



COMP 512
Rice University
Spring 2015

Data-Flow Analysis

Dominators to Reaching Definitions

Copyright 2015, Keith D. Cooper & Linda Torczon, all rights reserved.

Students enrolled in Comp 512 at Rice University have explicit permission to make copies of these materials for their personal use.

Faculty from other educational institutions may use these materials for nonprofit educational purposes, provided this copyright notice is preserved

Citation numbers refer to entries in the Eac2e bibliography.

Data-flow Analysis



Definition

Data-flow analysis (DFA) is a collection of techniques for compile-time reasoning about the run-time flow of values

- We use the results of DFA to prove safety & identify opportunities
 - ◆ Not an end unto itself
- Almost always involves building a graph
 - ◆ Control-flow graph, call graph, or graphs derived from them
 - ◆ Sparse evaluation graphs to model flow of values *(efficiency)*
- Usually formulated as a set of *simultaneous equations*
 - ◆ Sets attached to nodes and edges
 - ◆ Often use sets with a lattice or semilattice structure

We have seen several data-flow problems.

Prior Examples



Computing LIVEOUT Sets

- Domain is the set of variable names in the procedure
- Data-flow equations define **LIVE** at the end of a block, **LIVEOUT**

Initialization:

$$\text{LIVEOUT}(n) = \emptyset, \forall n$$

Fixed-point equations:

$$\text{LIVEOUT}(b) = \bigcup_{s \in \text{succs}(b)} (\text{UEVAR}(b) \cup (\text{LIVEOUT}(b) \cap \overline{\text{VARKILL}(b)}))$$

LIVE is a backward-flow problem

where

UEVAR(b) is the set of names used in b before definition in b

VARKILL(b) is the set of names defined in b



Prior Examples

Computing AVAIL Sets

- Domain is the set of expressions computed in the procedure
- Data-flow equations are more complex

Initialization:

$$\begin{aligned} \text{AVAIL}(n_0) &= \emptyset \\ \text{AVAIL}(n) &= \emptyset, \forall n \neq n_0 \end{aligned}$$

Fixed-point equations:

$$\text{AVAIL}(b) = \bigcap_{x \in \text{pred}(b)} (\text{DEEXPR}(x) \cup (\text{AVAIL}(x) \cap \overline{\text{EXPRKILL}(b)}))$$

AVAIL is a forward-flow problem

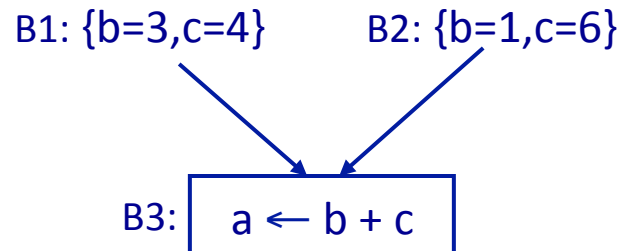
where

$\text{DEEXPR}(b)$ is the set of expressions defined in b and not subsequently killed in b
 $\text{EXPRKILL}(b)$ is the set of expressions killed in b because one or more operand is redefined in b



Prior Examples

Global constant propagation



Function “ f_3 ” models the effect of block B3

- $f_3(\{b=3,c=4\}) \Rightarrow \{a=7\}$
- $f_3(\{b=1,c=6\}) \Rightarrow \{a=7\}$
- $f_3(B1 \wedge B2) = f_3(\emptyset) \Rightarrow \{a=\perp\}$

Result depends on order of evaluation of the \wedge operation and application of f

- First condition in admissibility
 $\forall f \in F, \forall x,y \in L, f(x \wedge y) = f(x) \wedge f(y)$
- Constant propagation is not admissible
 - ◆ Kam & Ullman time bound does not hold
 - ◆ There are tight time bounds, however, based on lattice height
 - ◆ Require a variable-by-variable formulation ...

Because meet does not distribute over function application, constant propagation is not “admissible” in the Kam-Ullman sense.



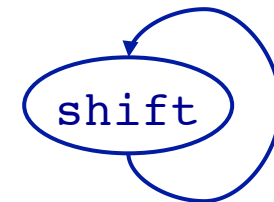
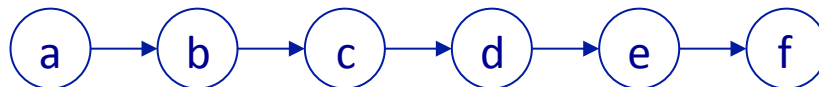
Prior Examples

Interprocedural May Modify sets

```
shift(a,b,c,d,e,f)
{
  local t;
  ...
  call shift(t,a,b,c,d,e);
  f = 1;
  ...
}
```

- Assume call-by-reference
- Compute the set of variables (in namespace of shift) that can be modified by a call to shift
- How long does it take?

- Iterations proportional to number of parameters
 - ◆ Not a function of the call graph
 - ◆ Can make example arbitrarily bad
- Proportional to length of chain of bindings...



Nothing to do with $d(G)$

Call-by-reference parameters plus recursion make the summary problems fail the Kam-Ullman “rapid” condition.



Proliferation of GDFAPs

In the late 1960's and the 1970's many data-flow problems were proposed

- GDFAP became the standard way to prove safety of a transformation
 - ◆ New transformation implied new GDFAP
 - ◆ Optimizing compilers spent a large fraction of compile time solving GDFAPs
 - ◆ Computers were relatively slow (1 – 10 MIPS) and small (16 to 32 MB)
 - ◆ Development of “frameworks” for GDFA
- Many papers showed a new GDFAP & a new transformation
 - ◆ Other applications arose for the GDFAP technology
 - ◆ See the papers on “DAVE” by Osterweil et al.

More GDFAPS: Very Busy Expressions



An expression e is very busy on exit from block b , iff e is evaluated & used along every path from b to n_f and evaluating e at the end of b would produce the same result as the next evaluation along those paths

The Plan

- Annotate each block n with a set $VERYBUSY(b)$ that contains expressions
 - ◆ Solve data-flow equations *(standard iterative solver)*
- If e is in $VERYBUSY(b)$, insert an evaluation at the end of n and eliminate the subsequent evaluations that it covers
 - ◆ If e is in $VERYBUSY(b)$ for successive blocks, want to insert it in the “right” block
 - ◆ Might be the last b (minimize register demand) or least frequently executed b (minimize dynamic number of evaluations) or ...
- This optimization aims, primarily, to reduce code space



More GDFAPS: Very Busy Expressions

Transformation: Hoisting

- e defined in every successor of b
- Evaluating e in b produces same result
- Saves code space, but shortens no path

Standard $f(x) = a \cup (b \cap c)$. \wedge is \cap .

Data-flow problem: Very Busy Expressions

$$\text{VERYBUSY}(b) = \bigcap_{s \in \text{succ}(b)} (\text{UEEXPR}(s) \cup (\text{VERYBUSY}(s) \cap \overline{\text{EXPRKILL}(s)}))$$

$$\text{VERYBUSY}(n_f) = \emptyset$$

- $\text{VERYBUSY}(b)$ contains expressions that are very busy at end of b
- $\text{UEEXPR}(b)$ is the set of expressions used before they are killed in b
- $\text{EXPRKILL}(b)$ is the set of expressions killed before they are used in b

VERYBUSY expressions is a **backward** flow problem

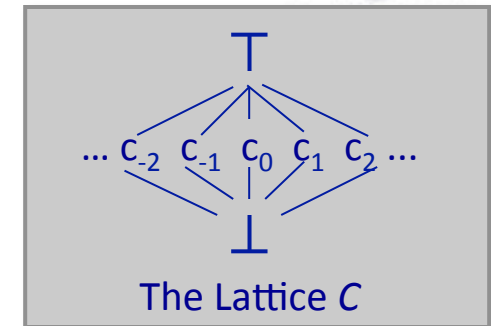
$$|\text{CONSTANTS}| = |\text{variables}|$$



More GDFAPS: Constant Propagation (Classic formulation)

Transformation: Global Constant Folding

- Along every path to p , v has same known value
- Specialize computation at p based on v 's value



Data-flow problem: Constant Propagation

Domain is the set of pairs $\langle v_i, c_i \rangle$ where v_i is a variable and $c_i \in C$

$$\text{CONSTANTS}(b) = \bigwedge_{p \in \text{preds}(b)} f_p(\text{CONSTANTS}(p))$$

- \bigwedge performs a pairwise meet on two sets of pairs
- $f_p(x)$ is a block specific function that models the effects of block p on the $\langle v_i, c_i \rangle$ pairs in x

Form of f is quite different than in the other GDFAPS that we have seen

Constant propagation is a **forward** flow problem



More GDFAPs: Constant Folding

Meet operation is more complex than we have already seen

- $c_1 \wedge c_2 = c_1$ if $c_1 = c_2$, else \perp (bottom & top as expected)

What about f_p ?

- If p has one statement then

$x \leftarrow y$ with $CONSTANTS(p) = \{\dots \langle x, l_1 \rangle, \dots \langle y, l_2 \rangle \dots\}$

then $f_p(CONSTANTS(p)) = CONSTANTS(p) - \langle x, l_1 \rangle + \langle x, l_2 \rangle$

$x \leftarrow y \text{ op } z$ with $CONSTANTS(p) = \{\dots \langle x, l_1 \rangle, \dots \langle y, l_2 \rangle \dots, \dots \langle z, l_3 \rangle \dots\}$

then $f_p(CONSTANTS(p)) = CONSTANTS(p) - \langle x, l_1 \rangle + \langle x, l_2 \text{ op } l_3 \rangle$

- If p has n statements then

$f_p(CONSTANTS(p)) = f_n(f_{n-1}(f_{n-2}(\dots f_2(f_1(CONSTANTS(p)))\dots)))$

where f_i is the function generated by the i^{th} statement in p

f_p does not fit into the mold of the functions in our Kam-Ullman rapid frameworks.

Building a Control-Flow Graph



The first step in almost any data-flow analysis is building a CFG

```
// find all the leaders, assume first op
// & block are numbered zero
next ← 0
leader[next] ← 0
for i ← 0 to n
  if op[i] is a jump
    then leader[next++] ← target(i)
  if op[i] is a branch then
    leader[next++] ← taken(i)
    leader[next++] ← not_taken(i)
// build all the blocks
for i ← 0 to next - 1
  j ← leader[i] + 1
  while j ≤ n and j ∉ leader
    j ← j + 1
  last[i] ← j - 1
```

If **target**, **taken**, or **not_taken** are ambiguous, then we must include all labeled ops as leaders.

Sources of ambiguous targets:

- ◆ Fall-through branch path
- ◆ Jump to a register

No Ambiguity In ILOC:

All branches in **ILOC** have two explicit targets. Branches and jumps target a label rather than a register.

In the original compiler, jump to register was followed with an advisory list of labels generated when the **ILOC** was generated.



Building a Control-Flow Graph

The first step in almost any data-flow analysis is building a CFG

```
// find all the leaders, assume first op
// & block are numbered zero
next ← 0
leader[next] ← 0
for i ← 0 to n
  if op[i] is a jump
    then leader[next++] ← target(i)
  if op[i] is a branch then
    leader[next++] ← taken(i)
    leader[next++] ← not_taken(i)
// build all the blocks
for i ← 0 to next - 1
  j ← leader[i] + 1
  while j ≤ n and j ∉ leader
    j ← j + 1
  last[i] ← j - 1
```

EXAMPLE

```
0      a ← 4
1      t1 ← a * 4
2  L1:  t2 ← t1 / c
3      if t2 < w then goto L2
4      m ← t1 * k
5      t3 ← m + i
6  L2:  h ← i
7      m ← t3 - h
8      if t3 ≥ 0 then goto L3
9      goto L1
10 L3:  halt
```

LEADER	0						
LAST							



Building a Control-Flow Graph

The first step in almost any data-flow analysis is building a CFG

```
// find all the leaders, assume first op
// & block are numbered zero
next ← 0
leader[next] ← 0
for i ← 0 to n
  if op[i] is a jump
    then leader[next++] ← target(i)
  if op[i] is a branch then
    leader[next++] ← taken(i)
    leader[next++] ← not_taken(i)
// build all the blocks
for i ← 0 to next - 1
  j ← leader[i] + 1
  while j ≤ n and j ∉ leader
    j ← j + 1
  last[i] ← j - 1
```

EXAMPLE

```
0      a ← 4
1      t1 ← a * 4
2  L1:  t2 ← t1 / c
3      if t2 < w then goto L2
4      m ← t1 * k
5      t3 ← m + i
6  L2:  h ← i
7      m ← t3 - h
8      if t3 ≥ 0 then goto L3
9      goto L1
10 L3:  halt
```

LEADER	0	6	4	9	10	2
LAST						

Building a Control-Flow Graph



The first step in almost any data-flow analysis is building a CFG

```
// find all the leaders, assume first op
// & block are numbered zero
next ← 0
leader[next] ← 0
for i ← 0 to n
  if op[i] is a jump
    then leader[next++] ← target(i)
  if op[i] is a branch then
    leader[next++] ← taken(i)
    leader[next++] ← not_taken(i)
// build all the blocks
for i ← 0 to next - 1
  j ← leader[i] + 1
  while j ≤ n and j ∉ leader
    j ← j + 1
  last[i] ← j - 1
```

EXAMPLE

```
0      a ← 4
1      t1 ← a * 4
2  L1:  t2 ← t1 / c
3      if t2 < w then goto L2
4      m ← t1 * k
5      t3 ← m + i
6  L2:  h ← i
7      m ← t3 - h
8      if t3 ≥ 0 then goto L3
9      goto L1
10 L3:  halt
```

LEADER	0	6	4	9	10	2
LAST	1	8	5	9	10	3

Building a Control-Flow Graph



The first step in almost any data-flow analysis is building a CFG

```
// find all the leaders, assume first op
// & block are numbered zero
next ← 0
leader[next] ← 0
for i ← 0 to n
  if op[i] is a jump
    then leader[next++] ← target(i)
  if op[i] is a branch then
    leader[next++] ← taken(i)
    leader[next++] ← not_taken(i)
// build all the blocks
for i ← 0 to next - 1
  j ← leader[i] + 1
  while j ≤ n and j ∉ leader
    j ← j + 1
  last[i] ← j - 1
```

EXAMPLE

```
0  a ← 4
1  t1 ← a * 4
2  L1: t2 ← t1 / c
3     if t2 < w then goto L2
4  m ← t1 * k
5  t3 ← m + i
6  L2: h ← i
7     m ← t3 - h
8     if t3 ≥ 0 then goto L3
9  goto L1
10 L3: halt
```

LEADER	0	6	4	9	10	2
LAST	1	8	5	9	10	3

Dominators



Definitions

In a flow graph, x dominates y if and only if every path from the entry of the control-flow graph to the node for y includes x

- By definition, x dominates x
- We associate a **DOM** set with each node
- $|\mathbf{DOM}(x)| \geq 1$

Immediate dominator

- For any node x , there must be a y in **DOM**(x) closest to x
 - ◆ Unless $x = n_0$, $x \neq \mathbf{IDOM}(x)$
- We call this y the immediate dominator of x
- As a matter of notation, we write this as **IDOM**(x)

*Original idea: R.T. Prosser. "Applications of Boolean matrices to the analysis of flow diagrams," *Proceedings of the Eastern Joint Computer Conference*, Spartan Books, New York, pages 133-138, 1959.*



Dominators

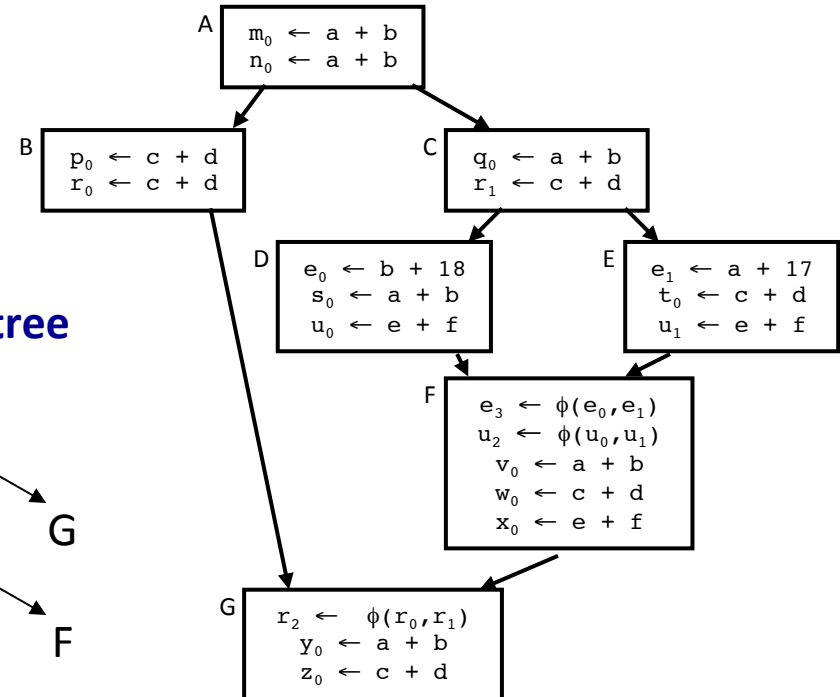
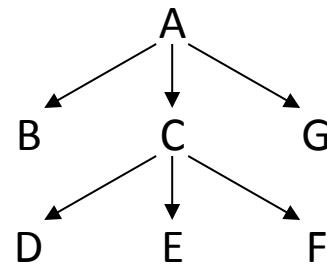
Dominators have many uses in analysis & transformation

- Finding loops
- Building SSA form
- Making code motion decisions

Dominator sets

Block	DOM	IDOM
A	A	—
B	A,B	A
C	A,C	A
D	A,C,D	C
E	A,C,E	C
F	A,C,F	C
G	A,G	A

Dominator tree





Computing Dominators

Critical first step in SSA construction and in DVNT

- A node n dominates m iff n is on every path from n_0 to m
 - ◆ Every node dominates itself
 - ◆ n 's immediate dominator is its closest dominator, $\text{IDOM}(n)^\dagger$

$$\text{DOM}(n_0) = \{n_0\}$$

$$\text{DOM}(n) = \{n\} \cup \left(\bigcap_{p \in \text{preds}(n)} \text{DOM}(p) \right)$$

Initially, $\text{DOM}(n) = N$,
 $\forall n \neq n_0$.
Can do better.

Computing DOM

- These simultaneous set equations the data-flow problem
 - ◆ The simplest equations we have seen
 - ◆ Transfer function is the identity function
- Equations have a unique fixed point solution
- An iterative fixed-point algorithm will solve them quickly



Round-robin Iterative Algorithm for DOM

```
DOM( $n_0$ )  $\leftarrow$   $n_0$ 
for  $x \leftarrow n_1$  to  $n_n$ 
    DOM( $x$ )  $\leftarrow$  { all nodes in graph }
change  $\leftarrow$  true
while (change)
    change  $\leftarrow$  false
    for  $x \leftarrow n_0$  to  $n_n$ 
        TEMP  $\leftarrow$  {  $x$  }  $\cup$  ( $\bigcap_{y \in \text{pred}(x)}$  DOM( $y$ ))
        if DOM( $x$ )  $\neq$  TEMP then
            change  $\leftarrow$  true
            DOM( $x$ )  $\leftarrow$  TEMP
```

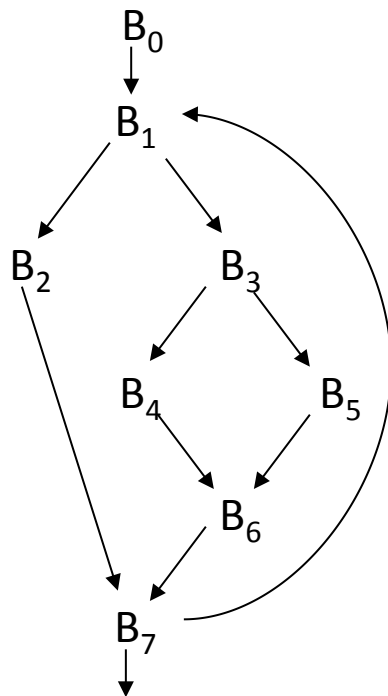
Termination

- Makes sweeps over the nodes
- Halts when some sweep produces no change

DOM Example



Progress of iterative solution for DOM



Flow Graph

Iteration	DOM(<i>n</i>)								
	0	1	2	3	4	5	6	7	
0	0	<i>N</i>	<i>N</i>	<i>N</i>	<i>N</i>	<i>N</i>	<i>N</i>	<i>N</i>	
1	0	0,1	0,1,2	0,1,3	0,1,3,4	0,1,3,5	0,1,3,6	0,1,7	
2	0	0,1	0,1,2	0,1,3	0,1,3,4	0,1,3,5	0,1,3,6	0,1,7	

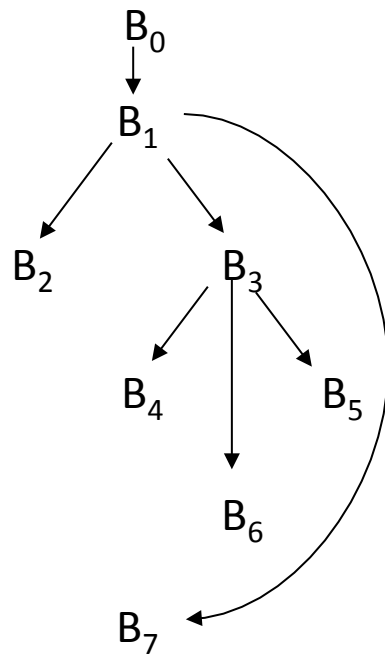
Example

If we have time, the last three slides show how to use **DOM** to improve **SVN**



Progress of iterative solution for DOM

Iteration	DOM(n)							
	0	1	2	3	4	5	6	7
0	0	N	N	N	N	N	N	N
1	0	0,1	0,1,2	0,1,3	0,1,3,4	0,1,3,5	0,1,3,6	0,1,7
2	0	0,1	0,1,2	0,1,3	0,1,3,4	0,1,3,5	0,1,3,6	0,1,7



Dominance Tree

Results of iterative solution for DOM

	0	1	2	3	4	5	6	7
DOM	0	0,1	0,1,2	0,1,3	0,1,3,4	0,1,3,5	0,1,3,6	0,1,7
IDOM	0	0	1	1	3	3	3	1

There are asymptotically faster algorithms.

With the right data structures, the iterative algorithm can be made extremely fast (competitive out to 10,000 or 20,000 nodes)

See Cooper, Harvey, & Kennedy [100], or § 9.5.2 in EaC2e.

Proliferation of GDFAPs



In the late 1960's and the 1970's many data-flow problems were proposed

- GDFAP became the standard way to prove safety of a transformation
 - ◆ New transformation implied new GDFAP
 - ◆ Optimizing compilers spent a large fraction of compile time solving GDFAPs
 - ◆ Computers were relatively slow (1 – 10 MIPS) and small (16 to 32 MB)
 - ◆ Development of “frameworks” for GDFA
- As transformations proliferated, need for a new paradigm emerged
 - ◆ One GDFAP that could be used for multiple transformations
 - ◆ Simplify the implementation
 - ◆ Reduce the time spent in analysis
 - ◆ The result was the development of information chains

In truth, the story is not that simple. Information chains did not arise overnight in response to an excessive number of GDFAPS; however, by the late 1980's they were being used to replace individual GDFAPs.

Information Chains



A tuple that connects 2 data-flow events is a *chain*

- Chains express data-flow relationships directly
- Chains provide a graphical representation
- Chains jump across unrelated code, simplifying search

event \cong definition
or use

We can build chains efficiently

Four interesting types of chain

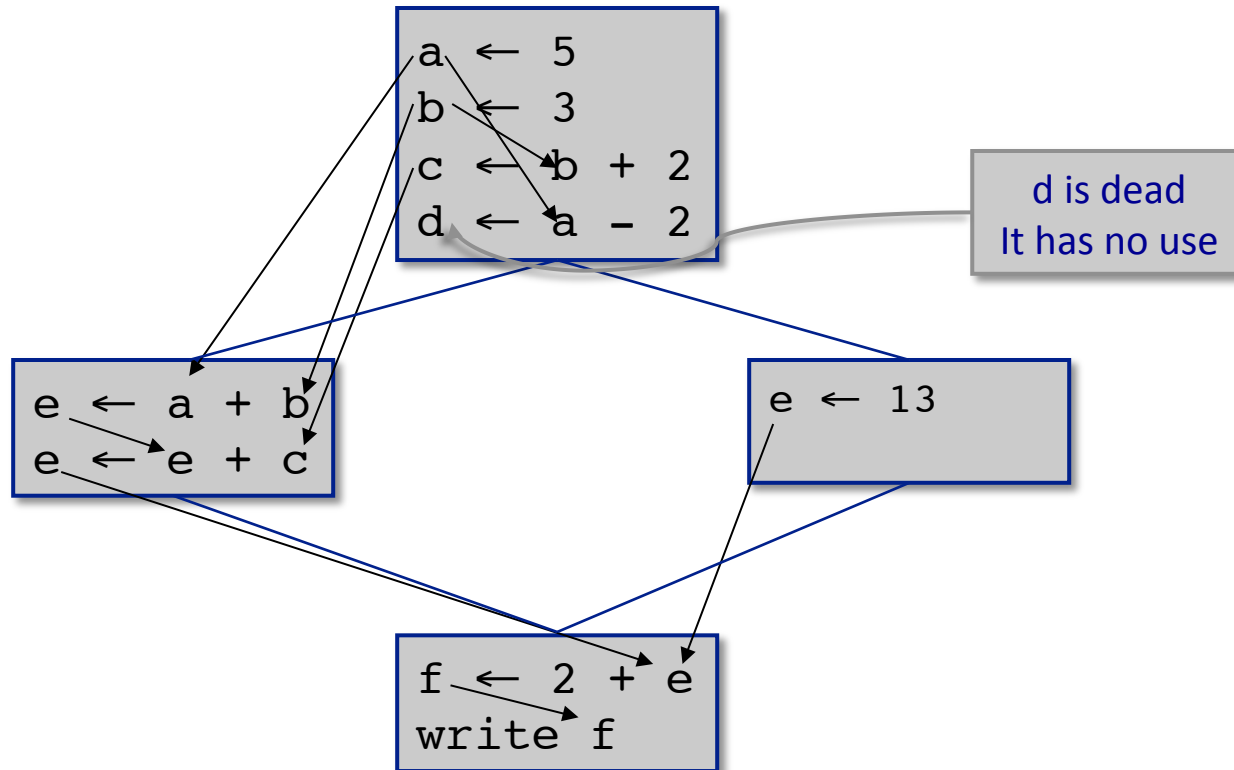
Source	Sink	Dependence Type
DEF	USE	<i>true, flow</i>
USE	DEF	<i>anti</i>
DEF	DEF	<i>output</i>
USE	USE	<i>input</i>

DEF-USE chains are the most common

Information Chains



Example



DEF-USE Chains

DEF-USE Chains form a sparse evaluation graph that we can use in many transformations.

For example, a DEF with no reachable use is *dead*.

Notation



Assume that, \forall operation i and each variable v ,

- $\text{DEFS}(v,i)$ is the set of operations that may have defined v most recently before i , along some path in the CFG
- $\text{USES}(v,i)$ is the set of operations that may use the value of v computed at i , along some path in the CFG

$$x \in \text{DEFS}(A,y) \Leftrightarrow y \in \text{USES}(A,x)$$

To construct DEF-USE chains, we solve *reaching definitions* (YAGDFAP)

- A definition d of some variable v reaches an operation i if and only if i reads v and there is a v -clear path from d to i
 - ♦ v -clear \Rightarrow no definition of v on the path
- Prior definition of v in same block $\Rightarrow |\text{DEFS}(v,i)| = 1$
- No prior definition $\Rightarrow |\text{DEFS}(v,i)| \geq 1$

Domain is |definitions|, same as number of operations



Reaching Definitions

The equations

$$\text{REACHES}(b) = \emptyset, \forall n \in N$$

$$\text{REACHES}(b) = \bigcup_{p \in \text{preds}(b)} (\text{DEDEF}(p) \cup (\text{REACHES}(p) \cap \overline{\text{DEFKILL}(p)}))$$

Form of f is same as in LIVE

- $\text{REACHES}(b)$ is the set of definitions that reach block b
- $\text{DEDEF}(b)$ is the set of definitions in n that reach the end of b
- $\text{DEFKILL}(b)$ is the set of defs obscured by a new def in b

Computing $\text{REACHES}(b)$

- Use any data-flow method (rapid framework)
- Zadeck shows a simple linear-time algorithm

F.K. Zadeck, "Incremental data-flow analysis in a structured program editor," *Proceedings of the SIGPLAN 84 Conf. on Compiler Construction*, June, 1984, pages 132-143.



Building DEFS Sets

The Plan

1. Find basic blocks & build the CFG
2. \forall block b , compute **REACHES**(b)
3. \forall block b , \forall operation i , \forall referenced name v ,
Set **DEFS**(v,i) according to the earlier rule
If there is a prior definition, d , of v in b
$$\mathbf{DEFS}(v,i) \leftarrow d$$
Otherwise
$$\mathbf{DEFS}(v,i) \leftarrow \{d \mid d \text{ defines } v \ \& \ d \in \mathbf{REACHES}(b)\}$$

To build USES

- Invert **DEFS**, or
- Solve *reachable uses*, and use the analogous construction

Building DEF-USE Chains



Miscellany

- Domain of **REACHES** is the set of definitions $(\propto |operations|)$
- Domain of **DEFS & USES** is also definitions
- Need a compact representation of **DEFS & USES**

Detecting Anomalies

- $DEFS(v,i) = \emptyset$ implies use of a never initialized variable
- Add a definition for each v to n_o to detect larger set of anomalies
 - ◆ If initial def $\in DEFS(v,i)$ then \exists a path to i with no initialization

NEXT LECTURE: using information chains & moving into SSA

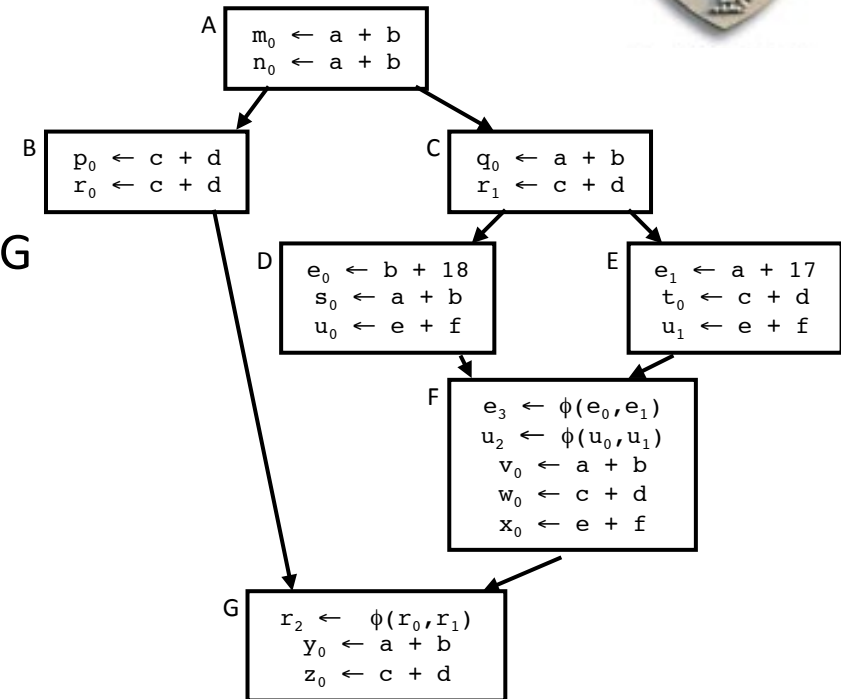
Back to Redundancy Elimination

Dominators Can Improve Superlocal Value Numbering



SVN did not help with blocks F or G

- Multiple predecessors
- Must decide what facts hold in F and in G
 - ◆ For G, combine B & F?
 - ◆ Merging state is expensive
 - ◆ Fall back on what's known
- Can use table from **IDOM**(x) to start x
 - ◆ Use C for F and A for G
 - ◆ Imposes a **DOM**-based application order



Leads to Dominator VN Technique (**DVNT**)

Dominator Value Numbering



The DVNT Algorithm

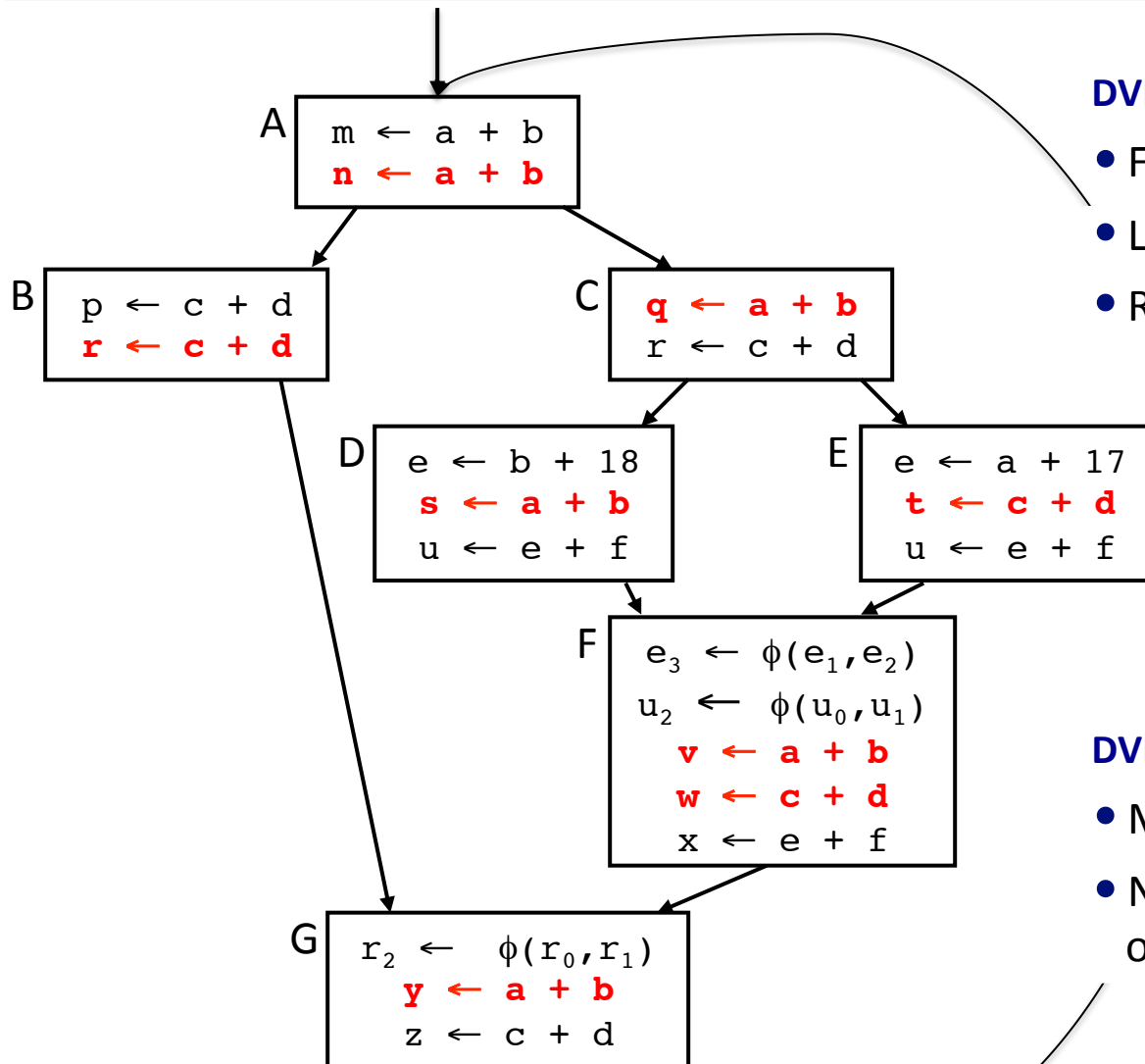
- Use superlocal algorithm on extended basic blocks
 - ◆ Retain use of scoped hash tables
 - ◆ Need to use the **SSA** name space to avoid bookkeeping headaches
- Start each node with table from its **IDOM**
 - ◆ **DVNT** generalizes the superlocal algorithm
- No values flow along back edges (*i.e.*, around loops)
- Constant folding, algebraic identities as before

Larger scope leads to (*potentially*) better results

- ◆ **LVN** + **SVN** + good start for **EBBs** missed by **SVN**



Dominator Value Numbering



DVNT advantages

- Find more redundancy
- Little additional cost
- Retains *online* character

DVNT shortcomings

- Misses some opportunities
- No loop-carried redundancies or constants