



软件分析与架构设计

指针分析

何冬杰
重庆大学

指针分析 (Pointer Analysis)

□ 指针分析计算指向关系: $Pts: \text{Var} \rightarrow \mathcal{P}(\text{OBJs})$

- $Pts(x)$: 确定指针变量 x 运行时的可能取值 (*may - points - to*)
- **Var**: 程序中的指针变量; **OBJs**: allocation site 抽象, 以表示从同一个分配点创建的所有运行时对象
- *Must-points-to* 关系不在考虑范围内
- 指向关系可以用来计算别名 (alias) 关系
 - $\text{Alias}(x, y)$ if and only if $Pts(x) \cap Pts(y) \neq \emptyset$

➤ 计算收集语义

□ 限定需要分析的3地址IR指令:

- 暂时只支持流不敏感 (flow-insensitive) 分析, 不考虑 goto 语句
- 暂时只涉及过程内 (intra-procedural) 分析, 不考虑 call 语句

```
Inst ::=  $x := y \mid x := \&p \mid x := *p \mid *p := x \mid x := y.f \mid x.f := y \mid p := q \mid x := \text{malloc}()$ 
```

指针分析 (Pointer Analysis)

□例子:

1. `q := malloc() // 01`

2. `p := malloc() // 02`

3. `p := q`

4. `r := &p`

5. `s := malloc() // 03`

6. `*r := s`

7. `t := &s`

8. `u := *t`

Strong Update

Pts(q) = {O1}

Pts(p) = {O2}

Pts(p) = {O1}

Pts(r) = {&p}

Pts(s) = {O3}

Pts(p) = {O3}

Pts(t) = {&s}

Pts(u) = {O3}

Weak Update

Pts(q) = {O1}

Pts(p) = {O2}

Pts(p) = {O1, O2}

Pts(r) = {&p}

Pts(s) = {O3}

Pts(p) = {O1, O2, O3}

Pts(t) = {&s}

Pts(u) = {O3}

流不敏感分析不区分指令执行顺序, 只有弱更新, 没有强更新

问题: 如何设计一个指向分析算法?

收集语义 (Collecting Semantics)

The **collecting semantics** requires us to know each execution of the program, assuming a (possibly infinite) trace for each run.

□ Infinite Domain value

$x \mapsto \{-\infty, \dots, -2, -1, 0, 1, 2, \dots, +\infty\}$

```
1. Read x
2. y = 2 * x
3. z = 0
4. if y > 0 goto 6
5. z = 2 * x - y
6. ...
```

□ Infinite execution path

```
1. read x
2. t = 0
3. if x < 0 goto 7
4. x = x - 1
5. t = t + x
6. goto 3
7. print t
```

循环次数取决于输入x

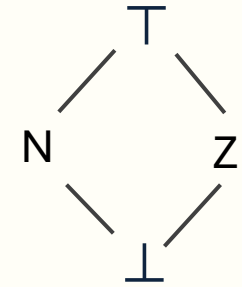
□ 需要对数据域和控制流等进行抽象

抽象域 (Abstraction)

□ 零分析 (Zero Analysis)

➤ 具体域: $\mathbb{Z} : \{-\infty, \dots, -2, -1, 0, 1, 2, \dots, +\infty\}$

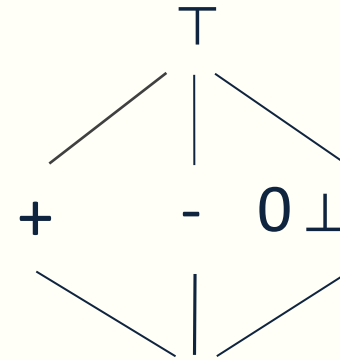
➤ 抽象域: $\{\perp, N, Z, \top\}$



□ 符号分析 (Sign Analysis)

➤ 具体域: $\mathbb{Z} : \{-\infty, \dots, -2, -1, 0, 1, 2, \dots, +\infty\}$

➤ 抽象域: $\{\perp, +, -, 0, \top\}$

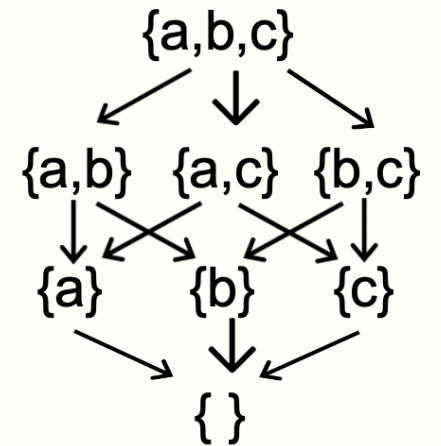


□ 指针分析 (Pointer Analysis)

➤ 具体域: 运行时堆对象或栈上位置

➤ 抽象域: $\&t$, $alloca$, 或 $new A$ 对应行号表示

○ Allocation site abstraction



□ 抽象域元素往往构成格 (lattice)

偏序 (Partial orders)

Given a set S , a partial order \sqsubseteq is a binary relation on S that satisfies:

- reflexivity: $\forall x \in S: x \sqsubseteq x$
- transitivity: $\forall x, y, z \in S: x \sqsubseteq y \wedge y \sqsubseteq z \rightarrow x \sqsubseteq z$
- anti-symmetry: $\forall x, y \in S: x \sqsubseteq y \wedge y \sqsubseteq x \rightarrow x = y$

□ (S, \sqsubseteq) 称为偏序集 (poset)

- 一个集合的幂集和 \subseteq (subset) 构成一个偏序集
- 偏序集中的两个元素可能无法比较, e.g., 右图中的 $\{a\}$ 和 $\{b\}$

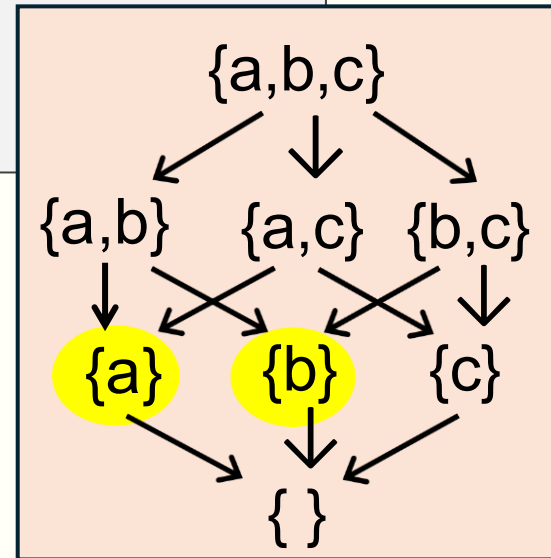
□ Upper and Lower bounds

Let $X \subseteq S$ be a subset, we say that $y \in S$ is an **upper bound** ($X \sqsubseteq y$) when $\forall x \in X: x \sqsubseteq y$;

We say that $y \in S$ is a **lower bound** ($y \sqsubseteq X$) when $\forall x \in X: y \sqsubseteq x$;

A **least upper bound** $\sqcup X$ is defined by $X \sqsubseteq \sqcup X \wedge \forall y \in S: X \sqsubseteq y \Rightarrow \sqcup X \sqsubseteq y$;

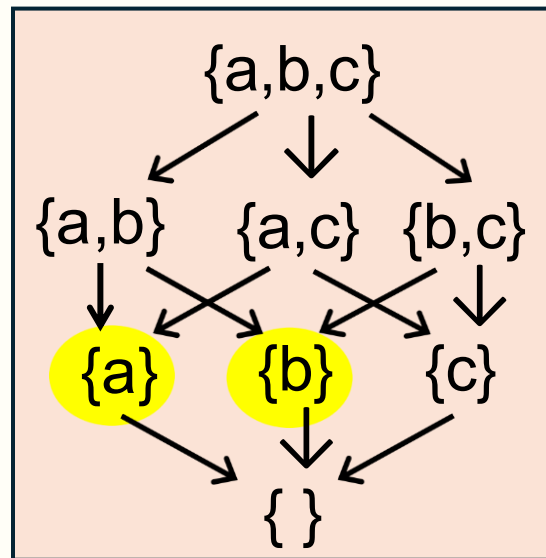
A **greatest lower bound** $\sqcap X$ is defined by $\sqcap X \sqsubseteq X \wedge \forall y \in S: y \sqsubseteq X \Rightarrow y \sqsubseteq \sqcap X$.



格 (Lattices)

□格与半格

- Given a poset (S, \sqsubseteq) , $\forall a, b \in S$,
- if $a \sqcup b$ exists, then (S, \sqsubseteq) is called a join semilattice;
- if $a \sqcap b$ exists, then (S, \sqsubseteq) is called a meet semilattice;
- if both exist, then (S, \sqsubseteq) is called a lattice.

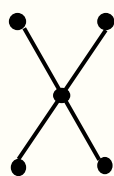


完全格

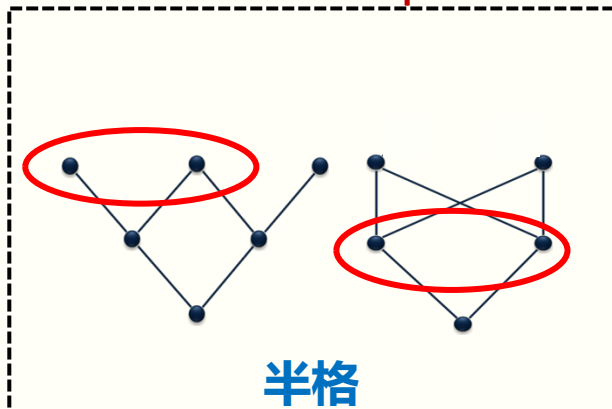
□完全格

- If $\forall X \subseteq S$, $\sqcap X$ and $\sqcup X$ exist, (S, \sqsubseteq) is called a complete lattice.
- A complete lattice must have a unique largest element $\top = \sqcup S$ and a unique smallest element $\perp = \sqcap S$
- A finite lattice is complete if \top and \perp exist (一般有限格都是完全格)

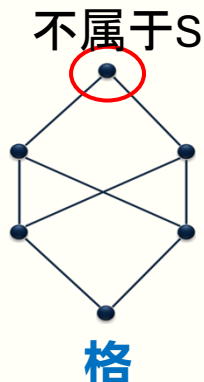
□举例:



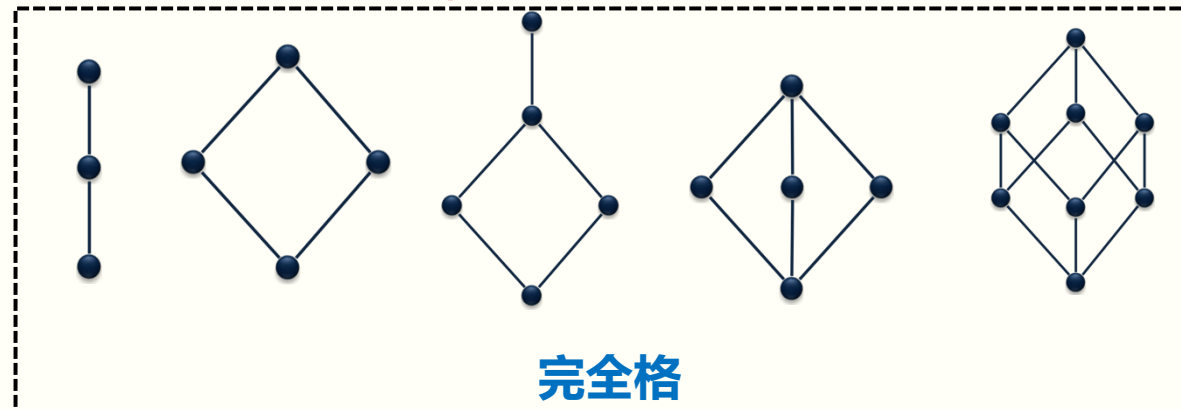
非格



半格



格



完全格

格 (Lattices)

□ **height: the length of the longest path in the lattice**

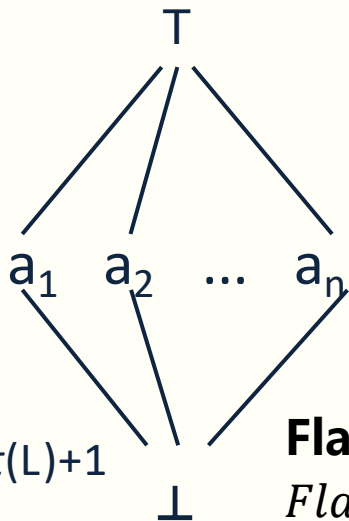
➤ For complete lattice: the length from \perp to \top

□ **Powerset lattice:**

➤ Powerset of every finite set A defines a complete lattice

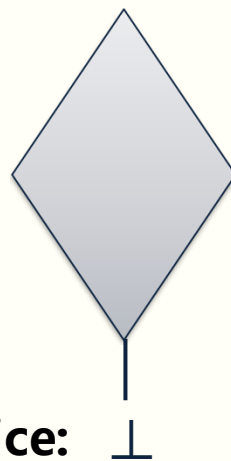
➤ $(\mathcal{P}(A), \subseteq), \perp = \emptyset, \top = A, x \sqcup y = x \cup y, x \sqcap y = x \cap y$

□ **Flat lattice and Lift lattice**



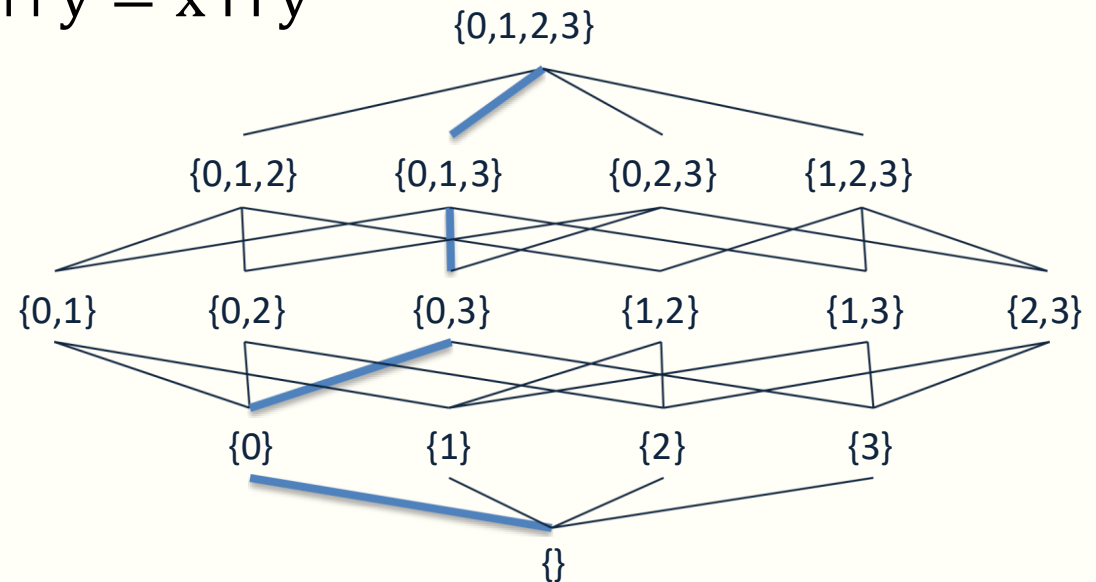
Lift lattice:

$$height(lift(L)) = height(L) + 1$$



Flat lattice:

$Flat(A)$, where A is a set, $height(flat(A)) = 2$



for $A = \{0,1,2,3\}$

格 (Lattices)

□ Product lattice:

- $L_1 \times L_2 \times \cdots \times L_n = \{(x_1, x_2, \dots, x_n) \mid x_i \in L_i\}$
- $L_i (1 \leq i \leq n)$ is a complete lattice
- \sqsubseteq, \sqcup , and \sqcap are computed pointwise
- $\text{Height}(L_1 \times L_2 \times \cdots \times L_n) = \sum_{i=1}^n \text{Height}(L_i)$

□ Map Lattice:

- $A \rightarrow L = \{[a_1 \mapsto x_1, a_2 \mapsto x_2, \dots] \mid a_i \in A \wedge x_i \in L\}$
- \sqsubseteq, \sqcup , and \sqcap are computed pointwise
- $\text{Height}(A \rightarrow L) = |A| \cdot \text{Height}(L)$

□ 例子：指针分析

- A is the set of program variables
- L is the powerset of allocation site abstractions (objects)

Andersen's Analysis

□ Key idea: cast as a constraint-solving problem

- One subset constraint per instruction
- Domain abstraction and function abstraction

$$\frac{}{\llbracket n: p := \text{malloc}() \rrbracket \hookrightarrow l_n \in p} \text{ malloc}$$

$$\frac{}{\llbracket p := \&x \rrbracket \hookrightarrow l_x \in p} \text{ address-of}$$

$$\frac{}{\llbracket p := q \rrbracket \hookrightarrow p \supseteq q} \text{ copy}$$

$$\frac{}{\llbracket *p := q \rrbracket \hookrightarrow *p \supseteq q} \text{ assign}$$

$$\frac{}{\llbracket p := *q \rrbracket \hookrightarrow p \supseteq *q} \text{ dereference}$$

- The constraint solver can give the most precise solution.

- Why? (optional)

- Constraints are equivalent to abstract functions

$$e \in v \equiv Pts(v) = Pts(v) \cup \{e\}$$

$$x \supseteq y \equiv Pts(x) = Pts(x) \cup Pts(y)$$

- Abstract function is monotone
 - Abstract domain is finite
 - Objects are finite
 - $(\text{Var} \rightarrow \mathcal{P}(\text{Objects}))$ forms a map lattice
 - Abstract function work on this lattice from \perp will reach the **least fixed point** (e.g., the most precise solution)

Kleene's fixed-point theorem

□ Monotonicity

A function $f: L \rightarrow L$ (L is a lattice) is monotonic if $\forall x, y \in L,$
 $x \sqsubseteq y \rightarrow f(x) \sqsubseteq f(y)$

□ Fixed point

$x \in L$ is a fixed point of function $f: L \rightarrow L$ iff $f(x) = x$

□ Kleene's fixed-point theorem

In a complete lattice with finite height, every monotone function f has a *unique least fixed-point*:

$$lfp(f) = \bigsqcup_{i \geq 0} f^i(\perp)$$

Proof of existence and uniqueness

□ Existence

- Clearly, $\perp \sqsubseteq f(\perp)$
- Since f is monotone, we also have $f(\perp) \sqsubseteq f^2(\perp)$; By induction, $f^i(\perp) \sqsubseteq f^{i+1}(\perp)$
- This means that $\perp \sqsubseteq f(\perp) \sqsubseteq f^2(\perp) \sqsubseteq \dots \sqsubseteq f^i(\perp) \sqsubseteq \dots$ is an increasing chain
- L has finite height, so for some k : $f^k(\perp) = f^{k+1}(\perp)$
- If $x \sqsubseteq y$ then $x \cup y = y$. So $\text{lfp}(f) = f^k(\perp)$

□ Uniqueness of **least**

- Assume that x is another fixed-point: $x = f(x)$
- Clearly, $\perp \sqsubseteq x$
- By induction and monotonicity, $f^i(\perp) \sqsubseteq f^i(x) = x$
- In particular, $\text{lfp}(f) = f^k(\perp) \sqsubseteq x$, i.e. $\text{lfp}(f)$ is least
- Uniqueness then follows from anti-symmetry
 - Suppose x is another least fixed-point, we have $x \sqsubseteq \text{lfp}(f)$

Andersen's Analysis

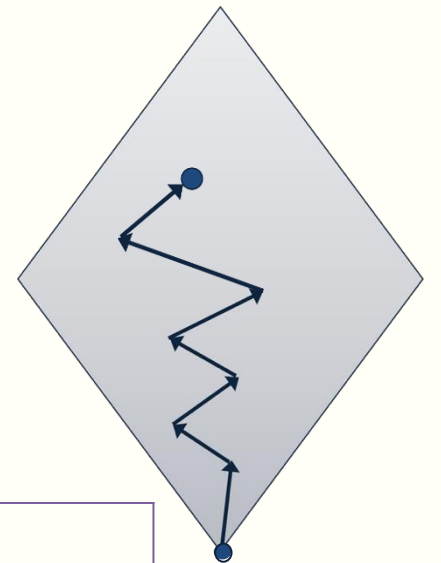
Let $f_1, f_2: L \rightarrow L$ (L is a lattice) be two monotonic functions, then their **composition** $f_1 \circ f_2$ is also **monotonic**.

Let $f: L \rightarrow L$ ($L = (\text{Var} \rightarrow \mathcal{P}(\text{OBJs}))$ is a map lattice) be the constraint solver of Andersen-style pointer analysis, obviously, f is the composition of many constraint-equivalent abstract functions, by Kleene's fixed-point theorem, the **solution of Andersen's analysis is $\text{lfp}(f)$** .

□ **The time complexity of $\text{lfp}(f)$ depends on:**

- the height of the lattice
- the cost of computing f
- the cost of testing equality

```
x = ⊥;  
do {  
  t = x;  
  x = f(x);  
} while (x ≠ t);
```



Implementation: [TIP/src/tip/solvers/FixpointSolvers.scala](https://github.com/TipToolbox/TipToolbox/blob/master/src/tip/solvers/FixpointSolvers.scala)

A more efficient implementation

□ The Cubic Framework ($O(n^3)$)

- A set of tokens $T = \{t_1, t_2, \dots, t_k\}$,
- A collection of constraint variables $V = \{x_1, \dots, x_n\}$
- A collection of constraints of these forms:

- $t \in x$;
- $x \subseteq y$; inclusion constraints
- $t \in x \implies y \subseteq z$; conditional constraints

□ Solution: reachability on a constraint graph

- Each variable is mapped to a node
- Each node has a bitvector in $\{0,1\}^k$, initially set to all 0's
- Each bit has a list of pairs of variables, modeling conditional constraints
- The edges model inclusion constraints

求解过程：传播+约束维护

□ $x.sol \subseteq T$:

- the set of tokens for x (the 1's in its bitvectors)

□ $x.succ \subseteq V$:

- the successors of x (the edges)

□ $x.cond(t) \subseteq V \times V$:

- the conditional constraints for x and t

□ $W \subseteq T \times V$:

- a worklist (initially empty)

- $t \in x$
addToken(t, x)
propagate()
- $x \subseteq y$
addEdge(x, y)
propagate()
- $t \in x \rightarrow y \subseteq z$
if $t \in x.sol$
 addEdge(y, z)
 propagate()
else
 add (y, z) to $x.cond(t)$

addToken(t, x):

```
if  $t \notin x.sol$   
  add  $t$  to  $x.sol$   
  add ( $t, x$ ) to  $W$ 
```

addEdge(x, y):

```
if  $x \neq y \wedge y \notin x.succ$   
  add  $y$  to  $x.succ$   
  for each  $t$  in  $x.sol$   
    addToken( $t, y$ )
```

propagate():

```
while  $W \neq \emptyset$   
  pick and remove ( $t, x$ ) from  $W$   
  for each ( $y, z$ ) in  $x.cond(t)$   
    addEdge( $y, z$ )  
  for each  $y$  in  $x.succ$   
    addToken( $t, y$ )
```

$O(n^3)$ 的时间复杂度, n 是变量数量

Andersen's Analysis

□使用Cubic框架求解

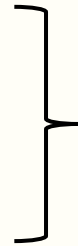
$$\frac{}{\llbracket n: p := \text{malloc}() \rrbracket \hookrightarrow l_n \in p} \text{ malloc}$$

$$\frac{}{\llbracket p := \&x \rrbracket \hookrightarrow l_x \in p} \text{ address-of}$$

$$\frac{}{\llbracket p := q \rrbracket \hookrightarrow p \supseteq q} \text{ copy}$$

$$\frac{}{\llbracket *p := q \rrbracket \hookrightarrow *p \supseteq q} \text{ assign}$$

$$\frac{}{\llbracket p := *q \rrbracket \hookrightarrow p \supseteq *q} \text{ dereference}$$



第一种形式: $t \in x$



第二种形式: $x \subseteq y$



第三种形式:

$$\&x \in p \Rightarrow q \subseteq x$$

$$\&y \in q \Rightarrow y \subseteq p$$

➤改进: SCC缩点、类型过滤 (eliminate type-incompatible Objects)

□理论复杂度 $O(n^3)$, 实际复杂度 $O(n^2)$ (Sridharan et al. [SAS,09])

Andersen's Analysis

□ 一些扩展

➤ 域敏感 (field sensitivity)

- 域不敏感: treats fields `f` as dereferences `*`
- 区分不同的域 (field), 在Java指针分析中很重要

$$\left. \begin{array}{l} \overline{\llbracket p := q.f \rrbracket} \hookrightarrow p \supseteq q.f \quad \text{field-read} \\ \overline{\llbracket p.f := q \rrbracket} \hookrightarrow p.f \supseteq q \quad \text{field-assign} \end{array} \right\} \longrightarrow$$

c/c++ 转化成第二种形式:

$$x \subseteq y$$

Java 转化成第三种形式(p.f相当于c/c++中的(*p).f):

$$0 \in q \Rightarrow 0.f \subseteq p$$

$$0 \in p \Rightarrow q \subseteq 0.f$$

➤ 支持函数调用

- 参数和返回值的传递建模成copy语句

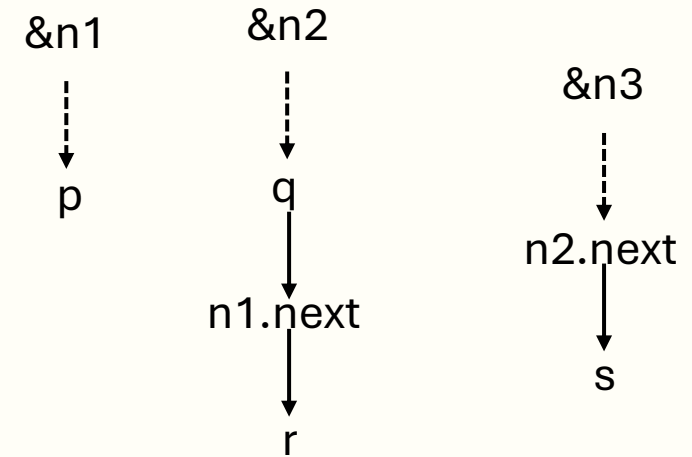
$$\begin{array}{l} r = m(a_1, \dots, a_n) \\ \dots \\ \text{return ret;} \end{array} \quad \left\{ \begin{array}{l} T \quad m(T \quad p_1, \dots, T \quad p_n) \\ \dots \\ \text{return ret;} \end{array} \right\} \quad \overline{\llbracket r = m(a_1, \dots, a_n) \rrbracket} \hookrightarrow p_1 \supseteq a_1, \dots, p_n \supseteq a_n, r \supseteq \text{ret} \quad \text{call}$$

Example Program

```
#include <stdio.h>
struct Node { struct Node *next; };
int main() {
    struct Node n1, n2, n3;
    struct Node *p = &n1;
    struct Node *q = &n2;
    n1.next = q;
    struct Node *r = p->next;
    n2.next = &n3;
    struct Node *s = r->next;
    return 0;
}
```

留了个更复杂的例子作为可选课后作业

Constraint Graph (so-called Pointer Assignment Graph)



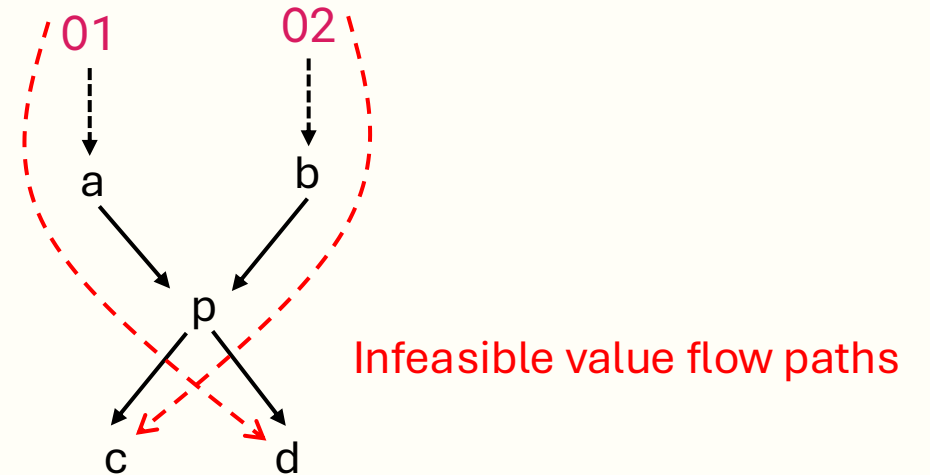
Vars	PTS	Vars	PTS
p	$\{\&n1\}$	$n2.next$	$\{\&n3\}$
q	$\{\&n2\}$	s	$\{\&n3\}$
$n1.next$	$\{\&n2\}$	r	$\{\&n2\}$

Context Sensitivity

```
1. #include <stdlib.h>
2. int *id(int *p) {
3.     return p;
4. }
5. int main() {
6.     int *a, *b, *c, *d;
7.     a = (int *)malloc(sizeof(int)); // 01
8.     c = id(a);
9.     b = (int *)malloc(sizeof(int)); // 02
10.    d = id(b);
11.    return 0;
12. }
```

- Andersen's analysis is context-insensitive
 - Joins information across callsites to same function
 - Loses precision due to modeling infeasible paths
 - Can we "remember" where to return?

Constraint Graph (so-called PAG)



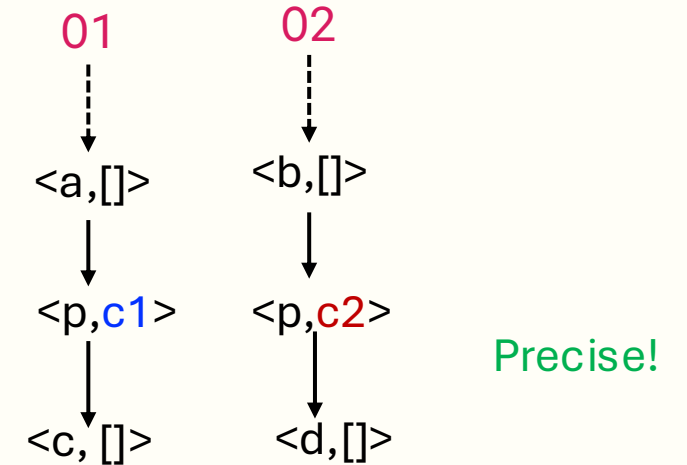
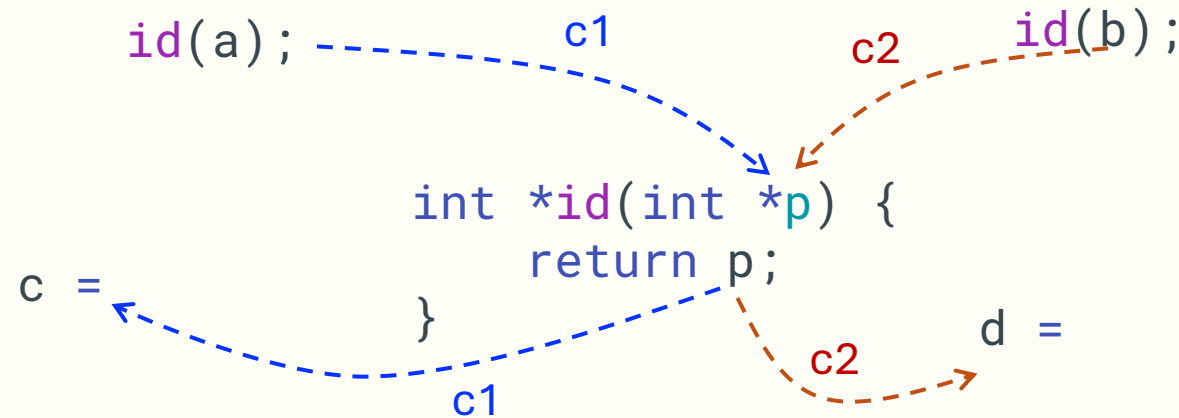
Vars	PTS	Vars	PTS
<code>a</code>	{O1}	<code>b</code>	{O2}
<code>p</code>	{O1, O2}		
<code>c</code>	{O1, O2}	<code>d</code>	{O1, O2}

Spurious Points-to Objects

Context Sensitivity

□ Key idea:

- Separate analyses for functions called in different "contexts"
- "context": an abstract concept that is statically definable; used to differentiate different calls to a function
- Contexts are used to decorate graph nodes



Vars	PTS	Vars	PTS
<code><a, []></code>	{O1}	<code><b, []></code>	{O2}
<code><p, c1></code>	{O1}	<code><p, c2></code>	{O2}
<code><c, []></code>	{O1}	<code><d, []></code>	{O2}

Types of Context Sensitivity

- ❑ No context sensitivity, []
- ❑ Callsite sensitivity
 - Call strings, $[l_1, \dots, l_n]$
 - l_i : the line label of the callsite
- ❑ Object sensitivity, $[o_1, \dots, o_n]$
 - o_i is an allocation site
- ❑ Type sensitivity, $[t_1, \dots, t_n]$
 - t_i is the type of an allocation object
- ❑ Value contexts
 - States (environment) at the given callsite
- ❑ k – limiting technique
- ❑ Unscalable when $k \geq 3$ in practice

```
1. int *id(int *p) {
2.     return p;
3. }
4. int *wid_1(int *p_1) {
5.     return id(p_1);
6. }
7. ...
8. int *wid_n(int *p_n) {
9.     return wid_(n-1)(p_n);
10. }
11. c = wid_n(a);
12. d = wid_n(b);
```

Two contexts for
differentiate id():

$[l_5, \dots, l_9, l_{11}]$

$[l_5, \dots, l_9, l_{12}]$

Context length: $n+1$

```
int factorial(int* n) {
    if (*n == 0) {
        return 1;
    } else {
        int t = *n;
        *n = t - 1;
        return t * factorial(n)
    }
}
int x = 100;
int y = factorial(&x);
```

Context length
could be infinite.

Steensgaard's Analysis

□ **Problem:** Quadratic-in-practice is still not ultra-scalable

➤ **Want a faster algorithm! Need ~LINEAR. How?**

□ **Challenge:**

➤ Solution space of pointer analysis (e.g. points-to sets) is $O(n^2)$

□ **Key idea:**

➤ Use constant-space per pointer.

➤ Merge aliases and alternates into the same equivalence class.

➤ p can point to q or r? Let's treat q and r as the same pseudo-var and merge everything we know about q and r.

➤ Points-to "sets" are basically singletons

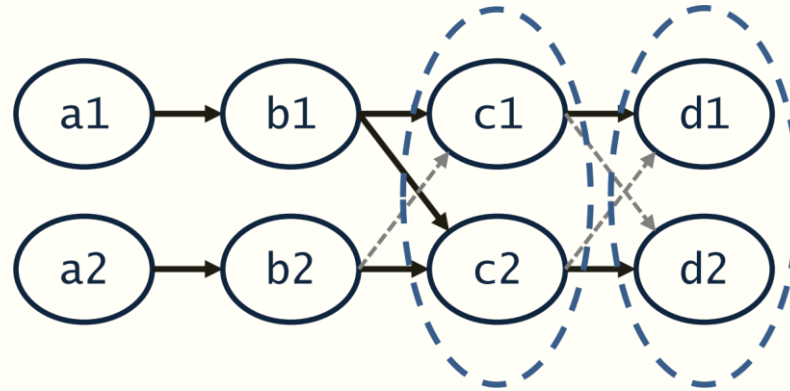
□ **Algorithm runs in $O(n * \alpha(n))$ using union-find structure**

➤ Almost linear, very scalable in practice (Millions of LoC)

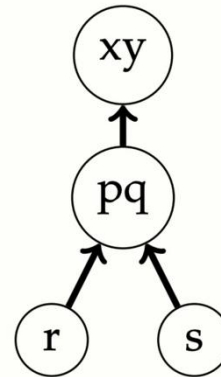
➤ **Yes, very imprecise!**

Steensgaard's Analysis-Examples

```
a1 = &b1;  
b1 = &c1;  
c1 = &d1;  
a2 = &b2;  
b2 = &c2;  
c2 = &d2;  
b1 = &c2;
```



```
1 :  $p := \&x$   
2 :  $r := \&p$   
3 :  $q := \&y$   
4 :  $s := \&q$   
5 :  $r := s$ 
```



Steensgaard's Analysis

$$\overline{\llbracket p := q \rrbracket} \hookrightarrow \text{join}(*p, *q) \text{ copy}$$
$$\overline{\llbracket p := \&x \rrbracket} \hookrightarrow \text{join}(*p, x) \text{ address-of}$$
$$\overline{\llbracket p := *q \rrbracket} \hookrightarrow \text{join}(*p, **q) \text{ dereference}$$
$$\overline{\llbracket *p := q \rrbracket} \hookrightarrow \text{join}(**p, *q) \text{ assign}$$

```
join(l1, l2)
  if (find(l1) == find(l2))
    return
  n1 ← *l1
  n2 ← *l2
  union(l1, l2)
  join(n1, n2)
```

- **Points-to “sets” are basically singletons**

时间复杂度 $O(n * \alpha(n))$ (almost linear), 空间复杂度 $O(n)$.

(Optional) 作业：指针分析应用

□用Andersen算法分析如下代码

```
#include <stdio.h>
struct Node { struct Node *next; int *data; };
int main() {
    int a, b, c; struct Node n1, n2, n3;
    int *pa = &a; int *pb = &b; int *pc = &c;
    struct Node *p = &n1;
    struct Node *q = &n2;
    struct Node *r = &n3;
    n1.next = &n2; q->next = r;
    p->data = pa; n2.data = &b
    struct Node *x = p->next;
    struct Node *y = x->next;
    int *d1 = p->data; int *d2 = q->data;
    r->data = pc;
    y = p->next;
    return 0;
}
```

□用Steensgaard算法分析

```
#include <stdlib.h>
struct Node { struct Node *next; };
int main() {
    struct Node *a, *b, *c, *tmp;
    a = (struct Node *)malloc(sizeof(struct
Node));
    b = (struct Node *)malloc(sizeof(struct
Node));
    c = (struct Node *)malloc(sizeof(struct
Node));
    tmp = a;
    tmp = b;
    a = c;
    a->next = b;
    b->next = c;
    return 0;
}
```

(Optional) 项目：熟悉指针分析实现

□Qilin指针分析框架

- <https://github.com/QilinPTA/Qilin>
- 支持多种指针分析优化技术

□任务：

- 选一组Java程序（3-4个即可），尝试用Qilin框架的技术进行指针分析
- 阅读源码，探究Qilin是怎么支持多种上下文敏感性的？
- 指出Qilin框架设计和实现上可以改进和优化的地方？
- 写一个1000字左右的阅读报告
- 语言要求言简意赅，不无病呻吟
- 可以使用示意图等辅助说明