

Why Global Dataflow Analysis?

Answer key questions at compile-time about the flow of values and other program properties over control-flow paths

Compiler fundamentals

What defs. of x reach a given use of x (and *vice-versa*)?

What $\{\langle \text{ptr}, \text{target} \rangle\}$ pairs are possible at each statement?

Scalar dataflow optimizations

Are any uses reached by a particular definition of x ?

Has an expression been computed on all incoming paths?

What is the innermost loop level at which a variable is defined?

Correctness and safety:

Is variable x defined on every path to a use of x ?

Is a pointer to a local variable live on exit from a procedure?

Parallel program optimization, program understanding, ...

Common Applications of Global Dataflow Analysis

Preliminary Analyses

- Pointer Analysis
- Detecting uninitialized variables
- Type inference
- Strength Reduction for Induction Variables

Static Computation Elimination

- Dead Code Elimination (DCE)
- Constant Propagation
- Copy Propagation

Redundancy Elimination

- Local Common Subexpression Elimination (CSE)
- Global Common Subexpression Elimination (GCSE)
- Loop-invariant Code Motion (LICM)
- Partial Redundancy Elimination (PRE)

Code Generation

- Liveness analysis for register allocation

Dataflow Analysis: Our Objectives

1. To distinguish different types of dataflow problems
 - may v. must*
 - forward v. backward*
 - intersection v. union*
2. To set up and solve the dataflow equations for a basic dataflow problem
3. To identify dataflow problems needed for a given optimization

Preliminary definitions

Value, Storage location, variable, pointer : *these should be familiar*

Alias or alias pair : Two different names for the same storage location

Reference : An occurrence of a name in a program statement

Use of a variable : A reference that *may read* the value of the variable.

Definition of a variable : A reference that *may store* a value into the storage location(s) named by the variable.

Examples: Assignment; FOR; input I/O

Unambiguous definition : *guaranteed* to store to X

must

Ambiguous definition : *may* store to X

may

Ambiguity comes from aliases, unpredictable side effects of procedure calls, arrays

Dataflow Analysis Basics

Point: A location in a basic block just before or after some statement.

Path: A path from p_1 to p_n is a sequence of points p_1, p_2, \dots, p_n such that (intuitively) some execution can visit these points in order.

[See section 9.2.1 of 2nd Ed. Dragon Book for formal definition]

Identifying Defs, Refs

	Examples
1	<code>X = Y + 1; // r_1_Y, d_1_X</code>
2	
3	<code>p = cond? &X : &Z; // d_3_p (what about X and Z?)</code>
4	<code>*p = Y + 1; // r_4_Y, d_4_X, d_4_Z</code>
5	
6	<code>A[i] = Y + 1; // r_6_Y, d_6_A</code>
7	<code>A[i+1] = Y + 1; // r_7_Y, d_7_A</code>
8	
9	<code>// On line 54: list->next = new ListNode(...);</code>
10	<code>list->next->val = list->val + 1; // r_10_H_54->val, d_10_H_54->val</code>

Principles of “naming” memory locations

- *Variable names* identify (sets of) memory locations
- Defs, refs apply to individual variables
- Arrays are usually named as a single variable
- Heap allocated objects can be named (i.e., treated as “dummy variables”) in different ways
 - Most common: H_k , k = line number of malloc/new

An Example Dataflow Problem: Reaching Definitions

Reaching Definitions

May or Must?

$\forall p$, compute REACH(p): the set of defs that reach point p .

Definition d reaches point p if there is a path from the point after d to p such that d is not *killed* along that path.

Kill of a Definition: A definition d of variable V is killed on a path if there is an unambiguous definition of V on that path.

Dataflow variables (for each block B)

- In (B) \equiv the set of defs that reach the point before first statement in B
- Out (B) \equiv the set of defs that reach the point after last statement in B
- Gen (B) \equiv the set of defs in B that are not killed in B .
- Kill (B) \equiv the set of all defs that are killed in B
(i.e., on the path from entry to exit of B , if def $d \notin B$;
or on the path from d to exit of B , if def $d \in B$).

The difference:

Gen (B), Kill (B) are local properties of block B alone.

In (B), Out (B) are global dataflow properties.

Dataflow Analysis for Reaching Definitions

Dataflow equations

$$In[B] = \bigcup_{p:p \rightarrow B} Out[p]$$

$$Out[B] = Gen[B] \cup (In[B] - Kill[B])$$

Dataflow algorithms

Goal: solve these $2n$ simultaneous dataflow equations ($n = \#$ basic blocks)

- Block-structured graph (no GOTO; no BREAK from loops):
 - bottom-up evaluation, one scope at a time
(section 10.5 of 1st Ed. Dragon Book)
- General flow-graphs:
 - iterative solution

Iterative Algorithm for Reaching Definitions

1. Initialize:

```

/* If there are globals or formals, in[s] ≠ φ */
in[B] = φ           ∀B
out[B] = gen[B]    ∀B

```

2. Iterate until Out[B] does not change:

```

do
  change = false
  for each block B do
    In[B] =  $\bigcup_{p:p \rightarrow B} Out[p]$ 
    oldout = Out[B]
    Out[B] = Gen[B]  $\cup$  (In[B] - Kill[B])
    if (oldout ≠ Out[B]) change = true
  end
while (change == true)

```

What is the algorithm doing?

```
1 (d0) X = ...
2     if (...) {
3         ...
4     } else {
5 (d1)     X = ...
6     }
7     ...
8 (d2) X = ...
9     ...
```

```
1 (d0) X = ...
2     while (...) {
3         ...
4         if (...) {
5 (d1)     X = ...
6         } else { ... }
7         ...
8     }
9     ...
```

Convergence of the Algorithm

OUT[B] must converge in a finite #iterations

- $Out[B]$ is finite $\forall B$
- $Out[B]$ never decreases for any B
 - Only KILL sets (constants) are ever subtracted from OUT sets
 - IN sets never decrease (if OUT sets never decrease)
 - But isn't that a circular argument?

Acyclic Property

- Definitions need propagate only over acyclic paths
- ⇐ Each block only adds $Gen[B]$, subtracts $Kill[B]$
 - $\cup, -$: only need to add, remove *once*
- ⇒ Must visit each block exactly once
- Also need one final iteration to check convergence

Section 10.10 of 1st Ed. Dragon Book has more details.

Efficient Orderings for Visiting Basic Blocks

[Assume *reducible* graphs for now \Rightarrow Cycles “formed by” back edges]

1. No back edges: 2 iterations
2. 1 back edge (on any acyclic path): 3 iterations
3. k back edges on an acyclic path: $k + 2$ iterations

Efficient Orderings for Visiting Basic Blocks

Goal: Propagate information as far as possible in each iteration

Postorder and Reverse Postorder

- Depth-first spanning tree (DFST): tree constructed by Depth-first Search
- DFST has 3 kinds of edges: *tree edges*, *cross-edges*, *up-edges*
- Graph excluding up-edges is acyclic (DAG)
- *Postorder* (on original graph) \equiv postorder traversal of resulting DAG

Properties of Reverse Postorder

- If $B_1 \rightarrow B_2$, then B_1 is visited before B_2 , except for up-edges of DFST.
- If CFG is reducible, up-edges are exactly the back edges!
- In any case, max. # number of up-edges on any acyclic path is never more than maximum loop nesting depth

Efficiency of the Algorithm

Rule-of-thumb: Typically 5 iterations or less!
(when dataflow information propagates only over acyclic paths)

Efficient dataflow ordering

- Use Reverse Postorder (RPO) for “forward” dataflow problems
 - Use Postorder (PO) for “backward” dataflow problems
- ⇒ Information propagates “as far as possible” in each iteration, until it reaches a “retreating” DFS edge. It flows across the retreating DFS edge in the next iteration.

Rule of thumb

- Knuth [1971]: Max. #up-edges on each acyclic path is typically 3 or fewer.

Available Expressions

Definitions

Available expressions: $x + y$ is available at point p if:

- (a) every path to p evaluates $x + y$
- (b) between the last such evaluation and p on each path, neither x nor y is modified.

Kill: Block B kills $x + y$ if it may assign to x or y , and it does not subsequently recompute $x + y$

Generate: Block B generates $x + y$ if it definitely evaluates $x + y$, and it does not subsequently modify x or y .

Dataflow variables:

Let \mathcal{U} = universal set of expressions in the program. Then:

$$\begin{aligned}
 in[B] &= \{\epsilon \in \mathcal{U} \mid \epsilon \text{ is avail at entry to } B\} \\
 out[B] &= \{\epsilon \in \mathcal{U} \mid \epsilon \text{ is avail at exit from } B\} \\
 e_gen[B] &= \{\epsilon \in \mathcal{U} \mid \epsilon \text{ is generated by } B\} \\
 e_kill[B] &= \{\epsilon \in \mathcal{U} \mid \epsilon \text{ is killed by } B\}
 \end{aligned}$$

Naming Expressions

		Examples
1	<code>a = x * y;</code>	<code>// eval e_1: x * y</code>
2	<code>b = x * y;</code>	<code>// eval e_1: x * y: redundant</code>
3	<code>x = 2;</code>	<code>// "kills" e_1</code>
4	<code>c = x * y;</code>	<code>// eval e_1: x * y</code>
5		
6	<code>if (...) { x=5; t= x+y; }</code>	<code>// eval e_2: x+y</code>
7	<code>else { x=9; t= x+y; }</code>	<code>// eval e_2: x+y</code>
8	<code>x = x+y;</code>	<code>// eval e_2: x+y: redundant!</code>
9		
10	<code>p = cond? &X : &Z;</code>	
11	<code>... = *p + 1;</code>	<code>// e_3: *p; e_4: X+1, e_5: Z+1</code>
12		<code>// e_4, e_5 may not be evaluated</code>
13	<code>... = X + 1;</code>	<code>// eval e_3: X+1 may not be redundant</code>

Dataflow Analysis for Available Expressions

Dataflow equations:

$$In[B] =$$

$$Out[B] =$$

Algorithm is identical to *Reaching Definitions* except:

- Confluence operator is \cap instead of \cup
- Algorithm must initialize sets as follows:

$$\begin{aligned} In[s] &= \phi \\ Out[s] &= e_gen[s] \\ Out[B] &= \mathcal{U} - e_kill[B] \end{aligned} \quad \forall B \neq s$$

Live Variables

Live Variables

Variable x is live at point p if x may be used along some path starting at p .

Dataflow variables

$$\begin{aligned}
 in[B] &= \{x \in \mathcal{V} \mid x \text{ is live at entry to } B\} \\
 out[B] &= \{x \in \mathcal{V} \mid x \text{ is live at exit from } B\} \\
 def[B] &= \{x \in \mathcal{V} \mid x \text{ is assigned in } B \text{ prior to use in } B\} \\
 use[B] &= \{x \in \mathcal{V} \mid x \text{ may be used in } B \text{ prior to being assigned in } B\}
 \end{aligned}$$

Dataflow equations

$$\begin{aligned}
 In[B] &= \\
 Out[B] &=
 \end{aligned}$$

General Approach to Dataflow Analysis

1. *Choose dataflow variables for problems of interest:*

$Gen(B) \equiv$ “information” generated in block B

$Kill(B) \equiv$ “information” killed in block B

$In(B), Out(B)$

2. *Set up dataflow equations*

Q. what is the transfer function for each block? E.g.,

$$Out[B] = Gen[B] \cup (In[B] - Kill[B])$$

Q. is it a forward vs. backward problem? E.g.,

$$In[B] = \bigcup_{p:p \rightarrow B} Out[p] \quad \text{or} \quad Out[B] = \bigcup_{s:B \rightarrow s} In[s]$$

Q. what is the “confluence” operator: \bigcup, \bigcap , other?

3. *Solve iteratively until convergence*

Postorder or Reverse Postorder

Def-Use and Use-Def Chains

Definitions

Use-Def chain or ud-chain: For each use u of a variable v , $\text{DEFS}(u)$ is the set of instructions that may have defined v last prior to u .

Def-Use chain or du-chain: For each def d of a variable v , $\text{USES}(d)$ is the set of instructions that may use the value of v computed at d

Note: $d \in \text{DEFS}(u)$ *iff* $u \in \text{USES}(d)$

Note: du-chains (or ud-chains) form a graph

Comparing with SSA

- Multiple defs reach each use, unlike SSA
- More edges in def-use graph than in SSA graph
- But fewer variable names, no ϕ functions

Computing and using du-chains and ud-chains

Construction

- Construct $DEFS(u)$ from the results of Reaching Definitions.
 - Then invert $DEFS$ to compute $USES$.
- ⇒ We can build chains very efficiently!

Some applications of chains:

- Building live ranges for graph-coloring register allocation
- Constant propagation
- Dead-code elimination
- Loop-invariant code motion