

gdfa: A Generic Data Flow Analyzer for GCC

(Version 1.1)

Uday Khedker

<http://cse.iitb.ac.in/~uday>

Department of Computer Science and Engineering

Indian Institute of Technology Bombay

Powai, Mumbai 400076 India.

February 25, 2009

Abstract

This document describes a *generic data flow analyzer* for *per function* (i.e., intraprocedural) *bit vector data flow analysis* in GCC 4.3.0. We call this infrastructure *gdfa*. The analyzers implemented using *gdfa* are called *pfbvdfa*. *gdfa* has been used to implement several bit vector data flow analyses.

1 Motivation

The design and implementation of *gdfa* is motivated by the following objectives:

- Demonstrating the practical significance of the following important generalization: Instead of implementing specific analyses directly, it is useful to implement a generic driver that is based on a carefully chosen set of abstractions. The task of implementing a particular analyzer then reduces to merely specifying the analysis by instantiating these abstractions to concrete values.
- Providing an easy to use and easy to extend data flow analysis infrastructure. The goal is to facilitate experimentation in terms of studying existing analyses, defining new analyses, and exploring different analysis algorithms.

Section 2 describes the specification mechanism of *gdfa* and shows how the resulting pass can be included in GCC 4.3.0. We illustrate it for the bit vector analyses implemented using *gdfa*. Section 3 describes the implementation of *gdfa*. This section also shows how local property computation can be driven by specifications. Finally Section 4 suggests some possible enhancements to *gdfa*.

In this document, we assume familiarity with data flow analysis and GCC internals. Section 5 point to further readings.

The source code of *gdfa* is available as a patch of the gcc directory for GCC-4.3.0 from the URL:

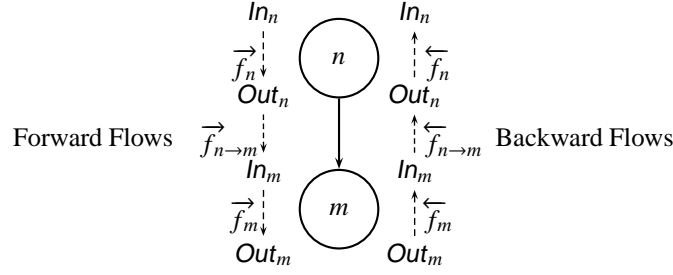


Figure 1: Associating flow functions with nodes and edges separately.

<http://www.cse.iitb.ac.in/uDAY/dfaBook>

Patches for later versions will be made available on this page whenever possible.

The code presented in this document is a slightly edited version of the original code. This was done to fit a page size constraints.

2 Specifying a Data Flow Analysis

In this section we look at how we can use the generic data flow analysis driver to implement a data flow analysis pass in GCC. The implemented pass has to be registered with the pass manager in GCC so that it can be executed by the compiler.

2.1 Generic Flow Functions and Data Flow Equations

Generic flow functions are defined in terms of flow functions illustrated in Figure 1. \vec{f} denotes a forward flow function whereas \overleftarrow{f} denotes a backward flow function. The subscripts used in flow function notation distinguish node flow functions from edge flow functions. Defining separate node and edge flow functions requires explicating In_n and Out_n rather than leaving one of them implicit. For forward unidirectional data flows, the forward flow functions associated with edges are identity functions and the backward node and edge flow functions compute \top . Analogous remarks hold for backward unidirectional data flows.

When separate flow functions are associated with nodes and edges, the generic data flow equations can be written as shown below.

$$In_n = \begin{cases} Bl_{Start} \sqcap \overleftarrow{f}_n(Out_n) & n = Start \\ \left(\prod_{m \in pred(n)} \vec{f}_{m \rightarrow n}(Out_m) \right) \sqcap \overleftarrow{f}_n(Out_n) & \text{otherwise} \end{cases} \quad (1)$$

$$Out_n = \begin{cases} Bl_{End} \sqcap \vec{f}_n(In_n) & n = End \\ \left(\prod_{m \in succ(n)} \overleftarrow{f}_{m \rightarrow n}(In_m) \right) \sqcap \vec{f}_n(In_n) & \text{otherwise} \end{cases} \quad (2)$$

where BI_{End} and BI_{Start} denotes boundary information for intraprocedural data flow analysis. These equations compute the *MFP* solution of an instance of a data flow framework.

2.2 Registering a Pass With the Pass Manager in GCC

gdfa works on the gimple version of the intermediate representation used by GCC. We have included *pfbvdfa* passes such that they are invoked by default when gcc is used for compiling a program. When gcc is built, this causes *pfbvdfa* passes to run on the entire source of gcc which consists of over a million lines of C code. This helps in ensuring that these do not cause any exception in the compilation sequence.

After constructing the gimple representation, gcc views the rest of the compilation as sequential execution of various passes. This is carried out by traversing a linked list whose nodes contain pointers to the entry functions of these passes. A pass is registered with the pass manager through the following steps:

- Instantiating a variable as an instance of `struct tree_opt_pass` in some file.
- Declaring this variable as an `extern` variable in header file `tree-pass.h`.
- Inserting this variable in the linked list of passes using the macro `NEXT_PASS` in function `init_optimization_passes` in file `passes.c`.
- Listing new file names in `gcc/Makefile.in` and configuring and building GCC.

Here is the declaration of `struct tree_opt_pass`. For convenience comments have been removed and are used in the explanation that follows.

```

0 struct tree_opt_pass
1 {
2     const char *name;
3     bool (*gate) (void);
4     unsigned int (*execute) (void);
5     struct tree_opt_pass *sub;
6     struct tree_opt_pass *next;
7     int static_pass_number;
8     unsigned int tv_id;
9     unsigned int properties_required;
10    unsigned int properties_provided;
11    unsigned int properties_destroyed;
12    unsigned int todo_flags_start;
13    unsigned int todo_flags_finish;
14    char letter;
15 };

```

The name of the pass (line 2) is used as a fragment of the dump file name. We have used the names like *gdfa_ave*. The gate function (line 3) is used to check whether

this pass and all its sub-passes should be executed or not. They are executed only if this function returns `true`. If no such checking is required, this function pointer can be `NULL`. The `execute` function (line 4) is entry function of the pass. If this function pointer is `NULL`, there should be sub-passes otherwise this pass does nothing. The return value tells `gcc` what more needs to be done. The variable `sub` (line 5) is a list of sub-passes that should be executed depending upon the gate predicate. If there are sub-passes that must be executed unconditionally, then they are listed in `next` (line 6). The static pass number (line 7) is used as a fragment of the dump file name. If it is specified as 0, the pass manager computes its value depending on the position of the pass. It is this that generated numbers 15, 16, 17, 18, and 19 for our data flow analyses. Variable `tv_id` is the variable that can be used as a time variable. The rest of the variables are self-explanatory. The last variable `letter` is used to annotate RTL code that is emitted.

We have registered available expressions analysis by creating a structure variable called `pass_gimple_pfbv_ave_dfa` as shown below.

```
struct tree_opt_pass pass_gimple_pfbv_ave_dfa =
{
  "gdfa_ave",          /* name */
  NULL,                /* gate */
  gimple_pfbv_ave_dfa, /* execute */
  NULL,               /* sub */
  NULL,               /* next */
  0,                  /* static_pass_number */
  0,                  /* tv_id */
  0,                  /* properties_required */
  0,                  /* properties_provided */
  0,                  /* properties_destroyed */
  0,                  /* todo_flags_start */
  0,                  /* todo_flags_finish */
  0                   /* letter */
};
```

This variable is declared as follows in file `tree-pass.h`

```
extern struct tree_opt_pass pass_gimple_pfbv_ave_dfa;
```

The next step in registering this pass is to include it in the list of passes. We show below the relevant code fragment from function `init_optimization_passes` in file `passes.c`:

```

NEXT_PASS (pass_build_cfg);
/* Intraprocedural dfa passes begin */
NEXT_PASS (pass_init_gimple_pfbvdfa);
NEXT_PASS (pass_gimple_pfbv_ave_dfa);
NEXT_PASS (pass_gimple_pfbv_pav_dfa);
NEXT_PASS (pass_gimple_pfbv_ant_dfa);
NEXT_PASS (pass_gimple_pfbv_lv_dfa);
NEXT_PASS (pass_gimple_pfbv_rd_dfa);
NEXT_PASS (pass_gimple_pfbv_pre_dfa);
/* Intraprocedural dfa passes end */

```

Finally, we need to include the new file names in the GCC build system. This is done by including the object file names and their dependencies in `Makefile.in` in the `gcc-4.3.0/gcc` directory.

2.3 Specifying Available Expressions Analysis

The specification mechanism supported by *gdfa* is simple and succinct. It follows the GCC mechanism of specification by using a `struct` as a hook and by requiring the user to create a variable by instantiating the members of the `struct` defined for the purpose.

For available expressions analysis, we define a variable called `gdfa_ave` which is of the type `struct gimple_pfbv_dfa_spec gdfa_ave`.

```

0 struct gimple_pfbv_dfa_spec gdfa_ave =
1 {
2     entity_expr,          /* entity          */
3     ONES,                 /* top_value       */
4     ZEROS,                /* entry_info      */
5     ONES,                 /* exit_info       */
6     FORWARD,             /* traversal_order  */
7     INTERSECTION,         /* confluence      */
8     entity_use,           /* gen_effect      */
9     down_exp,             /* gen_exposition  */
10    entity_mod,            /* kill_effect     */
11    any_where,             /* kill_exposition */
12    global_only,           /* preserved_dfi   */
13    identity_forward_edge_flow, /* forward_edge_flow */
14    stop_flow_along_edge,  /* backward_edge_flow */
15    forward_gen_kill_node_flow, /* forward_node_flow */
16    stop_flow_along_node  /* backward_node_flow */
17 };

```

Before we explain the above, we present the rest of the code required to complete

the specification.

```
18 pfbv_dfi ** AV_pfbv_dfi = NULL;
19
20 static unsigned int
21 gimple_pfbv_ave_dfa(void)
22 {
23
24     AV_pfbv_dfi = gdfa_driver(gdfa_ave);
25
26     return 0;
27 }
```

Nothing more is required for specifying available expressions analysis apart from registering it with the pass manager with function `gimple_pfbv_ave_dfa` as its entry point as described in Section 2.2. This function calls the *gdfa* driver passing the specification variable `gdfa_ave` as actual parameter. The data flow information computed by the driver is stored in a pointer to an array called `AV_pfbv_dfi`; each element of this array represents the data flow information for a basic block and is an instance of the following type defined by *gdfa*.

```
typedef struct pfbv_dfi
{
    dfvalue gen;
    dfvalue kill;
    dfvalue in;
    dfvalue out;
} pfbv_dfi;
```

The semantics expressed by `struct gimple_pfbv_dfa_spec gdfa_ave` is as described below: Line 2 declares that the relevant entities for this analysis are expressions (`entity_expr`). Line 3 specifies that \top is “all ONES” implying the universal set $\mathbb{E}xpr$. The specification “all ZEROS” on line 4 initializes the Bl_{Start} to \emptyset whereas ONES on line 5 renders Bl_{End} irrelevant because it is same as \top . Line 6 declares the direction of traversal to be FORWARD. Note that this is independent of the direction of flow and only influences the number of iterations. If we choose the direction of traversal as BACKWARD, the resulting data flow information will remain same except that it may take a much larger number of iterations. Line 7 declares the \sqcap to be \cap . Line 12 directs the driver to preserve only the global data flow information (*In* and *Out*); the driver can reclaim the space occupied by the local data flow information (*Gen* and *Kill*).

The most interesting elements of the specification are the specifications of local properties and flow functions:

- *Local property specification.*

Lines 8 to 11 define the *Gen* and *Kill* kill sets for a block. Observe that this mechanism closely follows the description in Section ??.

- Lines 8 and 9 say that when a downwards exposed (*down_exp*) use of an entity (*entity_use*) is found in a basic block, it is included in the *Gen* set of the block. From line 2 we know that the entity under consideration is an expression (*entity_expr*).
- Lines 10 and 11 say that when a modification of an entity (*entity_mod*) is found in a basic block, it is included in the *Kill* set of the block. This modification need not be upwards exposed or downwards exposed, it can appear *anywhere*.

This is possible because the *gdfa* driver is aware of the fact that the use of an entity could be affected by its modification and hence the notion of exposition of an entity is explicated in the specification.

- *Flow function specification.*

Lines 13 to 16 specify the flow functions for available expressions analysis as required by the generic data flow Equations (1) and (2).

- The forward edge flow function $\vec{f}_{n \rightarrow m}$ in available expressions analysis is identity (line 13).
- The forward node flow function \vec{f}_n is the conventional *Gen-Kill* function $f(X) = Gen \cup (In - Kill)$. This is specified by line 15.
- There is no backward flow i.e., \overleftarrow{f}_n and $\overleftarrow{f}_{n \rightarrow m}$ are \top . This is specified by lines 14 and 16.

All these functions are supported by *gdfa* and it is enough to associate the function pointers with appropriate functions.

When the nature of data flow is different from the default flows, it is also possible to write custom functions—we show how it is done for partial redundancy elimination.

2.4 Specifying Other Bit Vector Data Flow Analyses

Given the specification of available expressions analysis, it is easy to visualize specifications for other bit vector frameworks. We describe the required changes in the following:

- *Partially available expressions analysis.*

Confluence should be UNION, \top and Bl_{End} should be ZEROS.

- *Anticipable expressions analysis.*

In this case it is desirable, though not necessary, to choose the direction of traversal as BACKWARD. The exposition for *Gen* should be changed to *up_exp*. Bl_{Start}

should be ONES and Bl_{End} should be ZEROS. Flow functions would change as follows:

- forward edge flow function $\vec{f}_{n \rightarrow m}$ should be `stop_flow_along_edge`,
- forward node flow function \vec{f}_n should be `stop_flow_along_node`, and
- backward node flow function \overleftarrow{f}_n should be the default *Gen-Kill* function `backward_gen_kill_node_flow`.

- *Live variables analysis.*

This specification would be similar to that of anticipable expressions analysis except that the entity should be `entity_var`, confluence should be UNION, \top and Bl_{End} should be ZEROS.

- *Reaching definitions analysis.*

This is a forward data flow analysis similar to available expressions analysis except that the entity is `entity_defn`, confluence is UNION, and \top and Bl_{End} are ZEROS.

- *Partial redundancy elimination.*

Here it would useful to change the `gate` function to this pass to check that available expressions analysis and partially available expressions analysis has been performed.

The data flow equations for partial redundancy elimination are given below.

$$In_n = PavIn_n \cap (AntGen_n \cup (Out_n - Kill_n)) \cap \bigcap_{p \in pred(n)} (Out_p \cup AvOut_p) \quad (3)$$

$$Out_n = \begin{cases} Bl & n \text{ is } End \text{ block} \\ \bigcap_{s \in succ(n)} In_s & \text{otherwise} \end{cases} \quad (4)$$

The specification of data flow analysis would be similar to that of anticipable expressions analysis except that the node flow function in the equation for In_n would change. In particular, the forward edge flow function $\vec{f}_{n \rightarrow m}$ and the backward node flow function \overleftarrow{f}_n cannot be chosen from the default functions supported by *gdfa*. We define the required functions as shown below.

```

dfvalue
forward_edge_flow_pre(basic_block src, basic_block dest)
{
    dfvalue temp;

    temp = union_dfvalues (OUT(AV_pfbv_dfi,src),
                          CURRENT_OUT(src));

    return temp;
}

```

In this function, *src* and *dest* indicate the source and destination of an edge. Since this flow function is used in computing ln_n , *dest* represents n and *src* represents the given predecessor node p . Under the assumption that the data flow information of available expressions analysis is stored in the variable *AV_pfbv_dfi*, the term *OUT(AV_pfbv_dfi, src)* represents $AvOut_p$ whereas the Out_p is represented by the term *CURRENT_OUT(src)*. Thus this flow function computes $AvOut_p \cup Out_p$ for a given predecessor p .

The definition of backward node flow is similar to that of the default node flow except that we need to include the value of $Pavln_n$. This is easily achieved by the function defined below:

```

dfvalue
backward_node_flow_pre(basic_block bb)
{
    dfvalue temp1, temp2;

    temp1 = backward_gen_kill_node_flow(bb);

    temp2 = intersect_dfvalues (IN(PAV_pfbv_dfi,bb),
                              temp1);

    if (temp1)
        free_dfvalue_space(temp1);

    return temp2;
}

```

Here *bb* is the current node n . The default backward node flow function is used to compute the data flow information in the variable *temp1*. Under the assumption that the data flow information of partially available expressions analysis is stored in the variable *PAV_pfbv_dfi*, the term *IN(PAV_pfbv_dfi, bb)* represents $Pavln_n$. All that further needs to be done is to intersect them.

This completes the specification of partial redundancy elimination.

3 Implementing *gdfa*

We describe the implementation in terms of the specification primitives, interface with GCC, the generic functions for global property computation, and generic functions for local property computation.

3.1 Specification Primitives

The main data structure used for specification is:

```
0 struct gimple_pfbv_dfa_spec
1 {
2     entity_name          entity;
3     initial_value        top_value_spec;
4     initial_value        entry_info;
5     initial_value        exit_info;
6     traversal_direction   traversal_order;
7     meet_operation        confluence;
8     entity_manipulation   gen_effect;
9     entity_occurrence     gen_exposition;
10    entity_manipulation    kill_effect;
11    entity_occurrence      kill_exposition;
12    dfi_to_be_preserved    preserved_dfi;
13
14    dfvalue (*forward_edge_flow)(basic_block src,
15                                basic_block dest);
16    dfvalue (*backward_edge_flow)(basic_block src,
17                                  basic_block dest);
18    dfvalue (*forward_node_flow)(basic_block bb);
19    dfvalue (*backward_node_flow)(basic_block bb);
20
21 };
```

The types appearing on lines 2 to 12 are defined as enumerated types with the following possible values.

Enumerated Type	Possible Values
entity_name	entity_expr, entity_var, entity_defn
initial_value	ONES, ZEROS
traversal_direction	FORWARD, BACKWARD, BIDIRECTIONAL
meet_operation	UNION, INTERSECTION
entity_manipulation	entity_use, entity_mod
entity_occurrence	up_exp, down_exp, any_where
dfi_to_be_preserved	all, global_only, no_value

The type `dfvalue` is just another name for the type `sbitmap` supported by GCC.

We have used a different name to allow for the possibility of extending *gdfa* to other kinds of data flow values.

The entry point of each data flow analysis invokes the driver with its specification. The driver creates space for current data flow values in current data flow analysis in a variable `current_pfbv_dfi` which is declared as shown below:

```
typedef struct pfbv_dfi
{
    dfvalue gen;
    dfvalue kill;
    dfvalue in;
    dfvalue out;
} pfbv_dfi;

pfbv_dfi ** current_pfbv_dfi ;
```

For a basic block `bb`, different members of the data flow information are accessed using the following macros:

Data flow variable	<code>current_pfbv_dfi</code>	Given dfi
<i>Gen</i>	<code>CURRENT_GEN(bb)</code>	<code>GEN(dfi, bb)</code>
<i>Kill</i>	<code>CURRENT_KILL(bb)</code>	<code>KILL(dfi, bb)</code>
<i>In</i>	<code>CURRENT_IN(bb)</code>	<code>IN(dfi, bb)</code>
<i>Out</i>	<code>CURRENT_OUT(bb)</code>	<code>OUT(dfi, bb)</code>

Now we can describe the default functions that can be assigned to the function pointers on lines 14 to 19 in `struct gimple_pfbv_dfa_spec`. Alternatively, the users can define their own functions which have the same interface. The default functions supported by *gdfa* are:

Function	Returned value
<code>identity_forward_edge_flow(src, dest)</code>	<code>CURRENT_OUT(src)</code>
<code>identity_backward_edge_flow(src, dest)</code>	<code>CURRENT_IN(dest)</code>
<code>stop_flow_along_edge(src, dest)</code>	<code>top_value</code>
<code>identity_forward_node_flow(bb)</code>	<code>CURRENT_IN(bb)</code>
<code>identity_backward_node_flow(bb)</code>	<code>CURRENT_OUT(bb)</code>
<code>stop_flow_along_node(bb)</code>	<code>top_value</code>
<code>forward_gen_kill_node_flow(bb)</code>	$CURRENT_GEN(bb) \cup (CURRENT_IN(bb) - CURRENT_KILL(bb))$
<code>backward_gen_kill_node_flow(bb)</code>	$CURRENT_GEN(bb) \cup (CURRENT_OUT(bb) - CURRENT_KILL(bb))$

where `top_value` is of the type `initial_value` and is constructed based on the value of `top_value_spec` (line 3 in `struct gimple_pfbv_dfa_spec`).

This completes the description of the specification primitives.

3.2 Interface with GCC

The top level interface of *gdfa* with GCC is through the pass manager as described in Section 2.2. At the lower level, *gdfa* uses the support provided by GCC for traversals over CFGs, basic blocks etc.; discovering relevant features of statements, expressions, variables etc.; constructing and manipulating data flow values; and printing entities appearing in statements.

Traversal Over CFG and Basic Blocks

In a round robin iterative traversal, the basic blocks in a CFG are usually visited in the order of along control flow or against the order of control flow. In GCC, this is achieved as follows:

```
basic_block bb;

FOR_EACH_BB_FWD(ENTRY_BLOCK_PTR)
{
    /* process bb */
}
FOR_EACH_BB_BKD(EXIT_BLOCK_PTR)
{
    /* process bb */
}
```

In the above code, `basic_block` is a type supported by GCC. `ENTRY_BLOCK_PTR` and `EXIT_BLOCK_PTR` point to `ENTRY` and `EXIT` blocks of the current function being compiled. These macros have been defined by GCC. The two other macros used above are defined as follows:

```
#define FOR_EACH_BB_FWD(entry_bb) \
    for(bb=entry_bb->next_bb; \
        bb->next_bb!=NULL; \
        bb=bb->next_bb)
#define FOR_EACH_BB_BKD(exit_bb) \
    for(bb=exit_bb->prev_bb; \
        bb->prev_bb!=NULL; \
        bb=bb->prev_bb)
```

Given a basic block `bb`, its predecessor and successor blocks are traversed using an `edge_iterator` variable, an `edge` variable, and the macro `FOR_EACH_EDGE` as described below. All these are directly supported by GCC.

```

edge_iterator ei ;
edge e ;
basic_block succ_bb, pred_bb;

FOR_EACH_EDGE(e,ei,bb->preds)
{
    pred_bb = e->src;
    /* process the predecessor pred_bb */
}
FOR_EACH_EDGE(e,ei,bb->succs)
{
    succ_bb = e->dest;
    /* process successor succ_bb */
}

```

A statement is of the type `tree`. Further, all entities appearing in a statement are also of the type `tree`. All statements in a basic block can be traversed using a `block_statement_iterator` variable.

```

basic_block bb;
block_stmt_iterator bsi;
tree stmt;

FOR_EACH_STMT_FWD
{
    stmt = bsi_stmt(bsi);
    /* process stmt */
}
FOR_EACH_STMT_BKD
{
    stmt = bsi_stmt(bsi);
    /* process stmt */
}

```

The macros used in the above code are defined as follows:

```

#define FOR_EACH_STMT_FWD          \
    for(bsi=bsi_start(bb);         \
        !bsi_end_p(bsi);           \
        bsi_next(&bsi))

#define FOR_EACH_STMT_BKD          \
    for(bsi=bsi_last(bb);           \
        bsi.tsi.ptr!=NULL;         \
        bsi_prev(&bsi))

```

Discovering the Entities in a Statement

Statements can be of many types but only a few types are relevant to local data flow analysis. The lvalue and rvalue of a given statement `stmt` are of the type `tree` and are extracted as shown below:

```
tree expr=NULL, lval=NULL;

switch(TREE_CODE(stmt))
{ case COND_EXPR:
    expr = TREE_OPERAND(stmt,0);
    break;
  case MODIFY_EXPR:
    lval = TREE_OPERAND(stmt,0);
    expr = TREE_OPERAND(stmt,1);
  case GIMPLE_MODIFY_STMT:
    lval = GIMPLE_STMT_OPERAND(stmt,0);
    expr = GIMPLE_STMT_OPERAND(stmt,1);
    break;
  default:
    break;
}
```

The operands of relevant expressions are extracted as shown below:

```
tree op0=NULL, op1=NULL;

switch(TREE_CODE(expr))
{ case MULT_EXPR:
  case PLUS_EXPR:
  case MINUS_EXPR:
  case LT_EXPR:
  case LE_EXPR:
  case GT_EXPR:
  case GE_EXPR:
  case NE_EXPR:
  case EQ_EXPR:
    op1 = TREE_OPERAND(stmt,1);
    op0 = TREE_OPERAND(stmt,0);
    break;
  default:
    break;
}
```

Observe that this covers the set of expressions that is currently supported by *gdfa*.

Clearly, extending this set is easy.

Local variables are discovered by traversing `cfun->unexpanded_var_list` using `TREE_VALUE` and `TREE_CHAIN` macros supported by GCC. Here `cfun` represents the current function being compiled.

```
tree var, list;

list = cfun->unexpanded_var_list;
while (list)
{   var = TREE_VALUE (list);
    /* process variables */
    list = TREE_CHAIN(list);
}
```

Discovering definitions is easy: A statement with `TREE_CODE` as `MODIFY_EXPR` or `GIMPLE_MODIFY_STMT` is detected as a definition.

Constructing and Manipulating Data Flow Values

We define the type `dfvalue` as follows:

```
typedef sbitmap dfvalue;
```

`sbitmap` is a type supported by GCC to represent sets. We use the following `sbitmap` functions to construct and manipulate bitmaps. Note that these functions are not directly used in *gdfa*. Instead, *gdfa* code calls `dfvalue` functions that are defined in terms of these functions.

Name of the Function	Action
<code>sbitmap_equal(v_a, v_b)</code>	is v_a equal to v_b ?
<code>sbitmap_a_and_b(t, v_a, v_b)</code>	$t = v_a \cap v_b$
<code>sbitmap_union_of_diff(t, v_a, v_b, v_c)</code>	$t = v_a \cup (v_b - v_c)$
<code>sbitmap_a_or_b(t, v_a, v_b)</code>	$t = v_a \cup v_b$
<code>sbitmap_ones(v)</code>	set every bit in v to 1
<code>sbitmap_zero(v)</code>	set every bit in v to 0
<code>sbitmap_alloc(n)</code>	allocate a bitmap of n bits
<code>sbitmap_free(v)</code>	free the space occupied by v

Facilities for Printing Entities

We use the function `dump_sbitmap` to print bitmaps. For printing a statement, the function `print_generic_stmt` is used whereas function `print_generic_expr` prints an expression `expr`.

3.3 The Preparatory Pass

Before the *gdfa* driver is invoked, some preparatory work has to be performed by an earlier pass. The top level function of this pass is:

```
static unsigned int
init_gimple_pfbvdfa_execute (void)
{
    local_var_count=0;
    local_expr_count=0;
    number_of_nodes = n_basic_blocks+2;

    assign_indices_to_var();
    assign_indices_to_exprs();
    assign_indices_to_defns();

    dfs_ordered_basic_blocks = NULL;
    dfs_numbering_of_bb();

    return 0;
}
```

Function `assign_indices_to_var` assigns a unique index to each local variable by traversing `cfun->unexpanded_var_list` as explained in Section 3.2. These indices represent the bit position of a local variable. This requires adding an `integer` field to the `tree` data structure. The variables which are not interesting are assigned index -1.

Function `assign_indices_to_exprs` assigns a unique index to each expression whose operands are restricted to constants and variables that have been assigned a valid index. These indices represent the bit position of relevant expressions. Other expressions are assigned index -1.

Unlike local variables, there is no ready list of expressions. Hence the function `assign_indices_to_exprs` traverses the CFG visiting each statement and examining the expressions appearing in relevant statements. If the expression used in a statement qualifies as a local expression, it is first checked whether an index has already been assigned to it. This could happen because an expression could appear multiple times in a program.

Function `assign_indices_to_defns` assigns a unique index to each statement that is a definition.

Finally, function `dfs_numbering_of_bb` performs depth first numbering of the blocks in a CFG.

3.4 Local Data Flow Analysis

In production compilers, implementing global data flow analyzers is much easier compared to implementing local data flow analyzers. This is because local data flow anal-

ysis has to deal with the lower level intricate details of the intermediate representation and intermediate representation are the most complex data structures in practical compilers. Global data flow analyzers are insulated from these lower level details; they just need to know CFGs in terms of basic blocks. Thus most data flow analysis engines require the local property computation to be implemented by the user of the engine.

This situation can change considerably if we view local data flow analysis as a special case of global data flow analysis. The objective of local data flow analysis is to compute Gen_n and $Kill_n$ of a block n . This computation can be performed by traversing statements in block n in a manner similar to traversing blocks in a CFG. The only difference is that statements in a block cannot have multiple predecessors or successors.

The way In_{Start} (or Out_{End}) is computed by incorporating the effect of blocks in a CFG, Gen_n and $Kill_n$ can also be computed by incorporating the effects of individual statements in block n . The effect of statement s can be defined in terms of Gen_s and $Kill_s$. However, we need to overcome the following conceptual difficulty: When we compute Gen_n for block n , Gen_s of a statement s must be added to the cumulative effect of the statements processed so far. However, when we compute $Kill_n$, $Kill_s$ of statement s should be *added* to the cumulative effect instead of being removed. This deviates from the normal meaning of *Kill* which represents the entities to be removed.

We overcome this conceptual difficulty by renaming Gen_s and $Kill_s$ as Add_s and $Remove_s$ respectively. Now local data flow analysis does not depend on knowing whether the data flow property being computed is Gen_n or $Kill_n$. Given a local property specification such as below:

```
typedef struct lop_specs
{
    entity_name entity;
    entity_manipulation stmt_effect;
    entity_occurrence exposition;
} lp_specs;
```

Local data flow analysis searches for the effect of a given statement specified through `stmt_effect` and stores it in `add_entities`. If the specified `stmt_effect` is `entity_use`, the entities that qualify for `entity_mod` are stored in the variable `remove_entities`. Depending upon the `exposition`, the final decision of removal is taken.

Thus computation of Gen_n and $Kill_n$ depends upon setting up a variable of the type `lp_specs` and the solving the following recurrence

$$\text{accumulated_entities} = (\text{accumulated_entities} - \text{remove_entities}) \cup \text{add_entities}$$

Function `effect_of_a_statement` performs the above computation for a given statement. It is called by the top level function `local_dfa_of_bb`. The relevant code

fragment for downwards exposed entities is:

```
FOR_EACH_STMT_FWD
{
    stmt = bsi_stmt(bsi);
    accumulated_entities = effect_of_a_statement(lps_given,
                                                stmt, accumulated_entities);
}
```

For upwards exposed entities, the accumulation is against the control flow and the above traversal is performed using the macro `FOR_EACH_STMT_BKD`.

The main limitation of this approach is that it requires independent traversal of a basic block for computing *Gen* and *Kill*. However, by using a slightly more complicated data structure that passes both *Gen* and *Kill* to function `local_dfa_of_bb`, will solve this problem. The other limitation is that due to the generality, there are many checks that are done in the underlying functions. There are two possible solutions to this problem of efficiency:

- This is used as a rapid prototyping tool for a given data flow analysis. Once the details are fixed, one could spend time writing a more efficient data flow analyzer.
- Instead of interpreting the specifications, a program can generate a customized C code that is compiled with GCC source.

3.5 Global Data Flow Analysis

As observed earlier, implementation of global data flow analyzer is much simpler once local data flow analysis and interface with the underlying compiler infrastructure is in place. The fact that *gdfa* use generic data flow Equations (1) and (2) makes it possible to execute a wide variety of specifications without having to know the name of a particular analysis being performed. In other words, *gdfa* driver is not a collection of data flow analysis implementations but is capable of executing any specification within the limits of the possible values of specification primitives.

At the top level, the *gdfa* driver needs to perform the following tasks:

- Create special values like \top , BI_{Start} , and BI_{End} .
- Create space for data flow values
- Perform local data flow analysis
- Select flow functions
- Perform global data flow analysis

Function `gdfa_driver` performs the above tasks:

```

0 pfbv_dfi **
1 gdfa_driver(struct gimple_pfbv_dfa_spec dfa_spec)
2 {
3     if (find_entity_size(dfa_spec) == 0)
4         return NULL;
5     initialize_special_values(dfa_spec);
6     create_dfi_space();
7     traversal_order = dfa_spec.traversal_order;
8     confluence = dfa_spec.confluence;
9
10    local_dfa(dfa_spec);
11
12    forward_edge_flow = dfa_spec.forward_edge_flow;
13    backward_edge_flow = dfa_spec.backward_edge_flow;
14    forward_node_flow = dfa_spec.forward_node_flow;
15    backward_node_flow = dfa_spec.backward_node_flow;
16
17    perform_pfbvdfa();
18
19    preserve_dfi(dfa_spec.preserved_dfi);
20    return current_pfbv_dfi;
21 }

```

Lines 12 to 15 select the flow functions from the specifications. Below we show the code fragment of function `perform_pfbvdfa` when the direction of traversal is `FORWARD`.

```

do
{   iteration_number++;
    change = false;
    FOR_EACH_BB_IN_SPECIFIED_TRAVERSAL_ORDER
    {   bb = VARRAY_BB(dfs_ordered_basic_blocks, visit_bb);
        if(bb)
        {   if (traversal_order == FORWARD)
            {   change_at_in = compute_in_info(bb);
                change_at_out = compute_out_info(bb);
                change = change || change_at_out || change_at_in;
            }
            else /* compute in the opposite order */
            {
            }
        }
    }
} while(change);

```

The main code fragment of function `compute_in_info` is as shown below. It calls function `backward_node_flow` which is extracted from the specification.

```

if (!bb->preds)
    temp = combine(entry_info, backward_node_flow(bb));
else
    temp = combine(combined_forward_edge_flow(bb),
                  backward_node_flow(bb));
old = CURRENT_IN(bb);
change = is_new_info(temp,old);

if (change)
{
    CURRENT_IN(bb) = temp;
    if (old)
        free_dfvalue_space(old);
}
return change;

```

Function `combined_forward_edge_flow` computes the following term

$$\prod_{p \in \text{pred}(n)} \vec{f}_{p \rightarrow n}(\text{Out}_p)$$

Its main code fragment is as shown below. It calls function `forward_edge_flow` which is extracted from the specification.

```

edge_vec = bb->preds;
temp = make_initialized_dfvalue(top_value_spec);

if (forward_edge_flow == &stop_flow_along_edge)
    return temp;

FOR_EACH_EDGE(e,ei,edge_vec)
{
    pred_bb = e->src;
    new = combine(temp, forward_edge_flow(pred_bb,bb));
    if (temp)
        free_dfvalue_space(temp);
    temp = new;
}
return temp;

```

The code sequence corresponding to function `compute_out_info` is an exact dual of the above code sequence. This completes the description of generic global data flow analysis in *gdfa*.

4 Extending the Generic Data Flow Analyzer *gdfa*

Many extensions and enhancements of *gdfa* are possible. We suggest some of them by dividing them into the following categories.

- *Extensions that do not require changing the architecture of gdfa.*
 - Include space and time measurement of analyses.
 - Consider scalar formal parameters for analysis.
 - Support a work list based driver.
 - Extend *gdfa* to support other entities such as statements (e.g., for data flow analysis based program slicing), and basic blocks (e.g., for data flow analysis based dominator computation). Both these problems are bit vector problems.
 - Improve the implementation of *gdfa* to make it more space and time efficient. This may require compromising on the simplicity of the implementation but generality should not be compromised.
- *Extensions that may require minor changes to the architecture of gdfa.*
 - Implement incremental data flow analysis and measure its effectiveness by invoking in just before gimple is expanded into RTL.
This would require a variant of a work list based driver.
 - Explore the possibility of extending *gdfa* to the data flow frameworks where data flow information can be represented using bit vectors but the frameworks are not bit vector frameworks because they are non-separable e.g., faint variables analysis, possibly undefined variables, analysis, strongly live variables analysis.
This would require changing the local data flow analysis. One possible option is using matrix based local property computation. The other option is to treat a statement as an independent basic block.
- *Extensions that may require major changes to the architecture of gdfa.*
 - Extend *gdfa* to non-separable frameworks in which data flow information cannot be represented by bit vectors e.g., constant propagation, signs analysis, points-to analysis, alias analysis, heap reference analysis etc. Although the main driver would remain same, this would require making fundamental changes to the architecture.
 - Extend *gdfa* to support some variant of context and flow sensitive interprocedural data flow analysis.

5 Further Readings

Most texts on compilers discuss data flow analysis in varying lengths [1, 2, 3, 7, 8, 11]. Some of them discuss details [1, 2, 8]. A useful introductory chapter is by Khedker [5]. An advanced treatment of data flow analysis can be found in the books by Hecht [4], Muchnick and Jones [9], F. Nielson, H. R. Nielson and Hankin [10], and by Khedker, Sanyal, and Karkare [6].

We list below some useful resources for learning about GCC:

- GCC Internals
<http://gcc.gnu.org/onlinedocs/gccint.html>
This is the official internals document which exhaustively describes most details and is a part of the documentation distributed with the compiler code.
- GCC Internals documents developed at IIT Bombay
<http://www.cse.iitb.ac.in/grc/>
This is the website of *GCC Resource Center* at IIT Bombay. It hosts the GCC documents developed at IIT Bombay.
- The GCC Wiki
<http://gcc.gnu.org/wiki/>
The official GCC Wiki pages where the various aspects of GCC, including some description of the internals, are being developed by the GCC developers and others.
- The GCC Internals workshop held at IIT Bombay
<http://www.cse.iitb.ac.in/~uday/gcc-workshop/>
This workshop that focused mainly on the machine descriptions was held at IIT Bombay in June 2007. The slides and some associated software is available on the Downloads page of the workshop.
- The GCC on Wikipedia
http://en.wikipedia.org/wiki/GNU_Compiler_Collection
- The GCC Internals on Wikipedia
http://en.wikibooks.org/wiki/GNU_C_Compiler_Internals

6 Copyright

gdfa and this manual is a copyrighted material of the GCC Resource Center, Department of Computer Science and Engineering, Indian Institute of Technology Bombay. This material may be distributed only subject to the terms and conditions set forth in

- the GNU GPL v 2.0 (<http://www.gnu.org/licenses/gpl.html>) or later, for the source code of *gdfa*, and

- the GNU FDL v1.2 (<http://www.gnu.org/licenses/fdl.html>) or later, for the documentation.

In particular, the original content of these documents, when used in your work, must be clearly marked as “Copyright ©2008 GCC Resource Center, Department of Computer Science and Engineering, Indian Institute of Technology Bombay”. The documents and the source code has been provided for the sole and exclusive purpose of disseminating information. You are free to download them, but neither GCC Resource Center, nor Indian Institute of Technology Bombay, nor any person related to them are in any way responsible for anything you do with it. In other words, the documents are provided as is. In case you are interested in redistribution or republishing of the gdga source code or it’s manual, whole or in part, either modified or unmodified, and you have questions, please contact the author.

References

- [1] A. V. Aho, M. S. Lam, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2006.
- [2] A. W. Appel and M. Ginsburg. *Modern Compiler Implementation in C*. Cambridge University Press, 1998.
- [3] D. Grune, H. E. Bal, C. J. H. Jacobs, and K. G. Langendoen. *Modern Compiler Design*. John Wiley & Sons, 2000.
- [4] M. S. Hecht. *Flow Analysis of Computer Programs*. Elsevier North-Holland Inc., 1977.
- [5] U. P. Khedker. Data flow analysis. In Y. N. Srikant and Priti Shankar, editors, *The Compiler Design Handbook: Optimizations & Machine Code Generation*. CRC Press, USA, 2002.
- [6] U. P. Khedker, A. Sanyal, and B. Karkare. *Data Flow Analysis: Theory and Practice*. CRC Press (Taylor and Francis Group), 2009. (Under publication).
- [7] R. Morgan. *Building an Optimizing Compiler*. Butterworth-Hienemann, 1998.
- [8] S. S. Muchnick. *Advanced Compiler Design and Implementation*. Morgan Kaufmann Publishing Co., 1997.
- [9] S. S. Muchnick and N. D. Jones. *Program Flow Analysis : Theory and Applications*. Prentice-Hall Inc., 1981.
- [10] F. Nielson, H. R. Nielson, and C. Hankin. *Principles of Program Analysis*. Springer-Verlag, 1998.
- [11] R. Wilhelm and D. Maurer. *Compiler Design*. Addison-Wesley, 1995.