            //=> { function (x) {..} }
```

Abstract Evaluation (OCFA):

```
//  $\forall X, Y. (X \rightarrow Y) \rightarrow (X \rightarrow Y)$   
// The identity function (for functions).  
function id(f) {  
  function (x) { f(x) }  
};  
  
id(sqr)(5); //=> { sqr(5) }  
id(sin)(7) //=> { sin(7), sin(5), sqr(5), sqr(7) }
```

Abstract Evaluation (OCFA):

$av : \text{Exp} \times \text{Cache} \rightarrow \text{Cache}$

$av(x^\ell) = C(x) \subseteq C(\ell)$

$av(\text{function}^\ell(x) \{E\}) = \{\text{function}(x) \{E\}\} \subseteq C(\ell)$

$av(E^{\ell_1}; E^{\ell_2}) = av(E^{\ell_1}); av(E^{\ell_2})$

$av(E^{\ell_1}(E^{\ell_2})^\ell) = av(E^{\ell_1}); av(E^{\ell_2});$

$\forall(\text{function}(x) \{E^{\ell_0}\}) \in C(\ell_1) :$

$C(\ell_2) \subseteq C(x);$

$av(E^{\ell_0});$

$C(\ell_0) \subseteq C(\ell)$

The idea of *k*CFA:

Use calling contexts to make distinctions during analysis.

```
(function (f) { f( $E_1$ )1; f( $E_2$ )2 })
```

```
(function (x) {  $E_0$  })
```

E_0 will be analyzed twice:

- ▶ once with x meaning E_1 ,
- ▶ once with x meaning E_2 .

Abstract Evaluation (1CFA):

```
//  $\forall X, Y. (X \rightarrow Y) \rightarrow (X \rightarrow Y)$ 
// The identity function (for functions).
function id(f) {
  function (x) { f(x) }
};

id(sqr)1;           //=> {⟨function (x) {..}, [f ↦ 1]⟩ }
id(sin)2           //=> {⟨function (x) {..}, [f ↦ 2]⟩ }
                    //   C(f,1) = { sqr }
                    //   C(f,2) = { sin }
```

Abstract Evaluation (1CFA):

```
//  $\forall X, Y. (X \rightarrow Y) \rightarrow (X \rightarrow Y)$   
// The identity function (for functions).  
function id(f) {  
  function (x) { f(x) }  
};  
  
id(sqr)1(5)3; //=> { sqr(5) }  
id(sin)2(7)4 //=> { sin(7) }
```

Abstract Evaluation (*kCFA*):

Cache C maps locations ℓ and variables x with contours δ to sets of “abstract” *closures*, i.e. function, contour environment pairs.

$$\begin{aligned} C(\ell, \delta) &= \{ \langle \text{function } (x) \{ \dots \}, \phi_0 \rangle, \\ &\quad \langle \text{function } (y) \{ \dots \}, \phi_1 \rangle \} \\ C(x, \delta) &= \{ \langle \text{function } (x) \{ \dots \}, \phi_0 \rangle, \\ &\quad \langle \text{function } (y) \{ \dots \}, \phi_1 \rangle \} \end{aligned}$$

Soundness:

- ▶ if E^ℓ evaluates in context δ to $\langle f, \rho \rangle$, then $\langle f, \phi \rangle \in C(\ell, \delta)$, where $\rho \sqsubseteq \phi$,
- ▶ if x is bound in context δ to $\langle f, \rho \rangle$, then $\langle f, \phi \rangle \in C(x, \delta)$, where $\rho \sqsubseteq \phi$.

Instrumented Evaluation (∞ CFA):

$av : \text{Exp} \times \text{CEnv} \times \text{Contour} \times \text{Cache} \rightarrow \text{Cache}$

$$av(x^\ell, \phi, \delta) = C(x, \phi(x)) \subseteq C(\ell, \delta)$$

$$av(\text{function}^\ell(x) \{E\}, \phi, \delta) = \{\langle \text{function}(x) \{E\}, \phi \rangle\} \subseteq C(\ell, \delta)$$

$$av(E^{\ell_1}; E^{\ell_2}, \phi, \delta) = av(E^{\ell_1}, \phi, \delta); av(E^{\ell_2}, \phi, \delta)$$

$$av(E^{\ell_1} (E^{\ell_2})^\ell, \phi, \delta) = av(E^{\ell_1}, \phi, \delta); av(E^{\ell_2}, \phi, \delta);$$

$$\forall (\langle \text{function}(x) \{E^{\ell_0}\}, \phi_0 \rangle) \in C(\ell_1, \delta) :$$

$$C(\ell_2, \delta) \subseteq C(x, \delta\ell);$$

$$av(E^{\ell_0}, \phi_0[x \mapsto \delta\ell], \delta\ell);$$

$$C(\ell_0, \delta\ell) \subseteq C(\ell, \delta)$$

Abstract Evaluation (kCFA):

$av : \text{Exp} \times \text{CEnv} \times \text{Contour} \times \text{Cache} \rightarrow \text{Cache}$

$$av(x^\ell, \phi, \delta) = C(x, \phi(x)) \subseteq C(\ell, \delta)$$

$$av(\text{function}^\ell(x) \{E\}, \phi, \delta) = \{\langle \text{function}(x) \{E\}, \phi \rangle\} \subseteq C(\ell, \delta)$$

$$av(E^{\ell_1}; E^{\ell_2}, \phi, \delta) = av(E^{\ell_1}, \phi, \delta); av(E^{\ell_2}, \phi, \delta)$$

$$av(E^{\ell_1} (E^{\ell_2})^\ell, \phi, \delta) = av(E^{\ell_1}, \phi, \delta); av(E^{\ell_2}, \phi, \delta);$$

$$\forall (\langle \text{function}(x) \{E^{\ell_0}\}, \phi_0 \rangle) \in C(\ell_1, \delta) :$$

$$C(\ell_2, \delta) \subseteq C(x, \llbracket \delta \ell \rrbracket_k);$$

$$av(E^{\ell_0}, \phi_0[x \mapsto \llbracket \delta \ell \rrbracket_k], \llbracket \delta \ell \rrbracket_k);$$

$$C(\ell_0, \llbracket \delta \ell \rrbracket_k) \subseteq C(\ell, \delta)$$

“It did not take long to discover that the basic analysis, for any $k > 0$, was intractably slow for large programs.

In the ensuing years, researchers expended a great deal of effort deriving clever ways to tame the cost of the analysis.”

Olin Shivers, Higher-order control-flow analysis in retrospect:
Lessons learned, lessons abandoned (2004)

Part II

Complexity

- ▶ 0CFA: Linearity and Evaluation
- ▶ k CFA: Non-Linearity and Approximation

Theorem. OCFA is in PTIME.

Proof: C grows monotonically. Its size is $O(n^2)$. □

$$C(\ell_0) = \{\text{function}^{\ell_i} \{..\}, \dots, \text{function}^{\ell_j} \{..\}\}$$

⋮

$$C(\ell_n) = \{\text{function}^{\ell_i} \{..\}, \dots, \text{function}^{\ell_j} \{..\}\}$$

$$C(x_p) = \{\text{function}^{\ell_i} \{..\}, \dots, \text{function}^{\ell_j} \{..\}\}$$

⋮

$$C(x_q) = \{\text{function}^{\ell_i} \{..\}, \dots, \text{function}^{\ell_j} \{..\}\}$$

Theorem. 0CFA is PTIME-hard.

Proof (sketch):

- ▶ Analysis and evaluation coincide for linear programs.
- ▶ Linear programs can express circuits.
- ▶ The circuit value problem is PTIME-hard.



function (x) {... x ...} // x occurs *exactly* once

Observe, in a *linear* program:

- ▶ each function can be applied to at most one argument,
- ▶ each variable can be bound to at most one value.

Therefore:

Analysis of a linear program coincides exactly with its evaluation.

Theorem. 1CFA is in EXPTIME.

Proof: C grows monotonically. Its size is $O(2^n)$.



Theorem. 1CFA is EXPTIME-hard.

Proof: ...

Hardness proof of k CFA relies on two insights:

1. Terms are approximated by an exponential number of closures.
2. *Approximation* engenders reevaluation which provides *power*.

Many closures can flow to a single program point:

```
function (z) { z(x1)(x2)... (xn) }
```

- ▶ n free variables
- ▶ an *exponential* number of environments map these variables to contours of length 1 in 1CFA.

```
(function (f0) { f0(0); f0(1) })  
(function (x0) {  
  function (z) { z(x0) } })
```

```
⟨function (z) { z(x0) }, [x0 ↦ 0]⟩  
⟨function (z) { z(x0) }, [x0 ↦ 1]⟩
```

```

(function (f0) { f0(0); f0(1) })
(function (x0) {
  (function (f1) { f1(0); f1(1) })
  (function (x1) {
    function (z) { z(x0)(x1) }}}))

```

```

⟨function (z) { z(x0)(x1) }, [x0 ↦ 0, x1 ↦ 0]⟩
⟨function (z) { z(x0)(x1) }, [x0 ↦ 1, x1 ↦ 0]⟩
⟨function (z) { z(x0)(x1) }, [x0 ↦ 0, x1 ↦ 1]⟩
⟨function (z) { z(x0)(x1) }, [x0 ↦ 1, x1 ↦ 1]⟩

```

The rest is just using this idea to construct a TM simulator that runs a machine for an exponential number of steps.

(See ICFP'08 for details.)

k -CFA		EXPTIME
	⋮	
2-CFA		EXPTIME
1-CFA		EXPTIME
0-CFA		PTIME
Sub0-CFA		PTIME
Simple closure analysis		PTIME
	⋮	

Part III

Paradox

A unifying view of functions vs. objects:

A function is just an object with an `apply` method.

Both data- and control-flow are recursively intertwined:

$$n(o) \{ o.m(x) \}$$

It is easy to translate the construction to Java?

```
(function (f) { f(0); f(1) })  
(function (x) {  
  function (z) { z(x) }})
```

```
new Fun() { apply(f) { f.apply(0);f.apply(1) }}.apply(  
new Fun() { apply(x) {  
  new Fun() { apply(z) { z.apply(x) }}}})
```

Object-oriented variants of *k*CFA are widely used...

...and they are implemented in Datalog!

PTIME \subsetneq EXPTIME

Part IV

Resolution

Traditionally, Java programs are analyzed in lifted form:

```
new Fun() { apply(f) { f.apply(0);f.apply(1) }}.apply(
new Fun() { apply(x) {
    new Fun() { apply(z) { z.apply(x) }}}})
```

```
class Fun1 { apply(f) { f.apply(0); f.apply(1) }};
class Fun2 { apply(x) { new Fun3(x) }};
class Fun3 {
    x;
    Fun3(x) { this.x = x };
    apply(z) { z.apply(x) }
};
new Fun1().apply(new Fun2())
```

But now objects are *explicitly* closed.

Consider:

```
class Fun3 {  
  x;  
  Fun3(x) { this.x = x };  
  apply(z) { z.apply(x) }  
};
```

```
new Fun3(p);
```

```
new Fun3(q);
```

Creates two closures over one variable.

Now consider:

```
class FunN {  
  x0; ... xn;  
  FunN { this.x0 = x0; ... this.xn = xn };  
  apply(z) { z.apply(x0)...apply(xn) }  
};
```

```
new Fun3(p1, ..., pn);
```

```
new Fun3(q1, ..., qn)
```

Creates 2 closures over n variables; we don't consider 2^n closures.

Theorem. *k*CFA of lifted object-oriented programs is in PTIME.

Part V

Exploitation

Key idea: simulate lifting with the representation of closures.

- ▶ $\text{Exp} \times (\text{Var} \rightarrow \text{Contour})$
- ▶ $\text{Exp} \times \text{Contour}$

At call sites, copy values of closures into the store.

*m*CFA PTIME

⋮

0CFA PTIME

Speed: Pathological-case benchmark

Terms	$k = 1$	$m = 1$	poly., $k=1$
69	ϵ	ϵ	ϵ
123	ϵ	ϵ	ϵ
231	46 s	ϵ	2 s
447	∞	3 s	5 s
879	∞	48 s	1 m 8 s
1743	∞	51 m	∞

Speed and precision:

Prog/Terms	$k = 1$	$m = 1$	poly., $k=1$
eta / 49	ϵ 7	ϵ 7	ϵ 3
map / 157	ϵ 8	ϵ 8	ϵ 8
sat / 223	∞ -	ϵ 12	1s 12
regex / 1015	4s 25	3s 25	14s 25
scm2java / 2318	5s 86	3s 86	3s 79
interp / 4289	5s 123	4s 123	9s 123
scm2c / 6219	179s 136	143s 136	157s 131

Column 1: analysis time

Column 2: inlinings enabled

What to take home with you:

- ▶ Linearity subverts approximation
- ▶ Closures make *k*CFA hard
- ▶ Lifting makes object-oriented *k*CFA easy
- ▶ Flat closures make functional *m*CFA easy
- ▶ *m*CFA v *k*CFA: similar precision, lower cost

