


# Inter-procedural Control Flow Analysis



Using Constraint-based  
Approach

# The Dynamic Dispatch Problem

---

- Which function is called by  $p(x)$ ?

```
int myFunc ( int (*p)(int), ...)  
{  
    .....  
    return p(x);  
}
```

- P is a function pointer. What function could p point to (what is the value of p)?
- P is a function parameter, so the value of p is unknown unless inter-procedural dataflow analysis is performed
  - But inter-procedural data-flow requires an inter-procedural control flow graph (or a call graph)
- The problem is relevant for
  - Imperative languages that allow functions as parameters
  - Object oriented languages and functional languages

# Inter-procedural Control flow Analysis

---

## □ Example code

```
int f (int (*x)(int) { return x(1); }  
int g (int y) { return y + 2; }  
int h (int z) { return z + 3; }  
int main() {  
    return f(g) + f(h);  
}
```

- For each function call, what functions may be invoked?

# Defining the Analysis

---

- What is the domain of analysis
  - What is the solution space?
    - What could be the values for each function pointer expression?
- Specification of the analysis
  - How to compute the solution?
    - how to accommodate the information flow from function definitions to function invocations
- Well-definedness of the analysis
  - What are the properties of the solution space?
  - Does it compute a solution?
  - Does the algorithm terminate?
  - Is the solution precise?

# Specification of Domain

---

- What is the solution?
  - For each expression in the program, could it have a function pointer value? If yes, what functions may it point to? (if no, the solution is  $\emptyset$ )
  - Must keep track of the values of variables (especially function parameters)
- To represent the solution, label each expression within the program, compute
  - An abstract cache (C) so that for each expression e,
    - C(e) contains the set of function values e may have
  - An abstract environment (P) so that for each variable x,
    - P(x) contains the set of function values x may have

# The Input Language

---

## □ Assume a small functional language

```
e ::= c           // constant values
    | x           // variable reference
    | fun f x => e0 // function with name f, parameter x, and body e0
    | e1 e2       // invoking function e1 with argument e2
    | if e0 then e1 else e2 //if e0 is true, return e1, else return e2
    | let x = e1 in e2 // introduce local variable x=e1 in e2
```

## □ Why functional language?

- Functions are first-class objects; allow nested functions/scopes
  - Can be used to model virtual functions in object-oriented programming
- Dataflow is explicit (a single symbolic value for each variable). No variable is ever modified
  - For imperative programming languages, perform global data-flow analysis / build SSA

# Example Code and Control-flow Analysis Solution

- Example code

((fun f x => x) (fun g y => y))

- Labels: 1: x;

2: (fun f x => x)

3: y;

4: (fun g y => y)

5: ((fun f x => x) (fun g y => y))

- Example CFA solution (guesses of the (C,P) mappings)

1	{fun g y => y}
2	{fun f x => x}
3	∅
4	{fun g y => y}
5	{fun g y => y}
x	{ fun g y => y}
y	∅
f	{fun f x => x}
g	{ fun g y => y}

# Solution Space of CFA

---

## □ Formally

- Abstract values:  $Val = \text{Power}(\text{Term})$

- Each term is a function definition in the form  $(\text{fun } f \ x \Rightarrow e_0)$

- Abstract environment:  $Env = \text{Var} \rightarrow Val$

- Var: the set of all variables (including function parameters)

- Abstract cache:  $Cache = \text{Label} \rightarrow Val$

- Label: the set of labels (expressions)

- Each solution: a pair of  $(P, C) \subseteq (Env, Cache)$



# Specification of CFA

---

- What properties must be satisfied by  $(P,C)$  to be a correct/acceptable solution?
  - $(C,P) \models e$  means that  $(C,P)$  is an acceptable Control Flow Analysis Solution for the expression  $e$ 
    - $(C,P) \models c$ 

Arbitrary solutions are acceptable for a constant value  $c$
    - $(C,P) \models (x)\ell$  iff  $P(x) \subseteq C(\ell)$ 

The solution for an variable must be a subset of the solution for its label (each variable has a single value through each of its lifetime)
    - $(C,P) \models (\text{fun } f \ x \Rightarrow (e_0)\ell_0)\ell_1$  iff  $(C,P) \models (e_0)\ell_0$  and
$$\{\text{fun } f \ x \Rightarrow e_0\} \subseteq C(\ell_1) \text{ and } \{\text{fun } f \ x \Rightarrow e_0\} \subseteq P(f)$$

The solution for a function definition(abstraction) label must include the function definition(abstraction)

# Specification of CFA (2)

---

- Function invocation (application)
  - $(C,P) \models ((e1)\ell_1 (e2)\ell_2)\ell_3$  iff  $(C,P) \models (e1)\ell_1$ ,  $(C,P) \models (e2)\ell_2$ , and  
 $\forall (\text{fun } f \ x \Rightarrow (e0)\ell_0) \in C(\ell_1) : (C,P) \models (e0)\ell_0$ ,  $C(\ell_2) \subseteq P(x)$  and  $C(\ell_0) \subseteq C(\ell_2)$ 
    - The solution for function parameter (x) must contain that of the invocation argument (e2);
    - The solution of the function invocation must contain that of the function body
- Local variables (nested scopes)
  - $(C,P) \models (\text{let } x = (e1)\ell_1 \text{ in } (e2)\ell_2)\ell_3$  iff  $(C,P) \models (e1)\ell_1$ ,  $(C,P) \models (e2)\ell_2$ ,  
 $C(\ell_1) \subseteq P(x)$  and  $C(\ell_2) \subseteq C(\ell_3)$ 
    - The solution for the local variable (x) must contain that of its defined value
    - The solution of the outer scope must contain that of the inner scope
- Conditionals
  - $(C,P) \models (\text{if } (e0)\ell_0 \text{ then } (e1)\ell_1 \text{ else } (e2)\ell_2)\ell_3$  iff  $(C,P) \models (e0)\ell_0$ ,  $(C,P) \models (e1)\ell_1$ ,  
 $(C,P) \models (e2)\ell_2$ , and  $C(\ell_2) \subseteq C(\ell_3)$  and  $C(\ell_1) \subseteq C(\ell_3)$ 
    - The solution of the outer scope must contain that of the inner scopes (both branches)

# Example Code and Control-flow Analysis Solution

## Example code

```
((fun f x => x) (fun g y => y))
```

### Labels: 1: x;

2: (fun f x => x)

3: y;

4: (fun g y => y)

5: ((fun f x => x) (fun g y => y))

## Example CFA solution (guesses of the (C,P) mappings). Are the valid?

	(C,P)	(C',P')
1	{fun g y => y}	{fun g y => y}
2	{fun f x => x}	{fun f x => x}
3	∅	∅
4	{fun g y => y}	{fun g y => y}
5	{fun g y => y}	{fun g y => y}
x	{fun g y => y}	∅
y	∅	∅
f	{fun f x => x}	{fun f x => x}
g	{fun g y => y}	{fun g y => y}

$(C,P) \models ((fun f x => x) (fun g y => y))$

$(C',P') \not\models ((fun f x => x) (fun g y => y))$

# Well-definedness of CFA Analysis

- Difficulty: Cannot build  $(C,P) \models e$  by structural induction on the expression  $e$

- E.g. function invocation (application)

$(C,P) \models ((e1)\ell_1 (e2)\ell_2)\ell_3$  iff  $(C,P) \models (e1)\ell_1$ ,  $(C,P) \models (e2)\ell_2$ , and

$\forall (\text{fun } f \ x \Rightarrow (e0)\ell_0) \in C(\ell_1)$ ,  $(C,P) \models (e0)\ell_0$ ,  $C(\ell_2) \subseteq P(x)$  and  $C(\ell_0) \subseteq C(\ell_2)$

- There is no guarantee that  $C(\ell_0)$  has been computed correctly before computing  $C(\ell_2)$

- Coinductive definition: the solution space includes all guesses of  $(C,P)$  that satisfy the specifications

- Must apply all constraints to iteratively modify the solutions until they become correct
- The best solution is the smallest one that satisfies all the constraints

# Correctness of Specification

---

- If there is a possible evaluation of the program such that the function at a call point evaluates to some function definition
  - then this definition has to be in the set of possible definitions computed by the analysis.
- Existence of solutions
  - Every expression accepts a least CFA solution

# Constraint based Analysis

---

- Syntax-directed analysis
  - Reformulate the analysis specification
    - Construct a finite set of constraints based on structural induction
    - Compute the least solution of the set of constraints
- Each constraint has the form
$$(\text{sol1} \subseteq \text{sol2}) \text{ or } (\{t\} \subseteq \text{sol}) \text{ or } (\{t\} \subseteq \text{sol1} \Rightarrow \text{sol2} \subseteq \text{sol3})$$
  - where
    - Each sol is either  $C(\ell)$  or  $P(x)$ 
      - $\ell$  is label,  $x$  is a variable
    - Each  $t$  is either  $(\text{fn } x \Rightarrow e0)$  or  $(\text{fun } f \ x \Rightarrow e0)$

# Constraint-based Analysis

- For each expression  $e$ , compute  $\text{Cond}[e]$ 
  - $\text{Cond}[c] = \emptyset$  // constants
  - $\text{Cond}[(x)l] = \{ P(x) \subseteq C(l) \}$  // variables
  - $\text{Cond}[(\text{fun } f \ x \Rightarrow e_0)l] = \text{Cond}[e_0] \cup$   
 $\{ \{ \text{fun } f \ x \Rightarrow e_0 \} \subseteq C(l) \} \cup \{ \{ \text{fun } f \ x \Rightarrow e_0 \} \subseteq P(f) \}$  // function def.
  - $\text{Cond}[(e_1)l_1 \ (e_2)l_2]l_3 = \text{Cond}[e_1] \cup \text{Cond}[e_2] \cup$   
 $\{ \{t\} \in C(l_1) \Rightarrow C(l_2) \subseteq P(x) \ \forall t = (\text{fun } f \ x \Rightarrow (e_0)l_0) \} \cup$   
 $\{ \{t\} \in C(l_1) \Rightarrow C(l_0) \subseteq C(l_3) \ \forall t = (\text{fun } f \ x \Rightarrow (e_0)l_0) \}$
  - $\text{Cond}[(\text{let } x = (e_1)l_1 \ \text{in } (e_2)l_2)l_3] =$   
 $\text{Cond}[e_1] \cup \text{Cond}[e_2] \cup \{C(l_1) \subseteq P(x)\} \cup \{C(l_2) \subseteq C(l_3)\}$
  - $\text{Cond}[(\text{if } (e_0)l_0 \ \text{then } (e_1)l_1 \ \text{else } (e_2)l_2)l_3] =$   
 $\text{Cond}[e_0] \cup \text{Cond}[e_1] \cup \text{Cond}[e_2] \cup \{C(l_2) \subseteq C(l_3)\} \cup \{C(l_2) \subseteq C(l_3)\}$

# Example: Constraint Construction

---

$\text{Cond}[\text{((fun f x } \Rightarrow \text{(x))1)2 (fun g y } \Rightarrow \text{(y)3)4 )5}]$   
= { {fun f x  $\Rightarrow$  (x)}  $\subseteq$  C(2), {fun f x  $\Rightarrow$  (x)}  $\subseteq$  P(f),  
P(x)  $\subseteq$  C(1),  
{fun g y  $\Rightarrow$  (y)}  $\subseteq$  C(4), {fun g y  $\Rightarrow$  (y)}  $\subseteq$  P(g),  
P(y)  $\subseteq$  C(3),  
{fun f x  $\Rightarrow$  (x)}  $\subseteq$  C(2)  $\Rightarrow$  C(4)  $\subseteq$  P(x),  
{fun f x  $\Rightarrow$  (x)}  $\subseteq$  C(2)  $\Rightarrow$  C(1)  $\subseteq$  C(5),  
{fun g y  $\Rightarrow$  (y)}  $\subseteq$  C(2)  $\Rightarrow$  C(4)  $\subseteq$  P(y),  
{fun g y  $\Rightarrow$  (y)}  $\subseteq$  C(2)  $\Rightarrow$  C(3)  $\subseteq$  C(5) }



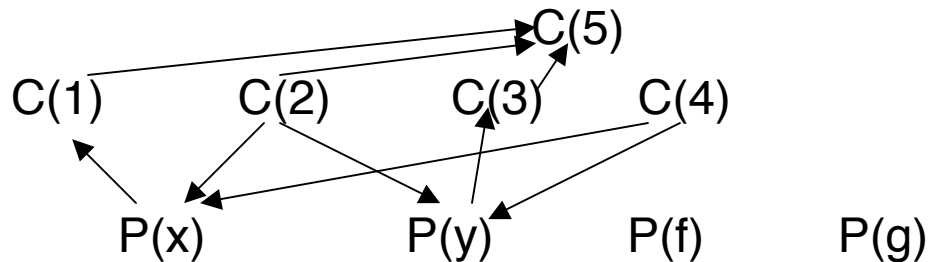
# Solving the constraints

---

- Input: a set of constraints for the entire program
- Output: the least solution (C,P) to the constraints
- Idea: equivalent to finding the least fixed point of a monotone function defined by the constraints
  - Straight-forward iterative algorithm has  $n^5$  cost, where  $n$  is the size of the program (expression)
  - A more sophisticated algorithm takes  $n^3$  complexity
- The graph-based algorithm
  - Build a graph where
    - Each node  $n$  corresponds to a unique  $C(l)$  or  $P(x) \Rightarrow \text{val}(n)$
    - Add an edge from node  $n1$  to  $n2$  if any change to  $\text{val}(n1)$  may require modifications to  $\text{val}(n2)$
  - Use a worklist to keep track of nodes to change

# Constraint Solving Algorithm (1)

- Define  $\text{add}(t, p) = \{ \text{if } (t \not\subseteq p) \{ p = p \cup t; \text{append}(p, \text{worklist}); \} \}$
- Step 1 Initialization
  - $\text{worklist} := \text{nil};$
  - for each label  $l$  (or variable  $x$ ) do
    - $\text{Val}[C(l)] = \text{nil}; \text{Edge}(C(l)) = \text{nil};$  (or  $\text{Val}[P(x)] = \text{nil}; \text{Edge}(P(x)) = \text{nil}$ )
- Step 2 Building the graph
  - for each  $cc$  in  $\text{Cond}[\text{program}]$  do
    - case  $cc$  of  $\{t\} \subseteq p$ :  $\text{add}(t, \text{Val}(p));$
    - $p1 \subseteq p2$ :  $\text{append}(cc, \text{Edge}[p1]);$
    - $\{t\} \subseteq p \Rightarrow p1 \subseteq p2$ :  $\text{append}(cc, \text{Edge}[p1]); \text{append}(cc, \text{Edge}[p]);$



# Constraint Solving Algorithm(2)

## □ Step 3 Iteration

- while worklist is not empty do

$q := \text{Remove-first}(W);$

    for each  $cc$  in  $\text{Edge}[q]$  do

        case  $cc$  of  $p1 \subseteq p2$ :  $\text{add}(p2, \text{Val}[p1]);$

$\{t\} \subseteq p \Rightarrow p1 \subseteq p2$ : if  $t \subseteq \text{Val}[p]$  then  $\text{add}(\text{val}(p1), p2);$

Val	Iteration 0	Propogate C(2)...	Propoage P(x)	Propoage C(1)
C(1)	$\emptyset$	$\emptyset$	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$
C(2)	$\{\text{fun } f \ x \Rightarrow x\}$	$\{\text{fun } f \ x \Rightarrow x\}$	$\{\text{fun } f \ x \Rightarrow x\}$	$\{\text{fun } f \ x \Rightarrow x\}$
C(3)	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
C(4)	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$
C(5)	$\emptyset$	$\emptyset$	$\emptyset$	$\{\text{fun } g \ y \Rightarrow y\}$
P(x)	$\emptyset$	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$
P(y)	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
P(f)	$\{\text{fun } f \ x \Rightarrow x\}$	$\{\text{fun } f \ x \Rightarrow x\}$	$\{\text{fun } f \ x \Rightarrow x\}$	$\{\text{fun } f \ x \Rightarrow x\}$
P(g)	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$	$\{\text{fun } g \ y \Rightarrow y\}$

# Summary

---

- Recording the solution of CFA analysis
  - for each label  $\ell$  (or variable  $x$ ) do
    - $C(\ell) = \text{Val}[C(\ell)]$  ( $P(x) = \text{Val}[P(x)]$ )
  
- Correctness and Termination
  - The worklist algorithm terminates and the result produced by the algorithm is the least solution to  $C[[e]]$ .
  - Complexity: The algorithm takes at most  $O(n^3)$  steps if the original expression  $e$  has size  $n$ .