

Introduction and Overview

Abstract

This assignment covers some common problems of compilers' middle-end, namely **Very Busy Expressions – VBE**, **Dominator Analysis – DA** and **Constant Propagation – CP**, addressing them with DFA – Data Flow Analysis – framework.

Summary

All the aforementioned problems can be framed as **data-flow problems**; solving them would allow us to drive new optimizations.

The table's In/Out notation is chosen under the assumption that data flow direction doesn't affect blocks' input-output naming, i.e. the **input** of a block is always its **entry point** and the **output** of a block is always its **exit point**, regardless of the problem-specific *direction*.

	Very Busy Expressions	Dominators Analysis	Constant Propagation
Domain	Set of expressions	Set of basic blocks	Set of tuples (var, lit)
Direction	Backward	Forward	Forward
Boundary Conditions	$In[BB_{Exit}] = \emptyset$	$Out[BB_{Entry}] = \{BB_{Entry}\}$	$Out[BB_{Entry}] = \emptyset$
Transfer Function	$In[BB_i] = Gen[BB_i] \cup (Out[BB_i] \setminus Kill[BB_i])$	$Out[BB_i] = In[BB_i] \cup \{BB_i\}$	$Out[BB_i] = Gen[BB_i] \cup (In[BB_i] \setminus Kill[BB_i])$
Meet Operator	\cap	\cap	\cap
Internal Points Init	$In[BB_i] = U \forall i \neq Exit$	$Out[BB_i] = U \forall i \neq BB_{Entry}$	$Out[BB_i] = U \forall i \neq BB_{Entry}$

Table 1: DFA Framework Summary Table

Very Busy Expressions

Problem Definition

DEF. An expression is **very busy** in a given point p **if and only if**, regardless of the path taken after p – i.e. for each path taken from p to *Exit* –, the expression is used before any of its operands is **redefined** (since, as a result, the expression is invalid).

Figuring out which expressions are **busy** in p , would enable us to hoist the expression itself to the point p .

DFA Solution

We defined the problem's **domain** as a **static set of expressions**, which is made up of all RHS expressions in the **CFG**.

The data flow direction is **backward**, since we only care about expressions that appear in the "future".

Set intersection is the **meet operator** that better fit the definition: " $[...] for each path from p to *Exit* [...]$ ".

Given that kind of meet operator, along with the need to keep the intersection with uncomputed interior points neutral, we opted for an initialization according to the following formula: $In[BB_i] = U \forall i \neq Exit$.

The analysis is kick-started by $In[BB_{Exit}] = \emptyset$, and equivalently $Out[BB_{Exit}] = \emptyset$.

Transfer Function's Sets Definition

Each basic block BB_i **inherits** a set of expressions from its **successor** BB_{i+1} ; this set makes up the output set (input from expressions-propagation POV) of block B_i :

$$Out[BB_i] = In[BB_{i+1}] \forall BB_i \in CFG : i \neq Exit$$

Also, each block **kills** the set of expressions invalidated by its internal definitions. An expression $E = OP1 \oplus OP2$ is killed by a basic block BB if either **OP1 or OP are defined inside the block**:

$$Kill[BB_i] = \{E = OP1 \oplus OP2 : OP1 \text{ redefined} \wedge OP2 \text{ redefined}\} \forall BB_i \in CFG$$

Lastly, the **Gen-set** is defined as the set of expressions that appear as right-hand-side of BB 's

internal definitions:

$$Gen[BB_i] = \{OP1_j \oplus OP2_j : E_j = OP1_j \oplus OP2_j \forall E_j \in BB_i\} \forall BB_i \in CFG$$

Illustrative Iteration

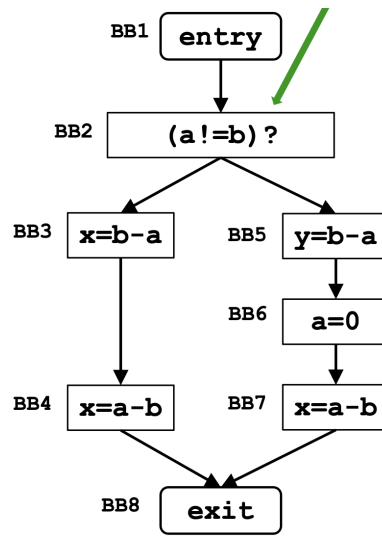


Figure 1: Reference CFG

Basic Block Name	Gen(BB)	Kill(BB)
BB2	\emptyset	\emptyset
BB3	b - a	\emptyset
BB4	a - b	\emptyset
BB5	b - a	\emptyset
BB6	\emptyset	b - a , a - b
BB7	a - b	\emptyset

Table 2: Gen-Kill Table

Node	In(BB)	Out(BB)
Entry	b - a	b - a
BB2	b - a	b - a
BB3	b - a , a - b	a - b
BB4	a - b	\emptyset
BB5	b - a	\emptyset
BB6	\emptyset	a - b
BB7	a - b	\emptyset
Exit	\emptyset	\emptyset

Table 3: Input-Output Table for Iteration 1

The execution stops after the first iteration since there are no back-edges in the CFG.

Dominator Analysis

Problem Definition

In a **CFG**, we say that a node X **dominates** another node Y if node X appears in every path of the graph that leads from the **Entry** block to block Y .

- We annotate each basic block B_i with a set **DOM** $[B_i]$
 - $B_i \in \text{DOM}[B_j]$, if and only if B_i dominates B_j .
- By definition, a **node dominates itself**
 - $B_i \in \text{DOM}[B_i]$

DFA Solution

We defined the domain of the problem as a **static set of blocks**, which are the **nodes** of the **CFG**.

The data flow direction is **forward**, since the concept of dominance is defined in the "past".

The **meet operator** can be inferred by the definition: ... *for each path from Entry to Y ...*, is equivalent to set intersection.

Given that kind of meet operator, along with the need to keep the **intersection** with uncomputed interior points **neutral**, interior points are initialized according to the following formula:
 $\text{Out}[BB_i] = U \ \forall i \neq \text{Entry}.$

The analysis is kick-started by $\text{Out}[BB_{\text{Entry}}] = \text{Entry}$, since according to the definition: *each node dominates itself*.

Transfer Function's Sets Definition

Each basic block BB_i **inherits** a set of expressions from its **predecessor** BB_{i-1} ; this set makes up the block's input set B_i :

$$\text{In}[BB_i] = \text{Out}[BB_{i-1}] \ \forall BB_i \in \text{CFG} : i \neq \text{Entry}$$

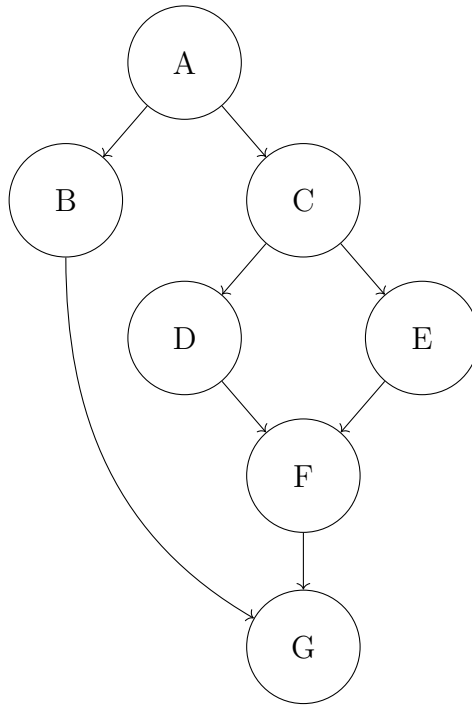
The **kill set is empty** because, considering each node BB_i , name conflicts don't show up. Each node **add** to the dominators' set **itself** without deleting other elements. Some elements could be removed only by applying the **intersection** between the **dominators sets coming from different paths**.

$$\text{Kill}[BB_i] = \emptyset \ \forall BB_i \in \text{CFG}$$

Lastly, as mentioned above, the **Gen-set** of the BB_i is defined as a set with only the BB_i since **each node dominates itself**:

$$Gen[BB_i] = \{BB_i\} \forall BB_i \in CFG$$

Illustrative Iteration

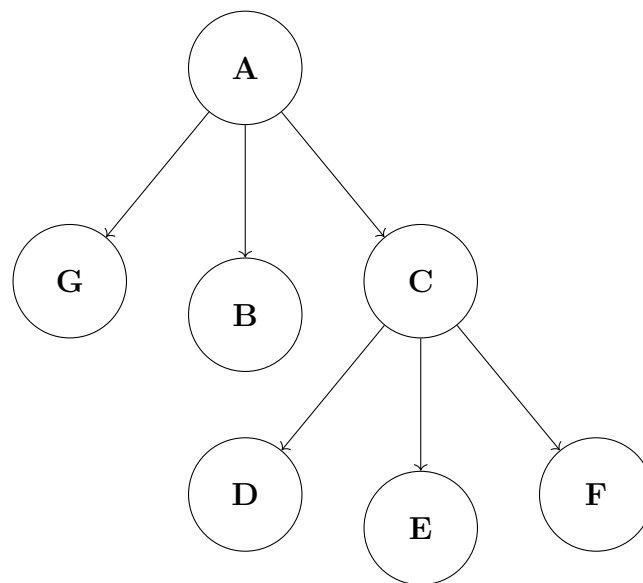


Node	In(BB)	Out(BB)
A	\emptyset	A
B	A	A , B
C	A	A , C
D	A , C	A , C , D
E	A , C	A , C , E
F	A , C	A , C , F
G	A	A , G

Table 4: Input-Output Table for Iteration 1

The execution stops after the first iteration since there are no back-edges in the CFG.

Dominator Tree



Constant Propagation

Problem Definition

The goal of **constant propagation** is to determine at which points in the program variables have a constant value.

The information to be computed for each node n of the Control Flow Graph (**CFG**) is a set of pairs of the form $\langle \text{variable}, \text{constant value} \rangle$.

If we have the pair $\langle x, c \rangle$ at node n , it means that x is guaranteed to have the value c every time n is reached during program execution.

DFA Solution

This solution works under the assumption that the transfer function can integrate some piece of logic. In particular, some sort of expression-solving mechanism is needed to compute expressions such as $k + 2$, expanding them to a constant value. For instance, $k + 2$ expression in BB3 of 2, should become 4 after passing through BB3's transfer function.

Consequently, the domain of this problem, defined as the set of tuples $\langle \text{variable}, \text{constant value} \rangle$, is variable-sized, thus making unfeasible a bit-vector implementation. This is due to the fact we have no prior knowledge about the number of tuples that could appear while algorithm is executing. As a matter of fact,

For a solution that admits a bit-vector implementation, have a look at our considerations in **Constant Propagation With Fixed-size Domain**.

We defined the domain of the problem as a set of definitions with literal or constant values as RHS, represented as tuples $\langle \text{variable}, \text{constant value} \rangle$.

The data flow is oriented forward because it's essential to reference past values of variables to ensure which ones remain constant.

The meet operator is equivalent to set intersection because a variable is deemed constant only if it's confirmed across every potential path.

Given that kind of meet operator, along with the need to keep the intersection with uncomputed interior points neutral, interior points are initialized according to the following formula:

$$Out[BB_i] = U \forall i \neq Entry.$$

The analysis is kick-started by $Out[BB_{Entry}] = \emptyset$.

The solution should use a dynamic size dictionary, because determining the value of a variable defined using operands with constant values requires computation. Each dictionary element's key consists of a pair $\langle \text{variable}, \text{constant value} \rangle$ associated to a boolean value representing whether the variable assumes a constant value at that program point.

Transfer Function's Sets Definition

The kill and gen set is defined using tuple of $\langle \mathbf{variable}, \mathbf{expression} \rangle$ to differentiate definition of the same variables which during the execution could have the same result. The Transfer Function will compute the result of the expression by inserting in the dictionary the tuple $\langle \mathbf{variable}, \mathbf{constant} \quad \mathbf{value} \rangle$. If in the point p of the program the input dictionary does not contain some of the variables in the expression which has to be computed, the tuple $\langle \mathbf{variable}, \mathbf{expression} \rangle$ is discarded from the output set.

Each basic block BB_i inherits a set of expressions from its predecessor BB_{i-1} ; this set makes up the block's input set B_i :

$$In[BB_i] = Out[BB_{i-1}] \quad \forall BB_i \in CFG : i \neq Entry$$

Also, each block kills the tuple $\langle \mathbf{variable}, \mathbf{expression} \rangle$ invalidated by its internal definitions. A tuple $\langle \mathbf{variable}, \mathbf{expression} \rangle$ is killed by a basic block BB if the variable is redefined inside the block:

$$Kill[BB_i] = \{ \langle \mathbf{variable}, \mathbf{expression} \rangle \mid \mathbf{variable} \text{ redefined} \} \quad \forall BB_i \in CFG$$

Lastly, the gen-set is defined as the set of tuple of $\langle \mathbf{variable}, \mathbf{expression} \rangle$ whenever the BB's internal definition of the variable is a *constant expression*, or rather a binary expression where one or both operands are variables whose constant value is known.

Let's define K as a set of constant values, and k as a constant value, $\forall k \in K$, we can declare an *Operand* with constant value as Op_k .

The **Gen-set** based on these assumptions is the following:

$$Gen[BB_i] = \{ \langle \mathbf{variable}, \mathbf{expression} \rangle \mid \forall E_j \in BB_i \wedge Op_k \quad \forall k \in K \} \quad \forall BB_i \in CFG$$

Illustrative Iteration

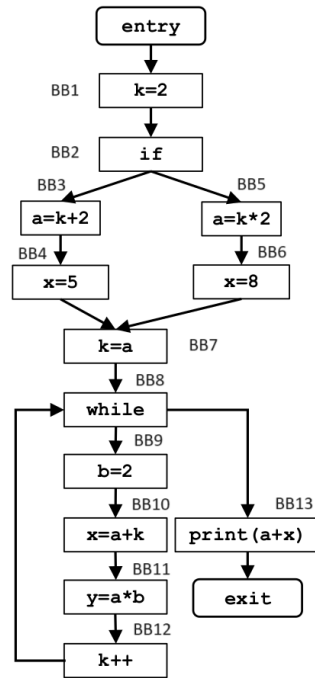


Figure 2: Reference CFG

Basic Block Name	Gen(BB)	Kill(BB)
BB1	$\langle k, 2 \rangle$	\emptyset
BB2	\emptyset	\emptyset
BB3	$\langle a, k + 2 \rangle$	$\langle a, k * 2 \rangle$
BB4	$\langle x, 5 \rangle$	$\langle x, 8 \rangle, \langle x, a + k \rangle$
BB5	$\langle a, k * 2 \rangle$	$\langle a, k + 2 \rangle$
BB6	$\langle x, 8 \rangle$	$\langle x, 5 \rangle, \langle x, a + k \rangle$
BB7	$\langle k, a \rangle$	$\langle k, 2 \rangle, \langle k, k++ \rangle$
BB8	\emptyset	\emptyset
BB9	$\langle b, 2 \rangle$	\emptyset
BB10	$\langle x, a + k \rangle$	$\langle x, 5 \rangle, \langle x, 8 \rangle$
BB11	$\langle y, a * b \rangle$	\emptyset
BB12	$\langle k, k++ \rangle$	$\langle k, 2 \rangle, \langle k, a \rangle$
BB13	\emptyset	\emptyset

Table 5: Gen-Kill Table

Iteration 1

Node	In(BB)	Out(BB)
Entry	\emptyset	\emptyset
BB1	\emptyset	$\langle k, 2 \rangle$
BB2	$\langle k, 2 \rangle$	$\langle k, 2 \rangle$
BB3	$\langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle$
BB4	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle, \langle x, 5 \rangle$
BB5	$\langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle$
BB6	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle, \langle x, 8 \rangle$
BB7	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 4 \rangle$
BB8	$\langle a, 4 \rangle, \langle k, 4 \rangle$	$\langle a, 4 \rangle, \langle k, 4 \rangle$
BB9	$\langle a, 4 \rangle, \langle k, 4 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle k, 4 \rangle$
BB10	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle k, 4 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle k, 4 \rangle, \langle x, 8 \rangle$
BB11	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle k, 4 \rangle, \langle x, 8 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle k, 4 \rangle, \langle x, 8 \rangle, \langle y, 8 \rangle$
BB12	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle k, 4 \rangle, \langle x, 8 \rangle, \langle y, 8 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle k, 5 \rangle, \langle x, 8 \rangle, \langle y, 8 \rangle$
BB13	$\langle a, 4 \rangle, \langle k, 4 \rangle$	$\langle a, 4 \rangle, \langle k, 4 \rangle$
Exit	$\langle a, 4 \rangle, \langle k, 4 \rangle$	

Table 6: Input-Output Table, Iteration 1

The execution continues after the first iteration since there are back-edges in the CFG.

Iteration 2

Node	In(BB)	Out(BB)
Entry	\emptyset	\emptyset
BB1	\emptyset	$\langle k, 2 \rangle$
BB2	$\langle k, 2 \rangle$	$\langle k, 2 \rangle$
BB3	$\langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle$
BB4	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle, \langle x, 5 \rangle$
BB5	$\langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle$
BB6	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle, \langle x, 8 \rangle$
BB7	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 4 \rangle$
BB8	$\langle a, 4 \rangle$	$\langle a, 4 \rangle$
BB9	$\langle a, 4 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle$
BB10	$\langle a, 4 \rangle, \langle b, 2 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle$
BB11	$\langle a, 4 \rangle, \langle b, 2 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle y, 8 \rangle$
BB12	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle y, 8 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle y, 8 \rangle$
BB13	$\langle a, 4 \rangle$	$\langle a, 4 \rangle$
Exit	$\langle a, 4 \rangle$	

Table 7: Input-Output Table, Iteration 2

In the BB10 the tuple of the gen set $\langle \mathbf{x}, \mathbf{a+k} \rangle$ is discarded because the variable k is not defined in the input dictionary.

Iteration 3

Node	In(BB)	Out(BB)
Entry	\emptyset	\emptyset
BB1	\emptyset	$\langle k, 2 \rangle$
BB2	$\langle k, 2 \rangle$	$\langle k, 2 \rangle$
BB3	$\langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle$
BB4	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle, \langle x, 5 \rangle$
BB5	$\langle k, 2 \rangle$	$\langle a, 4 \rangle; \langle k, 2 \rangle$
BB6	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 2 \rangle, \langle x, 8 \rangle$
BB7	$\langle a, 4 \rangle, \langle k, 2 \rangle$	$\langle a, 4 \rangle, \langle k, 4 \rangle$
BB8	$\langle a, 4 \rangle$	$\langle a, 4 \rangle$
BB9	$\langle a, 4 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle$
BB10	$\langle a, 4 \rangle, \langle b, 2 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle$
BB11	$\langle a, 4 \rangle, \langle b, 2 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle y, 8 \rangle$
BB12	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle y, 8 \rangle$	$\langle a, 4 \rangle, \langle b, 2 \rangle, \langle y, 8 \rangle$
BB13	$\langle a, 4 \rangle$	$\langle a, 4 \rangle$
Exit	$\langle a, 4 \rangle$	

Table 8: Input-Output Table, Iteration 3

The execution stops after the third iteration since there are no changes in the input and output sets.

Constant Propagation With Fixed-size Domain