




# Static Program Analysis

Xiangyu Zhang

The slides are compiled from  
Alex Aiken's  
Michael D. Ernst's  
Sorin Lerner's



# A Scary Outline

- ❑ Type-based analysis
- ❑ Data-flow analysis
- ❑ Abstract interpretation
- ❑ Theorem proving
- ❑ ...

# The Real Outline

- ❑ The essence of static program analysis
- ❑ The categorization of static program analysis
- ❑ Type-based analysis basics
- ❑ Data-flow analysis basics

# The Essence of Static Analysis

- ❑ Examine the program text (no execution)
- ❑ Build a model of the program state
  - An abstract of the run-time state
- ❑ Reason over the possible behaviors.
  - E.g. “run” the program over the abstract state

# The Essence of Static Analysis

Each program statement's **transfer function** indicates how it transforms the (abstract) state

Example: What is the transfer function for

```
y = x++;
```

?

# Selecting an abstract domain

Abstract domain: { even, odd, either }

$\langle x \text{ is odd}; y \text{ is odd} \rangle$

$y = x++;$

$\langle x \text{ is even}; y \text{ is odd} \rangle$

# Selecting an abstract domain

Abstract domain: { even, odd, either }

$\langle x \text{ is odd; } y \text{ is odd} \rangle$

$y = x++;$

$\langle x \text{ is even; } y \text{ is odd} \rangle$

Abstract domain: { prime, nonprime, anything }

$\langle x \text{ is prime; } y \text{ is prime} \rangle$

$y = x++;$

$\langle x \text{ is anything; } y \text{ is prime} \rangle$

# Selecting an abstract domain

Abstract domain: set of numbers, one set per variable

$\langle x = \{ 3, 5, 7 \}; y = \{ 9, 11, 13 \} \rangle$

$y = x++;$

$\langle x = \{ 4, 6, 8 \}; y = \{ 3, 5, 7 \} \rangle$



# Selecting an abstract domain

Abstract domain: set of numbers, one set per variable

$$\langle x = \{ 3, 5, 7 \}; y = \{ 9, 11, 13 \} \rangle$$
$$y = x++;$$
$$\langle x = \{ 4, 6, 8 \}; y = \{ 3, 5, 7 \} \rangle$$

Abstract domain: set of environments

- environment assigns a variable to a number

$$\langle x=3, y=11 \rangle, \langle x=5, y=9 \rangle, \langle x=7, y=1 \rangle$$
$$y = x++;$$
$$\langle x=4, y=3 \rangle, \langle x=6, y=5 \rangle, \langle x=8, y=7 \rangle$$

# Selecting an abstract domain

Abstract domain: set of numbers, one set per variable

$$\langle x = \{ 3, 5, 7 \}; y = \{ 9, 11, 13 \} \rangle$$

$y = x++;$

$$\langle x = \{ 4, 6, 8 \}; y = \{ 3, 5, 7 \} \rangle$$

Abstract domain: set of environments

- environment assigns a variable to a number

$$\langle x=3, y=11 \rangle, \langle x=5, y=9 \rangle, \langle x=7, y=1 \rangle$$

$y = x++;$

$$\langle x=4, y=3 \rangle, \langle x=6, y=5 \rangle, \langle x=8, y=7 \rangle$$

Abstract domain: symbolic expression per variable

$$\langle x_n = f_1(a_{n-1}, \dots, z_{n-1}); y_n = f_2(a_{n-1}, \dots, z_{n-1}) \rangle$$

$y = x++;$

$$\langle x_{n+1} = x_n + 1; y_{n+1} = x_n \rangle$$

# Categorization

- ❑ Flow sensitivity
- ❑ Context sensitivity.

# Flow Sensitivity

- ❑ Flow sensitive analyses
  - The order of statements matters
  - Need a control flow graph
  
- ❑ Flow insensitive analyses
  - The order of statements doesn't matter
  - Analysis is the same regardless of statement order

# Example Flow Insensitive Analysis

- What variables does a program modify?

$$G(x := e) = \{x\}$$

$$G(s_1; s_2) = G(s_1) \cup G(s_2)$$

- Note  $G(s_1; s_2) = G(s_2; s_1)$

# The Advantage

- Flow-sensitive analyses require a model of program state at each program point
  - E.g., liveness analysis, reaching definitions, ...
- Flow-insensitive analyses require only a single global state
  - E.g., for  $G$ , the set of all variables modified

# Notes on Flow Sensitivity

- ❑ Flow insensitive analyses seem weak, but:
- ❑ Flow sensitive analyses are hard to scale to very large programs
  - Additional cost: state size  $\times$  # of program points
- ❑ Beyond 1000's of lines of code, only flow insensitive analyses have been shown to scale (by Alex Aiken)

# Context-Sensitive Analysis

- What about analyzing across procedure boundaries?

```
Def f(x){...}
```

```
Def g(y){...f(a)...}
```

```
Def h(z){...f(b)...}
```

- Goal: Specialize analysis of **f** to take advantage of
    - **f** is called with **a** by **g**
    - **f** is called with **b** by **h**
-



# Flow Insensitive: Type-Based Analysis

# Outline

- A language
  - Lambda calculus
- Types
  - Type checking
  - Type inference
- Applications to software reliability
  - Representation analysis
    - ❖ Alias analysis and memory leak analysis.

# The Typed Lambda Calculus

- Lambda calculus
  - types are assigned to bound variables.
- Add integers, addition, if-then-else
- Note: Not every expression generated by this grammar is a properly typed term.

$e = x \mid \lambda x : \tau . e \mid e \ e \mid i \mid e + e \mid \text{if } e \ e \ e$

# Types

- ❑ Function types
- ❑ Integers
- ❑ Type variables
  - Stand for definite, but unknown, types

$$\tau = \alpha \mid \tau \rightarrow \tau \mid \text{int}$$

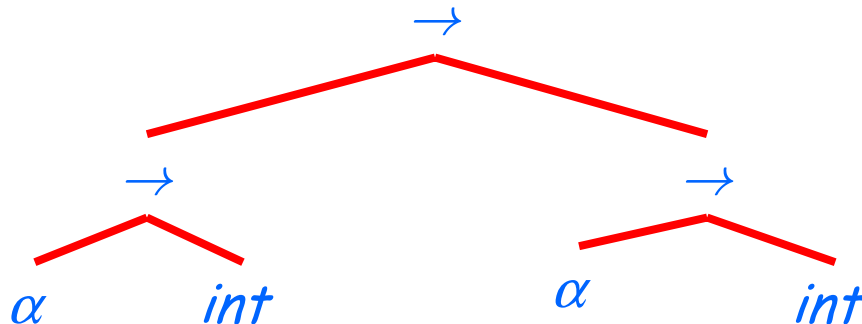
# Function Types

- Intuitively, a type  $\tau_1 \rightarrow \tau_2$  stands for the set of functions that map arguments of type  $\tau_1$  to results of type  $\tau_2$ .
- Placeholder for any other structured datatype
  - Lists
  - Trees
  - Arrays

# Types are Trees

- Types are terms
- Any term can be represented by a tree
  - The parse tree of the term
  - Tree representation is important in algorithms

$(\alpha \rightarrow \text{int}) \rightarrow \alpha \rightarrow \text{int}$



# Examples

- We write  $e:t$  for the statement “ $e$  has type  $t$ .”

$\lambda x. \alpha x. \alpha \rightarrow \alpha$

# Examples

- We write  $e:t$  for the statement “ $e$  has type  $t$ .”

$\lambda x. \alpha x. \alpha \rightarrow \alpha$

$\lambda x. \alpha \lambda y. \beta x. \alpha \rightarrow \beta \rightarrow \alpha$



# Examples

- We write  $e:t$  for the statement “ $e$  has type  $t$ .”

$\lambda x. \alpha x. \alpha \rightarrow \alpha$

$\lambda x. \alpha \lambda y. \beta x. \alpha \rightarrow \beta \rightarrow \alpha$

$\lambda f: \alpha \rightarrow \beta \lambda g: \beta \rightarrow \gamma \lambda x. \alpha (g (f x))$

# Examples

- We write  $e:t$  for the statement “ $e$  has type  $t$ .”

$$\lambda x. \alpha x. \alpha \rightarrow \alpha$$

$$\lambda x. \alpha \lambda y. \beta x. \alpha \rightarrow \beta \rightarrow \alpha$$

$$\lambda f: \alpha \rightarrow \beta. \lambda g: \beta \rightarrow \gamma. \lambda x. \alpha (g (f x)): (\alpha \rightarrow \beta) \rightarrow (\beta \rightarrow \gamma) \rightarrow \alpha \rightarrow \gamma$$

$$\lambda f: \alpha \rightarrow \beta \rightarrow \gamma. \lambda g: \alpha \rightarrow \beta. \lambda x. \alpha (f x) (g x): (\alpha \rightarrow \beta \rightarrow \gamma) \rightarrow (\alpha \rightarrow \beta) \rightarrow \alpha \rightarrow \gamma$$

# Type Environments

- ❑ To determine whether the types in an expression are correct we perform *type checking*.
- ❑ But we need types for free variables, too!
- ❑ A *type environment* is a function from variables to types. The syntax of environments is:

$$A = \emptyset \mid A, x : \tau$$

- ❑ The meaning is:

$$(A, x : \tau)(y) = \begin{array}{ll} \tau & \text{if } x = y \\ A(y) & \text{if } x \neq y \end{array}$$

# Type Checking Rules

- Type checking is done by structural induction.
  - One inference rule for each form
  - Assumptions contain types of free variables
  - A term is *well-typed* if  $\emptyset \mid e : \tau$

$$\frac{A(x) = \tau}{A \vdash x : \tau} \quad \frac{A, x : \tau \vdash e : \tau'}{A \vdash \lambda x : \tau. e : \tau \rightarrow \tau'} \quad \frac{A \vdash e_1 : \tau \rightarrow \tau' \quad A \vdash e_2 : \tau}{A \vdash e_1 e_2 : \tau'}$$
$$\frac{}{A \vdash i : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \text{int}}{A \vdash e_1 + e_2 : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \tau \quad A \vdash e_3 : \tau}{A \vdash \text{if } e_1 e_2 e_3 : \tau}$$

# Example

???

$$\emptyset \vdash \lambda x : \alpha. \lambda y : \beta. x : \alpha \rightarrow \beta \rightarrow \alpha$$

# Example

$$x:\alpha, y:\beta \vdash x:\alpha$$

$\frac{A(x)=\tau}{A \vdash x:\tau}$	$\frac{A, x:\tau \vdash e:\tau'}{A \vdash \lambda x:\tau. e:\tau \rightarrow \tau'}$	$\frac{A \vdash e_1:\tau \rightarrow \tau' \quad A \vdash e_2:\tau}{A \vdash e_1 e_2:\tau'}$
	$\frac{A \vdash e_1:\text{int} \quad A \vdash e_2:\text{int}}{A \vdash e_1 + e_2:\text{int}}$	$\frac{A \vdash e_1:\text{int} \quad A \vdash e_2:\tau \quad A \vdash e_3:\tau}{A \vdash \text{if } e_1 e_2 e_3:\tau}$

# Example

$$\frac{x:\alpha, y:\beta \vdash x:\alpha}{x:\alpha \vdash \lambda y:\beta. x:\beta \rightarrow \alpha}$$

$\frac{A(x)=\tau}{A \vdash x:\tau}$	$\frac{A, x:\tau \vdash e:\tau'}{A \vdash \lambda x:\tau. e:\tau \rightarrow \tau'}$	$\frac{A \vdash e_1:\tau \rightarrow \tau' \quad A \vdash e_2:\tau}{A \vdash e_1 e_2:\tau'}$
	$A \vdash e_1:\text{int}$	$A \vdash e_1:\text{int}$
	$A \vdash e_2:\text{int}$	$A \vdash e_2:\tau$
$A \vdash i:\text{int}$	$A \vdash e_1 + e_2:\text{int}$	$\frac{A \vdash e_3:\tau}{A \vdash \text{if } e_1 e_2 e_3:\tau}$

# Example

$$\frac{x:\alpha, y:\beta \vdash x:\alpha}{x:\alpha \vdash \lambda y:\beta. x:\beta \rightarrow \alpha}$$

---

$$\emptyset \vdash \lambda x:\alpha. \lambda y:\beta. x:\alpha \rightarrow \beta \rightarrow \alpha$$

$\frac{A(x)=\tau}{A \vdash x:\tau}$	$\frac{A, x:\tau \vdash e:\tau'}{A \vdash \lambda x:\tau. e:\tau \rightarrow \tau'}$	$\frac{A \vdash e_1:\tau \rightarrow \tau' \quad A \vdash e_2:\tau}{A \vdash e_1 e_2:\tau'}$
	$A \vdash e_1:\text{int}$	$A \vdash e_1:\text{int}$
	$A \vdash e_2:\text{int}$	$A \vdash e_2:\tau$
$A \vdash i:\text{int}$	$A \vdash e_1 + e_2:\text{int}$	$A \vdash e_3:\tau$
		$A \vdash \text{if } e_1 e_2 e_3:\tau$



# Not Straightforward

$$x:\alpha, y:\beta \vdash x:\alpha$$

$\frac{A(x)=\tau}{A \vdash x:\tau}$	$\frac{A, x:\tau \vdash e:\tau'}{A \vdash \lambda x:\tau. e:\tau \rightarrow \tau'}$	$\frac{A \vdash e_1:\tau \rightarrow \tau' \quad A \vdash e_2:\tau}{A \vdash e_1 e_2:\tau'}$
	$\frac{A \vdash e_1:\text{int} \quad A \vdash e_2:\text{int}}{A \vdash e_1 + e_2:\text{int}}$	$\frac{A \vdash e_1:\text{int} \quad A \vdash e_2:\tau \quad A \vdash e_3:\tau}{A \vdash \text{if } e_1 e_2 e_3:\tau}$

# Type Checking Algorithm

- There is a simple algorithm for type checking
- Observe that there is only one possible “shape” of the type derivation
  - only one inference rule applies to each form.

$$\frac{\frac{? \vdash x : ?}{? \vdash \lambda y : \beta . x : ?}}{\emptyset \vdash \lambda x : \alpha . \lambda y : \beta . x : ?}$$

# Algorithm (Cont.)

- Walk the proof tree from the root to the leaves, generating the correct environments.
- Assumptions are simply gathered from lambda abstractions.

$$\frac{\frac{x : \alpha, y : \beta \vdash x : ?}{x : \alpha \vdash \lambda y : \beta. x : ?}}{\emptyset \vdash \lambda x : \alpha. \lambda y : \beta. x : ?}$$

# Algorithm (Cont.)

- In a walk from the leaves to the root, calculate the type of each expression.
- The types are completely determined by the type environment and the types of subexpressions.

$$\frac{x : \alpha, y : \beta \vdash x : \alpha}{x : \alpha \vdash \lambda y : \beta. x : \beta \rightarrow \alpha}$$

---

$$\emptyset \vdash \lambda x : \alpha. \lambda y : \beta. x : \alpha \rightarrow \beta \rightarrow \alpha$$

# A Bigger Example

$$\frac{\frac{x : \alpha \rightarrow \alpha, y : \beta \vdash x : \alpha \rightarrow \alpha}{x : \alpha \rightarrow \alpha \vdash \lambda y : \beta. x : \beta \rightarrow \alpha \rightarrow \alpha}}{\emptyset \vdash \lambda x : \alpha \rightarrow \alpha. \lambda y : \beta. x : (\alpha \rightarrow \alpha) \rightarrow \beta \rightarrow \alpha \rightarrow \alpha} \quad \frac{z : \alpha \vdash z : \alpha}{\emptyset \vdash \lambda z : \alpha. z : \alpha \rightarrow \alpha}}{\emptyset \vdash (\lambda x : \alpha \rightarrow \alpha. \lambda y : \beta. x) \lambda z : \alpha. z : (\alpha \rightarrow \alpha) \rightarrow \beta \rightarrow \alpha \rightarrow \alpha}$$

# What Do Types Mean?

- **Thm.** If  $A \vdash e:\tau$  and  $e \rightarrow_{\beta}^* d$ , then  $A \vdash d:\tau$ 
  - Evaluation preserves types.
  
- This is the basis of a claim that there can be no runtime type errors
  - functions applied to data of the wrong type
    - ❖ Adding to a function
    - ❖ Using an integer as a function

# Type Inference

- The *type erasure* of  $e$  is  $e$  with all type information removed (i.e., the untyped term).
- Is an untyped term the erasure of some simply typed term? And what are the types?
- This is a *type inference* problem. We must infer, rather than check, the types.

# Type Inference

- recast the type rules in an equivalent form
- typing in the new rules reduces to a constraint satisfaction problem
- the constraint problem is solvable via term unification.



# New Rules

$$\frac{A(x) = \alpha_x}{A \vdash x : \alpha_x} \quad \frac{A, x : \alpha_x \vdash e : \tau}{A \vdash \lambda x. e : \alpha_x \rightarrow \tau} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 \rightarrow \beta}{A \vdash e_1 e_2 : \beta}$$

$$\frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 = \text{int}}{A \vdash i : \text{int}} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \text{int} \quad \tau_2 = \tau_3}{A \vdash \text{if } e_1 e_2 e_3 : \tau_2}$$

$$\frac{\frac{A(x) = \tau}{A \vdash x : \tau} \quad \frac{A, x : \tau \vdash e : \tau'}{A \vdash \lambda x : \tau. e : \tau \rightarrow \tau'} \quad \frac{A \vdash e_1 : \tau \rightarrow \tau' \quad A \vdash e_2 : \tau}{A \vdash e_1 e_2 : \tau}}{A \vdash i : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \tau}{A \vdash e_1 + e_2 : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \tau \quad A \vdash e_3 : \tau}{A \vdash \text{if } e_1 e_2 e_3 : \tau}$$

- Sidestep the problems by introducing explicit unknowns and constraints

# New Rules

$$\frac{A(x) = \alpha_x}{A \vdash x : \alpha_x} \quad \frac{A, x : \alpha_x \vdash e : \tau}{A \vdash \lambda x. e : \alpha_x \rightarrow \tau} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 \rightarrow \beta}{A \vdash e_1 e_2 : \beta}$$

$$\frac{}{A \vdash i : \text{int}} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 = \text{int}}{A \vdash e_1 + e_2 : \text{int}} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad A \vdash e_3 : \tau_3 \quad \tau_1 = \text{int} \quad \tau_2 = \tau_3}{A \vdash \text{if } e_1 e_2 e_3 : \tau_2}$$

$$\frac{A(x) = \tau \quad A, x : \tau \vdash e : \tau' \quad A \vdash e_1 : \tau \rightarrow \tau' \quad A \vdash e_2 : \tau}{A \vdash x : \tau \quad A \vdash \lambda x : \tau. e : \tau \rightarrow \tau' \quad A \vdash e_1 e_2 : \tau'}$$

$$\frac{}{A \vdash i : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \text{int}}{A \vdash e_1 + e_2 : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \tau \quad A \vdash e_3 : \tau}{A \vdash \text{if } e_1 e_2 e_3 : \tau}$$

- Type assumption for variable  $x$  is a fresh variable  $\alpha_x$

# New Rules

$$\frac{A(x) = \alpha_x}{A \vdash x : \alpha_x} \quad \frac{A, x : \alpha_x \vdash e : \tau}{A \vdash \lambda x. e : \alpha_x \rightarrow \tau} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 \rightarrow \beta}{A \vdash e_1 e_2 : \beta}$$

$$\frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 = \text{int}}{A \vdash i : \text{int}} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \text{int} \quad \tau_2 = \tau_3}{A \vdash \text{if } e_1 e_2 e_3 : \tau_2}$$

$$\frac{\frac{A(x) = \tau}{A \vdash x : \tau} \quad \frac{A, x : \tau \vdash e : \tau'}{A \vdash \lambda x. \tau. e : \tau \rightarrow \tau'} \quad \frac{A \vdash e_1 : \tau \rightarrow \tau' \quad A \vdash e_2 : \tau}{A \vdash e_1 e_2 : \tau'}}{\frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \tau}{A \vdash e_2 : \tau}} \quad \frac{A \vdash e_2 : \text{int} \quad A \vdash e_3 : \tau}{A \vdash \text{if } e_1 e_2 e_3 : \tau}}$$

- Hypotheses are all arbitrary
  - Can always complete a derivation, pending constraint resolution

# New Rules

$$\frac{A(x) = \alpha_x}{\vdash x : \alpha_x} \quad \frac{A, x : \alpha_x \vdash e : \tau}{A \vdash \lambda x. e : \alpha_x \rightarrow \tau} \quad \frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 \rightarrow \beta}{A \vdash e_1 e_2 : \beta}$$

$$\frac{A(x) = \tau}{A \vdash x : \tau} \quad \frac{A, x : \tau \vdash e : \tau'}{A \vdash \lambda x : \tau. e : \tau \rightarrow \tau'} \quad \frac{A \vdash e_1 : \tau \rightarrow \tau' \quad A \vdash e_2 : \tau}{A \vdash e_1 e_2 : \tau'}$$

$$\frac{}{A \vdash i : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \text{int}}{A \vdash e_1 + e_2 : \text{int}} \quad \frac{A \vdash e_1 : \text{int} \quad A \vdash e_2 : \tau \quad A \vdash e_3 : \tau}{A \vdash \text{if } e_1 e_2 e_3 : \tau}$$

$$\frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad \tau_1 = \tau_2 = \text{int}}{A \vdash e_1 + e_2 : \text{int}}$$

$$\frac{A \vdash e_1 : \tau_1 \quad A \vdash e_2 : \tau_2 \quad A \vdash e_3 : \tau_3 \quad \tau_1 = \text{int} \quad \tau_2 = \tau_3}{A \vdash \text{if } e_1 e_2 e_3 : \tau_2}$$

- Equality conditions represented as side constraints

# Solutions of Constraints

- The new rules generate a system of type equations.
- Intuitively, a solution of these equations gives a derivation.
- A solution is a substitution  $\text{Vars} \rightarrow \text{Types}$  such that the equations are satisfied.

# Example

$$\alpha = \beta \rightarrow \gamma$$

$$\alpha = \gamma \rightarrow \beta$$

$$\beta = \text{int}$$

- A solution is

$$\alpha = \text{int} \rightarrow \text{int}, \beta = \text{int}, \gamma = \text{int}$$

# Solving Type Equations

- Term equations are a unification problem.
  - Solvable in near-linear time using a union-find based algorithm.
  
- No solutions  $\alpha = T[\alpha]$  are permitted
  - The *occurs check*.
  - The check is omitted if we allow infinite types.

# Unification

- Four rules.
- If no inconsistency or occurs check violation found, system has a solution.
  - $\text{int} = x \rightarrow y$

$$\mathcal{S} \cup \{\alpha = \alpha\} \Rightarrow \mathcal{S}$$

$$\mathcal{S} \cup \{\alpha = \tau\} \Rightarrow \mathcal{S}[\tau/\alpha] \cup \{\alpha \equiv \tau\}$$

$$\mathcal{S} \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow \mathcal{S} \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$\mathcal{S} \cup \{\text{int} = \text{int}\} \Rightarrow \mathcal{S}$$



# Syntax

- We distinguish *solved* equations  $\alpha \cong \tau$
- Each rule manipulates only unsolved equations.

$$\mathcal{S} \cup \{\alpha = \alpha\} \Rightarrow \mathcal{S}$$

$$\mathcal{S} \cup \{\alpha = \tau\} \Rightarrow \mathcal{S}[\tau / \alpha] \cup \{\alpha \cong \tau\}$$

$$\mathcal{S} \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow \mathcal{S} \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$\mathcal{S} \cup \{\text{int} = \text{int}\} \Rightarrow \mathcal{S}$$

# Rules 1 and 4

- Rules 1 and 4 eliminate trivial constraints.
- Rule 1 is applied in preference to rule 2
  - the only such possible conflict

$$\mathcal{S} \cup \{\alpha = \alpha\} \Rightarrow \mathcal{S}$$

$$\mathcal{S} \cup \{\alpha = \tau\} \Rightarrow \mathcal{S}[\tau/\alpha] \cup \{\alpha \equiv \tau\}$$

$$\mathcal{S} \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow \mathcal{S} \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$\mathcal{S} \cup \{\text{int} = \text{int}\} \Rightarrow \mathcal{S}$$

# Rule 2

- Rule 2 eliminates a variable from all equations but one (which is marked as solved).
  - Note the variable is eliminated from all unsolved as well as solved equations

$$S \cup \{\alpha = \alpha\} \Rightarrow S$$

$$S \cup \{\alpha = \tau\} \Rightarrow S[\tau/\alpha] \cup \{\alpha \equiv \tau\}$$

$$S \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow S \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$S \cup \{\text{int} = \text{int}\} \Rightarrow S$$

# Rule 3

- Rule 3 applies structural equality to non-trivial terms.
- Note rule 4 is a degenerate case of rule 3 for a type constructor of arity zero.

$$\mathcal{S} \cup \{\alpha = \alpha\} \Rightarrow \mathcal{S}$$

$$\mathcal{S} \cup \{\alpha = \tau\} \Rightarrow \mathcal{S}[\tau / \alpha] \cup \{\alpha \equiv \tau\}$$

$$\mathcal{S} \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow \mathcal{S} \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$\mathcal{S} \cup \{\text{int} = \text{int}\} \Rightarrow \mathcal{S}$$

# Correctness

- Each rule preserves the set of solutions.
  - Rules 1 and 4 eliminate trivial constraints.
  - Rule 2 substitutes equals for equals.
  - Rule 3 is the definition of equality on function types.

$$\mathcal{S} \cup \{\alpha = \alpha\} \Rightarrow \mathcal{S}$$

$$\mathcal{S} \cup \{\alpha = \tau\} \Rightarrow \mathcal{S}[\tau / \alpha] \cup \{\alpha \equiv \tau\}$$

$$\mathcal{S} \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow \mathcal{S} \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$\mathcal{S} \cup \{\text{int} = \text{int}\} \Rightarrow \mathcal{S}$$

# Termination

- ❑ Rules 1 and 4 reduce the number of equations.
- ❑ Rule 2 reduces the number of variables in unsolved equations.
- ❑ Rule 3 decreases the height of terms.

$$\mathcal{S} \cup \{\alpha = \alpha\} \Rightarrow \mathcal{S}$$

$$\mathcal{S} \cup \{\alpha = \tau\} \Rightarrow \mathcal{S}[\tau/\alpha] \cup \{\alpha \equiv \tau\}$$

$$\mathcal{S} \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow \mathcal{S} \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$\mathcal{S} \cup \{\text{int} = \text{int}\} \Rightarrow \mathcal{S}$$

# Termination (Cont.)

- Rules 1, 3, and 4 always terminate
  - because terms must eventually be reduced to height 0.
- Eventually rule 2 is applied, reducing the number of variables.

$$S \cup \{\alpha = \alpha\} \Rightarrow S$$

$$S \cup \{\alpha = \tau\} \Rightarrow S[\tau/\alpha] \cup \{\alpha \equiv \tau\}$$

$$S \cup \{\tau_1 \rightarrow \tau_2 = \tau_3 \rightarrow \tau_4\} \Rightarrow S \cup \{\tau_1 = \tau_3, \tau_2 = \tau_4\}$$

$$S \cup \{\text{int} = \text{int}\} \Rightarrow S$$

# A Nitpick

- We really need one more operation.
- $\tau = \alpha$  should be flipped to  $\alpha = \tau$  if  $\tau$  is not a variable.
  - Needed to ensure rule 2 applies whenever possible.
  - We just assume equations are maintained in this “normal form”.



# Solutions

- The final system is a solution.
  - There is one equation  $\alpha \cong \tau$  for each variable.
  - This is a substitution with all the solutions of the original system
- Must also perform occurs check to guarantee there are no recursive constraints.

# Example

*rewrites*

$$\alpha = \beta \rightarrow \gamma, \alpha = \gamma \rightarrow \beta, \beta = \text{int}$$



# Example

*rewrites*

$$\alpha = \beta \rightarrow \gamma, \alpha = \gamma \rightarrow \beta, \beta = \text{int}$$

---

$$\alpha = \text{int} \rightarrow \gamma, \alpha = \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

# Example

*rewrites*

$$\alpha = \beta \rightarrow \gamma, \alpha = \gamma \rightarrow \beta, \beta = \text{int}$$

---

$$\alpha = \text{int} \rightarrow \gamma, \alpha = \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma \rightarrow \text{int} = \text{int} \rightarrow \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

# Example

*rewrites*

$$\alpha = \beta \rightarrow \gamma, \alpha = \gamma \rightarrow \beta, \beta = \text{int}$$

---

$$\alpha = \text{int} \rightarrow \gamma, \alpha = \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma \rightarrow \text{int} = \text{int} \rightarrow \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma = \text{int}, \text{int} = \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

# Example

*rewrites*

$$\alpha = \beta \rightarrow \gamma, \alpha = \gamma \rightarrow \beta, \beta = \text{int}$$

---

$$\alpha = \text{int} \rightarrow \gamma, \alpha = \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma \rightarrow \text{int} = \text{int} \rightarrow \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma = \text{int}, \text{int} = \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\text{int} = \text{int}, \gamma \cong \text{int}, \alpha \cong \text{int} \rightarrow \text{int}, \beta \cong \text{int}$$

# Example

*rewrites*

$$\alpha = \beta \rightarrow \gamma, \alpha = \gamma \rightarrow \beta, \beta = \text{int}$$

---

$$\alpha = \text{int} \rightarrow \gamma, \alpha = \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma \rightarrow \text{int} = \text{int} \rightarrow \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma = \text{int}, \text{int} = \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\text{int} = \text{int}, \gamma \cong \text{int}, \alpha \cong \text{int} \rightarrow \text{int}, \beta \cong \text{int}$$

---

$$\gamma \cong \text{int}, \alpha \cong \text{int} \rightarrow \text{int}, \beta \cong \text{int}$$

# An Example of Failure

$$\alpha = \beta \rightarrow \gamma, \alpha = \gamma \rightarrow (\beta \rightarrow \beta), \beta = \text{int}$$

$$\alpha = \text{int} \rightarrow \gamma, \alpha = \gamma \rightarrow (\text{int} \rightarrow \text{int}), \beta \cong \text{int}$$

$$\gamma \rightarrow (\text{int} \rightarrow \text{int}) = \text{int} \rightarrow \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

$$\gamma = \text{int}, \text{int} \rightarrow \text{int} = \gamma, \alpha \cong \gamma \rightarrow \text{int}, \beta \cong \text{int}$$

$$\text{int} \rightarrow \text{int} = \text{int}, \gamma \cong \text{int} \rightarrow \text{int}, \alpha \cong \text{int} \rightarrow \text{int}, \beta \cong \text{int}$$



# Notes

- ❑ The algorithm produces the *most general unifier* of the equations.
  - All solutions are preserved.
- ❑ Less general solutions are all substitution instances of the most general solution.
- ❑ There exists more efficient algorithm, amortized time complexity is close to linear

# Application – Treating Program Property as A Type

- INT, BOOL, and STRING are types, and
  - “ALLOCATED” and “FREED” can also be treated as types.

$$\frac{e_1:\tau \quad e_2:\tau}{e_1 = e_2:\tau}$$

For example,  $p=q$

# Uses

- Find bugs
  - Every equivalence class with a *malloc* should have a *free*
- Alias analysis
- Implemented for C in a tool Lackwit
  - O'Callahan & Jackson

# Where is Type Inference Strong?

- ❑ Handles data structures smoothly
- ❑ Works in infinite domains
  - Set of types is unlimited
- ❑ No forwards/backwards distinction
- ❑ Type polymorphism good fit for context sensitivity

# Where is Type Inference Weak?

- ❑ No flow sensitivity
  - Equality-based analysis only gets equivalence classes
- ❑ Context-sensitive analyses don't always scale
  - Type polymorphism can lead to exponential blowup in constraints

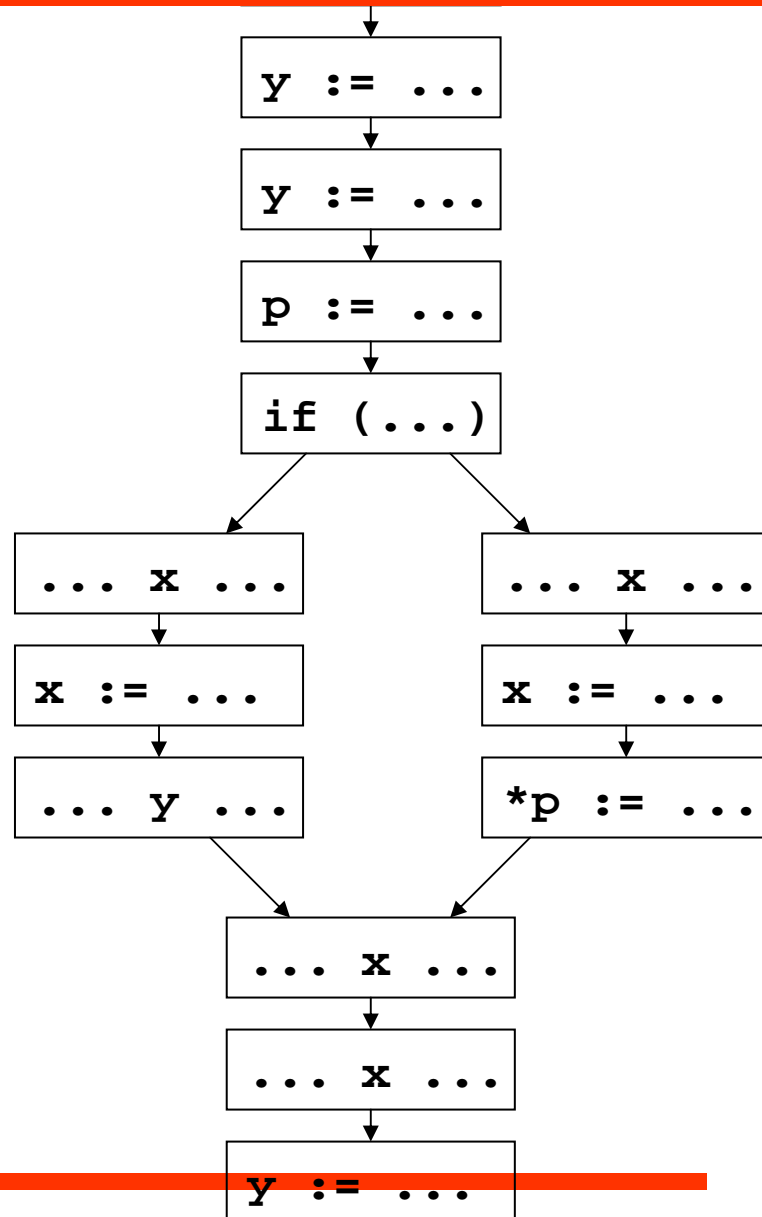
# Flow Sensitive: Data Flow Analysis

# An example DFA: reaching definitions

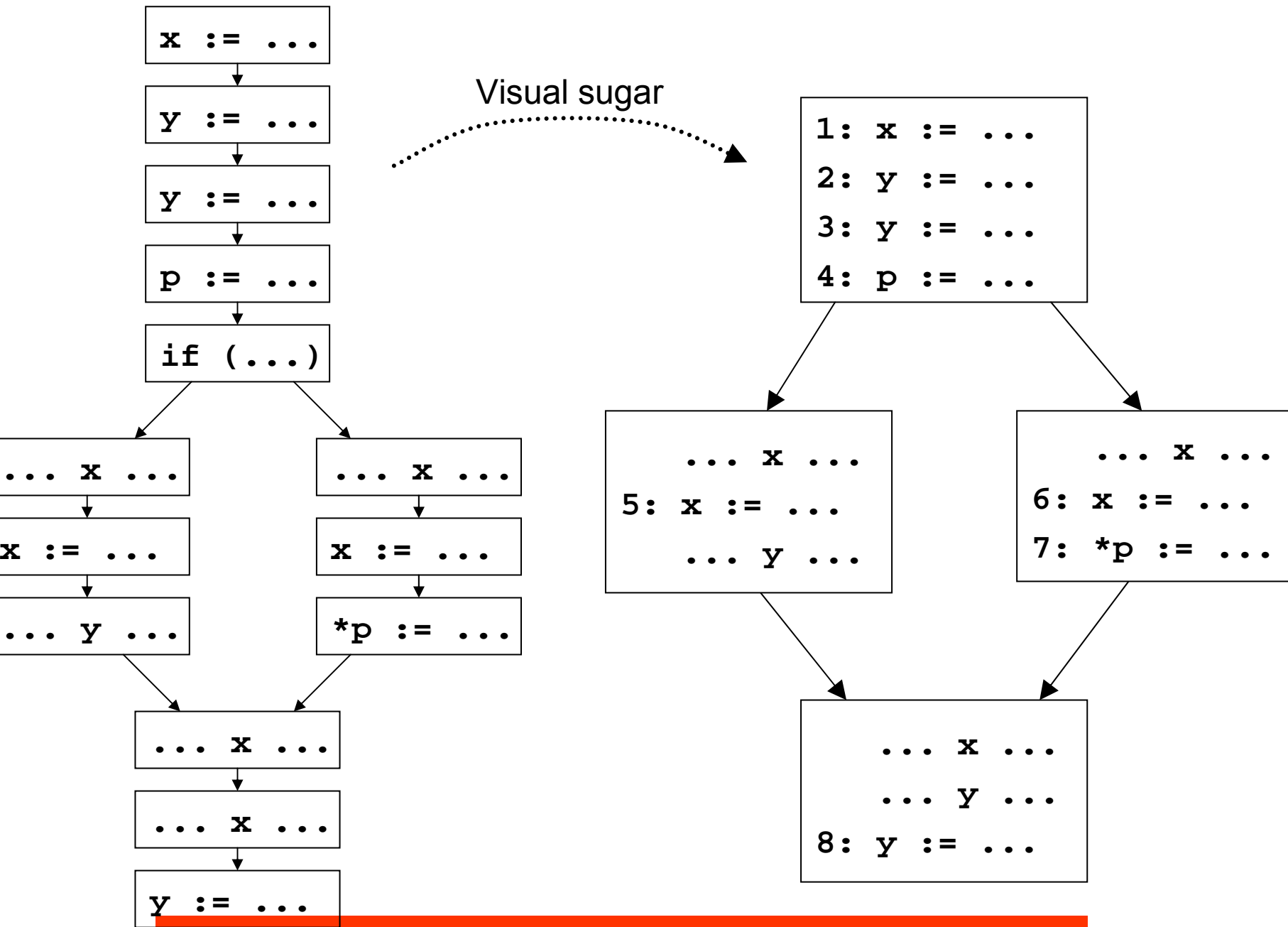
- ❑ For each use of a variable, determine what assignments could have set the value being read from the variable
- ❑ Information useful for:
  - performing constant and copy prop
  - detecting references to undefined variables
  - presenting “def/use chains” to the programmer
  - building other representations, like the program dependence graph
- ❑ Let’s try this out on an example

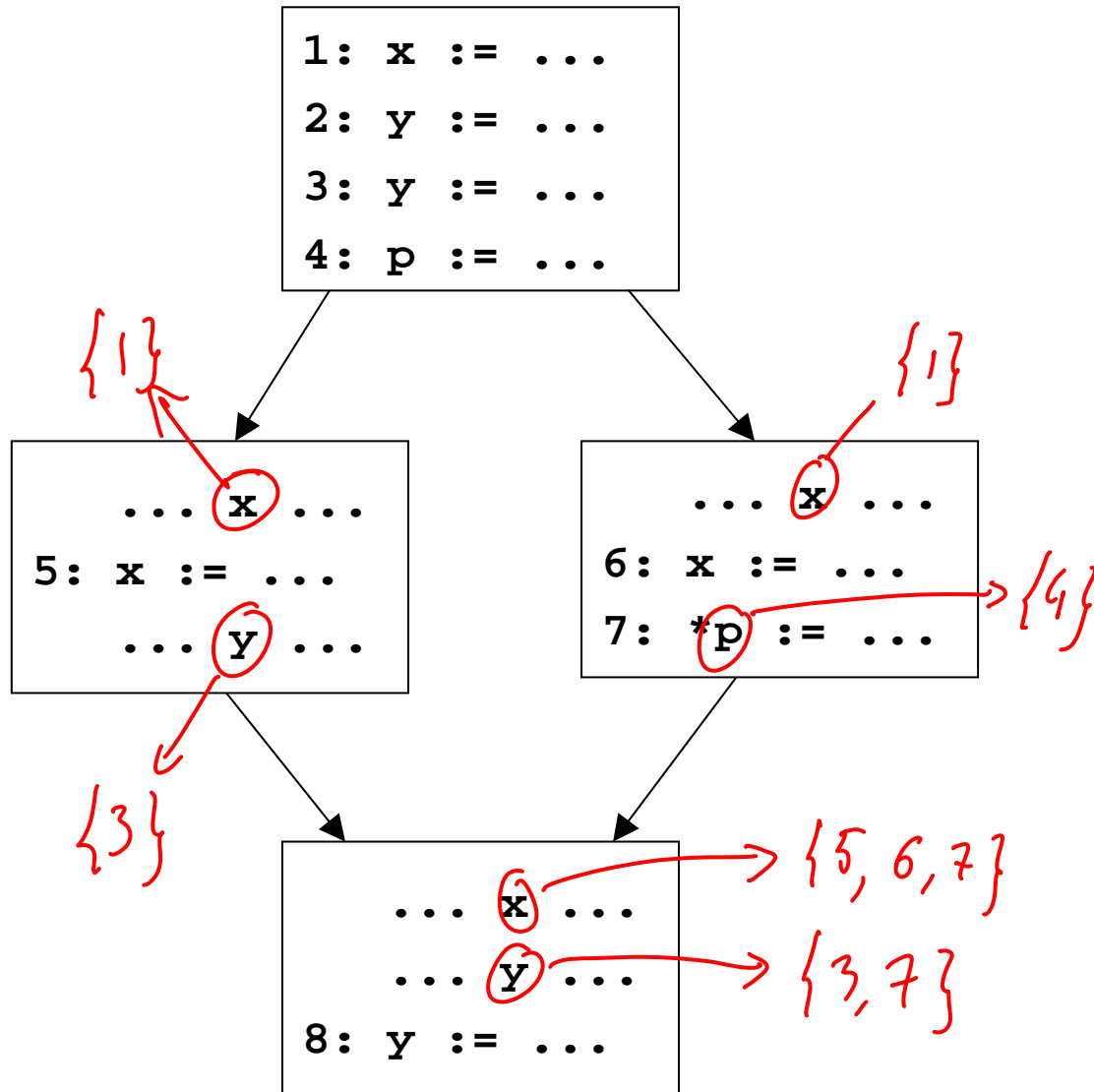
# Example CFG

```
x := ...
y := ...
y := ...
p := ...
if (...) {
  ... x ...
  x := ...
  ... y ...
}
else {
  ... x ...
  x := ...
  *p := ...
}
... x ...
... y ...
y := ...
```







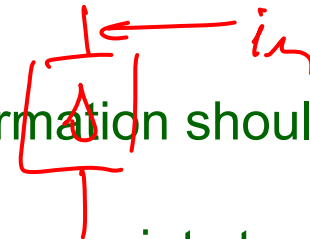


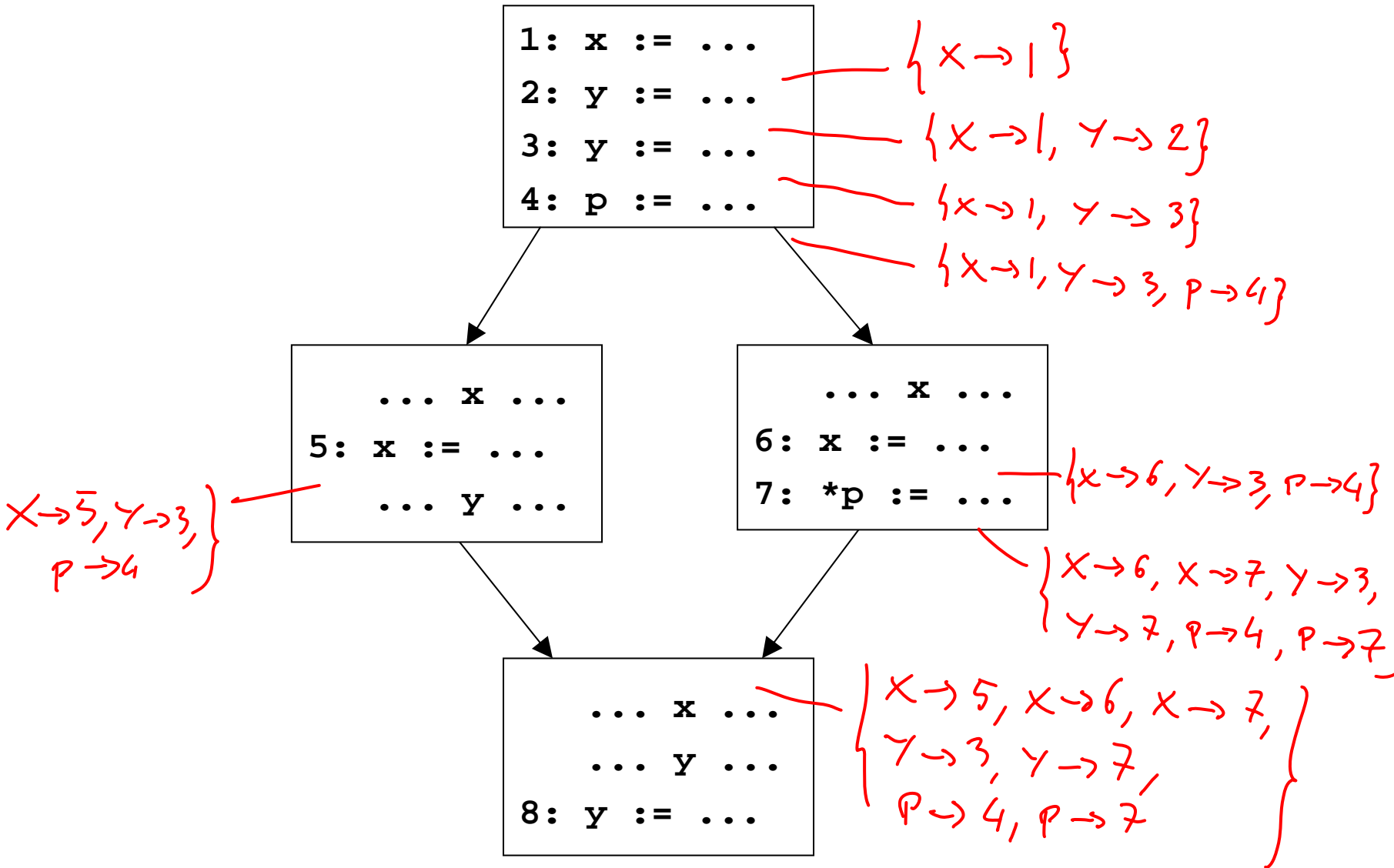
# Safety

- Safety:
  - can have more bindings than the “true” answer, but can’t miss any

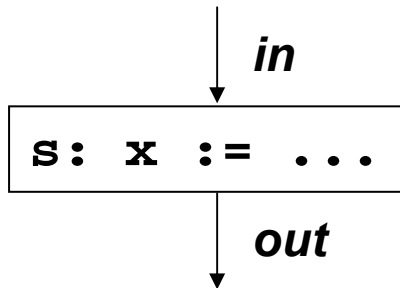
# Reaching definitions generalized

- ❑ Computed information at a program point is a set of  $\text{var} \rightarrow \text{stmt}$  bindings
  - eg:  $\{ x \rightarrow s_1, x \rightarrow s_2, y \rightarrow s_3 \}$
- ❑ How do we get the previous info we wanted?
  - if a var  $x$  is used in a stmt whose incoming info is  $in$ , then:  $\{ s \mid (x \rightarrow s) \in in \}$
- ❑ This is a common pattern
  - generalize the problem to define what information should be computed at each program point
  - use the computed information at the program points to get the original info we wanted

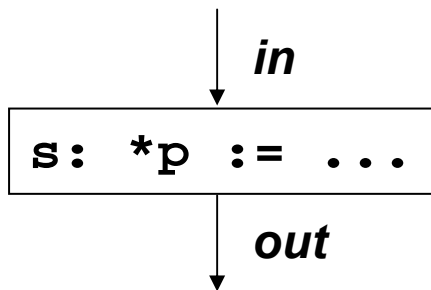




# Constraints for reaching definitions

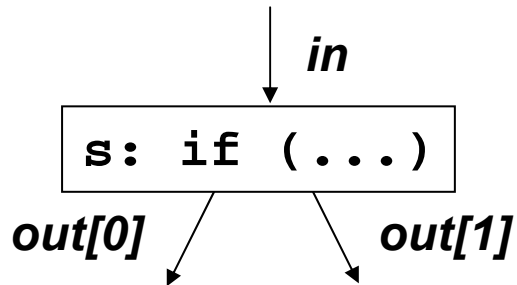


$$\text{out} = \text{in} - \{x \rightarrow s' \mid s' \in \text{stmts}\} \cup \{x \rightarrow s\}$$



$$\begin{aligned} \text{out} = \text{in} - \{x \rightarrow s' \mid x \in \text{must-point-to}(p) \wedge \\ s' \in \text{stmts}\} \\ \cup \{x \rightarrow s \mid x \in \text{may-point-to}(p)\} \end{aligned}$$

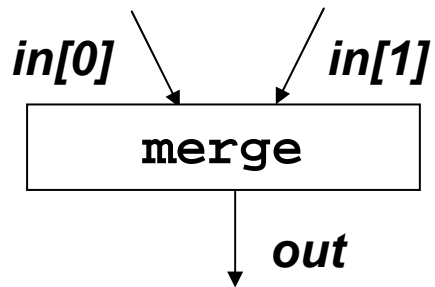
# Constraints for reaching definitions



$$out[0] = in \wedge$$

$$out[0] = in$$

more generally:  $\forall i . out[i] = in$



$$out = in[0] \cup in[1]$$

more generally:  $out = \bigcup_i in[i]$

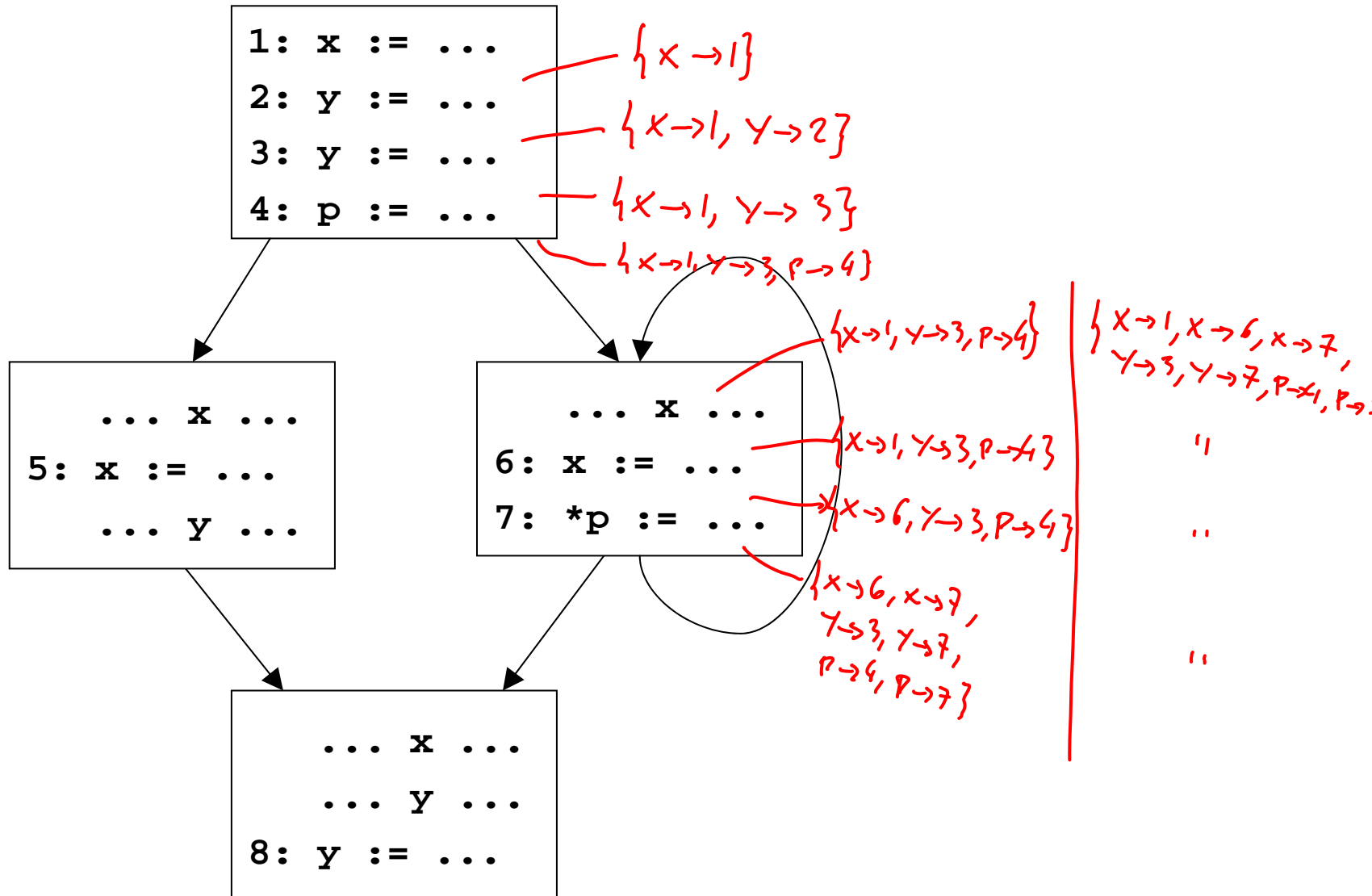
# Flow functions

- The constraint for a statement kind  $s$  often have the form:  $out = F_s(in)$
- $F_s$  is called a ~~flow~~ function
  - other names for it: dataflow function, transfer function
- Given information  $in$  before statement  $s$ ,  $F_s(in)$  returns information after statement  $s$



# The Problem of Loops

- ❑ If there is no loop, the topological order can be adopted to evaluate transfer functions of statements.
- ❑ What if loops?



# Solution: iterate!

- ❑ Initialize all sets to the empty
- ❑ Store all nodes onto a worklist
- ❑ while worklist is not empty:
  - remove node  $n$  from worklist
  - apply flow function for node  $n$
  - update the appropriate set, and add nodes whose inputs have changed back onto worklist

# Termination

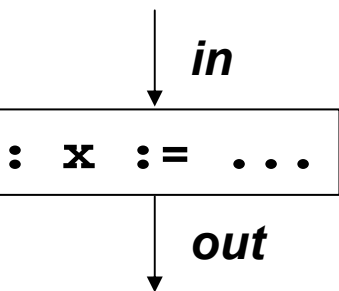
- How do we know the algorithm terminates?
- Because
  - operations are *monotonic*
  - the domain is finite

# Monotonicity

- Operation  $f$  is monotonic if

$$X \subseteq Y \Rightarrow f(x) \subseteq f(y)$$

- We require that all operations be monotonic
  - Easy to check for the set operations
  - Easy to check for all transfer functions; recall:



$$out = in - \{ x \rightarrow s' \mid s' \in stmts \} \cup \{ x \rightarrow s \}$$

# Termination again

- ❑ To see the algorithm terminates
  - All variables start empty
  - Variables and rhs's only increase with each update
  - Sets can only grow to a max finite size
  
- ❑ Together, these imply termination

# What Else In DFA

- ❑ May vs. must
- ❑ Backward vs. Forward
- ❑ Lattice
  - Mere goal: help prove the termination of the analysis
  - To show the domain is finite (has finite height)

# Where is Dataflow Analysis Useful?

- ❑ Best for flow-sensitive, context-insensitive, distributive problems on small pieces of code
  - E.g., the examples we've seen and many others
- ❑ Extremely efficient algorithms are known
  - Use different representation than control-flow graph, but not fundamentally different



# Where is Dataflow Analysis Weak?

- Lots of places

# Data Structures

- ❑ Not good at analyzing data structures
- ❑ Works well for atomic values
  - Labels, constants, variable names
- ❑ Not easily extended to arrays, lists, trees, etc.

# The Heap

- ❑ Good at analyzing flow of values in local variables
- ❑ No notion of the heap in traditional dataflow applications
  - Aliasing

# Context Sensitivity

- Standard dataflow techniques for handling context sensitivity don't scale well

# Flow Sensitivity (Beyond Procedures)

- ❑ Flow sensitive analyses are standard for analyzing single procedures
- ❑ Not used (or not aware of uses) for whole programs
  - Too expensive

# The Call Graph

- ❑ Dataflow analysis requires a call graph
  - Or something close
- ❑ Inadequate for higher-order programs
  - First class functions
  - Object-oriented languages with dynamic dispatch
- ❑ Call-graph hinders algorithmic efficiency

# Coming Back: The Essence of Static Analysis

- ❑ Examine the program text (no execution)
- ❑ Build a model of the program state
  - An abstract of the run-time state
- ❑ Reason over the possible behaviors.
  - E.g. “run” the program over the abstract state
- ❑ The property an analysis needs to promise is that it **TERMINATES**
  - Slogan of most researchers:

Finite Lattices + Monotonic Functions =  
Program Analysis

---

# Tips on Designing Analysis

- ❑ Program analysis is a formalization of INTUITIVE insights.
  - Type inference
  - Reaching definition
  - ...
- ❑ Steps
  - Look at the code (segment), gain insights;
  - More systematic: manually “runs” through the code with your abstraction.
  - Works? Good, lets do formalization.



# Next Lecture

- Dynamic Program Analysis