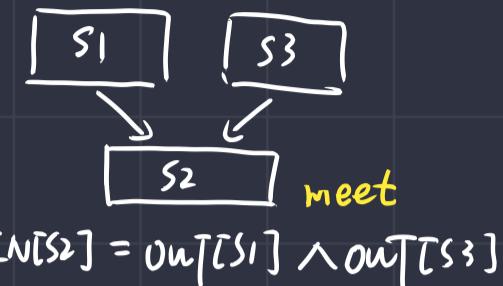
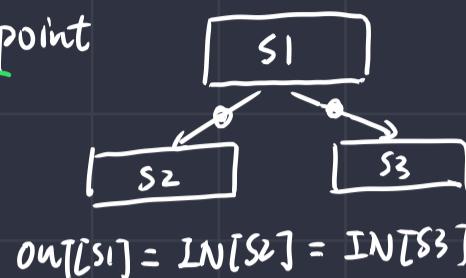
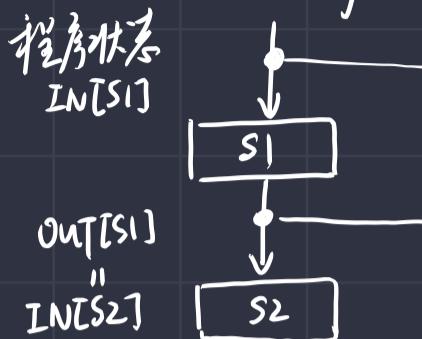


# 软件分析

## 一. 数据流分析应用

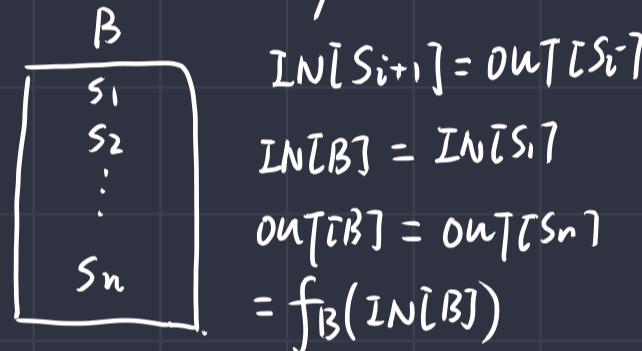
### 1. Preliminaries of Data Flow Analysis



PP 关联一个数据流值，代表程序状态的抽象

{ Forward Analysis  $IN[S] \downarrow$ , if  $OUT[S] = f_S(IN[S])$

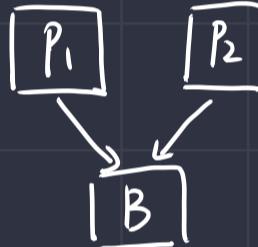
. Backward Analysis,  $IN[S] = f_S(OUT[S])$



$$IN[B_i] = IN[B_{i-1}]$$

$$OUT[B] = OUT[S_n]$$

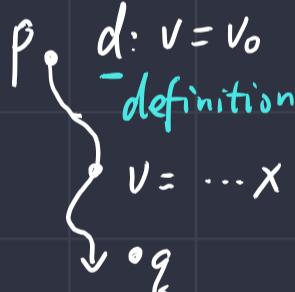
$$= f_B(IN[B])$$



$$IN[B] = \bigwedge_{P \text{ a predecessor of } B} OUT[P]$$

$$\Rightarrow f_B = f_{S_n} \circ \dots \circ f_{S_1}$$

### 2. Reaching Definitions Analysis



Reach 条件: ①  $P \rightarrow q$  存在路径  
② d 不被 killed (v 不被重新定义)

bit vectors:  $D_1, D_2, \dots, D_n$  (definition)  
 $0, 0, \dots, 1$  (0 表示不能 reach)

transfer function

$$B | \underline{D: v = x \text{ op } y} \rightarrow \text{kill 所有定义 } v \text{ 的地方} \Rightarrow OUT[B] = \underline{\text{gen}_B \cup (IN[B] - kill_B)}$$

$$IN[B] = \bigcup_{P \text{ a predecessor of } B} OUT[P]$$

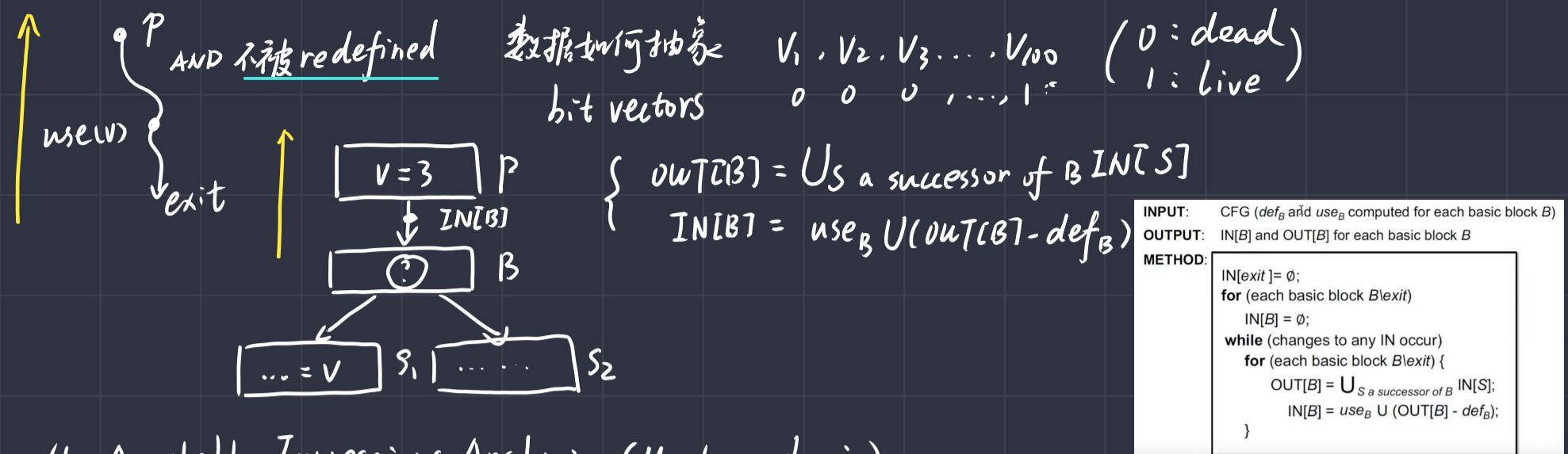
迭代算法:

```

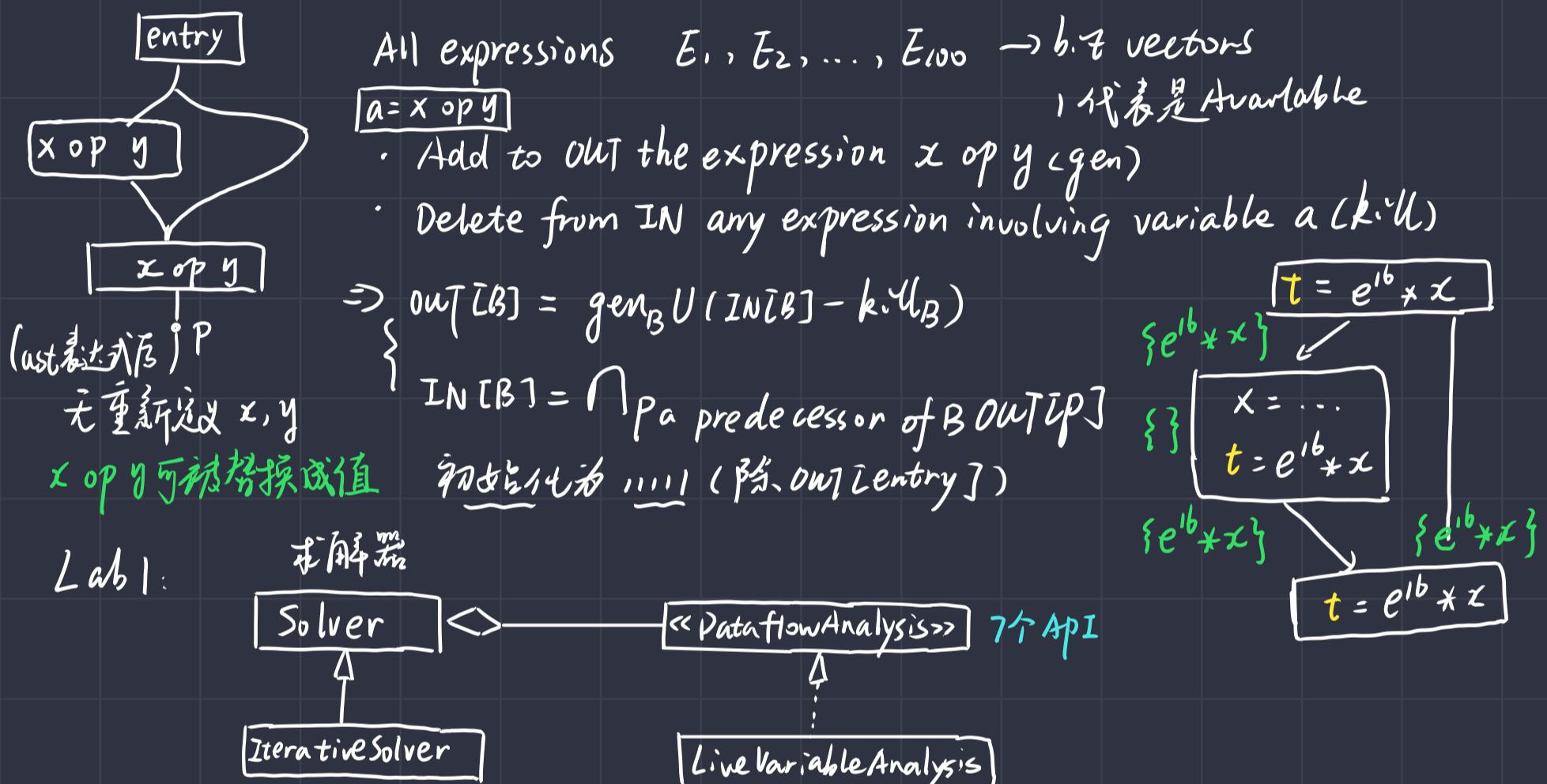
    OUT[entry] = φ;
    for (each basic block B \ entry)
        OUT[B] = φ; 初始化
        while (changes to any OUT occur)
            for (each basic block \ entry) {
                IN[B] =  $\bigcup_{P \text{ a predecessor of } B} OUT[P]$ ;
                OUT[B] =  $\text{gen}_B \cup (IN[B] - kill_B)$ 
            }
    
```

应用: 检测半初始化

### 3. Live Variables Analysis 语义变量分析



### 4. Available Expressions Analysis (Must analysis)



## 二. 数据流分析基础

mode

1. CFG 有  $k$  个 statement. domain of the values in data flow analysis is  $V^k$   
定义  $k$  元组  $(\text{out}[n_1], \text{out}[n_2], \dots, \text{out}[n_k]) \rightarrow V_1 \times V_2 \times \dots \times V_k \Rightarrow V^k$

每次迭代 as a Function  $F: V^k \rightarrow V^k$

$$\text{init} \rightarrow ( \perp, \perp, \dots, \perp ) = X_0$$

$$\text{iter1} \rightarrow (V_1^1, V_2^1, \dots, V_k^1) = X_1 = F(X_0)$$

$$\text{iter2} \rightarrow (V_1^2, V_2^2, \dots, V_k^2) = X_2 = F(X_1)$$

$$\text{iteri} \rightarrow (V_1^i, V_2^i, \dots, V_k^i) = X_i = F(X_{i-1})$$

$$\text{iteri+1} \rightarrow (V_1^{i+1}, V_2^{i+1}, \dots, V_k^{i+1}) = X_{i+1} = F(X_i)$$

$$\because X_i = X_{i+1}$$

$$\therefore X_i = X_{i+1} = F(X_i) \text{ 不成立.}$$

## 2. 偏序 Partial Order

poset(偏序集) as a pair  $(P, \sqsubseteq)$

e.g.



- 集合 偏序关系  $\sqsubseteq$  {
- 1) 自反性  $\forall x \in P, x \sqsubseteq x$  Reflexivity
  - 2) 反对称性  $\forall x, y \in P, x \sqsubseteq y \wedge y \sqsubseteq x \Rightarrow x = y$  Antisymmetry
  - 3) 传递性  $\forall x, y, z \in P, x \sqsubseteq y, y \sqsubseteq z \Rightarrow x \sqsubseteq z$  Transitivity

3. 上界和下界. Given a poset  $(P, \sqsubseteq)$ , 子集  $S \subseteq P$ ,

$u \in P$  是 upper bound of  $S$ . if  $\forall x \in S, x \sqsubseteq u$ . Similarly.

$l \in P$  是 lower bound of  $S$ . if  $\forall x \in S, l \sqsubseteq x$ .

最小上界 least upper bound (lub 或 join)  $\Rightarrow \sqcup S$  对且上界  $u$ .  $\sqcup S \sqsubseteq u$

最大下界 least greatest lower bound (glb 或 meet)  $\Rightarrow \sqcap S$ . ...  $l \sqsubseteq \sqcap S$

若集合只有2个元素.  $\sqcup S \Rightarrow a \sqcup b$ .  $\sqcap S \Rightarrow a \sqcap b$

性质: 偏序集不一定有 lub 和 glb. 若有则唯一

4. Lattice 格. 对 poset  $(P, \sqsubseteq)$ .  $\forall a, b \in P$ . if  $a \sqcup b$  and  $a \sqcap b$  存在. 那么  $(P, \sqsubseteq)$  为格

半格: 只存在  $a \sqcup b \rightarrow$  join 半格 ; 只存在  $a \sqcap b$ . meet 半格

△ 全格  $\forall S \subseteq P$ .  $\sqcup S$  and  $\sqcap S$  exist  $\rightarrow$  存在集合所有属性

(Complete Lattice)  $\begin{cases} \text{最大元素 } \top = \sqcup P \text{ called top} \\ \text{最小元素 } \perp = \sqcap P \text{ called bottom} \end{cases}$   $\forall$  有界的格  $\rightarrow$  全格

Product Lattice.  $L_1 = (P_1, \sqsubseteq_1), L_2 = (P_2, \sqsubseteq_2) \dots$ , 对  $\forall i$ ,  $(P_i, \sqsubseteq_i)$  has  $\sqcup_i$  and  $\sqcap_i$

例  $L^n = (P, \sqsubseteq) : P = P_1 \times P_2 \times \dots \times P_n$ ,  $(x_1, x_2, \dots, x_n) \sqsubseteq (y_1, y_2, \dots, y_n) \Leftrightarrow (x_1 \sqsubseteq y_1) \wedge (x_2 \sqsubseteq y_2) \dots$

## 5. Data Flow Analysis's Framework via Lattice

$(D, L, F)$   
方向 Lattice transfer functions from  $V$  to  $V$

$\sqcup, \sqcap, \sqsubseteq$

6. Lattice 上的函数  $f: L \rightarrow L$

单调性  $f$  is monotonic. if  $\forall x, y \in L, x \sqsubseteq y \Rightarrow f(x) \sqsubseteq f(y)$

Fixed-Point Theorem 不动点定理

Given a complete lattice  $(L, \sqsubseteq)$  if

(1)  $f: L \rightarrow L$  is monotonic (2)  $L$  is finite (有限)

则能找到 the least fixed point (最小不动点) 通过迭代

$f(\perp), f(f(\perp)), \dots, f^k(\perp)$  until a fixed point is reached

greatest fixed point (最大不动点)

$f(T), f(f(T)), \dots, f^k(T)$  until a fixed point is reached

证明: By the definition of  $\perp$  and  $f: L \rightarrow L$

$$\perp \sqsubseteq f(\perp)$$

As  $f$  is monotonic.  $f(\perp) \sqsubseteq f(f(\perp)) = f^2(\perp)$

Similarly,  $\perp \sqsubseteq f(\perp) \sqsubseteq f^2(\perp) \sqsubseteq \dots \sqsubseteq f^i(\perp)$

As  $L$  is finite, for some  $k$   $f^{\text{Fix}} = f^k(\perp) = f^{k+1}(\perp)$ , 不动点存在

Assume we have another fixed point  $x$ , i.e.,  $x = f(x)$ , we have  $\perp \sqsubseteq x$

数学归纳法

$$f(\perp) \sqsubseteq f(x)$$

Assume  $f^i(\perp) \sqsubseteq f^i(x)$ , as  $f$  is monotonic  $f^{i+1}(\perp) \sqsubseteq f^{i+1}(x)$

Thus by induction,  $f^i(\perp) \sqsubseteq f^i(x)$

Thus  $f^i(\perp) \sqsubseteq f^i(x) = x$ ,  $f^{\text{Fix}} = f^k(\perp) \sqsubseteq x$ . 得证. (唯一不存)

7. 应用到 DFA.  $\{ f_i: L \rightarrow L \text{ monotonic}$

$\cup/\cap L \times L \rightarrow L \text{ monotonic}$  Proof:  $\forall x, y, z \in L. x \sqsubseteq y \Rightarrow x \cup z \sqsubseteq y \cup z$

$y \sqsubseteq y \cup z$  由  $x \sqsubseteq y$  得  $x \sqsubseteq y \cup z$  又  $z \sqsubseteq y \cup z$   $h=3$

$x \cup z$  是最上界  $\Rightarrow x \cup z \sqsubseteq y \cup z$

$T \{a, b, c\}$

$\{a, b\} \downarrow \{a, c\} \rightarrow \{b, c\}$

When will the Algorithm Reach the Fixed Point?

$k$ -tuple, 最坏迭代  $h \cdot k$  次

8. May and Must Analysis - a Lattice View

All definitions may reach

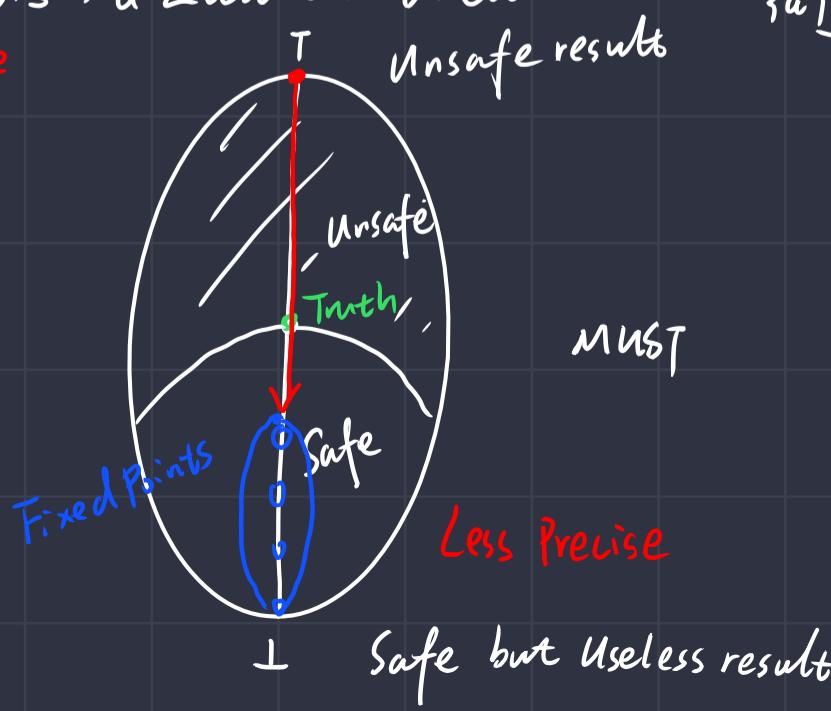


Safe  $\rightarrow$  Less Precise

Fixed Points  
Least Fixed Point

MAY

No definitions can reach  $\perp$  Unsafe result



MUST

Less Precise



## 9 Meet-Over-All-Paths Solution (MOP)

$$MOP[S_i] = \bigcup_{P \in \text{Paths from Entry to } S_i} F_p(\text{OUT}[entry])$$

A Path P from Entry to  $S_i$

$$\text{Ours} = F(x \cup y)$$

$$MOP = F(x) \cup F(y)$$

$$\Rightarrow MOP \subseteq \text{Ours}$$

Bit-vector or Gen/k.4

distributive

$\exists F$  is distributive  $F(x \cup y) = F(x) \cup F(y)$ ,  $MOP = \text{Ours}$

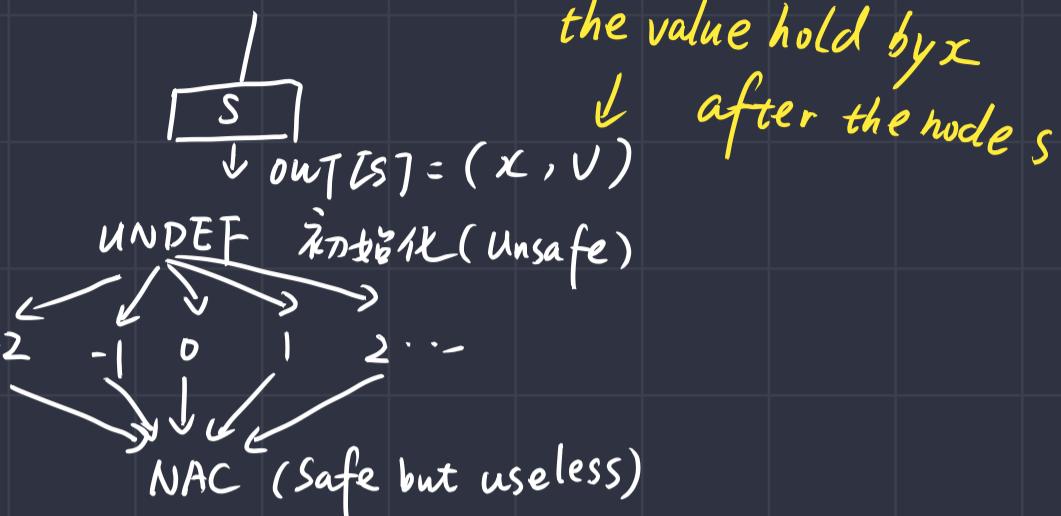
## 10 Constant Propagation 常量传播

$P \rightarrow \{x\}$ ,  $x$  是否保证持有常数

D, L, F

Forward

• 值域: ...



### Meet Operator $\sqcap$

$$NAC \sqcap v = NAC$$

UNDEF  $\sqcap v = v$  — Uninitialized variables are not focus in our constant program analysis.

$$c \sqcap c = c$$

$$c_1 \sqcap c_2 = NAC$$

/ 通配符

### Transfer Function $S: x = \dots, F: \text{OUT}[S] = \text{gen} \cup (\text{IN}[S] - \{(x, -)\})$

$$S: x = c; \quad \text{gen} = \{(x, c)\}$$

$$S: x = y; \quad \text{gen} = \{(x, \text{val}(y))\}$$

$$S: x = y \text{ op } z; \quad \text{gen} = \{(x, f(y, z))\}$$

$$f(y, z) = \begin{cases} \text{val}(y) \text{ op } \text{val}(z) & // \text{val}(y), \text{val}(z) \text{ 都是常量} \\ NAC & // \text{val}(y) 或 \text{val}(z) 是 NAC \\ UNDEF & // \text{Otherwise} \end{cases}$$

不是 distributed

## 11 Worklist Algorithm (an optimization of Iterative Algorithm)

初始化

Worklist  $\leftarrow$  all basic block

while (Worklist is not empty)

pick a basic block  $B$  from Worklist

$$\text{old-OUT} = \text{OUT}[B]$$

$$\text{IN}[B] = \bigcup_{P \text{ a predecessor of } B} \text{OUT}[P];$$

$$\text{OUT}[B] = \text{gen}_B \cup (\text{IN}[B] - k.4_B);$$

$$\text{if } (\text{old-OUT} \neq \text{OUT}[B])$$

Add all successors of  $B$  to Worklist

SPFA 类似思想



## 5. Interprocedural Control-Flow Graph (ICFG)

ICFG = 各方法的CFGs + Call edges + Return edges

```
void foo() {
    bar(...); // call site
    int n=3; // return site
}
```

返回点，紧跟call site

## 6. Interprocedural Data-Flow Analysis

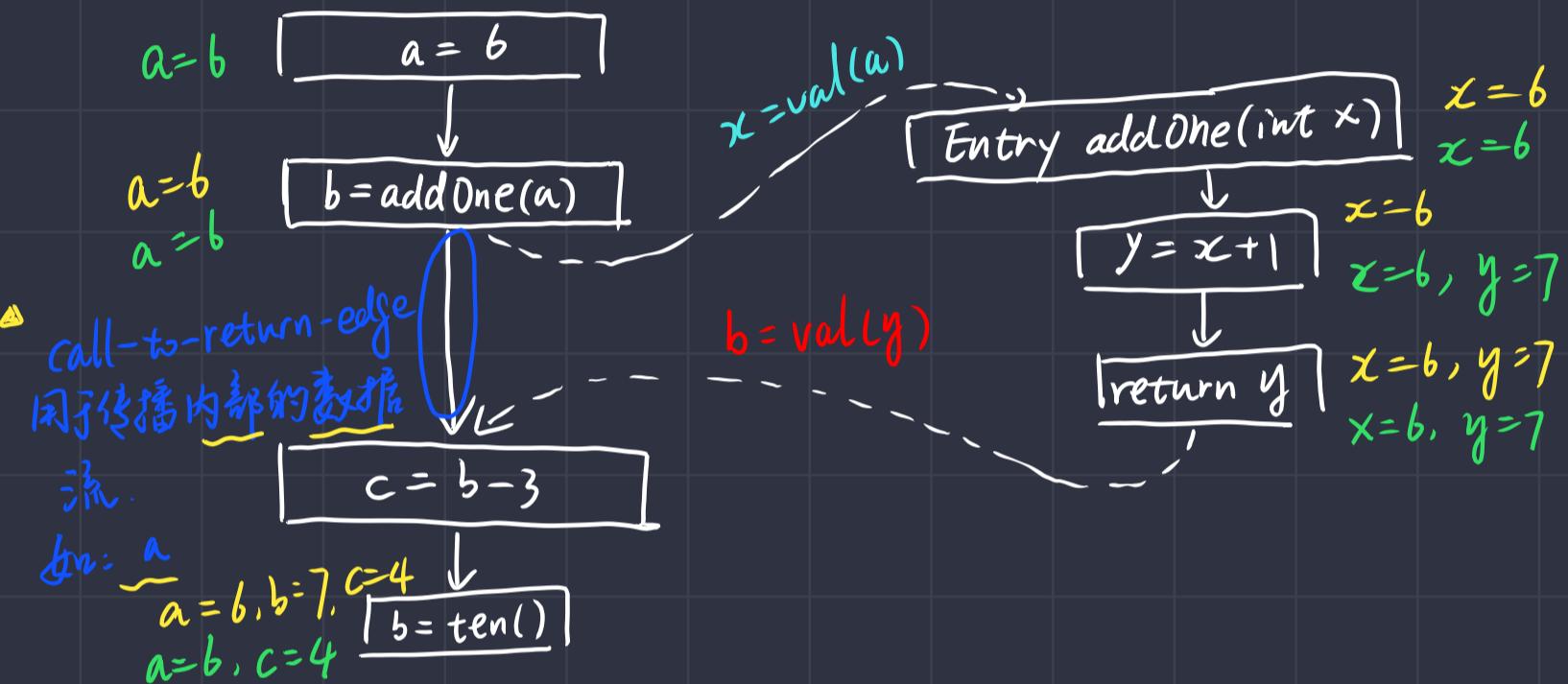
与过程中相比，还需要 edge transfer

{ call edge transfer  
Return edge transfer

过程间常量传播

{ call edge transfer 传递参数值  
Return edge transfer 传递返回值

⊕ b = addOne(a); kill掉 b (LHS)



## 四. 指针分析，回答指针能指向哪些内存地址. 对象

别名 aliases :  $p = \text{new } C();$   $p$  和  $q$  指向同一个对象

四个要素

1. Heap Abstraction 堆抽象 隔离静态分析对象数据，保证可终止

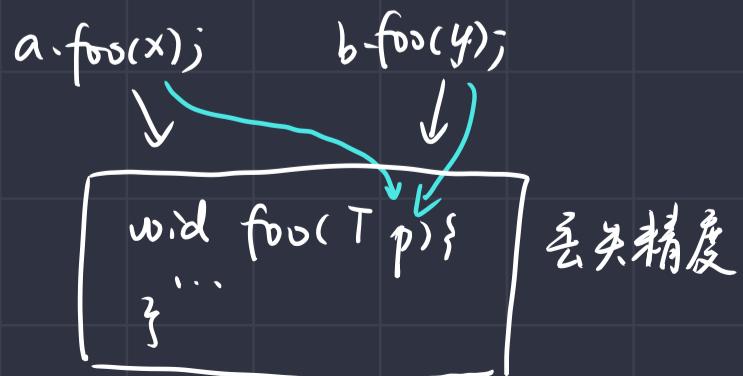
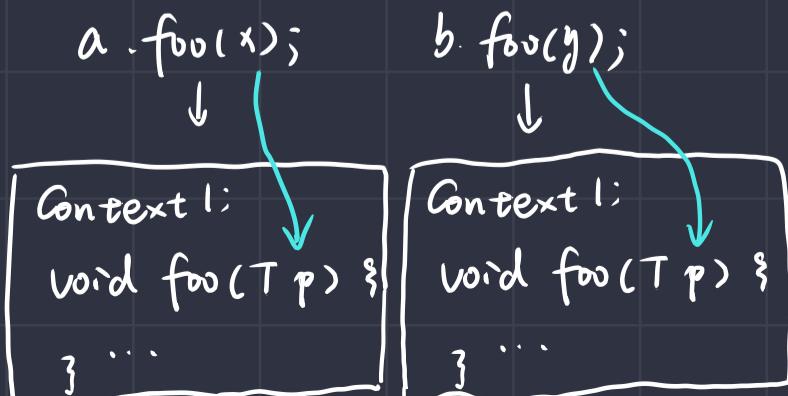
分支: Allocation-Site Abstraction ✓ 一个抽象对象代表创建点创建的所有对象

```
1 for (i=0; i<3; ++i) {
2     a = new A(); // O2 (创建点)
3     ...
4 }
```

$O_2, i=0$   
 $O_2, i=1 \rightarrow O_2$   
 $O_2, i=2$

2. Context Sensitivity 上下文敏感 ✓

上下文不敏感 ✓



### 3. Flow Sensitivity 流敏感

```

1 | C = new C(); → C → {O1}
2 | C.f = "x"; → C → {O1}
3 | S = C.f; → O1.f → {"x"}
4 | C.f = "y"; → S → {"x"}

```

非流敏感: 忽视语句顺序

C → {O<sub>1</sub>}  
O<sub>1</sub>.f → {"y"}  
S → {"z"}

✓ 主流做法 (OO 程序分析)

C → {O<sub>1</sub>}  
O<sub>1</sub>.f → {"x", "y"}  
S → {"x", "y"}

false positive

### 4. Analysis Scope 因该应该分析程序哪一部分

Whole-program ✓ Demand-driven (需求驱动)

```

1 X = new A();
2 y = x;
3 y.foo();
4 z = new T();
5 z.bar();

```

X → {O<sub>1</sub>}  
y → {O<sub>1</sub>}  
z → {O<sub>4</sub>}

若只对 Z 行感兴趣

Z → {O<sub>4</sub>}

指针分析关注的语句

Pointers in Java

Local Variable : x  
Static field : C.f      global variable  
Instance field : x.f      建模成 object with a field f  
Array element : array[i]      忽略下标, array[.] ⇒ array.arr  
建模成单个 field

Pointer-Affecting Statements :

New	x = new T()
Assign	x = y
Store	x.f = y
Load	y = x.f
Call	r = x.k(a, ...)

x.f.g.h = y;

↓  
t<sub>1</sub> = x.f;  
t<sub>2</sub> = t<sub>1</sub>.g;  
t<sub>2</sub>.h = y;

三地址码

### 五、指针分析基础

#### 1. Domains and Notations

Variable :  $x, y \in V$

Fields :  $f, g \in F$

Objects :  $O_i, O_j \in O$

Instance fields :  $O_i.f, O_j.g \in O \times F$

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

powerset of O

O 的幂集

pt = Pointer → P(O)

pt(p) : points-to set of p

## 2. Rules 规则

Kind	Statement	Rule
New	$i = \text{new } T()$ 开始	$o_i \in \text{pt}(x)$ 把 $o_i$ 加入工的指针集
Assign	$x = y$	$\frac{o_i \in \text{pt}(y)}{o_i \in \text{pt}(x)}$ 前提条件 Premises
Store	$x.f = y$	$\frac{o_i \in \text{pt}(x), o_j \in \text{pt}(y)}{o_j \in \text{pt}(o_i.f)}$ Conclusion
Load	$y = x.f$	$\frac{o_i \in \text{pt}(x), o_j \in \text{pt}(o_i.f)}{o_j \in \text{pt}(y)}$

形成一系列约束关系

关键: when  $\text{pt}(x)$  changed, 要传播改变部分给  $x$  的相关指针  $\Rightarrow$  图

## 3. Pointer Flow Graph (PFG) 指针流图(有向图)

Node: Pointer =  $V \cup (O \times F)$  a variable or a field of an abstract object

Edges: Pointer  $\times$  Pointer  $x \rightarrow y$   $x$  指向 object may flows to  $y$  指向 object

通过规则按红色箭头指示.

e.g.  $a = b; \quad ①$   
 $c.f = a; \quad ②$   
 $d = c; \quad ③$  隐依赖  
 $c.f = d; \quad ④$   
 $e = d.f; \quad ⑤$   
 $o_i \in \text{pt}(c), o_i \in \text{pt}(d)$



PFG 和 关系传播  
相互依赖. 动态增加

$\Rightarrow$  传递闭包  
transitive closure

## 4. 算法

### Pointer Analysis: Algorithms

Solve( $S$ )

```

WL = [], PFG = {}
foreach  $i: x = \text{new } T()$   $\in S$  do
    add  $\langle x, \{o_i\} \rangle$  to WL 初始化
foreach  $x = y \in S$  do
    AddEdge(y, x)

while WL is not empty do
    remove  $\langle n, pts \rangle$  from WL
     $\Delta = pts - pt(n)$  避免冗余处理
    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)
  
```

main algorithm

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

要分析语句的集合

$WL \triangleq \langle \text{Pointer}, PLO \rangle^*$

$\langle n, pts \rangle : pts$  需要加入  $pt(n)$

## 5. 指针分析处理方法调用

Kind	Statement	Rule
Call	$L: r = x.k(a_1, \dots, a_n)$	$O_i \in pt(x), m = dispatch(O_i, k)$ $O_{n-j} \in pt(a_j), 1 \leq j \leq n$ $O_v \in pt(m_{ret})$ $\Delta O_i \in pt(m_{this})$ $O_{n-j} \in pt(m_{pj}), 1 \leq j \leq n$ $O_v \in pt(r)$ <p style="margin-left: 20px;">↑ 传递值 PFG: <math>r \leftarrow m_{ret}</math></p> <p style="margin-left: 20px;">↑ 定义 <math>a_j \rightarrow m_{pj}</math></p>

## 算法: Algorithms

```

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( $x, o_i$ )**  $o_i$  是流到  $x$  的新对象  
**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( $m_{ret}, r$ )

115

解目标方法

## 六. 上下文敏感指针分析 Context-Sensitive Pointer Analysis

C.I. (上下文不敏感)



C.S. call-site sensitivity  $id \rightarrow [1] \text{ 和 } [2]$

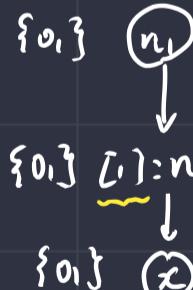
⇒ 不唯一，速度快

$x = id(n_1);$   
 $y = id(n_2);$   
 $int i = x.get();$   
 $Number id(Number n);$   
 $}$

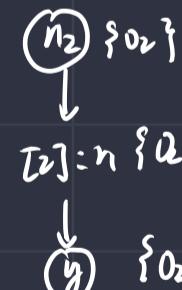
### 1. Cloning-Based Context Sensitivity 方法、变量加上下文

方法、变量加上下文

Cl.m



Cl.v



共享

### 2. Context-Sensitive Heap 上下文敏感堆区抽象对象

### 3. Domains and Notations

Context  $c, c', c'' \in C$  上下文

Context-Sensitive method:  $c:m \in C \times M$

Context-Sensitive variables:  $c:x, c':y \in C \times V$

Context-Sensitive objects:  $c:o_i, c':o_j \in C \times O$

Fields:  $f, g \in F$

Instance fields:  $c: o_i.f, c': o_j.g \in C \times O \times F$

Context-Sensitive pointers:  $CSpointer = (C \times V) \cup (C \times O \times F)$

$pt: CSpointer \rightarrow P(C \times O)$

## 4. KPI Rules

Kind	Statement	Rule (under context $c$ )	PFG edges
New	$i: x = \text{new } T()$	$c: o_i \in \text{pt}(c: x)$	
Assign	$x = y$	$c': o_i \in \text{pt}(c: y)$	$c: x \leftarrow c: y$
Store	$x.f = y$	$c': o_i \in \text{pt}(c: x), c'': o_j \in \text{pt}(c: y)$	$c'': o_j \in \text{pt}(c': o_i.f)$
Load	$y = x.f$	$c': o_i \in \text{pt}(c: x), o_j \in \text{pt}(c': o_i.f)$	$c'': o_j \in \text{pt}(c: y)$
Call	$l: r = x.k(a_1, \dots, a_n)$	$c': o_i \in \text{pt}(c: x), m = \text{Dispatch}(o_i, k), c^t = \text{Select}(c, l, c': o_i, m)$ $c'': o_n \in \text{pt}(c: a_j), 1 \leq j \leq n$ $c''' : o_v \in \text{pt}(c^t : m_{\text{ret}})$	目标方法上下文 $c': o_i \in \text{pt}(c: x)$ $c'': o_n \in \text{pt}(c^t : m_{pj}), 1 \leq j \leq n$ $c''' : o_v \in \text{pt}(c^t : m_{\text{ret}})$

## 5. PFG with C.S.

Nodes:  $\text{CSPointer} = (C \times V) \cup (C \times O \times F)$

6. 算法: Edges: CSPointer × CSPointer

## C.S. Pointer Analysis: Algorithm

```

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( $n, \Delta$ ) § C.1.6 扩散-凝聚
    if  $n$  represents a variable  $c: x$  then
        foreach  $c': o_i \in \Delta$  do
            foreach  $x.f = y \in S$  do
                AddEdge( $c: y, c': o_i.f$ )
            foreach  $y = x.f \in S$  do
                AddEdge( $c': o_i.f, c: y$ )
        ProcessCall( $c: x, c': o_i$ )
    
```

$S$  Set of **reachable** statements  
 $S_m$  Set of **statements in method  $m$**   
 $RM$  Set of **C.S. reachable** methods  
 $CG$  **C.S. call graph edges**

```

AddReachable( $c: m$ )
if  $c: m \notin RM$  then
    add  $c: m$  to  $RM$ 
 $S \cup= S_m$ 
foreach  $i: x = \text{new } T() \in S_m$  do
    add  $\langle c: x, \{c: o_i\} \rangle$  to WL
foreach  $x = y \in S_m$  do
    AddEdge( $c: y, c: x$ ) § C.1.5 扩散-凝聚

```

```

ProcessCall( $c: x, c': o_i$ )
foreach  $l: r = x.k(a_1, \dots, a_n) \in S$  do
     $m = \text{Dispatch}(o_i, k)$ 
     $c^t = \text{Select}(c, l, c': o_i)$ 
    add  $\langle c^t : m_{this}, \{c': o_i\} \rangle$  to WL
    if  $c: l \rightarrow c^t : m \notin CG$  then
        add  $c: l \rightarrow c^t : m$  to  $CG$ 
    AddReachable( $c^t : m$ )
    foreach parameter  $p_i$  of  $m$  do
        AddEdge( $c: a_i, c^t : p_i$ )
        AddEdge( $c^t : m_{\text{ret}}, c: r$ )

```

$c: z \rightarrow c^t : T.\text{foo}(A, A) \in CG$

## 7. Context Sensitivity Variants



1) Call-site sensitivity  $\text{Select}(c, l, c': o_i, m) = [l', \dots, l'', l]$   
where  $c = [l', \dots, l'']$

Also called call-string sensitivity or k-CFA

适用？限制上下文长度以调用链末尾 k 个调用点 (usually  $k \leq 3$ ) 常用方法  $k=2$   
堆  $k=1$

2) Object Sensitivity  $\text{Select}(c, l, c': o_i, m) = [o_j, \dots, o_k, o_i]$   
where  $c' = [o_j, \dots, o_k]$

实际表现 (OO 语言) Object Sensitivity 表现更好

3) Type Sensitivity  $\text{Select}(c, l, c': o_i, m) = [t', \dots, t'', \text{InType}(o_i)]$   
where  $c' = [t', \dots, t'']$

class X {  
void main() {  
 Y y1 = new Y(); // 下方: TX  
 }

牺牲精度，换取速度

## 7. 轨迹分析 - 安全 Security

### 1. 信息流安全 Information Flow Security

互斥 { Access Control 访问控制 }  
Information Flow Security 线关心信息如何传播

1) 信息流 variable  $x$  信息传递到  $y \Rightarrow$  信息流  $x \rightarrow y$

$$\boxed{\begin{array}{l} a = x; \\ b.f = a; \\ y = b.f; \end{array}}$$

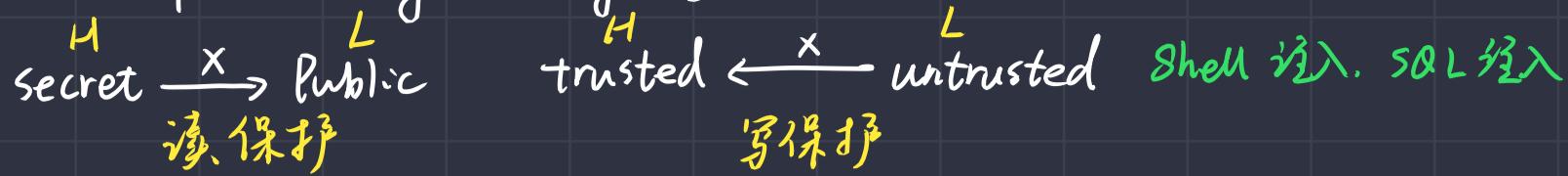
2) Security Levels (Classes) 分级 {  
H high security  
L low security  
modeled as lattice  $L \leq H$  可以是复杂的格

### 3) Information Flow Policy

Noninterference Policy :  $\text{Secret}(H) \xrightarrow{*} \text{Public}(L)$

$$X_L = Y_H \times \quad X_L = Y_L + Z_H \times$$

## 2. Confidentiality and Integrity



一致性的意义 Integrity: 正确性. 完整性. 一致性

## 3. Explicit Flows and Covert Channels

$$x_H = y_H$$

$$x_L = y_H$$

$$x_L = y_H + z_H$$

隐式流 implicit flow:

$\text{secret}_H = \text{getSecret}();$

```
if ( $\text{secret}_H < 0$ )
  public_L = 1;
else
  public_L = 0;
```

// 有信息流, 无数据流

Does Secret information Leak?

```
while ( $\text{secret}_H < 0$ ) { ... };
```

Leak, we can conclude that secret is negative if the program does not terminate

Termination channels

```
if ( $\text{secret}_H < 0$ )
  for (int i = 0; i < 1000000; ++i) { ... };
```

Leak, we can conclude that secret is negative if the program execution spends more time

Timing channels

```
if ( $\text{secret}_H < 0$ )
  throw new Exception("...");
```

Leak, we can conclude that secret is negative if we observe the exception

Exceptions

```
int sa_H[] = getSecretArray();
sa_H[secret_H] = 0;
```

Leak, exception may reveal that secret is negative

Covert/Hidden Channels

## 4. Taint Analysis 污点分析 同位素标记法

把程序数据分为两类 { Data of interest tainted data  
Other data untainted data }

两类问题  $\Rightarrow$  { Confidentiality  
Integrity }

Source  $\rightarrow$  sink Can tainted flow to a sink?

$\Rightarrow$  Which tainted data a pointer (at a sink) can point to?

拓展指针分析域

Variable:  $x, y \in V$

Fields:  $f, g \in F$

Objects:  $O_i, O_j \in O$

Instance fields:  $O_i.f, O_j.g \in O \times F$

Pointers:  $\text{Pointer} = V \cup (O \times F)$

Tainted data:  $t_i, t_j \in T \subset O$   $i$  表示 call site  $i$

Inputs: Sources: a set of source methods

Sinks: a set of sink methods

Outputs: TaintFlows

元素  $(t_i, m)$  表示 tainted data from call site  $i$  may flow to sink method  $m$

Rules: Call

## ① 处理 Source

Call  $l : r = x \cdot k(a_1, \dots, a_n)$

$$\frac{l \rightarrow m \in CG \\ m \in Sources}{t_l \in pt(r)}$$

(传播规则和指针分析完全一致)

## ② 处理 Sink

Call  $l : r = x \cdot k(a_1, \dots, a_n)$

$$\frac{l \rightarrow m \in CG \\ m \in Sinks \\ \exists i, 1 \leq i \leq n : t_j \in pt(a_i)}{\langle t_j, m \rangle \in TaintFlows}$$

## 八. 基于 Datalog 的程序分析

高级语言 (C++ Java) VS 声明式语言 (SQL)

1. Datalog: Prolog 子集、声明式逻辑编程语言

= Data + Log 无副作用、无控制流、无函数、非图灵完备

1) • Predicate  $\rightarrow$  table. 关系.

• fact: Age("Xiaoming", 18) 存在于表中

• Atoms E.g. Age("Xiaoming", 18), age  $\geq 18$ , Age(person, age)

• Rules:  $H \leftarrow \underbrace{B_1, B_2, \dots, B_n}_{\text{atom atom}} \quad \text{逻辑与}$

person	age
Xiaoming	18
Alan	16

例程: Adult(person)  $\leftarrow$

Age(person, age),  
age  $\geq 18$ .

2) Predicate 命令  $\begin{cases} EDB & \text{原始的} \\ IDB & \text{推导产生的} \end{cases}$

$H \leftarrow \underbrace{B_1, B_2, B_3}_{IDB} \quad EDB \text{ 或 } IDB$

3) 表达逻辑或 1. 写两条 Rules 2. SportFan(person)  $\leftarrow$

Hobby(person, "jogging");  
Hobby(person, "Swimming").

表达逻辑非 !  $B_2(X_2, X_4)$

e.g. 找补考学生: MakeUpExam Std(student)  $\leftarrow$   
Student(student),  
! PassedStd(student).

4) 支持递归: e.g. 可达性

Reach(from, to)  $\leftarrow$   
Edge(from, to).

Reach(from, to)  $\leftarrow$   
Reach(from, node),  
Edge(node, to).

rule is safe,  $\forall x \exists y$  且只出现在一个 non-negated relational atom 产生“无害”产生

不能出现  $A(x) \leftarrow B(y), x > y$        $A(x) \leftarrow B(y), !C(x, y)$

Recursion and Negation     $A(x) \leftarrow B(x), !A(x)$

## 2. 通用 Datalog 的指针分析

### Datalog Model for Pointer Analysis

Kind	Statement
New	$i: x = \text{new } T()$
Assign	$x = y$
Store	$x.f = y$
Load	$y = x.f$

EDB

$\text{New}(x : V, o : O)$

$\text{Assign}(x : V, y : V)$

$\text{Store}(x : V, f : F, y : V)$

$\text{Load}(y : V, x : V, f : F)$

IDB

$\text{VarPointsTo}(v : V, o : O)$

e.g., fact  $\text{VarPointsTo}(x, o_i)$  represents  $o_i \in pt(x)$

Variables: V  
Fields: F  
Objects: O

$\text{FieldPointsTo}(oi : O, f : V, oj : O)$

e.g., fact  $\text{FieldPointsTo}(o_i, f, o_j)$  represents  $o_j \in pt(o_i.f)$

$\frac{x = \text{new } T()}{o_i \in pt(x)}$

$\text{VarPointsTo}(x, o) \leftarrow \text{New}(x, o).$

$\frac{x = y}{\begin{array}{l} o_i \in pt(y) \\ o_i \in pt(x) \end{array}}$

$\text{VarPointsTo}(x, o) \leftarrow \text{Assign}(x, y)$   
 $\text{VarPointsTo}(y, o)$ .

Datalog 处理方法调用 Label

EDB:  $VCall(l: S, x: V, k: M)$

$\text{ThisVar}(m: M, this: V)$

$Dispatch(o: O, k: M, m: M)$

Method

IDB:  $\text{Reachable}(m: M)$

$\text{CallGraph}(l: S, m: M)$

Rule

$o_i \in pt(x), m = dispatch(o_i, k)$

$o_n \in pt(a_j), 1 \leq j \leq n$

$o_v \in pt(m_{ret})$

实参值

$\Delta o_i \in pt(m_{this})$

形参 PFG

$o_n \in pt(m_{pj}), 1 \leq j \leq n$

连接

$o_v \in pt(r)$

$a_j \rightarrow m_{pj}$

参数

参数

PFG:  $r \leftarrow m_{ret}$

①  $\text{VarPointsTo}(this, o),$

$\text{Reachable}(m),$

$\text{CallGraph}(l, m) \leftarrow$

$VCall(l, x, k),$

$\text{VarPointsTo}(x, o),$

$\text{Dispatch}(o, k, m),$

$\text{ThisVar}(m, this).$

EDB:  $\text{Argument}(l: S, i: N, a: V)$

②  $\text{VarPointsTo}(pi, o) \leftarrow$

$\text{CallGraph}(l, m),$

$\text{Argument}(l, i, ai),$

$\text{Parameter}(m, i, pi),$

$\text{VarPointsTo}(ai, o)$

Parameter(m: M, i: N, p: V)

EDB: MethodReturn(m: M, ret: V)

CallReturn(l: S, r: V) 接收返回值变量

③ VarPointsTo(r, o) ←

callGraph(L, m),

Method(m, ret),

VarPointTo(ret, o),

CallReturn(l, r).

3. Datalog 语义分析 → PPT

九. CFL-Reachability and IDFS 程序分析框架

动态运行时不属执行到

Infeasible Paths: Paths in CFG that do not correspond to actual executions

Realizable Paths: The paths in which "returns" are matched with corresponding "calls"

补充(回忆) 上下文无关文法  $S \rightarrow aSb$  左边可被右边替换无论 S 在哪里  
 $S \rightarrow \epsilon$  ("不吃肉")

上下文有关文法:  $zS \rightarrow zaSb$

1. CFL-Reachability

括号匹配

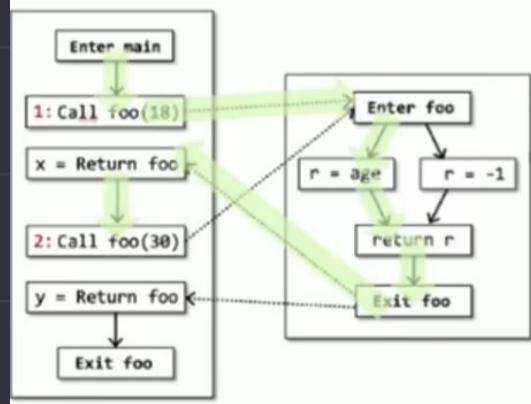
CFL-Reachability

A path is considered to connect two nodes A and B, or B is reachable from A, only if the concatenation of the labels on the edges of the path is a word in a specified context-free language.

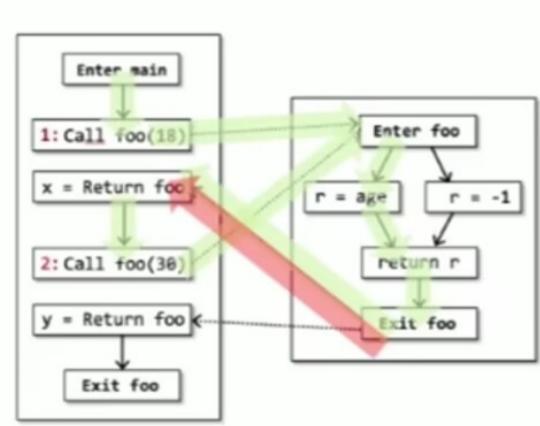
- A valid sentence in language L must follow L's grammar.
- A context-free language is a language generated by a **context-free grammar (CFG)**.

L(realizable):

realizable → matched realizable  
 $\rightarrow (\_ \text{realizable})_1$   
 $\rightarrow \epsilon$   
matched → ( $\_ \text{matched}$ ) $_1$   
 $\rightarrow e$   
 $\rightarrow \epsilon$   
 $\rightarrow \text{matched matched}$



$e(1eee)_1e \in L(\text{realizable})$



交错路径:

$e(1eee)_1e(2eee)_1$

$\notin L(\text{realizable})$

2. IDFS : A Program Analysis Framework via Graph Reachability (无传播 → 图可达性)

用于过程间数据流分析 → (Interprocedual, Finite, Distribute, Subset Problem)

MOP

⇒ MRP Solution

TPP  $F(x \cap y) = F(x) \cap F(y)$

复杂MOP:  $p f_p = f_n \circ \dots \circ f_2 \circ f_1$

$MOP_n = \bigcup_{p \in \text{Paths}(\text{start}, n)} p f_p (\perp)$

$$MRP_n = \bigcup_{p \in RPaths(start, n)} Pf_p(\perp) \quad MRP_n \subseteq mOP_n$$

IDFS

## Overview of IFDS

Given a program P, and a dataflow-analysis problem Q

- Build a **supergraph  $G^*$**  for P and define **flow functions** for edges in  $G^*$  based on Q
- Build **exploded supergraph  $G^{\#}$**  for P by transforming flow functions to **representation relations (graphs)**
- Q can be solved as graph reachability problems (find out MRP solutions) via applying **Tabulation algorithm** on  $G^{\#}$

$$\lambda S. \{a\} : 1S.S \\ \begin{matrix} o & o & o \\ o & a & b \\ o & o & o \end{matrix}$$

$$\lambda S. (S - \{a\}) \cup \{b\} \\ \begin{matrix} o & a & b & c \\ o & o & o & o \\ o & o & o & o \end{matrix}$$

(有点记不出来了，可参考PPT或笔记网站)

2(D+1)结点的图

representation relation of flow function  $f$ :

$$R_f \subseteq (D \cup O) \times (D \cup O)$$

$$R_f = \{(o, o)\}$$

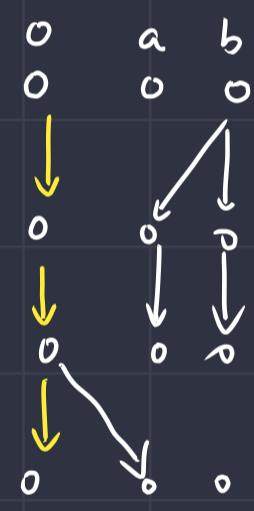
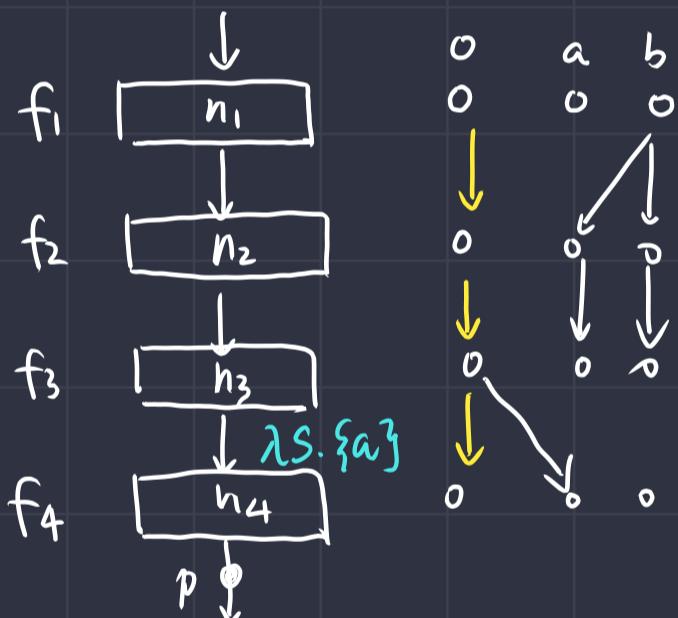
$$\cup \{(o, y) \mid y \in f(\phi)\}$$

$$\cup \{(x, y) \mid y \notin f(\phi) \text{ and } y \in f(\{x\})\}$$

Edge:  $o \rightarrow o$

Edge:  $o \rightarrow d_i$

Edge:  $d_i \rightarrow d_j$



### 3. IFDS 中的 Distributivity (可分配性)

IFDS 不能直接用于常量传播、指针分析

Distributivity:  $F(x \wedge y) = F(x) \wedge F(y)$ , transfer function 只能处理一个值

结论：分析需要依赖多个输入  $\Rightarrow$  不能用 IFDS

### 4. Soundness (2015年提出)

1. Soundness  $\rightarrow$  保真估计 所有程序行为(无论什么输入、运行行为)

Hard Language Features 静态分析难以分析的语言特征.

- 例如：
- Java 及其 native code
  - JavaScript eval、DOM
  - C/C++ Pointer arithmetic

可能指向的内存不确定

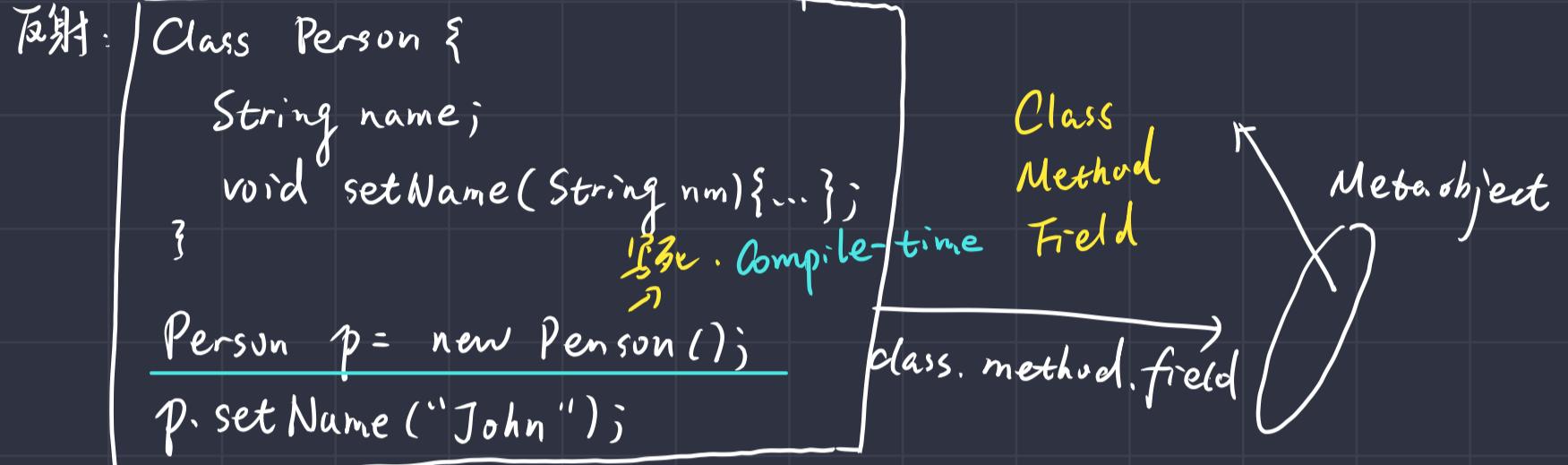
函数指针

## 2. Soundness

soundly analysis typically means analysis mostly sound, with well-identified unsound treatments to hard/specific language features.

## 3. Hard language features in Java 分析

### 1) Java Reflection 分析的噩梦 - 公开对象的问题



代码:

```
Class c = Class.forName("Person");
```

Method m = c.getMethod("setName", ...); m.invoke(a, ...);

Object p = c.newInstance(); 等价 Run-time

反射有利于程序的松耦合

分析反射: String Constant analysis + Pointer Analysis 常量分析

↓  
使用中推导: 利用反射语义和Java类型系统

### 2) Native Code

```
private native void writeBytes(byte b[], int off, int len)
throws IOException. 声明. 实现由 C/C++ 完成
```

JNI Java Native Interface (A function module of JVM)

方法: 手动建模

```
System.arraycopy(src, sp, dest, dp, ln);
```

↓  
for (int i=0; i<ln; i++)  
 dest[dp+i] = src[sp+i];