

Interprocedural Data Flow Analysis

Aditya V. Nori
MSR India

Intraprocedural Analysis

proceduralless program \leadsto flow graph

Assumption: All paths in the graph represent actual executions of the program.

!

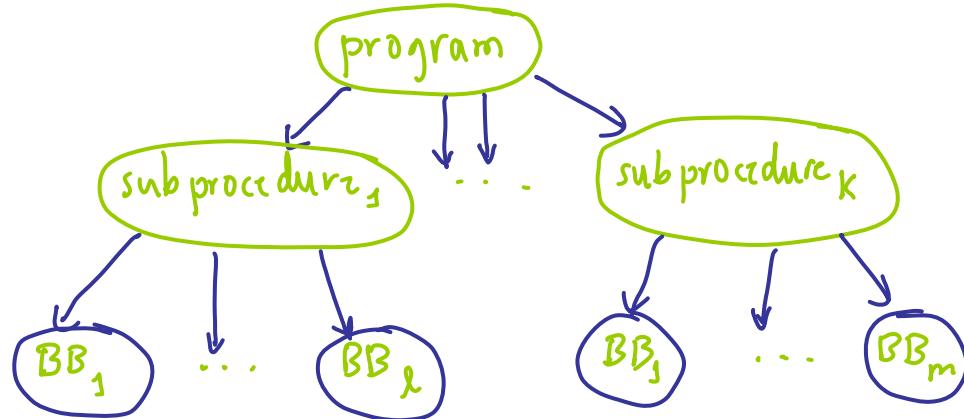
Not true!

But this model is popular:

- Relatively simple structure enables the development of simple algorithms for program analysis.
- Isolation of feasible paths from infeasible paths is an undecidable problem.

Reference: Flow Analysis of Computer Programs - Hecht 1977.

Notation



A basic block is a maximal single entry multiexit sequence of code.

Each procedure p is assumed to have a single entry block r_p and a single exit block e_p .

The flow graph G_p of p is a rooted directed graph whose nodes are the basic blocks of p , with root r_p . There is an edge (m, n) in G_p if there is a direct transfer of control from basic block m to basic block n .

Notation

Let $G = (N, E, r)$ be a directed graph.

$\uparrow \quad \uparrow \quad \uparrow$
nodes edges root

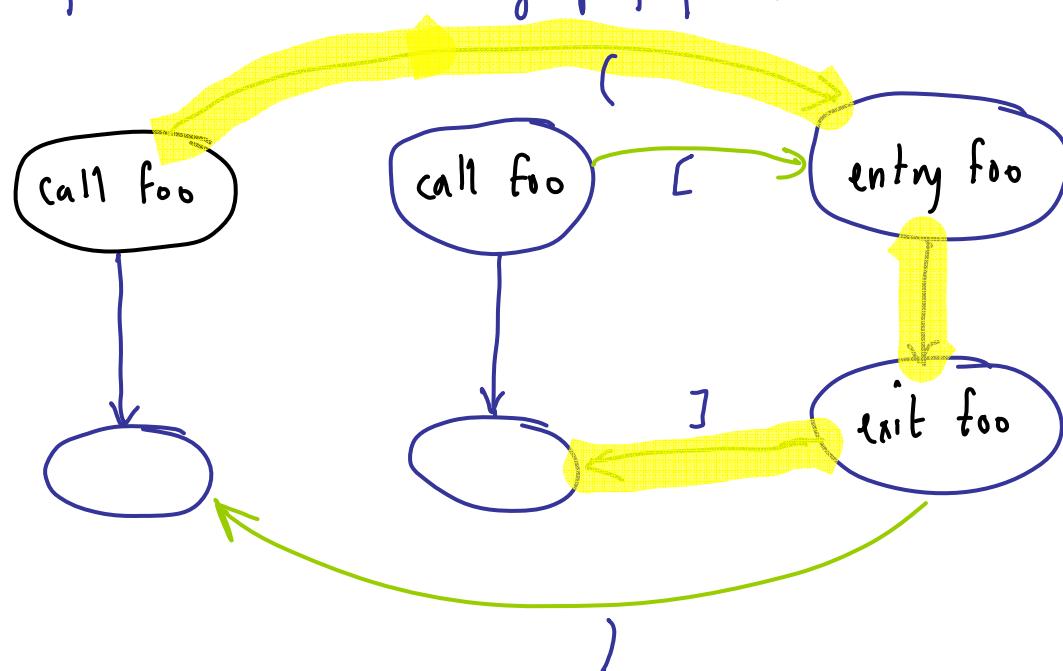
A path p in G is a sequence of nodes in N (n_1, n_2, \dots, n_k) st. for each $1 \leq j < k$, $(n_j, n_{j+1}) \in E$.

The length of p is the no. of edges along p , denoted $\text{length}(p)$.

$\text{path}_G(m, n) =$ set of all paths in G leading from m to n .

Interprocedural Analysis

- (A) Functional Approach: Views procedures as collections of structured program blocks and aims to establish I/O relations for each block.
- (B) Call string Approach: Information is tagged with an encoded history of procedure calls during propagation.



A global data flow frame is a pair (L, F) , where L is a semilattice of data and $F: L \rightarrow L$ is a space of functions. F is closed under composition and meet and is monotone.

(L, F) is a distributive framework if:

$$\forall f \in F \text{ and } x, y \in L, f(x \wedge y) = f(x) \wedge f(y)$$

Data propagation equations: (Given (L, F) and G)

$$x_r = 0$$

$$x_n = \bigwedge_{(m,n) \in E} f_{(m,n)}(x_m) \quad n \in N - \{r\}$$



Note: every edge (m, n) of G is associated with a propagation fn $f_{(m,n)}$ which represents change of data as control transfers from m to n .

Available Expression Analysis

Aim: For each program point, determine which expressions have already been computed, and not later modified, on all paths to a program point.

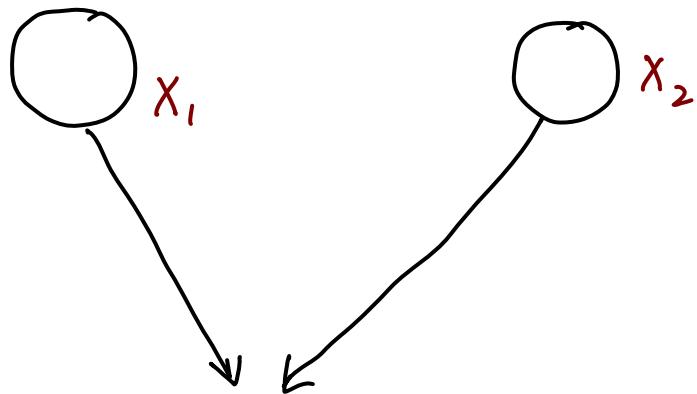
Example:

```
x := a + b;  
y := a * b;  
while y > a + b  
do  
    a := a + 1;  
    x := a + b;  
od
```



```
x := a + b;  
y := a * b;  
while y > x  
do  
    a := a + 1;  
    x := a + b;  
od
```

Basic Idea



$$M = X_1 \cap X_2$$

$x := a$

$x = M \setminus \{ \text{expressions with an } x \}$

kill

gen

$\cup \{ \text{subexpressions of } a \text{ with an } x \}$

$$L = (\mathcal{P}(\text{Exps}), \sqsupseteq, \sqcap, \text{Exps}) \rightarrow \text{semilattice}$$

The least upper bound operator is \sqcap .

The least element \perp is Exps .

L satisfies the ascending chain condition.

F is the set of transfer functions (gen & kill).

Exercise: Show that this is a distributive dfa framework.

Examples: [ASU].

Meet Over all Paths (MOP) :

$$y_n = \bigwedge \{ f_p(0) : \text{path}_G(r, n) \}$$

For each path $p = (n_1, n_2, \dots, n_k)$, $f_p = f_{(n_{k-1}, n_k)} \circ f_{(n_{k-2}, n_{k-1})} \circ \dots \circ f_{(n_1, n_2)}$

If p is empty, $f_p = \text{id}_L$.

Computing y_n is undecidable in general.

[Kildall]: If (L, F) is distributive, then $x_n = y_n, \forall n \in N$.

[Kam]: If (L, F) is monotone, then $x_n \leq y_n, \forall n \in N$.

Interprocedural Flow Graphs

$G = \bigcup \{ G_p : p \text{ is a procedure in the program} \}$, where

$$G_p = (N_p, E_p, r_p) \text{ s.t.}$$

- N_p : set of all basic blocks in p .
- $E_p = E_p^0 \cup E_p^1$: set of edges of G_p .
 - $(m, n) \in E_p^0 \rightarrow$ there can be a direct flow of control from m to n .
 - $(m, n) \in E_p^1 \rightarrow m$ is a call block and n follows that block.

$G^* = (N^*, E^*, r_1)$, where $N^* = \bigcup_p N_p$, and $E^* = E^0 \cup E^1$,

$E^0 = \bigcup_p E_p^0$ ← defined earlier.

$(m, n) \in E^1 \rightarrow$ m is a call block and n is the called proc. entry (call edge)
 \leftarrow m is an exit block of some proc p and n follows the
call to p (return edge).

The call edge (m, r_p) and a return edge (r_q, n) are said to correspond to each other if $p = q$ and $(m, n) \in E_s^1$ for some procedure s .

Note: All paths in G^* are not feasible!

$\text{IVP}(r_1, n)$: set of all interprocedurally valid paths in G^* , from r_1 to n .

$q \in \text{path}_{G^*}(r_1, n) \in \text{IVP}(r_1, n)$ if $q|_{E^1}$ is proper.

(a) $q|_{E^1}$ is proper if it has no return edges.

(b) If $q|_{E^1}$ contains return edges \rightarrow they are matched by corresponding calls.

$\bigcup_n \text{IVP}(r_1, n)$ is a CFL.

$\text{path}_{G^*}(r_1, n)$ is regular

A path $q \in \text{IVP}_o(r_p, n)$ iff $q_1 = q|_{E^1}$ is complete



- The null tuple is complete.
- q_1 is complete if
 - (a) it is the concatenation of two complete tuples, or
 - (b) it starts with a call edge and terminates with a corresponding return edge, and the rest of the tuple is complete.

Example

Main program

```

read a, b;
t := a * b;
call p;
t := a * b;
print t;
stop;

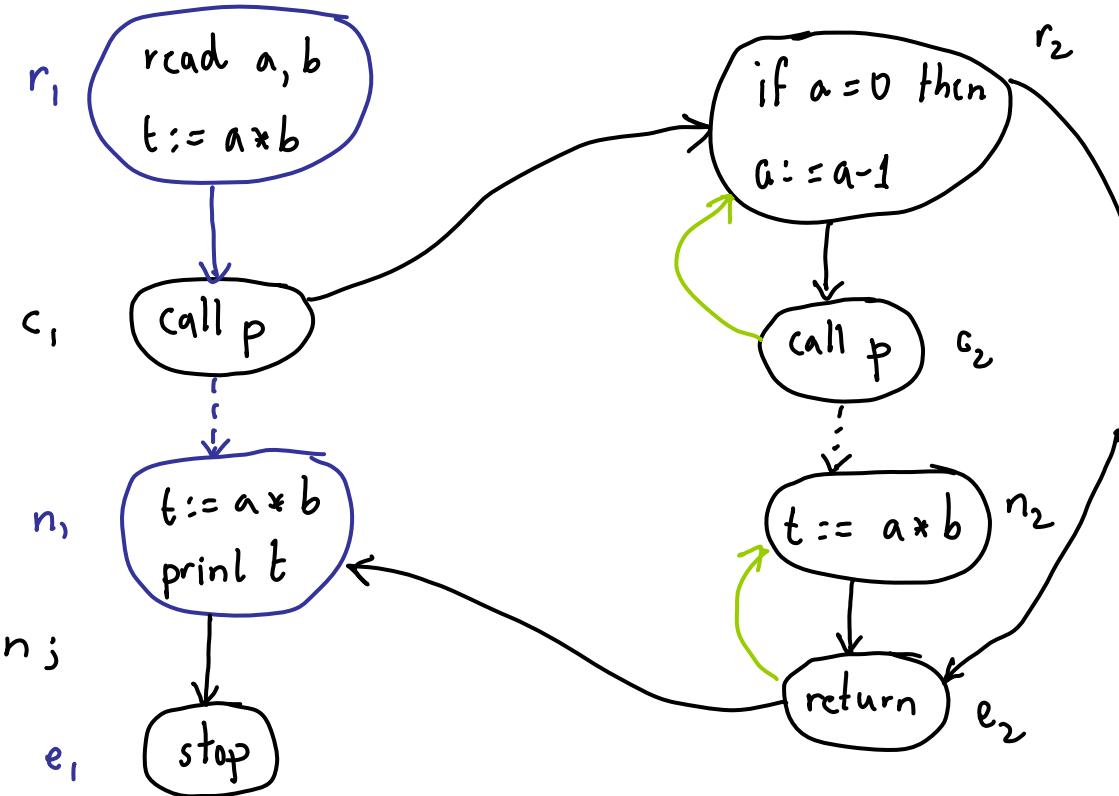
```

procedure p

```

if a = 0 then return;
else
  a := a - 1;
  call p;
  t := a * b;
endif;
return;

```



$$q_1 = (r_1, c_1, r_2, c_2, r_2, e_2, n_2, e_2, n_1, e_1) \in IVP(r_1, e_1)$$

$$q_2 = (r_1, c_1, r_2, c_2, r_2, e_2, n_1, e_1) \notin IVP(r_1, e_1)$$

$$q_3 = (r_2, c_2, r_2, e_2, n_2) \in IVP_0(r_2, n_2)$$

$$q_4 = (r_2, c_2, r_2, c_2, e_2, n_2) \notin IVP_0(r_2, n_2)$$

Lemma: Let $n \in N^*$ and $q \in IVP(r_1, n)$. Then there exist procedures p_1, p_2, \dots, p_j , where p_1 is the main program and p_j the procedure containing n , and calls c_1, \dots, c_{j-1} such that for each $i < j$, c_i is in p_i and calls p_{i+1} , and q can be represented as

$$q = q_1 \sqcap (c_1, r_{p_2}) \sqcap q_2 \sqcap (c_2, r_{p_3}) \sqcap \dots \sqcap (c_{j-1}, r_{p_j}) \sqcap q_j$$

where for each $i < j$, $q_i \in IVP_0(r_{p_i}, c_i)$ and $q_j \in IVP_0(r_{p_j}, n)$. Conversely, any path which admits such a decomposition is in $IVP(r_1, n)$ and this decomposition is unique.

Proof: Simple induction on path length.

(Reading exercise!)

A Functional Approach to Interprocedural Analysis

Let (L, F) be a distributive flow data flow framework for G .

Phase I : $\phi_{(r_p, r_p)} \leq id_L$ for each procedure p .

$$\phi_{(r_p, n)} = \bigwedge_{(m, n) \in E_p} (h_{(m, n)} \circ \phi_{(r_p, m)}) \quad \forall n \in N_p.$$

where

$$h_{(m, n)} = \begin{cases} f_{(m, n)} & \text{if } (m, n) \in E_p^0 \\ \phi_{(r_q, e_q)} & \text{if } (m, n) \in E_p^{-1} \text{ and } m \text{ calls procedure } q. \end{cases}$$

Define an ordering over F as follows:

$$\forall g_1, g_2 \in F, \quad g_1 \geq g_2 \iff g_1(x) \geq g_2(x) \quad \forall x \in L.$$

$$\left. \begin{array}{l} \text{Let } \phi_{(r_p, r_p)}^o = \text{id}_L \text{ for each procedure } p. \\ \phi_{(r_p, n)}^o = f_{SL} \text{ for each } n \in N_p - \{r_p\}. \end{array} \right\} \text{initialization}$$

Compute the maximum fixed point for this set of equations.

Phase II For each basic block n , let $x_n \in L$ denote the info available at the start of n .

$$x_{r_{\text{main}}} = 0 \in L$$

$$x_{r_p} = \bigwedge \{ \phi_{(r_q, c)}(x_{r_q}) : q \text{ is a procedure and } c \text{ is a call to } p \text{ in } q \}$$

for each procedure p .

$$x_n = \phi_{(r_p, n)}(x_{r_p}) \text{ for each procedure } p, \text{ and } n \in N_p - \{r_p\}.$$

Example - Available expressions

$$L = \{0, 1, 2\}$$

$$F = \{0, 1, id_L, f_2\}$$

a * b is available

Phase I : $\phi_{(r_1, r_1)} = \phi_{(r_2, r_2)} = id_L$

$$\phi_{(r_1, c_1)} = 1 \circ \phi_{(r_1, r_1)}$$

$$\phi_{(r_1, n_1)} = \phi_{(r_2, e_2)} \circ \phi_{(r_1, c_1)}$$

$$\phi_{(r_1, e_1)} = 1 \circ \phi_{(r_1, n_1)}$$

$$\phi_{(r_2, c_2)} = 0 \circ \phi_{(r_2, r_2)}$$

$$\phi_{(r_2, n_2)} = \phi_{(r_2, e_2)} \circ \phi_{(r_2, c_2)}$$

$$\phi_{(r_2, e_2)} = [id_L \circ \phi_{(r_2, r_2)}] \wedge [1 \circ \phi_{(r_2, n_2)}]$$

<u>Function</u>	<u>Initial value</u>	<u>It 1</u>	<u>It 2</u>	<u>It 3</u>	
$\phi_{(r_1, r_1)}$	id_L	id_L	id	id	fixed point!
$\phi_{(r_1, g)}$	f_n	$\underline{1}$	$\underline{1}$	$\underline{1}$	
$\phi_{(r_1, n_1)}$	f_n	f_n	$\underline{1}$	$\underline{1}$	
$\phi_{(r_1, u)}$	f_n	f_n	$\underline{1}$	$\underline{1}$	
$\phi_{(r_2, r_2)}$	id_L	id_L	id_L	id_L	
$\phi_{(r_2, g_2)}$	f_n	$\underline{0}$	0	0	
$\phi_{(r_2, n_2)}$	f_n	f_n	0	0	
$\phi_{(r_2, e_2)}$	f_n	id_L	id_L	id_L	

Phase II : $x_{r_1} = 0$

$$\begin{aligned}x_{r_2} &= \phi_{(r_1, e_1)}(x_{r_1}) \wedge \phi_{(r_2, e_2)}(x_{r_2}) \\&= \underline{1}(x_{r_1}) \wedge \underline{0}(x_{r_2})\end{aligned}$$

After two iterations : $x_{r_1} = x_{r_2} = 0$

Complete solution : $x_{r_1} = x_{r_2} = x_{c_2} = x_{n_2} = x_{e_2} = 0$

$$x_{c_1} = x_{n_1} = x_{e_1} = 1$$

$\Rightarrow a * b$ is available at the start of n_1 .

Interprocedural MOP

$$\Psi_n = \bigwedge \{ f_q : q \in IVP(r_{\text{main}}, n) \} \in F \text{ for each } n \in N^*.$$

$$y_n = \Psi_n(0) \text{ for each } n \in N^*.$$

* this is the MOP solution.

Lemma: Let $n \in N_p$ for some procedure p . Then

$$\phi_{(r_p, n)} = \bigwedge \{ f_q : q \in IVP_0(r_p, n) \}$$

Proof:

Step 1: Prove by induction on $i \rightarrow \phi_{(r_p, n)}^i \geq \bigwedge \{ f_q : q \in IVP_0(r_p, n) \}.$

Basis: $i=0 \rightarrow$ (a) $n=r_p \rightarrow \phi_{(r_p, r_p)}^0 = id_L \geq \bigwedge \{ f_q : q \in IVP_0(r_p, r_p) \} = id_L$

 (b) $n \neq r_p \rightarrow \phi_{(r_p, n)}^0 = f_{-n} \geq f, \forall f \in F$

Inductive hypothesis: Assume $\phi_{(r_p, n)}^i \geq \Lambda \{ f_q : q \in IVP_0(r_p, n) \}$ for some i .

$$\begin{aligned}
\phi_{(r_p, n)}^{i+1} &= \Lambda_{(m, n) \in E_p} (h_{(m, n)} \circ \phi_{(r_p, m)}^i) \\
&\geq \Lambda_{(m, n) \in E_p} (h_{(m, n)} \circ \Lambda \{ f_q : q \in IVP_0(r_p, m) \}) \text{ for each procedure } p \\
&\quad \text{and } n \in N - \{r_p\} \\
&= \Lambda_{(m, n) \in E_p^0} (f_{(m, n)} \circ \Lambda \{ f_q : q \in IVP_0(r_p, m) \}) \\
&\quad \Lambda_{(m, n) \in E_p^1} (\phi_{(r_{p'}, e_{p'})} \circ \Lambda \{ f_q : q \in IVP_0(r_p, m) \}) \\
&\quad m \text{ calls } p' \\
&\geq \Lambda_{(m, n) \in E_p^0} (\Lambda \{ f_{q|| (m, n)} : q \in IVP_0(r_p, m) \}) \Lambda \\
&\quad \Lambda_{(m, n) \in E_p^1} (\Lambda \{ f_{q'} : q' \in IVP_0(r_{p'}, e_{p'}) \} \circ \Lambda \{ f_q : q \in IVP_0(r_p, m) \})
\end{aligned}$$

$$\begin{aligned}
&= \bigwedge_{(m,n) \in E_p} \left(\bigwedge \left\{ f_{q \parallel (m,n)} : q \in IVP_0(r_p, m) \right\} \right) \wedge \\
&\quad \bigwedge_{\substack{(m,n) \in E_p \\ m \text{ calls } p'}} \left(\bigwedge \left\{ f_{q \parallel (m, r_{p'}) \parallel q' \parallel (e_{p'}, n)} : q \in IVP_0(r_p, m), q' \in IVP_0(r_{p'}, e_{p'}) \right\} \right) \\
&\geq \left\{ f_q : q \in IVP_0(r_p, n) \right\}.
\end{aligned}$$

Step 2: For each $q \in IVP_0(r_p, n)$, show that $f_q \geq \phi_{(r_p, n)}$.

Proof by induction on the length of q .

Basis: If $\text{length}(q) = 0$, then $n = r_p \rightarrow f_q = \phi_{(r_p, r_p)} = \text{id}_L$.

Inductive hypothesis: Assume that $f_q \geq \phi_{(r_p, n)}$. $\forall p, n$ and $\forall q \in IVP_0(r_p, n)$ s.t. $\text{length}(q) \leq k$.

Induction: Let p, q, n be such that $\text{length}(q) = k+1$. Let (m, n) be the last edge in $q \rightarrow q = q_1 \parallel (m, n)$.

If $(m, n) \in E_p^o \Rightarrow q_1 \in IVP_o(r_p, m)$ and $\text{length}(q_1) \leq k$.

$$\therefore f_{q_1} \geq \phi_{(r_p, m)} \text{ and } f_{q_1} = f_{(m, n)} \circ f_{q_1} \geq h_{(m, n)} \circ \phi_{(r_p, m)} \geq \phi_{(r_p, n)}$$

If $(m, n) \in E^{'}, \text{ then } m = e_p,$ for some procedure p' .

$$\begin{aligned} \therefore q &= q_1 || (m_1, r_{p'}) || q_2 || (e_{p'}, n) \\ &\quad \uparrow \qquad \qquad \uparrow \qquad \qquad \downarrow (m_1, n) \in E_p^{'} \\ &\quad IVP_o(r_p, m_1) \qquad IVP_o(r_{p'}, e_{p'}) \end{aligned}$$

$$\text{Since } f_{(m_1, r_{p'})} = f_{(e_{p'}, n)} = \text{id}_L \Rightarrow f_q = f_{q_2} \circ f_{q_1},$$

$$\Rightarrow f_{q_1} \geq \underbrace{\phi_{(r_{p'}, e_{p'})} \circ \phi_{(r_p, m_1)}}_{\text{by IH.}} = h_{(m_1, n)} \circ \phi_{(r_p, m_1)} \geq \phi_{(r_p, n)}$$

For each basic block n , define

$$\chi_n = \bigwedge \{ \phi_{(r_{p_j}, n)} \circ \phi_{(r_{p_{j-1}}, c_j)} \circ \dots \circ \phi_{(r_{p_1}, c_1)} : p_1 = \text{main program},$$

p_j is the procedure containing n , and for each $i < j$, c_j is a call to p_{i+1} from $p_i\}$

$$z_n = \chi_n(0)$$

Thm: $\psi_n = \chi_n$ for each $n \in N^*$.

Proof: let $q \in IVP(r_{\text{main}}, n)$. Then $q = q_1 \parallel (c_1, r_{p_2}) \parallel q_2 \parallel \dots \parallel (c_{j-1}, r_{p_j}) \parallel q_j$

(path decomposition thm). $p_1 = \text{main program}$

$p_2, \dots, p_{j-1} = \text{procedures}$

$p_j = \text{procedure containing } n$

$\begin{cases} c_1, \dots, c_{j-1} = \text{calls s.t.} \\ \forall i < j, c_i \text{ is a call to} \\ p_{i+1} \text{ from } p_i, \end{cases}$

and $q_i \in IVP_o(r_{p_i}, c_i)$ and $q_j \in IVP_o(r_{p_j}, n)$

By the earlier lemma:

$$f_q = f_{q_j} \circ f_{q_{j-1}} \circ \dots \circ f_{q_1} \geq \phi_{(r_{p_j}, n)} \circ \phi_{(r_{p_{j-1}}, c_{j-1})} \circ \dots \circ \phi_{(r_{p_1}, c_1)} \geq x_n$$

Therefore $\psi_n \geq x_n$.

Conversely, let $p_1, \dots, p_j, c_1, \dots, c_{j-1}$ be defined as before.

$$\phi_{(r_{p_j}, n)} \circ \phi_{(r_{p_{j-1}}, c_{j-1})} \circ \dots \circ \phi_{(r_{p_1}, c_1)} = \bigwedge \{ f_{q_j} \circ f_{q_{j-1}} \circ \dots \circ f_{q_1};$$

$q_i \in IVP_0(r_{p_i}, c_i)$ for each $i < j$

and $q_j \in IVP_0(r_{p_j}, n) \}$

$$= \{ f_{q_1} \| (c_1, r_{p_1}) \| \dots \| (c_{j-1}, r_{p_{j-1}}) \| q_j : \dots \} = \mathbb{P}$$

Every path in \mathbb{P} belongs to $IVP(r_{\text{main}}, n)$

$$\Rightarrow \mathbb{P} \supseteq \bigwedge \{ f_q : q \in IVP(r_{\text{main}}, n) \}$$

Therefore $x_n \geq \psi_n \Rightarrow x_n = \psi_n \quad \forall n \in \mathbb{N}^*$.

Main Result

Thm: For each basic block $n \in N^*$, $x_n = y_n = z_n$.

Proof: $\because \underbrace{\psi_n = x_n \ \forall n \in N^*}_{\text{previous thm}} \Rightarrow y_n = z_n$. (Note: $y_n = \psi_n(0)$ and $z_n = x_n(b)$)

Claim: $x_{rp} = z_{rp}$, $\forall p \in \text{Procedures}$.

If this claim is true, then by

(a) $x_n = \phi_{(r_p, n)}(x_{rp})$ for each procedure p and $n \in N - \{r_p\}$.

(b) $x_n = \wedge \{ \phi_{(r_{pj}, n)} \circ \phi_{(r_{pj-1}, c_{j-1})} \circ \dots \circ \phi_{(r_p, c_j)} : \dots \}$

(c) $z_n = x_n(0)$



$x_n = z_n$ for each $n \in N^*$

Proof of the claim

Define a new flow graph $G_c = (N_c, E_c, r_c)$, where N_c is the set of all call and entry blocks in the program.

$$E_c = E_c^0 \cup E_c'$$

$(m, n) \in E_c^0 \iff m$ is the entry node of some procedure p and n is a call within p .

$(m, n) \in E_c' \iff m$ is a call to some procedure p and n is the entry of p .

Define a DFP for G_c by associating $g_{(m,n)} \in F$, for each $(m,n) \in E_c$.

$$g_{(m,n)} = \begin{cases} \phi_{(m,n)} & \text{if } (m,n) \in E_c^0 \\ id_L & \text{if } (m,n) \in E_c' \end{cases}$$

Iterative eqns for this DFP:

$$x_{r_{\text{main}}} = 0 \in L$$

$$x_{r_p} = \bigwedge \{ \phi_{(r_q, c)} (x_{r_c}) : q \text{ is a procedure and } c \text{ is a call to } p \text{ in } q \}$$

for each procedure p .

The MOP solution for the new problem:

$$\chi_n = \wedge \{ \phi_{(r_{pj}, n)} \circ \phi_{(r_{pj-1}, c_{j-1})} \circ \dots \circ \phi_{(r_p, c_1)} :$$

p_1 = main program, p_2, \dots, p_j are procedures }, $n \in \{\text{call, entry}\}$

$$z_n = \chi_n(0).$$

Since F is distributive, from Kildall's thm $u_{rp} = z_{rp}$ for each procedure p ,