

Compiler Techniques

Lecture 11: Dataflow Analysis (Forward)

Tianwei Zhang

Outline

- ▶ **Available Expression Analysis**
- ▶ **Reaching Definition Analysis**
- ▶ **Summary of Dataflow Analysis**

Outline

- ▶ **Available Expression Analysis**
- ▶ Reaching Definition Analysis
- ▶ Summary of Dataflow Analysis

Common Subexpression Elimination

- Identify expressions that have been computed before, and reuse previously computed result

The value of $z*2$ has been computed here:

So we can replace the assignment with $z = t$

```
int r, t, x, y, z;  
x:= @param0; y:= @param1;  
z = 1;  
r = x;  
11: if z==y goto 13;  
    t = z*2;  
    if t>y goto 12;  
    r = r*r;  
    z = z*2;  
    goto 11;  
12:  
    r = r*x;  
    z = z+1;  
    goto 11;  
13: return r;
```

Correctness Conditions

- ▶ Two conditions need to be held in order to eliminate a common subexpression:
 1. The expression must have been computed previously on *every* possible execution path, not just on one
 2. None of the variables involved in computing the expression may have been updated in the meantime



```
if z > 0 goto l
x = y + z
l: r = y + z
```



```
x = y + z
y = y + 1
r = y + z
```

Available Expressions

- ▶ These conditions can be formularized as the concept of an available expression in terms of the CGF
- ▶ An expression e is available at a point p if
 1. e is computed on every path from ENTRY to p , and
 2. no variable in e is overwritten between the computation of e and p
- ▶ Flow sets:
 - ▶ $\text{in}_A(n)$: the set of available expressions before n
 - ▶ $\text{out}_A(n)$: the set of available expression after n
- ▶ Goal of available expression analysis: compute $\text{in}_A(n)$ and $\text{out}_A(n)$ for every CFG node n
 - ▶ The set of possible expressions is infinite, but we only need to consider expressions that actually occur in the method (which is a finite set)
 - ▶ For simplicity, we only consider expressions computed from local variables using arithmetic operators

Transfer Functions

- ▶ If we already know $\text{in}_A(n)$, it is easy to compute $\text{out}_A(n)$:
 - ▶ An expression e is available after n if
 - (1) Node n computes e , or
 - (2) e is available before n and n does not write to any variable in e
- ▶ More denotations:
 - ▶ $\text{vars}(e)$: the set of (local) variables in an expression e
 - ▶ $\text{comp}(n)$: the set of expressions computed by n
 - ▶ $\text{def}(n)$: the set of variables node n writes to
- ▶ **Transfer function** for $\text{out}_A(n)$:
$$\text{out}_A(n) = \text{in}_A(n) \setminus \{e \mid \text{vars}(e) \cap \text{def}(n) \neq \emptyset\} \cup \text{comp}(n)$$
- ▶ This is **forward flow analysis**.

Transfer Functions

- ▶ How do we get $\text{in}_A(n)$? There are two cases:
 - ▶ Node n is the Entry node. Then no expression is available before a method:

$$\text{in}_A(\text{Entry}) = \emptyset$$

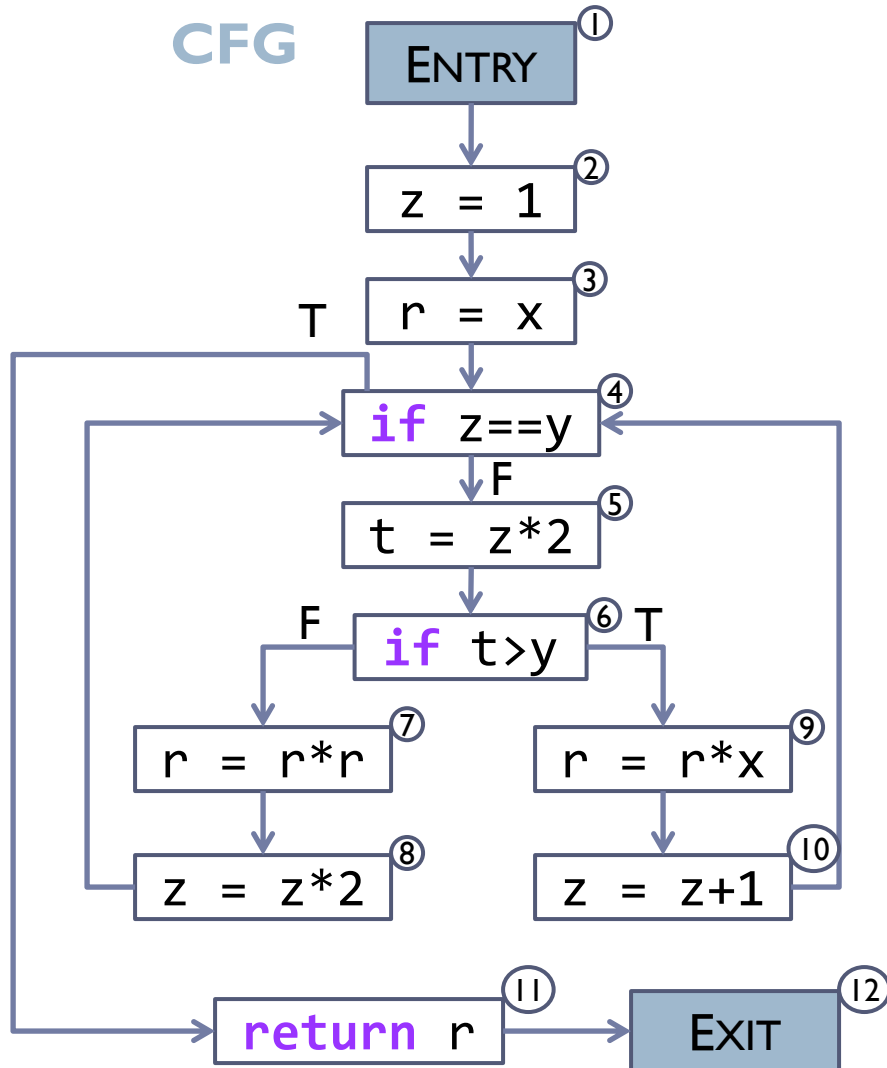
- ▶ Node n has at least one predecessor node.
 - ▶ $\text{pred}(n)$: the set of predecessor nodes of n .
 - ▶ An expression e is available before n if it is available after every predecessor of n

$$\text{in}_A(n) = \cap \{ \text{out}_A(m) \mid m \in \text{pred}(n) \}$$

- ▶ This is **must analysis**, as intersection is used to combine results from predecessor nodes.

Available Expressions Example

CFG



Transfer Functions

$$\text{out}_A(1) = \text{in}_A(1)$$

$$\text{out}_A(2) = \text{in}_A(2) \setminus \{ z*2, z+1 \}$$

$$\text{out}_A(3) = \text{in}_A(3) \setminus \{ r*r, r*x \}$$

$$\text{out}_A(4) = \text{in}_A(4)$$

$$\text{out}_A(5) = \text{in}_A(5) \cup \{ z*2 \}$$

$$\text{out}_A(6) = \text{in}_A(6)$$

$$\text{out}_A(7) = \text{in}_A(7) \setminus \{ r*r, r*x \} \cup \{ \}$$

$$\text{out}_A(8) = \text{in}_A(8) \setminus \{ z*2, z+1 \} \cup \{ \}$$

$$\text{out}_A(9) = \text{in}_A(9) \setminus \{ r*r, r*x \} \cup \{ \}$$

$$\text{out}_A(10) = \text{in}_A(10) \setminus \{ z*2, z+1 \} \cup \{ \}$$

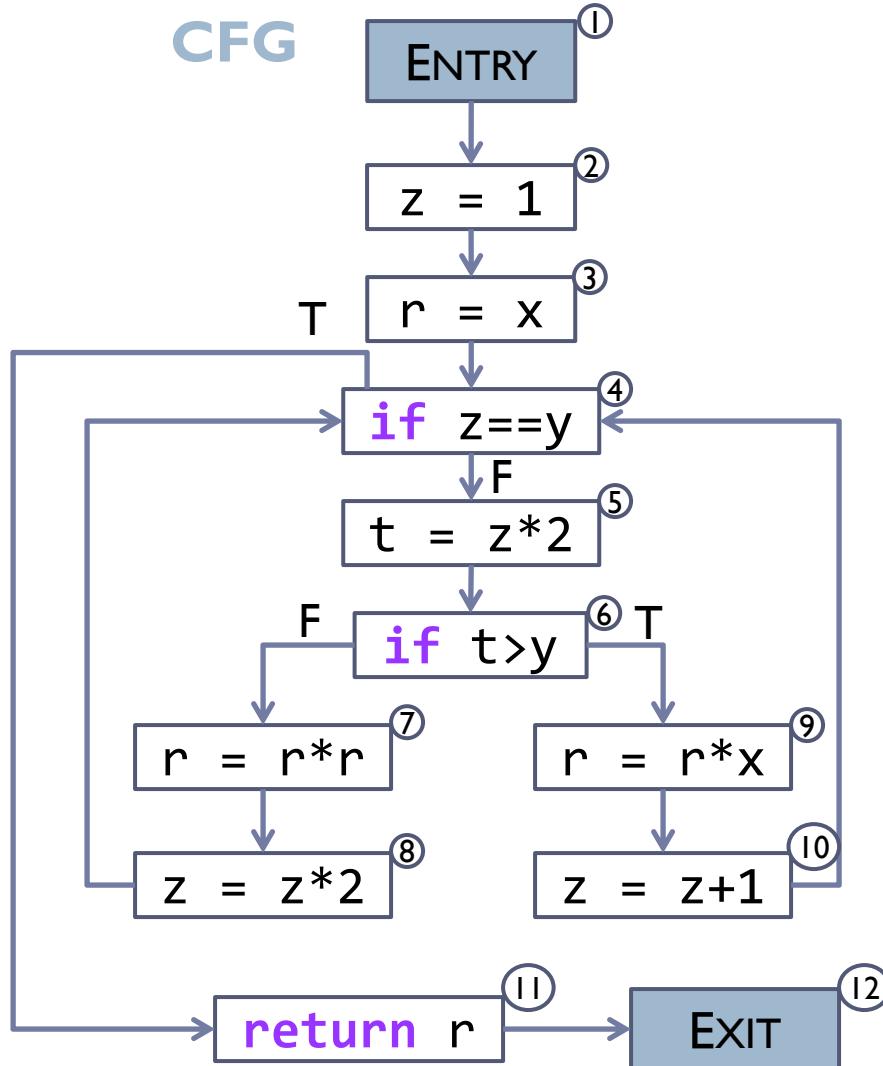
$$\text{out}_A(11) = \text{in}_A(11)$$

$$\text{out}_A(12) = \text{in}_A(12)$$

$r*r$ should not be included in $\text{comp}(7)$ as the value of r has been changed immediately by the statement in node 7. The same can be said for node 8, 9, and 10.

Available Expressions Example

CFG



Transfer Functions

$$\text{in}_A(1) = \emptyset$$

$$\text{in}_A(2) = \text{out}_A(1)$$

$$\text{in}_A(3) = \text{out}_A(2)$$

$$\text{in}_A(4) = \text{out}_A(3) \cap \text{out}_A(8) \cap \text{out}_A(10)$$

$$\text{in}_A(5) = \text{out}_A(4)$$

$$\text{in}_A(6) = \text{out}_A(5)$$

$$\text{in}_A(7) = \text{out}_A(6)$$

$$\text{in}_A(8) = \text{out}_A(7)$$

$$\text{in}_A(9) = \text{out}_A(6)$$

$$\text{in}_A(10) = \text{out}_A(9)$$

$$\text{in}_A(11) = \text{out}_A(4)$$

$$\text{in}_A(12) = \text{out}_A(11)$$



Method 1: Iterative Solution

- ▶ The equation system can be solved by iteration:
 1. Set $\text{out}_A(n) = \text{in}_A(n) = U$ for all nodes n where U is the set of *all* expressions in the method, since we use intersection to compute $\text{in}_A(n)$. Exception: $\text{out}_A(\text{ENTRY}) = \text{in}_A(\text{ENTRY}) = \emptyset$
 2. For every node n , recompute $\text{in}_A(n)$ based on the values we have computed in the previous iteration and then recompute $\text{out}_A(n)$ from $\text{in}_A(n)$
 3. Keep doing step 2 until the values do not change any further

Simplifying Transfer Equations

- ▶ To make it easier to solve the equations, we substitute away $\text{in}_A(n)$ equations, so we only have to solve for $\text{out}_A(n)$

$$\begin{aligned}\text{out}_A(1) &= \text{in}_A(1) \\ \text{out}_A(2) &= \text{in}_A(2) \setminus \{z*2, z+1\} \\ \text{out}_A(3) &= \text{in}_A(3) \setminus \{r*r, r*x\} \\ \text{out}_A(4) &= \text{in}_A(4) \\ \text{out}_A(5) &= \text{in}_A(5) \cup \{z*2\} \\ \text{out}_A(6) &= \text{in}_A(6) \\ \text{out}_A(7) &= \text{in}_A(7) \setminus \{r*r, r*x\} \cup \{\} \\ \text{out}_A(8) &= \text{in}_A(8) \setminus \{z*2, z+1\} \cup \{\} \\ \text{out}_A(9) &= \text{in}_A(9) \setminus \{r*r, r*x\} \cup \{\} \\ \text{out}_A(10) &= \text{in}_A(10) \setminus \{z*2, z+1\} \cup \{\} \\ \text{out}_A(11) &= \text{in}_A(11) \\ \text{out}_A(12) &= \text{in}_A(12)\end{aligned}$$
$$\begin{aligned}\text{in}_A(1) &= \emptyset \\ \text{in}_A(2) &= \text{out}_A(1) \\ \text{in}_A(3) &= \text{out}_A(2) \\ \text{in}_A(4) &= \text{out}_A(3) \cap \text{out}_A(8) \cap \text{out}_A(10) \\ \text{in}_A(5) &= \text{out}_A(4) \\ \text{in}_A(6) &= \text{out}_A(5) \\ \text{in}_A(7) &= \text{out}_A(6) \\ \text{in}_A(8) &= \text{out}_A(7) \\ \text{in}_A(9) &= \text{out}_A(6) \\ \text{in}_A(10) &= \text{out}_A(9) \\ \text{in}_A(11) &= \text{out}_A(4) \\ \text{in}_A(12) &= \text{out}_A(11)\end{aligned}$$

$$\begin{aligned}\text{out}_A(1) &= \emptyset \\ \text{out}_A(2) &= \text{out}_A(1) \setminus \{z*2, z+1\} \\ \text{out}_A(3) &= \text{out}_A(2) \setminus \{r*r, r*x\} \\ \text{out}_A(4) &= \text{out}_A(3) \cap \text{out}_A(8) \cap \text{out}_A(10) \\ \text{out}_A(5) &= \text{out}_A(4) \cup \{z*2\} \\ \text{out}_A(6) &= \text{out}_A(5) \\ \text{out}_A(7) &= \text{out}_A(6) \setminus \{r*r, r*x\} \\ \text{out}_A(8) &= \text{out}_A(7) \setminus \{z*2, z+1\} \\ \text{out}_A(9) &= \text{out}_A(6) \setminus \{r*r, r*x\} \\ \text{out}_A(10) &= \text{out}_A(9) \setminus \{z*2, z+1\} \\ \text{out}_A(11) &= \text{out}_A(4) \\ \text{out}_A(12) &= \text{out}_A(11)\end{aligned}$$

Method 2: Worklist Algorithm (Forward)

- ▶ $\text{out}_A(n)$ can only change if $\text{in}_A(n)$ changes, and $\text{in}_A(n)$ changes only if $\text{out}_A(m)$ changes for some predecessor m of n
- ▶ We can avoid unnecessary recomputation by keeping a *worklist* of nodes for which $\text{in}_A(n)$ has changed:

```
worklist = [ all nodes ]
while worklist != empty do
     $m = \text{removeFirst}(\text{worklist})$ 
    recompute  $\text{out}_A(m)$ 
    if  $\text{out}_A(m)$  has changed then
        for each successor  $n$  of  $m$ 
            compute  $\text{in}_A(n)$ 
            if  $\text{in}_A(n)$  has changed then
                put  $n$  into worklist (if not already in worklist)
```

Worklist Algorithm Example (Forward)

- ▶ Available expression analysis using the worklist algorithm.
 - ▶ We initialize the worklist to contain all nodes in the CFG to ensure that each node is evaluated at least once
 - ▶ These nodes can be in an arbitrary order. For forward flow analysis, we sort these nodes as the normal order in the worklist, to reduce number of iterations.
 - ▶ For all nodes, $\text{out}_A(n)$ and $\text{in}_A(n)$ are initialized to $U = \{z*2, r*r, r*x, z+1\}$, the set of all expressions in the method, except $\text{in}_A(l) = \emptyset$.

Example

Worklist

1, ..., 12

out_A(m) &

out_A(1) = \emptyset

2, ..., 12

out_A(2) = \emptyset

3, ..., 12

out_A(3) = \emptyset

4, ..., 12

out_A(4) = \emptyset

5, ..., 12

out_A(5) = {z*2}

6, ..., 12

out_A(6) = {z*2}

7, ..., 12

out_A(7) = {z*2}

8, ..., 12

out_A(8) = \emptyset

9, ..., 12

out_A(9) = {z*2}

10, ..., 12

out_A(10) = \emptyset

11, 12

out_A(11) = \emptyset

12

out_A(12) = \emptyset

in_A(n)

in_A(2) = \emptyset

in_A(3) = \emptyset

in_A(4) = \emptyset

in_A(5) = \emptyset

in_A(11) = \emptyset

in_A(6) = {z*2}

in_A(7) = {z*2}

in_A(9) = {z*2}

in_A(8) = {z*2}

in_A(4) = \emptyset

(no change)

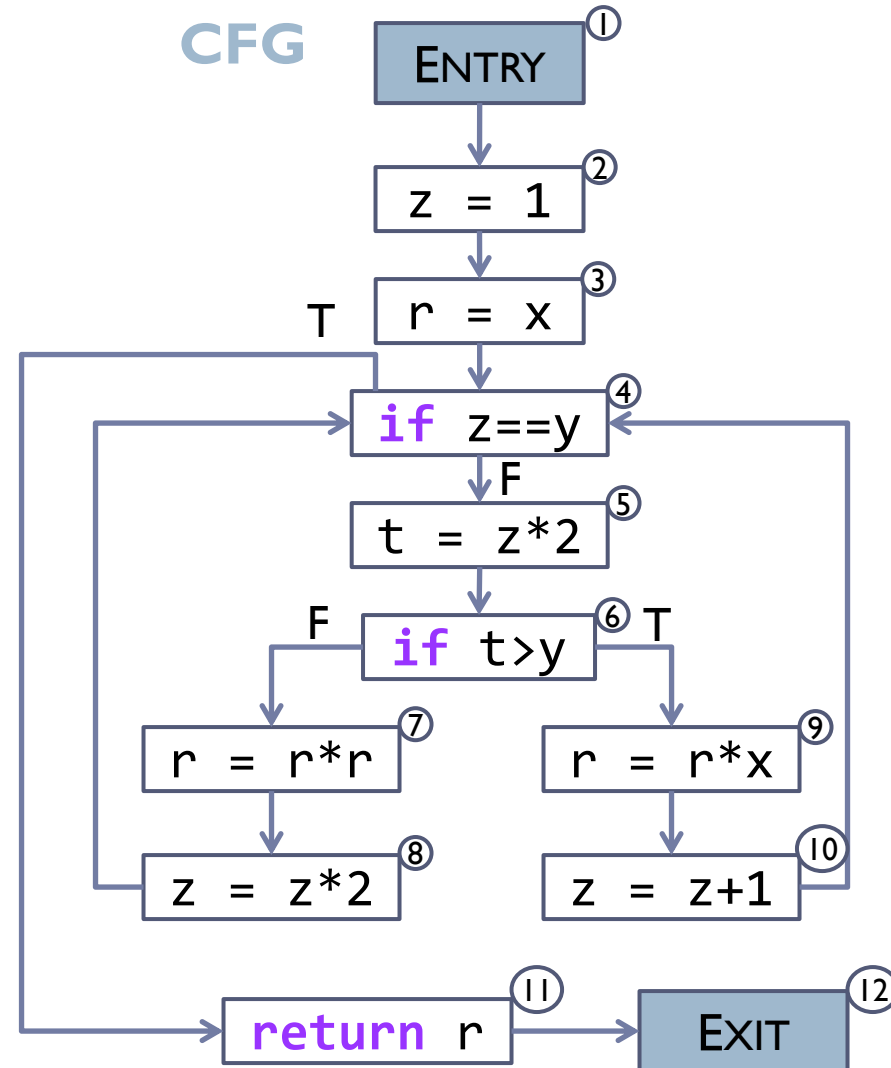
in_A(10) = {z*2}

in_A(4) = \emptyset

(no change)

in_A(12) = \emptyset

CFG

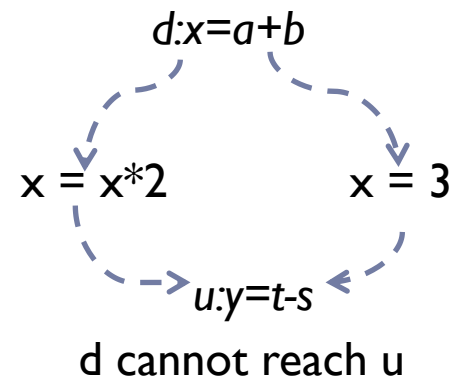
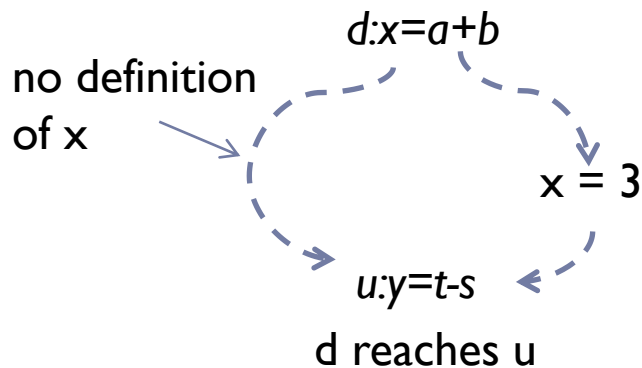


Outline

- ▶ Available Expression Analysis
- ▶ **Reaching Definition Analysis**
- ▶ Summary of Dataflow Analysis

Reaching Definition

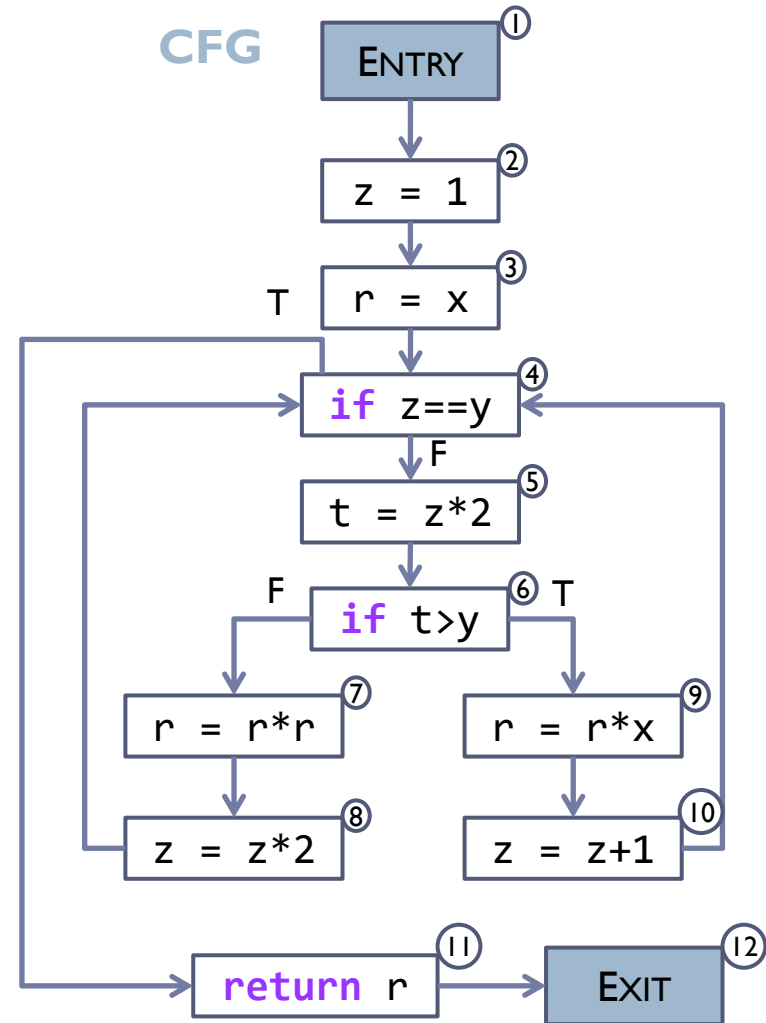
- ▶ Every assignment is a definition
- ▶ Definition of a variable x at a program point d reaches point u if there exists a path from d to u such that x is not re-defined on this path.
- ▶ Problem statement:
 - ▶ For each point in the program, determine if each definition in the program reaches the point



Reaching Definition

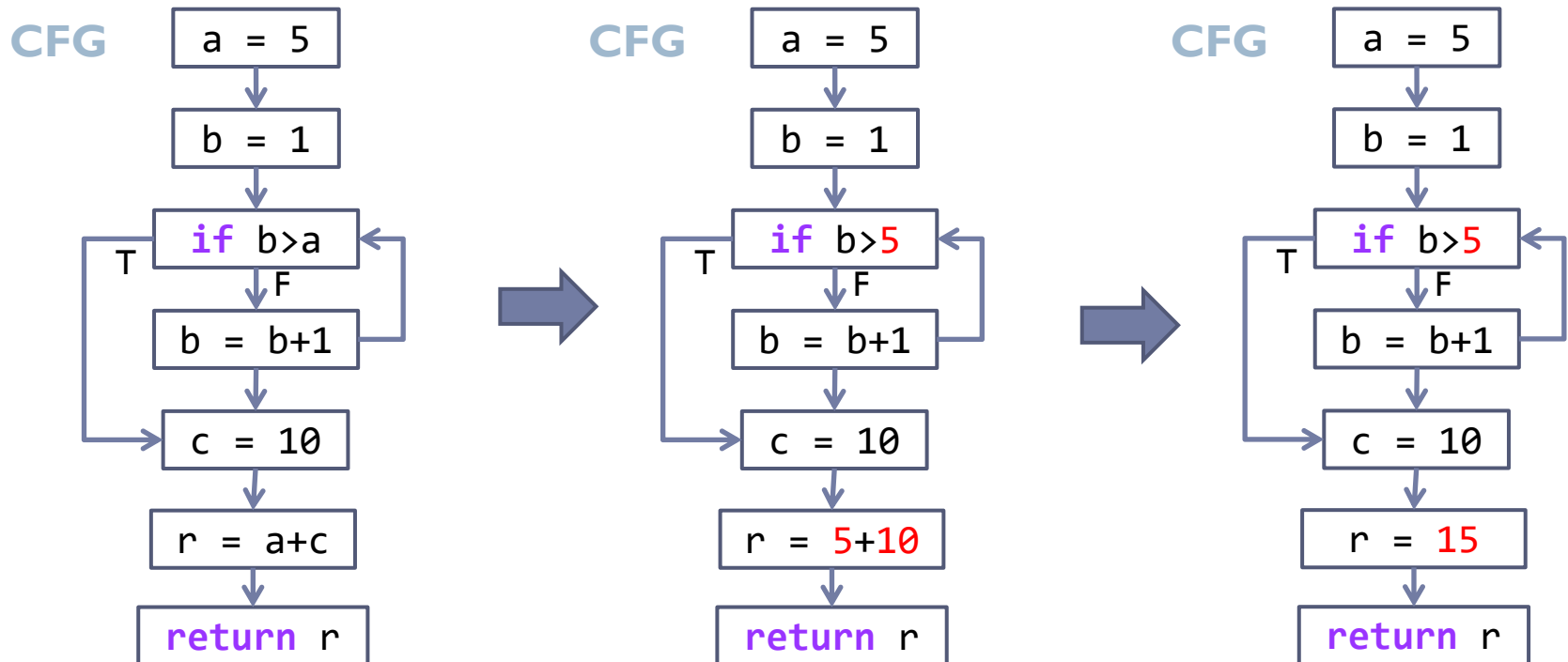
- ▶ d_2 can reach: {3,4,5,6,7,9,11}
- ▶ d_3 can reach: {4,5,6,11}
- ▶ d_5 can reach: {6,7,8,9,10,4,11}
- ▶ d_7 can reach: {8,4,5,6,11}
- ▶ d_8 can reach: {4,5,6,7,9,11}
- ▶ d_9 can reach: {10,4,5,6,11}
- ▶ d_{10} can reach: {4,5,6,7,9,11}

- ▶ Node 3: { d_2 }
- ▶ Node 4: { $d_2, d_3, d_5, d_7, d_8, d_9, d_{10}$ }
- ▶ Node 5: { $d_2, d_3, d_7, d_8, d_9, d_{10}$ }
- ▶



Application: Constant Propagation and Folding

- ▶ Given a node $n: t = a \text{ op } b$:
 - ▶ If definition $e: a = c$ (where c is a constant) reaches n and no other definition of a reaches n , then replace n with $e: t = c \text{ op } b$
 - ▶ If a and b are constant, compute v as $a \text{ op } b$ and replace n as $t = v$



Reaching Definition Analysis

- ▶ A definition e (at node m) can reach node n if
 1. There exists a path p from m to n , and
 2. The variable of e is not redefined along this path
- ▶ Flow sets:
 - ▶ $\text{in}_R(n)$: the set of definitions that can reach before n
 - ▶ $\text{out}_R(n)$: the set of definitions that can reach after n
- ▶ Goal of reaching definition analysis: compute $\text{in}_R(n)$ and $\text{out}_R(n)$ for every CFG node n
 - ▶ We only need to consider definitions that actually occur in the method

Transfer Functions

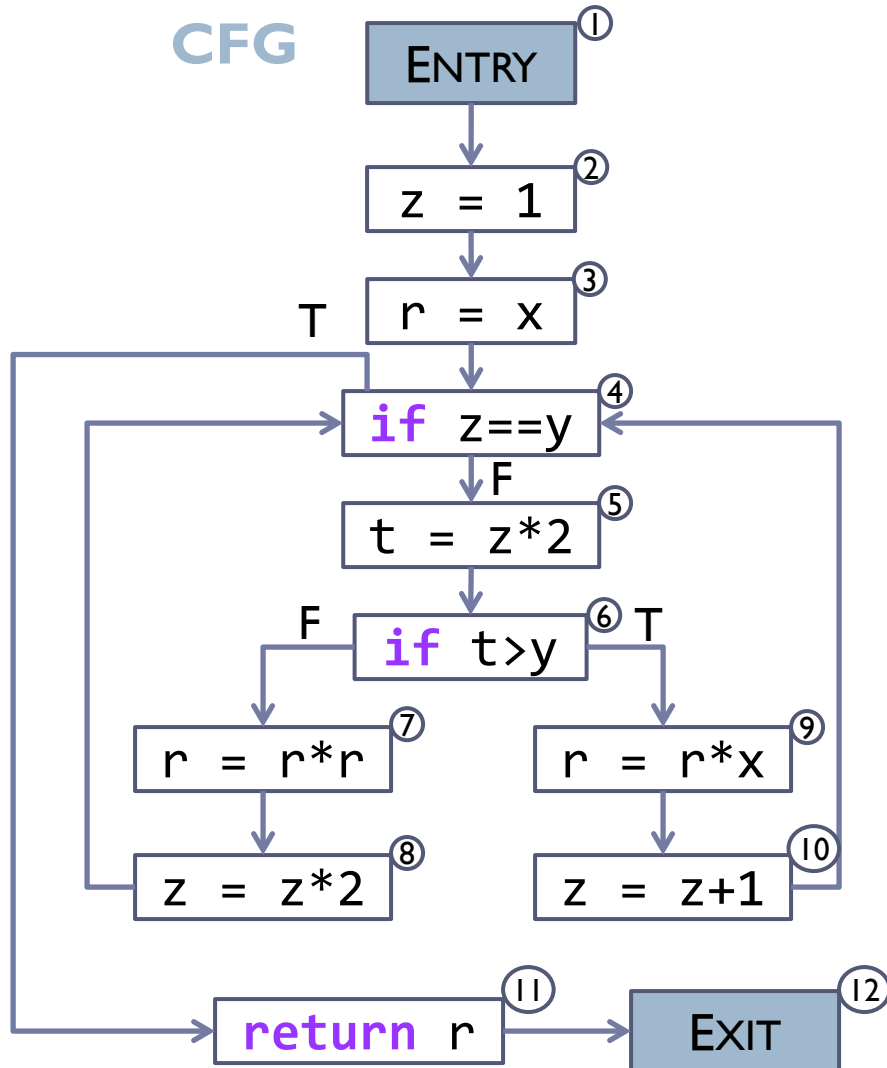
- ▶ If we already know $\text{in}_R(n)$, it is easy to compute $\text{out}_R(n)$:
 - ▶ A definition e can reach after n if
 - (1) Node n gives e , or
 - (2) e reaches before n and n does not re-define the variable in e
- ▶ More denotations:
 - ▶ $\text{gen}(n)$: the set of definitions given by n
 - ▶ $\text{kill}(n)$: the set of definitions whose variables are written by node n
- ▶ **Transfer function** for $\text{in}_R(n)$:
$$\text{out}_R(n) = \text{in}_R(n) \setminus \text{kill}(n) \cup \text{gen}(n)$$
- ▶ This is **forward flow analysis**.

Transfer Functions

- ▶ How do we get $\text{in}_R(n)$? There are two cases:
 - ▶ Node n is the Entry node. Then no definition can reach before a method:
$$\text{in}_R(\text{Entry}) = \emptyset$$
 - ▶ Node n has at least one predecessor node.
 - ▶ $\text{pred}(n)$: the set of predecessor nodes of n .
 - ▶ A definition e can reach before n if it can reach any predecessor of n
- $$\text{in}_R(n) = \bigcup \{ \text{out}_R(m) \mid m \in \text{pred}(n) \}$$
- ▶ This is **may analysis**, as union is used to combine results from predecessor nodes.

Reaching Definition Example

CFG



Transfer Functions

$$\text{out}_A(1) = \text{in}_A(1)$$

$$\text{out}_A(2) = \text{in}_A(2) \setminus \{d_8, d_{10}\} \cup \{d_2\}$$

$$\text{out}_A(3) = \text{in}_A(3) \setminus \{d_7, d_9\} \cup \{d_3\}$$

$$\text{out}_A(4) = \text{in}_A(4)$$

$$\text{out}_A(5) = \text{in}_A(5) \setminus \{\} \cup \{d_5\}$$

$$\text{out}_A(6) = \text{in}_A(6)$$

$$\text{out}_A(7) = \text{in}_A(7) \setminus \{d_3, d_9\} \cup \{d_7\}$$

$$\text{out}_A(8) = \text{in}_A(8) \setminus \{d_2, d_{10}\} \cup \{d_8\}$$

$$\text{out}_A(9) = \text{in}_A(9) \setminus \{d_3, d_7\} \cup \{d_9\}$$

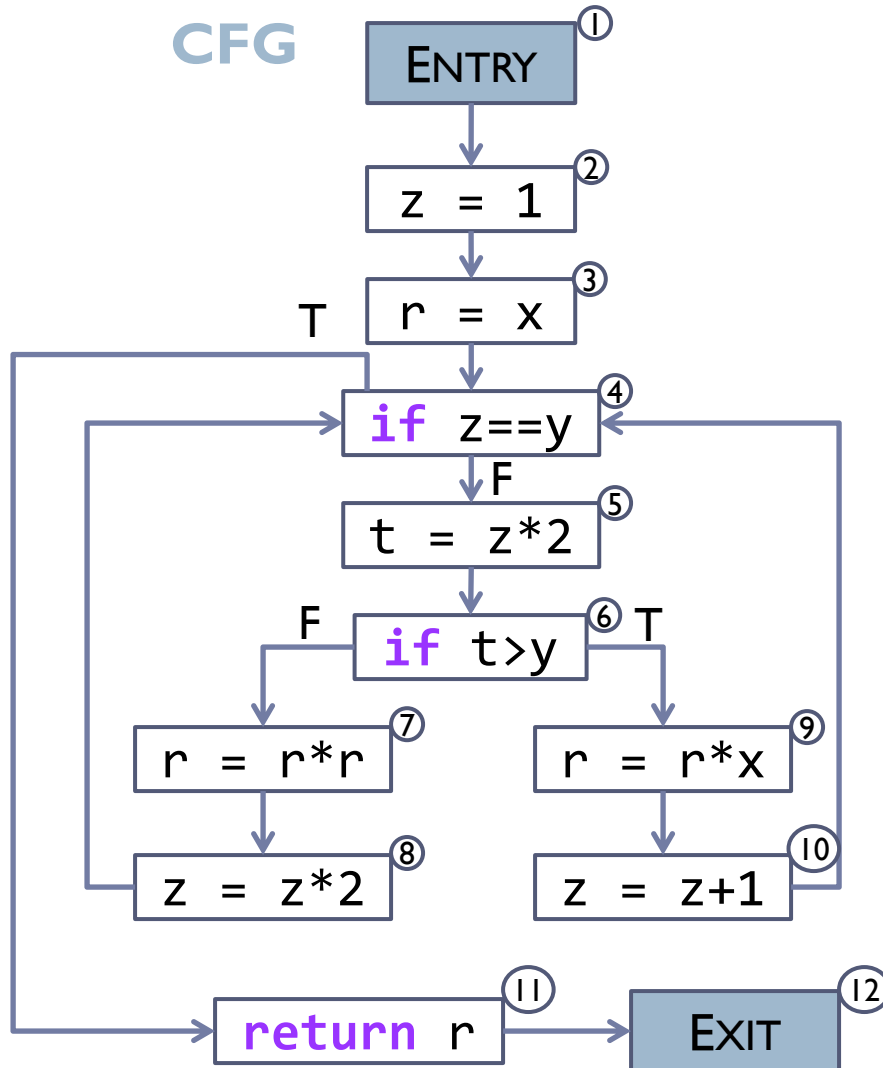
$$\text{out}_A(10) = \text{in}_A(10) \setminus \{d_2, d_8\} \cup \{d_{10}\}$$

$$\text{out}_A(11) = \text{in}_A(11)$$

$$\text{out}_A(12) = \text{in}_A(12)$$

Reaching Definition Example

CFG



Transfer Functions

$$\text{in}_A(1) = \emptyset$$

$$\text{in}_A(2) = \text{out}_A(1)$$

$$\text{in}_A(3) = \text{out}_A(2)$$

$$\text{in}_A(4) = \text{out}_A(3) \cup \text{out}_A(8) \cup \text{out}_A(10)$$

$$\text{in}_A(5) = \text{out}_A(4)$$

$$\text{in}_A(6) = \text{out}_A(5)$$

$$\text{in}_A(7) = \text{out}_A(6)$$

$$\text{in}_A(8) = \text{out}_A(7)$$

$$\text{in}_A(9) = \text{out}_A(6)$$

$$\text{in}_A(10) = \text{out}_A(9)$$

$$\text{in}_A(11) = \text{out}_A(4)$$

$$\text{in}_A(12) = \text{out}_A(11)$$

Solving Transfer Functions

- ▶ The transfer functions can be solved in the same way as available expression analysis, using iterative or worklist algorithms.
- ▶ Initialization:
 - ▶ $\text{out}_R(n) = \text{in}_R(n) = \emptyset$

Outline

- ▶ Available Expression Analysis
- ▶ Reaching Definition Analysis
- ▶ **Summary of Dataflow Analysis**

Categorization of Dataflow Analysis

- ▶ Forward analysis: the set after the node depends on the set before the node
- ▶ Backward analysis: the set before the node depends on the set after the node
- ▶ Must analysis: intersection is used to combine results from predecessor/successor nodes.
- ▶ May analysis: union is used to combine results from predecessor/successor nodes.

	Forward	Backward
Must	Available Expression	Very Busy Expression
May	Reaching Definition	Liveness

Optimization Purposes

- ▶ Liveness:
 - ▶ Analyzing **variables**. Used for **register allocation**.
- ▶ Very busy expression:
 - ▶ Analyzing **expressions**. Used for **code hoisting**.
- ▶ Available expression
 - ▶ Analyzing **expressions**. Used for **common subexpression elimination**.
- ▶ Reaching definition:
 - ▶ Analyzing **definitions**. Used for **constant propagation and folding**

	Forward	Backward
Must	Available Expression	Very Busy Expression
May	Reaching Definition	Liveness

Transfer Functions

- ▶ Two flow sets $\text{in}(n)$ and $\text{out}(n)$ for before and after every node n in the CFG
 - ▶ Liveness: set of variables
 - ▶ Very busy expression: set of expressions
 - ▶ Available expression: set of expressions
 - ▶ Reaching definition: set of definitions
- ▶ For every node, there is a function that computes either $\text{in}(n)$ from $\text{out}(n)$ (backward analysis), or $\text{out}(n)$ from $\text{in}(n)$ (forward analysis)
- ▶ There is also another function that computes either $\text{out}(n)$ from $\text{in}(m)$ for the successors m of n (backward analysis), or $\text{in}(n)$ from $\text{out}(m)$ for the predecessors m of n (forward analysis)

Worklist Algorithm

Backward analysis

```
worklist = [ all nodes ]
while worklist != empty do
     $m = \text{removeFirst}(\text{worklist})$ 
    recompute  $\text{in}(m)$ 
    if  $\text{in}(m)$  has changed then
        for each predecessor  $n$  of  $m$ 
            compute  $\text{out}(n)$ 
            if  $\text{out}(n)$  has changed then
                put  $n$  into worklist (if not
                    already in worklist)
```

Forward analysis

```
worklist = [ all nodes ]
while worklist != empty do
     $m = \text{removeFirst}(\text{worklist})$ 
    recompute  $\text{out}(m)$ 
    if  $\text{out}(m)$  has changed then
        for each successor  $n$  of  $m$ 
            compute  $\text{in}(n)$ 
            if  $\text{in}(n)$  has changed then
                put  $n$  into worklist (if not
                    already in worklist)
```

Worklist Algorithm Initialization

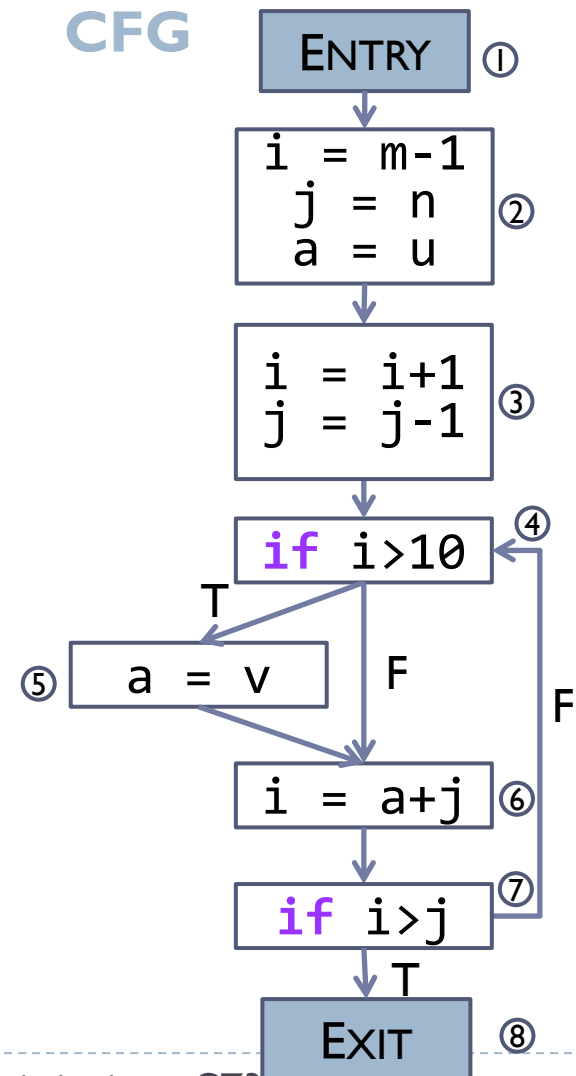
- ▶ For all nodes, if the join operator is union, then $\text{out}(n) = \text{in}(n) = \emptyset$, otherwise U
 - ▶ Exception: $\text{in}(\text{ENTRY})$ (forward analysis) or $\text{out}(\text{EXIT})$ (backward analysis) which are defined by the analysis
- ▶ We initialise the worklist to contain all nodes in the CFG to ensure that each node is evaluated at least once
- ▶ The order in which the worklist is initialized is important:
 - ▶ As far as possible the nodes should be ordered so that a node is only evaluated after all its predecessors or successors have updated their values.
 - ▶ Generally, we initialize the worklist in a normal order for forward analysis, and reversed order for backward analysis. But this can be difficult when the CFG has cycles

For a New Dataflow Analysis Problem

- ▶ Direction: forward or backward?
- ▶ Two sets of transfer function
 - ▶ May or must?
 - ▶ Sets of definition
- ▶ Initial values
 - ▶ Special case for ENTRY or EXIT node

Control Flow Graph with Basic Blocks

- ▶ In these two lectures, we build control flow graph with each node as one instruction.
- ▶ Actually a node can be a basic block, which consists of a sequence of instructions without branches
- ▶ Data flow analysis methods can be performed in the same way. The definition of each set (out, in, use, def, comp, etc.) needs consider all the instructions in the same block
- ▶ Analysis at the basic block level can improve the performance



Conflicts between Optimisations

- ▶ Sometimes different optimisations may be in conflict with each other
 - ▶ e.g., common subexpression elimination increases the live range of local variables (*i.e.*, the number of nodes where they are live)
 - ▶ This makes register allocation harder, so the performance gain of avoiding recomputation may be lost because more memory accesses are needed
- ▶ Sometimes the order of optimisations may also affect the final results.
- ▶ Deciding which optimisations to apply and when to apply them is difficult, and there is no general recipe