软件分析与架构设计

# 收集语义和指针分析

**何冬杰**

**重庆大学**

# 三地址码 (Three-Address)

## ❑Review SIMP Syntax

$$
\begin{array}{llll}
S & ::= & x := a & \\
& | & \text{skip} & \\
& | & S_1; S_2 & \\
& | & \text{if } b \text{ then } S_1 \text{ else } S_2 & \\
& | & \text{while } b \text{ do } S &
\end{array}
\qquad
\begin{array}{lll}
b & ::= & \text{true} \\
& | & \text{false} \\
& | & \text{not } b \\
& | & b_1 \; op_b \; b_2 \\
& | & a_1 \; op_r \; a_2
\end{array}
\qquad
\begin{array}{lll}
a & ::= & x \\
& | & n \\
& | & a_1 \; op_a \; a_2 \\
\end{array}
\qquad
\begin{array}{lll}
op_b & ::= & \text{and} \mid \text{or} \\
op_r & ::= & < \; \mid \; \leqslant \; \mid \; = \\
& & \mid \; > \; \mid \; \geqslant \\
op_a & ::= & + \mid - \mid * \mid /
\end{array}
$$

## ❑三地址码：一种编译器常用中间表示（IR）形式

➤3-address Syntax:

➤解决AST难以追踪程序执行时的数据流和控制流问题

$Inst ::= \; x := n \mid x := y \mid x := y \; op \; z \mid \text{goto } n \mid \text{if } x \; op_r \; 0 \text{ goto } n$

$op_a ::= \; + \mid - \mid * \mid / \mid \dots$

$op_r ::= \; < \mid \; \leq \; \mid \; = \; \mid > \; \mid \; \geq \; \mid \; \dots$

$P \in Num \rightarrow Inst$

*Inst*: instructions
*Num*: number literals
$x, y \in Var$: variables

更多3地址指令，后续会用

$Inst ::= \cdots \mid x := f(y) \mid \text{return } x \mid x := y.m(z) \mid \text{read } x \mid \text{print } x \mid x := \&p \mid x :=* p \mid * p := x \mid x := y.f \mid x.f := y \mid \text{halt}$

# 三地址码 (Three-Address)

w = x * y + z

if b then S1 else S2

while b do S

3地址码

```
1. t = x * y
2. w = t + z
```

3地址码

```
1. if b goto 4
2. S2
3. goto 5
4. S1
5. ...
```

3地址码

```
1. if !b then goto 3
2. S
3. ...
```

❑ **AST到3地址的转换是直接的、简单的**
  ➢ Compiler Explorer：https://godbolt.org
  ➢ (可选)作业：clang如何将AST转成LLVM-IR?

# An Abstract Machine of 3-Address IR

❑ **Configuration (state):**
  ➢ Environment
  ➢ Program counter
  $$c \in E \times \mathbb{N}$$

❑ **Program:**
  ➢ Maps labels to instructions
  $$P \in \mathbb{N} \to Inst$$

❑ **Transition or Execution:**
  ➢ Small-step Semantics

$$P \vdash \langle E, n \rangle \rightsquigarrow \langle E', n' \rangle$$

$$Inst ::= \ x := n \mid x := y \mid x := y \ op \ z \mid \text{goto} \ n \mid \text{if} \ x \ op_r \ 0 \ \text{goto} \ n$$

$$\frac{P(n) = x := m}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E[x \mapsto m], n+1 \rangle} \ \textit{step-const}$$

$$\frac{P[n] = x := y}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E[x \mapsto E(y)], n+1 \rangle} \ \textit{step-copy}$$

$$\frac{P(n) = x := y \ op \ z \quad E(y) \ \mathbf{op} \ E(z) = m}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E[x \mapsto m], n+1 \rangle} \ \textit{step-arith}$$

$$\frac{P(n) = \text{goto} \ m}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, m \rangle} \ \textit{step-goto}$$

$$\frac{P(n) = \text{if} \ x \ op_r \ 0 \ \textbf{goto} \ m \quad E(x) \ \mathbf{op_r} \ 0 = true}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, m \rangle} \ \textit{step-iftrue}$$

$$\frac{P(n) = \text{if} \ x \ op_r \ 0 \ \textbf{goto} \ m \quad E(x) \ \mathbf{op_r} \ 0 = false}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, n+1 \rangle} \ \textit{step-iffalse}$$

# 收集语义 (Collecting Semantics)

The **collecting semantics** is used to define a collection of all the states (defined by operational semantics) of the program at a given point.

❑ **零分析 (zero analysis):**
- ➢ Determine whether the value of a variable is 0
- ➢ Collecting semantics: $\mathbb{E} = \{E \mid \langle E, n \rangle \text{ is a state}\}$
- ➢ Abstract function: $\alpha_{ZA}(x) = \{i \mid \exists E \in \mathbb{E}, \text{such that } E(x) = i \}$
- ➢ Given a variable x : if $\alpha_{ZA}(x) = \{0\}$, the value of x is zero

```
1. x = 10
2. y = 2 * x
3. z = 0
4. if y > 0 goto 6
5. z = 2 * x - y
6. ...
```

$x \mapsto \{10\}$
$x \mapsto \{10\}, y \mapsto \{20\}$
$x \mapsto \{10\}, y \mapsto \{20\}, z \mapsto \{0\}$
$x \mapsto \{10\}, y \mapsto \{20\}, z \mapsto \{0\}$
$x \mapsto \{10\}, y \mapsto \{20\}, z \mapsto \{0\}$
$x \mapsto \{10\}, y \mapsto \{20\}, z \mapsto \{0\}$

# 收集语义 (Collecting Semantics)

The **collecting semantics** is used to define a collection of all the states (defined by operational semantics) of the program at a given point.

☐ **达到定值 (reaching definitions)**
- ➤ Extending environment $E$ (called $E_{\text{RD}}$): $E_{RD} \in Var \to \mathbb{Z} \times \mathbb{N}$
- ➤ $x \mapsto v, n$: indicating that x was last defined as $v$ at the location $n$

Collecting Semantics for Reaching definitions

$$\frac{P[n] = x := m}{P \vdash E, n \rightsquigarrow E[x \mapsto m, n], n+1} \; step\text{-}const$$

$$\frac{P(n) = \mathbf{goto} \; m}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, m \rangle} \; step\text{-}goto$$

$$\frac{P[n] = x := y}{P \vdash E, n \rightsquigarrow E[x \mapsto E[y], n], n+1} \; step\text{-}copy$$

$$\frac{P(n) = \mathbf{if} \; x \; op_r \; 0 \; \mathbf{goto} \; m \quad E(x) \; \mathbf{op_r} \; 0 = true}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, m \rangle} \; step\text{-}iftrue$$

$$\frac{P[n] = x := y \; op \; z \quad E[y] \; \mathbf{op} \; E[z] = m}{P \vdash E, n \rightsquigarrow E[x \mapsto m, n], n+1} \; step\text{-}arith$$

$$\frac{P(n) = \mathbf{if} \; x \; op_r \; 0 \; \mathbf{goto} \; m \quad E(x) \; \mathbf{op_r} \; 0 = false}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, n+1 \rangle} \; step\text{-}iffalse$$

New rules

```
1.  x = 3
2.  t = x - 2
3.  if t > 0 goto 5
4.  x = 4
5.  y = x + 5
```

$$\alpha(E, 4) = \{2,4\}$$
$$\alpha(E, 5) = \{1,2,4,5\}$$

➤ Abstract function: $\alpha_{RD}(E_{RD}, n) = \{m \mid \exists x \in domain(E_{RD}) \text{ such that } E_{RD}(x) = i, m\}$

$m$处定义的变量可到达节点$n$

# 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**

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

```
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
```
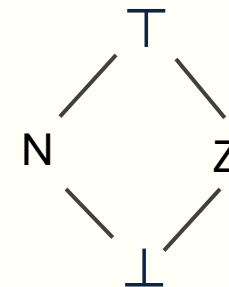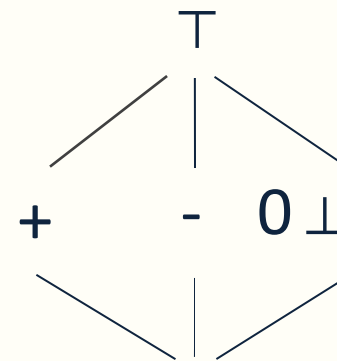
循环次数取决于输入x

❑ 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：$\mathbb{Z} : \{-\infty, \cdots, -2, -1, 0, 1, 2, \cdots + \infty\}$
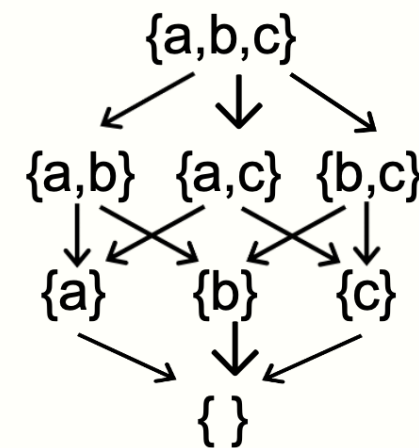  ➢抽象域：$\{\perp, N, Z, \top\}$

❑符号分析 **(Sign Analysis)**
  ➢具体域：$\mathbb{Z} : \{-\infty, \cdots, -2, -1, 0, 1, 2, \cdots + \infty\}$
  ➢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
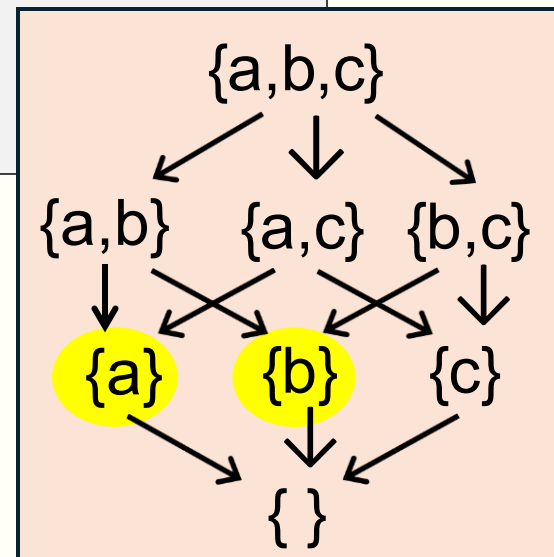- ➢ 偏序集中的两个元素可能无法比较,e.g., 右图中的{a}和{b}

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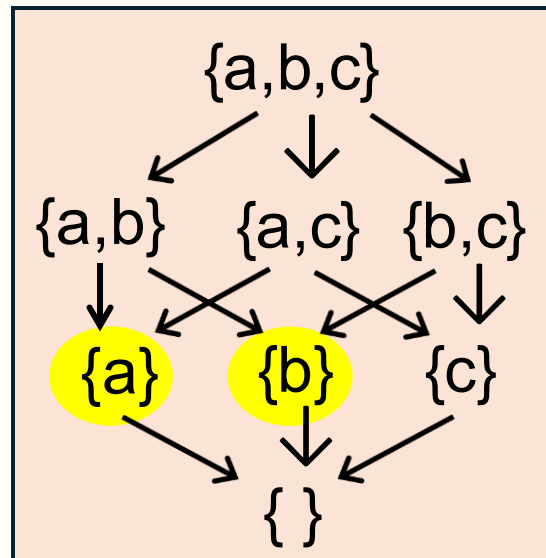
A ***greatest lower bound*** $\sqcap X$ is defined by $\sqcap X \sqsubseteq X \wedge \forall\, y \in S:\ y \sqsubseteq X \Rightarrow y \sqsubseteq \sqcap X$.

# 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;
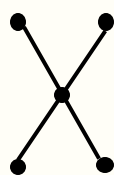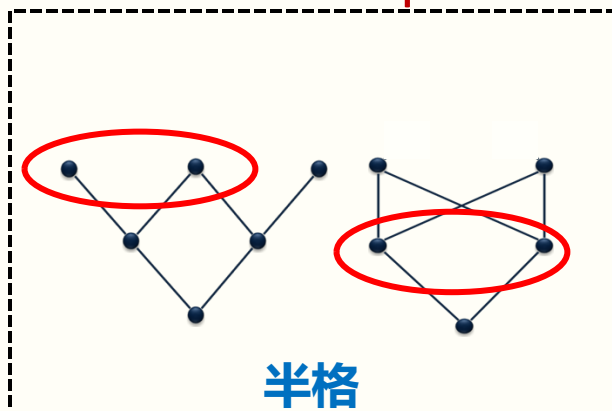➢ if both exist, then $(S, \sqsubseteq)$ is called a lattice.

## ❑完全格
➢ 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 $\bot = \sqcap S$
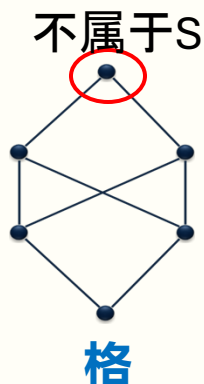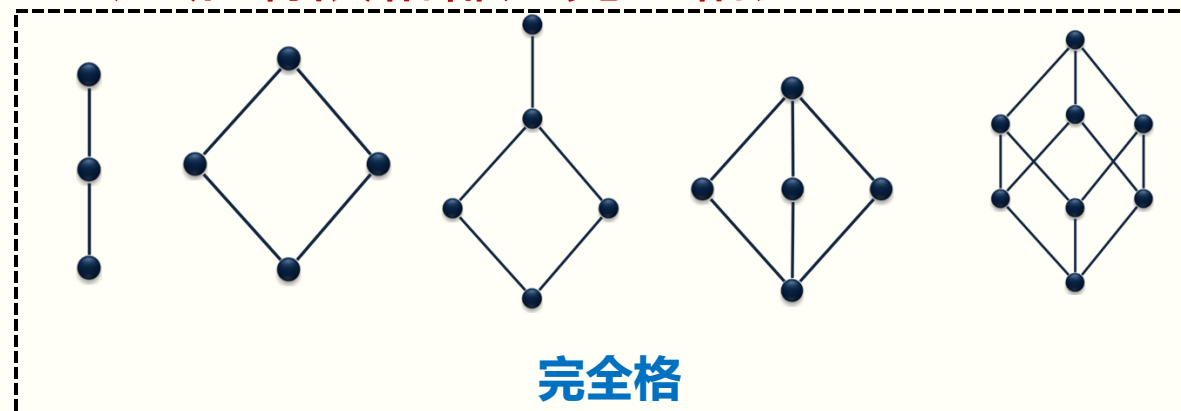➢ A finite lattice is complete if $\top$ and $\bot$ exist （一般有限格都是完全格）

## ❑举例:

{a,b,c}

{a,b}  {a,c}  {b,c}

{a}  {b}  {c}

{ }

完全格

不属于S

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 $\bot$ to $\top$

❑**Powerset lattice:**
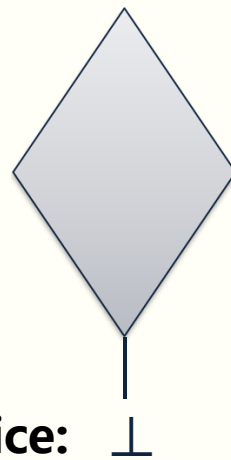  ➢Powerset of every finite set $A$ defines a complete lattice
  ➢$(\mathcal{P}(A), \subseteq), \bot = \emptyset, \top = A, x \sqcup y = x \cup y, x \sqcap y = x \cap y$
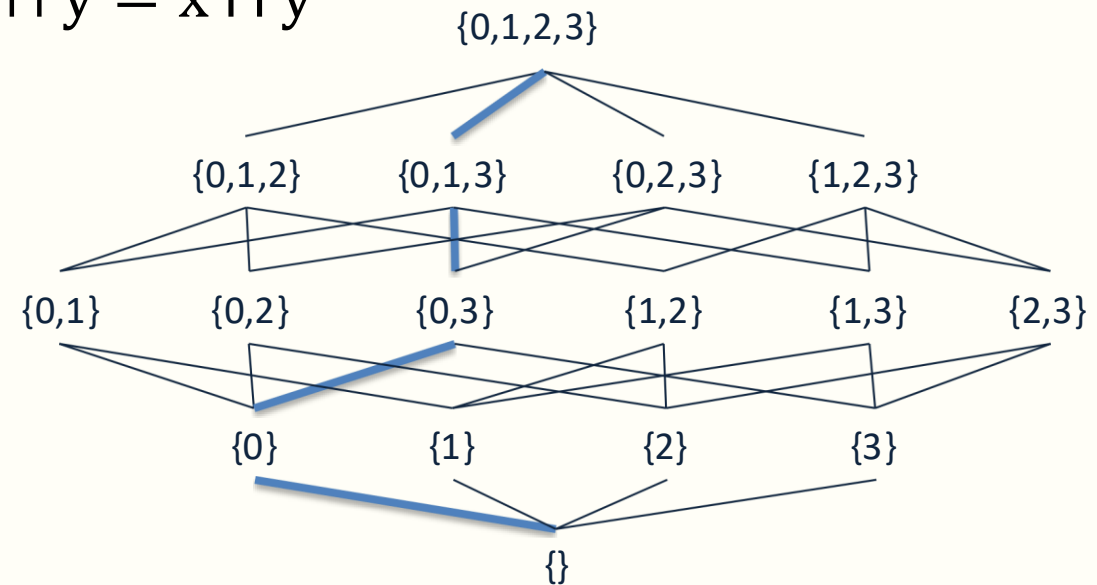
❑ **Flat lattice and Lift lattice**



**Lift lattice:**
$height(lift(L)) = height(L)+1$

**Flat lattice:**  $\bot$
$Flat(A)$, where $A$ is a set, $height(flat(A)) = 2$

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

# 格 (Lattices)

❑**Product lattice:**
- ➢$L_1 \times L_2 \times \cdots \times L_n = \{(x_1, x_2, \cdots, x_n) \mid x_i \in L_i\}$
- ➢$L_i (1 \leq i \leq n)$ 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:**
- ➢$\boldsymbol{A \to L = \{[a_1 \mapsto x_1, a_2 \mapsto x_2, \cdots] \mid a_i \in A \wedge x_i \in L\}}$
- ➢$\sqsubseteq, \sqcup, and \ \sqcap \ are\ computed\ pointwise$
- ➢$\text{Height}(A \to L) = |A| \cdot \text{Height}(L)$

❑**例子：指针分析**
- ➢$A$ is the set of program variables
- ➢$L$ is the powerset of allocation site abstractions (objects)

# 指针分析 (Pointer Analysis)

❑ **指针分析计算指向关系**： $\boldsymbol{Pts}$: Var → $\mathcal{P}$(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

 ➢ $Must\text{-}points\text{-}to$关系不在考虑范围内

 ➢ 指向关系可以用来计算别名(alias)关系

  o Alias(x, y) if and only if Pts(x) ∩ Pts(y) ≠ ∅

❑ **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

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

# 指针分析 (Pointer Analysis)

❑ 例子：

**Strong Update**           **Weak Update**

```
1. q := malloc() // O1
```
Pts(q) = {O1}        Pts(q) = {O1}

```
2. p := malloc() // O2
```
Pts(p) = {O2}        Pts(p) = {O2}

```
3. p := q
```
Pts(p) = {O1}        Pts(p) = {O1, O2}

```
4. r := &p
```
Pts(r) = {&p}        Pts(r) = {&p}

```
5. s := malloc() // O3
```
Pts(s) = {O3}        Pts(s) = {O3}

```
6. *r := s
```
Pts(p) = {O3}        Pts(p) = {O1, O2, O3}

```
7. t := &s
```
Pts(t) = {&s}        Pts(t) = {&s}

```
8. u := *t
```
Pts(u) = {O3}        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？

# Andersen's Analysis

## ❑Key idea: cast as a constraint-solving problem
  ➢One subset constraint per instruction
  ➢Domain abstraction and function abstraction

$$\overline{[\![n\colon p := malloc()]\!] \hookrightarrow l_n \in p} \quad malloc$$

$$\overline{[\![p := \&x]\!] \hookrightarrow l_x \in p} \quad address\text{-}of$$

$$\overline{[\![p := q]\!] \hookrightarrow p \sqsupseteq q} \quad copy$$

$$\overline{[\![*p := q]\!] \hookrightarrow *p \sqsupseteq q} \quad assign$$

$$\overline{[\![p := *q]\!] \hookrightarrow p \sqsupseteq *q} \quad dereference$$

- ▪ The constraint solver can give the most precise solution.
  - ▪ Why? (optional)
- ▪ Constraints are equivalent to abstract functions

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

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

# Kleene's fixed-point theorem

## ❑Monotonicity

A function f: L → 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*:

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

# 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 $\perp \sqsubseteq f(\perp) \sqsubseteq f^2(\perp) \sqsubseteq \ldots f^i(\perp) \ldots$ 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 $lfp(f) = f^k(\perp)$

## ❑ Uniqueness

- ➤ 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, $lfp(f) = f^k(\perp) \sqsubseteq x$, i.e. $lfp(f)$ is least
- ➤ Uniqueness then follows from anti-symmetry

# Andersen's Analysis
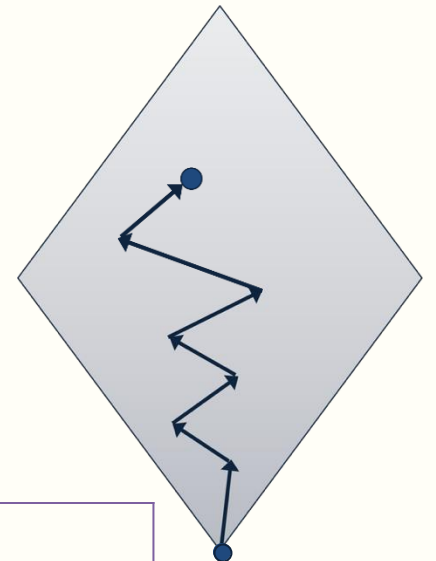
Let $f_1, f_2: L \to L$ (L is a lattice) be two monotonic functions, then their **composition** $f_1 \circ f_2$ is also **monotonic**.

Let $f: L \to L$ ($L = (Var \to \mathcal{P}(OBJs))$ 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 Ansersen's analysis is $lfp(f)$.

❑ **The time complexity of *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

# A more efficient implementation

❑**The Cubic Framework ($\boldsymbol{O(n^3)}$)**

➢A set of tokens $T = \{t_1, t_2, \cdots, t_k\}$,

➢A collection of constraint variables $V = \{x_1, \cdots, x_n\}$

➢A collection of constraints of these forms:

- $t \in x$ ;
- $x \subseteq y$ ; inclusion constraints
- $t \in x \Longrightarrow y \subseteq z$ ; 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)

$O(n^3)$的时间复杂度，$n$是变量数量

- $t \in x$
  addToken(t, x)
  propagate()

- $x \subseteq y$
  addEdge(x, y)
  propagate()

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

**addToken**(t, x):
  if t ∉ x.sol
    add t to x.sol
    add (t, x) to W

**addEdge**(x, y):
  if x ≠ y ∧ y ∉ x.succ
    add y to x.succ
    for each t in x.sol
      addToken(t, y)

**propagate**():
  while W ≠ ∅
    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)

# Andersen's Analysis

❑ **使用Cubic框架求解**

$$\overline{[\![n: p := malloc()]\!] \hookrightarrow l_n \in p}\ malloc$$

$$\overline{[\![p := \&x]\!] \hookrightarrow l_x \in p}\ address\text{-}of$$

$\Longrightarrow$ 第一种形式：$t \in x$

$$\overline{[\![p := q]\!] \hookrightarrow p \supseteq q}\ copy$$

$\Longrightarrow$ 第二种形式：$x \subseteq y$

$$\overline{[\![*p := q]\!] \hookrightarrow *p \supseteq q}\ assign$$

$$\overline{[\![p := *q]\!] \hookrightarrow p \supseteq *q}\ dereference$$

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

$$\frac{}{[\![p := q.f]\!] \hookrightarrow p \supseteq q.f} \textit{field-read}$$

$$\frac{}{[\![p.f := q]\!] \hookrightarrow p.f \supseteq q} \textit{field-assign}$$

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)：
$$O \in q \Rightarrow O.f \subseteq p$$
$$O \in p \Rightarrow q \subseteq O.f$$

### ➢支持函数调用
- 参数和返回值的传递建模成copy语句

```
r = m(a1, ..., an)      T m(T p1, ... T pn) {
                        ...
                        return ret;
                        }
```

$$\frac{}{[\![r = m(a1, \cdots, an]\!] \hookrightarrow p1 \supseteq a1, \ldots, pn \supseteq an, r \supseteq ret} \textit{call}$$

# Example Program

```c
#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

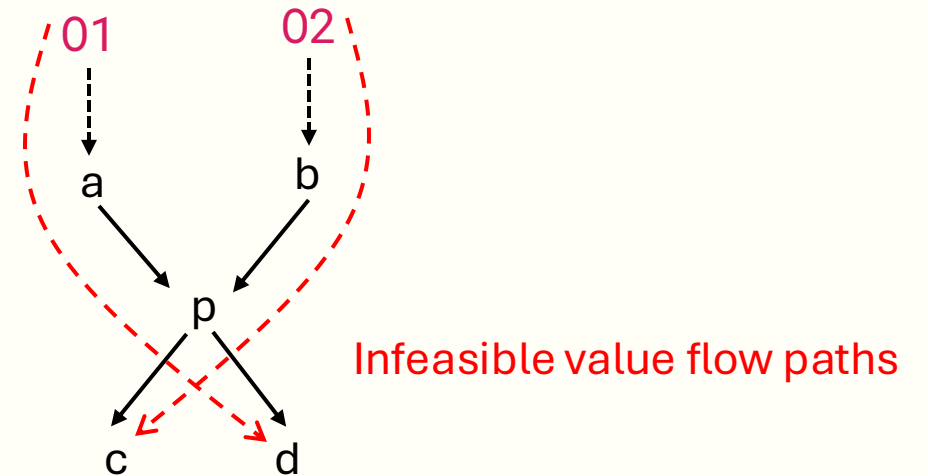**Constraint Graph** (so-called Pointer Assignment Graph)



| 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)); // O1
8.    c = id(a);
9.    b = (int *)malloc(sizeof(int)); // O2
10.   d = id(b);
11.    return 0;
12. }
```

**Constraint Graph** (so-called PAG)



Infeasible value flow paths

| Vars | PTS | Vars | PTS |
|------|-----|------|-----|
| a | {O1} | b | {O2} |
| p | {O1, O2} | | |
| c | {O1, O2} | d | {O1, O2} |

Spurious Points-to Objects

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

# 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



```
id(a);            c1        c2        id(b);

              int *id(int *p) {
                  return p;
c =           }
                                      d =
```

O1                  O2

&lt;a,[]&gt;             &lt;b,[]&gt;

&lt;p,c1&gt;            &lt;p,c2&gt;          Precise!

&lt;c, []&gt;           &lt;d,[]&gt;

| Vars | PTS | Vars | PTS |
|------|-----|------|-----|
| &lt;a, []&gt; | {O1} | &lt;b, []&gt; | {O2} |
| &lt;p, c1&gt; | {O1} | &lt;p, c2&gt; | {O2} |
| &lt;c, []&gt; | {O1} | &lt;d, []&gt; | {O2} |

# Types of Context Sensitivity

❑ No context sensitivity, []

❑ Callsite sensitivity
- ➢ Call strings, $[l_1, \cdots, l_n]$
- ➢ $l_i$: the line label of the callsite

❑ Object sensitivity, $[o_1, \cdots, o_n]$
- ➢ $o_i$ is an allocation site

❑ Type sensitivity, $[t_1, \cdots, t_n]$
- ➢ $t_i$ 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():

[l5, ..., l9, l11]
[l5, ..., l9, l12]

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.
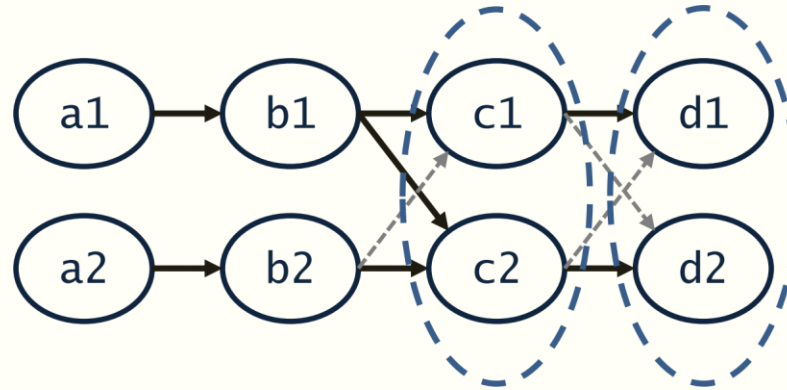  ➢Points-to "sets" are basically singletons
❑**Algorithm runs in** $O(n * \alpha(n))$ **using union-find structure**
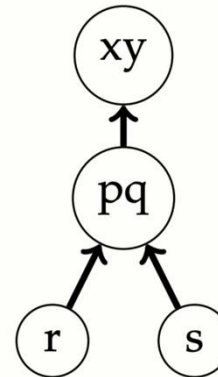  ➢Almost linear, very scalable in practice (Millions of LoC)
  ➢Yes, very imprecise!

# Steensgaard's Analysis-Examples

a1 = &b1;

b1 = &c1;

c1 = &d1;

a2 = &b2;

b2 = &c2;

c2 = &d2;
b1 = &c2;



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

# Steensgaard's Analysis

$$\frac{}{[\![p := q]\!] \hookrightarrow join(*p, *q)} \; copy$$

$$\frac{}{[\![p := \&x]\!] \hookrightarrow join(*p, x)} \; address\text{-}of$$

$$\frac{}{[\![p := *q]\!] \hookrightarrow join(*p, **q)} \; dereference$$

$$\frac{}{[\![*p := q]\!] \hookrightarrow join(**p, *q)} \; assign$$

```
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   ❏用Steensgaard算法分析

```c
#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;
}
```

```cpp
#include <cstdlib>
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) 作业：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？