

# 一、PTA的规则

我们先针对上节中提到的除涉及方法调用的 Pointer-affecting 语句进行分析，涉及方法调用的之后再回来谈。

## (1) 符号约定

Variables:

$x, y \in V$

$\rightarrow$ 一般属于某个实例

Fields:

$f, g \in F$

Objects:

$o_i, o_j \in O$

$\rightarrow$ 实例集

Instance fields:

$o_i.f, o_j.g \in O \times F$

Pointers:

Pointer =  $V \cup (O \times F)$

Points-to relations:

$pt : \text{Pointer} \rightarrow \mathcal{P}(O)$

- $\mathcal{P}(O)$  denotes the powerset of  $O$
- $pt(p)$  denotes the points-to set of  $p$

实际上我们要求解的就是  $pt$  映射，表征两类 Pointer 到具体 object 集的映射（也即  $O$  的某个子集）

## (2) Rules

Kind	Statement	Rule
New	$i : x = \text{new } T()$	$\overline{o_i \in pt(x)}$
Assign	$x = y$	$\frac{o_i \in pt(y)}{o_i \in pt(x)}$
Store	$x.f = y$	$\frac{o_i \in pt(x), o_j \in pt(y)}{o_j \in pt(o_i.f)}$
Load	$y = x.f$	$\frac{o_i \in pt(x), o_j \in pt(o_i.f)}{o_j \in pt(y)}$

① i:  $x = \text{new } T()$

没有条件，直接将新产生的对象  $o_i$  加入  $\text{pt}(x)$

②  $x = y$

假如对象  $o_i$  属于集合  $\text{pt}(y)$ ，就把它加入到  $\text{pt}(x)$  中

③  $x.f = y$

若  $o_i \in \text{pt}(x)$ ,  $o_j \in \text{pt}(y)$ , 则  $o_j \in \text{pt}(o_i.f)$

④  $y = x.f$

若  $o_i \in \text{pt}(x)$ ,  $o_j \in \text{pt}(o_i.f)$ , 则  $o_j \in \text{pt}(y)$

Kind	Rule	Illustration
New	$\overline{o_i \in \text{pt}(x)}$	
Assign	$\frac{o_i \in \text{pt}(y)}{o_i \in \text{pt}(x)}$	
Store	$\frac{o_i \in \text{pt}(x), o_j \in \text{pt}(y)}{o_j \in \text{pt}(o_i.f)}$	
Load	$\frac{o_i \in \text{pt}(x), o_j \in \text{pt}(o_i.f)}{o_j \in \text{pt}(y)}$	

## 二、如何进行PTA?

### 1. PTA的逻辑抽象

- Essentially, pointer analysis is to **propagate** points-to information among pointers (variables & fields)

Kind	Statement	Rule
New	$i: x = \text{new } T()$	$\frac{}{o_i \in pt(x)}$
Assign	$x = y$	$\frac{o_i \in pt(y)}{o_i \in pt(x)}$
Store	$x.f = y$	$\frac{o_i \in pt(x), o_j \in pt(y)}{o_j \in pt(o_i.f)}$
Load	$y = x.f$	$\frac{o_i \in pt(x), o_j \in pt(o_i.f)}{o_j \in pt(y)}$

一切信息的来源

主要就是在指针之间传播指向信息。

也即PTA实际上就是在节点间基于某种约束进行

Points-to Information 的传播

所谓的“闭包”，实际上就是由若干确定的一组规则约束，约束的变化一定可以传播到了集中。

对于约束结构我们可以使用图(PFG)来抽象

### 2. PFG的构成

Pointer flow graph of a program is a **directed graph** that expresses how objects flow among the pointers in the program.

(1) Nodes:  $V \cup (D \times F)$

PFG的节点由程序中出现的局部变量、实例

成员构成，也即各种 Pointers (分析对象)

2) Edges: Pointers × Pointers

边  $x \rightarrow y$  表征  $x$  的指向信息 可能会流向  $y$ ，其本质上代表着由 Pointer-affecting 谓句所确立的闭包约束。

- PFG edges are added according to the statements of the program and the corresponding rules

Kind	Statement	Rule	PFG Edge
New	$i: x = \text{new } T()$	$\overline{o_i \in pt(x)}$	N/A
Assign	$x = y$	$\frac{o_i \in pt(y)}{o_i \in pt(x)}$	$x \leftarrow y$
Store	$x.f = y$	$\frac{o_i \in pt(x), o_j \in pt(y)}{o_j \in pt(o_i.f)}$	$o_i.f \leftarrow y$
Load	$y = x.f$	$\frac{o_i \in pt(x), o_j \in pt(o_i.f)}{o_j \in pt(y)}$	$y \leftarrow o_i.f$

由程序很容易获得 PFG 的节点集，关键在于如何根据 Pointer-affecting 谓句导出边集。你可能会说这还不容易，根据节点集直接梭哈不就行了吗？但很不幸，PFG 中的节点由变量 ( $v$ ) 与实例成员 ( $o_i.f$ ) 构成，前者可以直接从谓句中提取，后者的的确立则需要  $x.f$  中  $x$  的指向信息，而这个信息在 PTA 之前我们是

不知道�！给定所需某例的指向信息，

Program edge很容易由 Pointer flow graph

( $a_i \in pt(c)$ ,  $a_i \in pt(d)$ ) 对应语句的规则  
求解得到

> Variable node

v

> Instance field node

$o_i.f$

→ j: b = new T();

a = b; ①

c.f = a; ②

d = c; ③

c.f = d; ④

e = d.f; ⑤

$pt(a) = \{o_j\}$  ①  $pt(b) = \{o_j\}$

②  $pt(o_i.f) = \{o_j\}$  ⑤

③  $pt(e) = \{o_j\}$  ④

⑤  $pt(d) = \{o_j\}$

With PFG, pointer analysis can be solved by computing **transitive closure** of the PFG

E.g., e is reachable from b on the PFG, which means that the objects pointed by b may flow to and also be pointed by e

因此 PFG 构建的难点，就在于涉及  $x.f$  的语句该如何处理？( $pt(x)$  随着分析的迭代而变化，则 PFG 也会随之动态变化，进而又会影响信息传递的迭代)

1. Build pointer flow graph (PFG)

Mutually dependent

2. Propagate points-to information on PFG

构建 PFG 与在 PFG 上传播信息是一个相互影响的整体的过程。

### 3. 差分法

Solve( $S$ )

$WL = [], PFG = \{\}$

**foreach**  $i: x = \text{new } T() \in S$  **do**  
    add  $\langle x, \{o_i\} \rangle$  to  $WL$

**foreach**  $x = y \in S$  **do**  
    AddEdge( $y, x$ )

**while**  $WL$  is not empty **do**  
    remove  $\langle n, pts \rangle$  from  $WL$

$\Delta = pts - pt(n)$

Propagate( $n, \Delta$ )

if  $n$  represents a variable  $x$  **then**

**foreach**  $o_i \in \Delta$  **do**  
        **foreach**  $x.f = y \in S$  **do**  
            AddEdge( $y, o_i.f$ )  
        **foreach**  $y = x.f \in S$  **do**  
            AddEdge( $o_i.f, y$ )

AddEdge( $s, t$ )

**if**  $s \rightarrow t \notin PFG$  **then**  
    add  $s \rightarrow t$  to  $PFG$

**if**  $pt(s)$  is not empty **then**  
    add  $\langle t, pt(s) \rangle$  to  $WL$

Propagate( $n, pts$ )

**if**  $pts$  is not empty **then**  
     $pt(n) \cup= pts$

**foreach**  $n \rightarrow s \in PFG$  **do**  
    add  $\langle s, pts \rangle$  to  $WL$

$S$  Set of statements of  
the input program

$WL$  Work list

$PFG$  Pointer flow graph

vi)  $S$  是包含程序中所有语句的无序集

Flow - Insensitive

vii) Worklist 存储待接收的指向信息  $pts$  及其对应的指针  $n$

- Worklist contains the points-to information to be processed
  - $WL \subseteq \langle \text{Pointer}, \mathcal{P}(O) \rangle^*$

- Each worklist entry  $\langle n, pts \rangle$  is a pair of pointer  $n$  and points-to set  $pts$ , which means that  $pts$  should be propagated to  $pt(n)$ 
  - E.g.,  $[(x, \{o_i\}), (y, \{o_j, o_k\}), (o_j.f, \{o_l\}) \dots]$

(3) Differential Propagation

注意到在  $WL$  中的待传播项进行处理时，

我们实际传播的信息是  $\Delta = pts - pt(n)$ , 也即原本的  $pt(n)$  中没有的指向信息。虽然后面进行 Propagate 操作时进不进行逐步去重不会对算法的正确性产生影响, 但是对于 Propagate 中的①归并  $pts$  至  $pt(n)$  中 ②传播  $pts$  至  $n$  的所有后继节点, 减少  $pts$  的规模无疑可以大大提高算法效率。

- **Differential propagation** is employed to avoid propagation and processing of redundant points-to information
- **Insight:** existing points-to information in  $pt(n)$  have already been propagated to  $n$ 's successors, and **no need** to be propagated again

### Solve( $S$ )

```

while WL is not empty do
    remove  $\langle n, pts \rangle$  from WL
     $\Delta = pts - pt(n)$ 
    Propagate( $n, \Delta$ )
    ...
  
```

$$pt(a) = \{o_1, o_2, o_3\}$$

$$pt(b) = \{o_1, o_3, o_5\}$$



In practice,  $\Delta$  is usually small compared with the original set, so propagating only the new points-to information ( $\Delta$ ) improves efficiency

Differential propagation:  
 $\Delta = pts - pt(c)$   
 $= \{o_1, o_3, o_5\} - \{o_1, o_2, o_3\}$   
 $= \{o_5\}$

PFG

64

### 4) Load/Store 语句

对于 New 语句, 我们简单地将初始化信息加入 WL 中即可; 对于 Assign 语句, 加条边即可。唯独在于涉及  $x.f$  的语句该如何处理? 我们的的想法很简单, 当指向信息传播到代表变量  $x$  的节点,

$n$ 时,  $n$ 指向的 Object 可能会发生变化, 因此对所有涉及  $x.f$  的语句, 我们都需要在  $\Delta$  中的所有  $O_i$  上增加对应的新边。

## ★ 离法总结

① 对于无序语句集  $S$ , 先处理其中的 New 语句与 Assign 语句:

New —— 添加初始指向信息进入 WL

Assign —— 增加新边(同时添加初始推广信息入 WL)

② 之后对 WL 中各项进行处理  $\langle n, pts \rangle$

A. 首先算出  $pts$  相较于  $n$  原有指向信息的增量  $\Delta$

B. 然后归并  $\Delta$  入  $pt(n)$  中(同时添加次级推广信息入 WL)

C. 最终如果  $n$  表征的变量  $x$  在  $S$  中涉及  $x.f$ , 在  $pt(n)$  新增加的可能性指向  $\Delta$  上, 根据相关语句增加新边(同时添加初始推广信息入 WL)

### 三、含有 Method Calls 的 PTA

回顾之前学到的 Inter-procedural Analysis，其关键点在于 Call Graph 的构建（可以基于此得到 ICFG，进而进行 IDFA），而构建 CG 的关键点又在于如何求解 call site  $\rightarrow$  target method（核心在于 virtual call 的求解），其中如何估计动态分发的起始点又是重中之重。

CHA：用 declared Type 估计起始点。

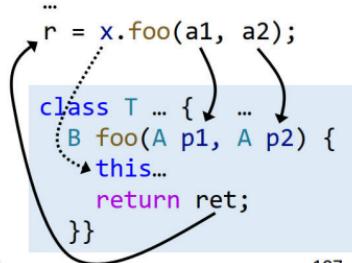
PTA：用  $pt(a)$  近似起始点，更准确。

#### (1) Rules

Kind	Statement	Rule	PFG Edge
Call	$L: r = x.k(a_1, \dots, a_n)$	$\begin{array}{l} \textcircled{1} o_i \in pt(x), m = \text{Dispatch}(o_i, k) \\ o_u \in pt(a_j), 1 \leq j \leq n \\ o_v \in pt(m_{ret}) \\ \hline \textcircled{2} o_i \in pt(m_{this}) \\ v_u \in pt(m_{pj}), 1 \leq j \leq n \\ o_v \in pt(r) \end{array}$	$a_1 \rightarrow m_{p1}$ $\dots$ $a_n \rightarrow m_{pn}$ $r \leftarrow m_{ret}$

- $\text{Dispatch}(o_i, k)$ : resolves the virtual dispatch of  $k$  on  $o_i$  to a target method (based on type of  $o_i$ )
- $m_{this}$ : this variable of  $m$
- $m_{pj}$ : the  $j$ -th parameter of  $m$
- $m_{ret}$ : the variable that holds the return value of  $m$

`C x = new T();`



从而处理 call site 处指向信息流动的规则  
实际上与方法调用的机制是十分相似的。

- ① 基于 Receive Object 的指向信息 + 方法签名确定其 target method.
- ② 确定目标方法 m 的 this 指针的指向。
- ③ 在 cs 处的参数与 m 处的参数间构建 PFG 边，使得前者的指向信息可以流向后者。
- ④ 在 m 的返回值与 cs 处的左值间构建新边，使得前者的指向信息可以流向后者。

Ps: 之所以②没有像③④一样构建新的 PFG 边  
使得信息传导过去，是由于 this 指针所指向的对象  
实际上就是 target method 所属类的实例，  
而 target method 所属的类在 Dispatch 时已经  
完全确定了（要不然怎么能说 target method 是  
确定的呢？）因此为了精确，我们直接传过去

Rule: Call

Why not add PFG edge  $x \rightarrow m_{this}$ ?

Kind	Statement	Rule	PFG Edge
Call	$L: r = x.k(a_1, \dots, a_n)$	$\begin{array}{l} o_i \in pt(x), o_i = Dispatch(o_i, k) \\ o_u \in pt(u), 1 \leq i \leq n \\ o_v \in pt(v_{new}) \\ o_w \in pt(m_{this}) \\ o_x \in pt(m_j), 1 \leq j \leq n \\ o_y \in pt(r) \end{array}$	$a1 \rightarrow m_{p1}$ $an \rightarrow m_{pn}$ $r \rightarrow m_{ret}$

Receiver object should only flow  
to this variable of the  
corresponding target method

PFG edge  $x \rightarrow m_{this}$  would  
introduce spurious points-to  
relations for this variable

With  $x \rightarrow m_{this}$











































































































































































































































































































































</

## (2) Interprocedural PTA (IPTA)

- Run **together** with call graph construction

```
void foo(A a) {
    ...
    b = a.bar();
    ...
}
```

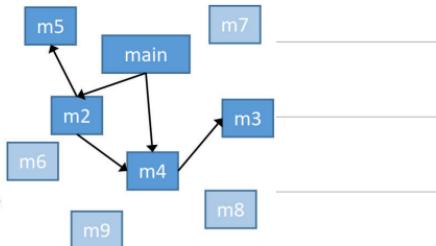
要知道a的  
pts才能构造  
Call Graph

Pointer analysis  
Mutually dependent  
Call graph construction

知道谁调用  
了foo才能对a  
进行操作分析

- Call graph forms a “reachable world”

- Entry methods (e.g., the main method) are reachable from the beginning
- The other reachable methods are gradually discovered during analysis
- Only **reachable** methods and statements are analyzed



实际上依旧是  
从入口方法 (main) 开始，不断扩张，  
分析可达世界。①只分析可达世界，效率自然更高  
②避免了不可达方法的干扰，因此精度更高。

Solve( $m^{entry}$ )

$WL = [ ]$ ,  $PFG = \{ \}$ ,  $S = \{ \}$ ,  $RM = \{ \}$ ,  $CG = \{ \}$

AddReachable( $m^{entry}$ )

while  $WL$  is not empty do

remove  $\langle n, pts \rangle$  from  $WL$

$\Delta = pts - pt(n)$

Propagate( $n, \Delta$ )

if  $n$  represents a variable  $x$  then

foreach  $o_i \in \Delta$  do

foreach  $x.f = y \in S$  do

AddEdge( $y, o_i.f$ )

foreach  $y = x.f \in S$  do

AddEdge( $o_i.f, y$ )

ProcessCall( $x, o_i$ )

$S$  Set of **reachable** statements

$S_m$  Set of statements in method  $m$

$RM$  Set of **reachable** methods

$CG$  Call graph edges

AddReachable( $m$ )

if  $m \notin RM$  then

add  $m$  to  $RM$

$S \cup= S_m$

foreach  $i: x = new T() \in S_m$  do

add  $\langle x, \{o_i\} \rangle$  to  $WL$

foreach  $x = y \in S_m$  do

AddEdge( $y, x$ )

ProcessCall( $x, o_i$ )

foreach  $l: r = x.k(a_1, \dots, a_n) \in S$  do

$m = Dispatch(o_i, k)$

add  $\langle m_{this}, \{o_i\} \rangle$  to  $WL$

if  $l \rightarrow m \notin CG$  then

add  $l \rightarrow m$  to  $CG$

AddReachable( $m$ )

foreach parameter  $p_i$  of  $m$  do

AddEdge( $a_i, p_i$ )

AddEdge( $m_{ret}, r$ )

