Static analysis of modern software systems: Taming control-flow

Matt Might University of Utah matt.might.net

ucombinator.org

Problem

Software fails.

Problem

- Software fails because we can't engineer it.
- We can't engineer what we can't predict.
- We can't predict the behavior of software.

Message

- Static analysis of modern software is hard!
- Control-flow analysis is the gatekeeper.
- Yet, precise control-flow analysis is possible.



Goal



Software "engineering"



Why we need software engineering

Security vulnerabilities

\$80 billion in cyber-crime each year.

FBI

Cost and cause of insecurity



Cost of cybercrime (\$80bn)



Source: CSI/FBI Survey 2007



Software bugs

Bugs cost U.S. economy \$60 billion annually.

NIST

Why bugs are bad

<exploding-rocket-video />











Source: Intel



2010

GHz



2010

GHz

Source: Intel



1968

2010



1968

2010



1968

2010

Tomorrow's software



Tomorrow's software



Bottom line

If we want software that is...

- ...more parallel,
- ...more correct,
- ...more secure,

then we need engineering.

Why can't we predict what software will do?

Why can't we predict what software will do?

Because Alan Turing said we can't.



"Thou shalt not write a program which determines whether a program halts."

Banned by corollary

- Will a program eventually do X?
- Will a program never do Y?

A "loop" hole

- Always answering "yes" or "no" is impossible.
- Answering "yes," "no" or "maybe" is allowed.











"Full employment theorem." - Appel
Why analyzing modern software is hard

animal.eat(food);

What is animal?

animal.eat(food);

What is food?

void process (Animal animal) {
food = world.gather();
animal.eat(food);

Who calls process?

void process (Animal animal) {
food = world.gather();
animal.eat(food);

What is world?

The control-flow problem

The control-flow problem

23

Control-flow

Data-flow

Gatekeeper

Before we can do anything interesting, we must bound interprocedural control-flow.

Essence of the problem

Value = Object

Essence of the problem

Value = Object = Class + Record

Essence of the problem

Value = Object = Class + Record \subseteq Code + Data

Which language is the paragon of value = code + data?

A-calculus

Assertion

If we can analyze λ -calculus expressions, we can analyze object-oriented programs.

λ -calculus (Church, 1928)

• Minimalist, universal language



- Minimalist, universal language
- Three expression types:

v [variable]

$\lambda v. e$ [anonymous function]



- Minimalist, universal language
- Three expression types:
 - v [variable]
 - $e_1(e_2)$ [function application]



- Minimalist, universal language
- Three expression types:
 - v [variable]
 - $e_1(e_2)$ [function application]
 - $\lambda v. e$ [anonymous function]



Lisp and Scheme

• $v \equiv V$

•
$$f(e) = (f e)$$

• $\lambda v.e = (lambda (v) e)$

λ -fortified

- Lisp
- SML
- Haskell
- Scala
- Java
- C#

- C++ (Boost)
- Python
- Ruby
- Smalltalk
- JavaScript
- PHP(!)

One rule: *β*-reduction

$(\lambda v.v^2)(3)$

One rule: *β*-reduction

Another interpretation

- Functions = Closures
- $Closure = \lambda \times Env$
- $Env = Var \rightarrow Value$
- Ex: $(\lambda x.x + z, [z \mapsto 1])$

Essence of the essence

- Value = Code + Data
- Closure = λ + Environment

Control-flow question

Given a call site f(x), what could f be?

Control-flow scenarios



Control-flow scenarios

let $f = \lambda z.z$ in f(x)

Control-flow scenarios

 $\lambda f.f(x)$

Control-flow analysis

A **control-flow analysis** conservatively approximates the procedures which may be invoked at a given call site.

Control-flow analysis

A **control-flow analysis** conservatively approximates the procedures which may be invoked at a given call site.

A value-flow analysis conservatively approximates the values to which an expression may evaluate.

Techniques for CFA

- Ad hoc techniques
- Constraint-solving
- Type-based analysis
- Abstract interpretation

Techniques for CFA

- Ad hoc techniques
 - Constraint-solving
- Type-based analysis
- Abstract interpretation

Techniques for CFA

- Ad hoc techniques
 - Constraint-solving
- Type-based analysis
- Abstract interpretation

Constraint-based OCFA

What is OCFA?

Lambda-flow analysis.
The OCFA approximation

- Value = Code x Data
- Closure = Lambda x Env
- Object = Class x Record

The OCFA approximation

- Value = Code
- Closure = Lambda
- Object = Class

0CFA



OCFA

 $\lambda v.e_{
m b}$ $e_1(e_2)$













 $\lambda v.e_{b} \in \text{FlowsTo}[e_{1}] \text{ and } val \in \text{FlowsTo}[e_{b}]$ $val \in \text{FlowsTo}[e_{1}(e_{2})]$

OCFA

 $\frac{\lambda v.e_{b} \in \text{FlowsTo}[e_{1}] \text{ and } val \in \text{FlowsTo}[e_{b}]}{val \in \text{FlowsTo}[e_{1}(e_{2})]}$

OCFA

 $\lambda v.e_{\rm b} \in {\rm FlowsTo}[\lambda v.e_{\rm b}]$

 $\frac{\lambda v.e_{b} \in \text{FlowsTo}[e_{1}] \text{ and } val \in \text{FlowsTo}[e_{b}]}{val \in \text{FlowsTo}[e_{1}(e_{2})]}$

 $\frac{\lambda v.e_{b} \in \text{FlowsTo}[e_{1}] \text{ and } val \in \text{FlowsTo}[e_{2}]}{val \in \text{FlowsTo}[v]}$

OCFA (Palsberg, 1995)

 $\{\lambda v.e_{\rm b}\} \subseteq {\rm FlowsTo}[\lambda v.e_{\rm b}]$

 $\frac{\lambda v.e_{b} \in \text{FlowsTo}[e_{1}]}{\text{FlowsTo}[e_{b}] \subseteq \text{FlowsTo}[e_{1}(e_{2})]}$

 $\frac{\lambda v.e_{b} \in \text{FlowsTo}[e_{1}]}{\text{FlowsTo}[e_{2}] \subseteq \text{FlowsTo}[v]}$

OCFA (Palsberg, 1995)

 $\{\lambda v.e_{\rm b}\} \subseteq {\rm FlowsTo}[\lambda v.e_{\rm b}]$

 $\frac{\lambda v.e_{b} \in \text{FlowsTo}[e_{1}]}{\text{FlowsTo}[e_{b}] \subseteq \text{FlowsTo}[e_{1}(e_{2})]}$

 $\frac{\lambda v.e_{b} \in \text{FlowsTo}[e_{1}]}{\text{FlowsTo}[e_{2}] \subseteq \text{FlowsTo}[v]}$

OCFA (Palsberg, 1995)

 $\{\lambda v.e_{\rm b}\} \subseteq {\rm FlowsTo}[\lambda v.e_{\rm b}]$

 $\frac{\lambda v.e_{\rm b} \in \operatorname{FlowsTo}[e_1]}{\operatorname{FlowsTo}[e_{\rm b}] \subseteq \operatorname{FlowsTo}[e_1(e_2)]}$

+ Constraint Solver

 $\frac{\lambda v.e_{\rm b} \in \operatorname{FlowsTo}[e_1]}{\operatorname{FlowsTo}[e_2] \subseteq \operatorname{FlowsTo}[v]}$

= Control-flow analysis

Applications

- Classic data-flow analysis
- Data-flow optimizations
- Global register allocation
- Defunctionalization
- Static method resolution

- Global constant propagation
- Global copy propagation
- Loop detection/optimization
- Escape analysis
- Constant folding

A problem with the "CFA-first" approach

map f list

fireMissile(n)

map f list

| |

fireMissile(n) map f list















Solution platform: Small-step abstract interpretation







Small-step strategy

- Model program as infinite-state machine
- Approximate program with finite-state machine

Small-step machine

Small-step machine

• Convert program e into machine state s_0



Small-step machine

- Convert program e into machine state s_0
- Transition from state s_n to state s_{n+1}


Abstract machine















Theorem: The abstract simulates the concrete.

Abstraction



Example: Abstract graph

```
 (\text{letrec } ((\text{lp1 } (\lambda \text{ (i x)}) \\ (\text{if } (= 0 \text{ i}) \text{ x} \\ (\text{letrec } ((\text{lp2 } (\lambda \text{ (j f y)}) (\text{if } (= 0 \text{ j}) \\ (\text{lp1 } (- \text{ i 1}) \text{ y}) \\ (\text{lp2 } (- \text{ j 1}) \text{ f} \\ (\text{f y}))))) \\ (\text{lp2 } 10 (\lambda (\text{n}) (+ \text{n i})) \text{ x})))))
```

(lp1 10 0))

Example: Abstract graph



Small-step CFA for CPS

Continuation-passing style

$$v \in Var$$

 $f, e \in Exp = Var + Lam$
 $lam \in Lam ::= (\lambda (v_1 \dots v_n) call)$
 $call \in Call ::= (f e_1 \dots e_n)$

Concrete state-space

 $\varsigma \in \Sigma = \mathsf{Call} \times BEnv \times Store \times Time$ $\beta \in BEnv = \mathsf{Var} \rightharpoonup Addr$ $\sigma \in Store = Addr \rightharpoonup Clo$ $clo \in Clo = \mathsf{Lam} \times BEnv$ $a \in Addr \text{ is a set of addresses}$ $t \in Time \text{ is a set of time-stamps}$

Simpler option

 $\varsigma \in \Sigma = \text{Call} \times Env$ $\rho \in Env = \text{Var} \rightharpoonup Clo$ $clo \in Clo = \text{Lam} \times Env$

Concrete state-space

 $\varsigma \in \Sigma = \mathsf{Call} \times BEnv \times Store \times Time$ $\beta \in BEnv = \mathsf{Var} \rightharpoonup Addr$ $\sigma \in Store = Addr \rightharpoonup Clo$ $clo \in Clo = \mathsf{Lam} \times BEnv$ $a \in Addr \text{ is a set of addresses}$ $t \in Time \text{ is a set of time-stamps}$





Injector

$\mathcal{I}: \mathsf{Call} \to \Sigma$ $\mathcal{I}(call) = (call, [], [], t_0)$





Factored evaluator

 $\mathcal{E}(v,\beta,\sigma) = \sigma(\beta(v))$ $\mathcal{E}(lam,\beta,\sigma) = (lam,\beta)$

Concrete semantics

When $call = \llbracket (f e_1 \dots e_n) \rrbracket$:

 $(call, \beta, \sigma, t) \Rightarrow (call', \beta'', \sigma', t'),$ where $(lam, \beta') = \mathcal{E}(f, \beta, \sigma)$ $clo_i = \mathcal{E}(e_i, \beta, \sigma)$ $lam = \llbracket (\lambda (v_1 \dots v_n) call') \rrbracket$ t' = tick(call, t) $a_i = alloc(v_i, t')$ $\beta'' = \beta'[v_i \mapsto a_i]$ $\sigma' = \sigma[a_i \mapsto clo_i]$

Concrete semantics

When $call = [[(f e_1 ... e_n)]]$:

 $(call, \beta, \sigma, t) \Rightarrow (call', \beta'', \sigma', t'),$ where $(lam, \beta') = \mathcal{E}(f, \beta, \sigma)$ $clo_i = \mathcal{E}(e_i, \beta, \sigma)$ $lam = \llbracket (\lambda \ (v_1 \dots v_n) \ call') \rrbracket$ t' = tick(call, t) $a_i = alloc(v_i, t')$ $\beta'' = \beta'[v_i \mapsto a_i]$ $\sigma' = \sigma[a_i \mapsto clo_i]$

The easy solution

$Time = \mathbb{N}$ $Addr = \operatorname{Var} \times Time$

 $tick(_, t) = t + 1$ alloc(v, t) = (v, t)





Abstracting into a control-flow analysis

Abstract state-space $\hat{\varsigma} \in \hat{\Sigma} = \mathsf{Call} \times \widetilde{BEnv} \times \widetilde{Store} \times \widetilde{Time}$ $\hat{\beta} \in \widetilde{BEnv} = \operatorname{Var} \to \widetilde{Addr}$ $\hat{\sigma} \in \widehat{Store} = \widehat{Addr} \to \mathcal{P}\left(\widehat{Clo}\right)$ $clo \in Clo = Lam \times BEnv$ $\hat{a} \in Addr$ is a **finite** set of addresses $\hat{t} \in Time$ is a **finite** set of time-stamps

Abstract state-space

 $\hat{\varsigma} \in \hat{\Sigma} = \operatorname{Call} \times \widehat{BEnv} \times \widehat{Store} \times \widehat{Time}$ $\hat{\beta} \in \widehat{BEnv} = \operatorname{Var} \to \widehat{Addr}$ $\hat{\sigma} \in \widehat{Store} = \widehat{Addr} \to \mathcal{P}\left(\widehat{Clo}\right)$ $\widehat{clo} \in \widehat{Clo} = \operatorname{Lam} \times \widehat{BEnv}$ $\hat{a} \in \widehat{Addr} \text{ is a finite set of addresses}$ $\hat{t} \in \widehat{Time} \text{ is a finite set of time-stamps}$

 $\varsigma \in \Sigma = \mathsf{Call} \times BEnv \times Store \times Time$ $\beta \in BEnv = \mathsf{Var} \rightharpoonup Addr$ $\sigma \in Store = Addr \rightharpoonup Clo$ $clo \in Clo = \mathsf{Lam} \times BEnv$ $a \in Addr \text{ is a set of addresses}$ $t \in Time \text{ is a set of time-stamps}$

Abstract state-space

 $\hat{\varsigma} \in \hat{\Sigma} = \operatorname{Call} \times \widehat{BEnv} \times \widehat{Store} \times \widehat{Time}$ $\hat{\beta} \in \widehat{BEnv} = \operatorname{Var} \to \widehat{Addr}$ $\hat{\sigma} \in \widehat{Store} = \widehat{Addr} \to \mathcal{P}\left(\widehat{Clo}\right)$ $\widehat{clo} \in \widehat{Clo} = \operatorname{Lam} \times \widehat{BEnv}$ $\hat{a} \in \widehat{Addr} \text{ is a finite set of addresses}$ $\hat{t} \in \widehat{Time} \text{ is a finite set of time-stamps}$

 $\varsigma \in \Sigma = \mathsf{Call} \times BEnv \times Store \times Time$ $\beta \in BEnv = \mathsf{Var} \rightharpoonup Addr$ $\sigma \in Store = Addr \rightharpoonup Clo$ $clo \in Clo = \mathsf{Lam} \times BEnv$ $a \in Addr \text{ is a set of addresses}$ $t \in Time \text{ is a set of time-stamps}$

Is this state-space finite?

Non-recursive

 $\hat{\varsigma} \in \hat{\Sigma} = \mathsf{Call} \times \widehat{BEnv} \times \widehat{Store} \times \widehat{Time}$ $\hat{\sigma} \in \widehat{Store} = \widehat{Addr} \to \mathcal{P}\left(\widehat{Clo}\right)$ $\widehat{clo} \in \widehat{Clo} = \text{Lam} \times \widehat{BEnv}$ $\hat{\beta} \in \widetilde{BEnv} = \mathsf{Var} \to \widetilde{Addr}$ $\hat{a} \in Addr$ is a **finite** set of addresses $\hat{t} \in Time$ is a **finite** set of time-stamps

Non-recursive

 $\hat{\varsigma} \in \hat{\Sigma} = \mathsf{Call} \times \widehat{BEnv} \times \widehat{Store} \times \widehat{Time}$ $\hat{\sigma} \in \widehat{Store} = Addr \longrightarrow \mathcal{P}$ \widehat{Clo} $\widehat{clo} \in \widehat{Clo} = \operatorname{Lam} \times BEnv$ $\hat{\beta} \in BEnv = \operatorname{Var} \to Aadr$ $\hat{a} \in \widetilde{Addr}$ is a **finite** set of addresses $\hat{t} \in Time$ is a **finite** set of time-stamps







$\hat{\mathcal{I}}(call) = (call, [], [], \hat{t}_0)$




Abstract components

$(\leadsto) \subseteq \hat{\Sigma} \times \hat{\Sigma}$

 $\hat{\mathcal{E}}: \mathsf{Exp} \times \widehat{BEnv} \times \widehat{Store} \to \mathcal{P}\left(\widehat{Clo}\right)$

 $\hat{\mathcal{E}} : \mathsf{Exp} \times \widehat{BEnv} \times \widehat{Store} \to \mathcal{P}\left(\widehat{Clo}\right)$

 $\hat{\mathcal{E}}(v,\hat{\beta},\hat{\sigma}) = \hat{\sigma}(\hat{\beta}(v))$ $\hat{\mathcal{E}}(lam, \hat{\beta}, \hat{\sigma}) = \left\{ (lam, \hat{\beta}) \right\}$

Abstract semantics

When $call = \llbracket (f e_1 \dots e_n) \rrbracket$: $(call, \hat{\beta}, \hat{\sigma}, \hat{t}) \rightsquigarrow (call', \hat{\beta}'', \hat{\sigma}', \hat{t}'), \text{ where }$ $(lam, \hat{\beta}') \in \hat{\mathcal{E}}(f, \hat{\beta}, \hat{\sigma})$ $\hat{C}_i = \hat{\mathcal{E}}(e_i, \hat{\beta}, \hat{\sigma})$ $lam = \llbracket (\lambda \ (v_1 \dots v_n) \ call') \rrbracket$ $\hat{t}' = \hat{tick}(call, \hat{t})$ $\hat{a}_i = \widehat{alloc}(v_i, \hat{t}')$ $\hat{\beta}'' = \hat{\beta}' [v_i \mapsto \hat{a}_i]$ $\hat{\sigma}' = \hat{\sigma} \sqcup [\hat{a}_i \mapsto \hat{C}_i]$





Abstraction maps

$$\begin{split} \alpha(call,\beta,\sigma,t) &= (call,\alpha(\beta),\alpha(\sigma),\alpha(t)) \\ \alpha(\beta) &= \lambda v.\alpha(\beta(v)) \\ \alpha(\sigma) &= \lambda \hat{a}. \bigsqcup_{\alpha(a)=\hat{a}} \alpha(\sigma(a)) \\ \alpha(lam,\beta) &= \{(lam,\alpha(\beta))\} \\ \alpha(a) \text{ is set by parameter} \\ \alpha(t) \text{ is set by parameter} \end{split}$$

$$\hat{\sigma} \sqcup \hat{\sigma}' = \lambda \hat{a}.(\hat{\sigma}(\hat{a}) \cup \hat{\sigma}'(\hat{a}))$$







Theorem: If the concrete takes a step, then the abstract can take a matching step.





Application: Buffer-overflow checks



i < length(a)</pre>



i < length(a)</pre>



Beyond CFA

- Abstract G.C.
- Nondeterministic A.I.
- Frame-string analysis
- Abstract counting
- Logic-flow analysis
- Dependence analysis
- See matt.might.net

Thanks!