

A Knowledge-Based Approach to the Analysis of Loops¹

Salwa K. Abd-El-Hafiz and Victor R. Basili
Institute for Advanced Computer Studies
Department of Computer Science
University of Maryland
College Park, MD 20742, U.S.A.

March 3, 1994

Abstract

This paper presents a knowledge-based analysis approach which generates first order predicate logic annotations of loops. A classification of loops according to their complexity levels is presented. Based on this taxonomy, variations on the basic analysis approach that best fit each of the different classes are described. In general, mechanical annotation of loops is performed by first decomposing them using data flow analysis. This decomposition encapsulates closely related statements in events, which can be analyzed individually. Specifications of the resulting loop events are then obtained by utilizing patterns, called plans, stored in a knowledge base. Finally, a consistent and rigorous functional abstraction of the whole loop is synthesized from the specifications of its individual events. To test the analysis techniques and to assess their effectiveness, a case study was performed on a pre-existing program of reasonable size. Results concerning the analyzed loops and the plans designed for them are given.

Index terms: First order predicate logic, formal specifications, knowledge base, loops, program understanding, reverse engineering.

¹Research for this study was supported in part by the ONR grant N00014-87-k-0307 to the University of Maryland.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 3 |
| 2 | Definitions | 5 |
| 3 | A loop taxonomy | 6 |
| 4 | Analysis of flat loops | 6 |
| 4.1 | Normalization of the loop representation | 7 |
| 4.2 | Decomposition of the loop body | 7 |
| 4.3 | Formation of the loop events | 9 |
| 4.4 | A knowledge base of plans | 10 |
| 4.5 | Analysis of the events | 15 |
| 5 | Analysis of nested loops | 16 |
| 5.1 | Definitions | 18 |
| 5.2 | Analysis of inner loops and representation of their analysis results | 19 |
| 5.3 | Analysis of outer loops | 22 |
| 5.4 | Adaptation of inner loops specifications | 24 |
| 6 | Discussion | 27 |
| 7 | Case Study | 28 |
| 7.1 | Objectives | 28 |
| 7.2 | Method | 29 |
| 7.3 | Results and analysis | 29 |
| 8 | Conclusion | 33 |

1 Introduction

Program understanding plays an important role in nearly all software related tasks. It is vital to the maintenance and reuse activities which cannot be performed without a deep and correct understanding of the component to be maintained or reused. It is indispensable for improving the quality of software development activities such as code reviews, debugging, and some testing approaches all require programmers to read and understand programs. There has been considerable research on techniques and tools for analyzing and understanding computer programs. Within these efforts, substantial interest is usually directed towards the specific topic of analyzing loops. This interest stems mainly from inherent reasoning difficulties involving repeated program state modifications and the fact that loops have a major effect on program understandability [36].

To analyze loops and reason about their properties, some approaches define heuristics which can be used to guide a search for a loop invariant [19] or function [29]. However, heuristic techniques in general are not always useful. After applying the heuristics a considerable number of times, one may or may not succeed in finding a correct invariant or function. Other approaches focus on developing algorithmic techniques for finding the invariants or functions of specific simple classes of loops. The research performed by Basu and Misra [8], Dunlop and Basili [12], Katz and Manna [25], and Wegbreit [42] is representative of these loop analysis approaches. The advantage of these approaches is that they analyze loops through the use of formal, semantically sound, and unambiguous notation. Although some of them provide guidelines on how to mechanically generate loop invariants or functions, they were not actually used to implement automatic analysis systems. A different approach, which analyzes loops by mechanically decomposing them into smaller fragments, is adopted by Waters [41]. Even though Waters' approach does not address the issue of how to use this decomposition to mechanically annotate loops, it is especially interesting because of its practicality.

To analyze complete programs, the knowledge-based approaches utilize a knowledge base of expert-designed plans in providing intelligent analysis results. They are inspired by the cognitive studies [28, 37, 35] which suggest that the understanding process is a process in which programmers make use of stereotyped solutions to problems in making sophisticated high-level decisions about a program. These knowledge-based approaches are all implemented, to varying degrees, in automatic analysis systems. Some of these approaches are: graph-parsing [32, 43]; top-down analysis using the program's goals as input [23, 22, 30]; heuristic-based object-oriented recognition [15, 16]; transformation of a program into a semantically equivalent but more abstract form with the help of plans and transformation rules [27, 40]; and decomposition of a program into smaller more tractable parts using proper decomposition [17] or program slicing [18]. Even though these approaches demonstrate the feasibility and usefulness of the automation of program understanding, they lack some important features.

Most of the knowledge-based program analysis and understanding approaches produce program documentation which is, more or less, in the form of structured natural language text [9, 15, 16, 17, 23, 31, 27, 32, 43]. Such informal documentation gives expressive and intuitive descriptions of the code. However, there is no semantic basis that makes it possible to determine whether or not the documentation has the desired meaning. This lack of a firm semantic basis makes informal natural language documentation inherently ambiguous.

Some of the knowledge-based approaches rely on user-supplied information which might not

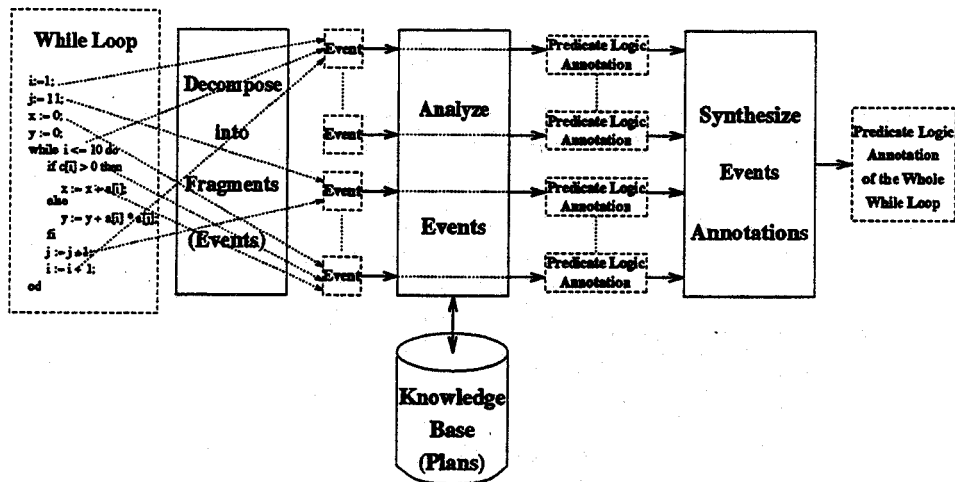


Figure 1: Overview of the analysis approach

be available at all times. For instance, the goals a program is supposed to achieve [23] or the transformation rules that are appropriate for analyzing a specific code fragment [40] are not always clear to the user. Others have difficulty in analyzing non-adjacent program statements [27]. In addition, a significant amount of program analysis and understanding research has used toy programs to validate proposed approaches. Realistic evaluations of these approaches, which give quantifiable results about recognizable and unrecognizable concepts in real and pre-existing programs, are needed. Such evaluations can also serve as a basis for empirical studies and future comparisons with other approaches [34].

To address the above mentioned drawbacks, we present a knowledge-based approach to the automation of program analysis. It combines and builds on the strengths of a practical program decomposition method [41], the axiomatic correctness notation [19], and the knowledge-based analysis approaches. It mechanically documents programs by generating first order predicate logic annotations of their loops. The advantages of predicate logic annotations are that they are unambiguous and have a sound mathematical basis. This allows correctness conditions to be stated and verified, if desired. Another advantage is that they can be used in assisting the formal development of software using such languages as VDM and Z [24, 39].

A family of analysis techniques has been developed and tailored to cover different levels of program complexity. This complexity is determined by classifying while loops along three dimensions. The first dimension focuses on the control computation part of the loop. The other two dimensions focus on the complexity of the loop condition and body. Based on this taxonomy, the analysis techniques which can be applied to the different loop classes are described.

In general, we annotate loops with predicate logic assertions in a step by step process as depicted in Figure 1. The analysis of a loop starts by decomposing it into fragments, called *events*. Each event encapsulates the loop parts which are closely related, with respect to data flow, and separates them from the rest of the loop. The resulting events are then analyzed, using plans stored in a knowledge base, to deduce their individual predicate logic annotations. Finally, the annotation of the whole loop is synthesized from the annotations of its events [2].

This study tests several hypotheses related to the presented analysis approach:

- A loop classification class is an indicator of its amenability to analysis.
- The loop decomposition and plan design methods can make the plans applicable in many different loops and, hence, increase their utilization.
- The analysis techniques can be automated.

To test the first two hypotheses and to characterize the practical limits of the analysis approach, a case study on a set of 77 loops in a pre-existing program for scheduling university courses has been performed. The approach was found to be effective. The program has 1400 executable lines of code and the loops analyzed have the usual programming language features such as pointers, procedure and function calls, and nested loops. To test the third hypothesis, a prototype tool, which annotates loops with predicate logic annotations, was developed.

Section 2 of this paper gives some of the definitions used. Section 3 introduces the loop taxonomy. Sections 4 and 5 describe the techniques used for analyzing flat and nested loops, respectively. Section 6 discusses the approach presented and highlights its advantages and limitations. Section 7 describes how the case study was performed and gives the results of the analysis. Finally, conclusions and future research directions are given in Section 8.

2 Definitions

We start by defining some of the notation used throughout this paper. First, we give the definitions related to the representation of while loops.

A *control-flow graph* is a directed graph that has one node for each simple statement and one node for each control predicate. There is an edge from node I to node J if an execution of J can immediately follow that for I [20].

Let the *abstract representation of the while loop* be *while B do S* where the condition *B* has no side effects and the statements *S* are representable by a single-entry single-exit control-flow graph. This representation abstracts from the syntax of the specific imperative programming language being used. Though the approach described here applies to all loops having this abstract representation, examples and illustrations are given using Pascal. Using this abstract representation, a *control variable* of the while loop is a variable that exists in the condition *B* and is modified in the body *S*.

Now, we give some definitions which introduce the language and terminology used in the analysis. A *concurrent assignment* is a statement in which several variables can be assigned simultaneously. We use the form $v_1, v_2, \dots, v_n := e_1, e_2, \dots, e_n$ to assign every *i*th expression from the right hand list to its corresponding *i*th variable from the left hand list [14, 29]. A *conditional assignment* is a set of one or more guarded concurrent assignments separated by commas ','. When the guard (i.e., the boolean expression), of a concurrent assignment is satisfied, the modifications performed on a variable are given by the concurrent assignment [14, 29].

Any variable assigned in a conditional assignment defines the *data flow out* of the statement. Any variable referenced by a conditional assignment defines the *data flow into* the statement. Two conditional assignments are said to be *circularly dependent* if some variable is responsible for data flow out of one statement and into the other, either directly or indirectly, and vice versa.

3 A loop taxonomy

To design the analysis techniques that best fit different levels of program complexity, we classify while loops along three dimensions. The first dimension focuses on the control computation part of the loop. The other two dimensions focus on the complexity of the loop condition and body. Along each dimension, a loop must belong to one of two complementary classes as shown in Table 1. In this classification, the loops in the middle column are expected to be more amenable to analysis than the corresponding ones in the right column.

| Dimension | Complementary classes | |
|----------------------------|------------------------|---------------------|
| 1. Control computation | Simple loop | General loop |
| 2. Complexity of condition | Noncomposite condition | Composite condition |
| 3. Complexity of body | Flat loop | Nested loop |

Table 1: The three dimensions used for classifying loops

Within the first dimension, we differentiate between simple and general loops. *Simple loops* have a behavior similar to that of *for* loops. They are defined by imposing two restrictions: the loop has a unique control variable, and the modification of the control variable does not depend on the values of other variables modified within the loop body. Loops which do not satisfy these conditions are called *general loops*.

Along the second dimension, the complexity of the loop condition can vary between two cases. In the *noncomposite* case, B is a logical expression that consists of one clause of the conjunctive normal form [33]. In the *composite* case, more than one clause exists. Along the third dimension, the complexity of the loop body varies between *flat* and *nested* loop structures. In flat loop structures, the loop body can not contain any other loop inside it which is not the case in nested structures.

4 Analysis of flat loops

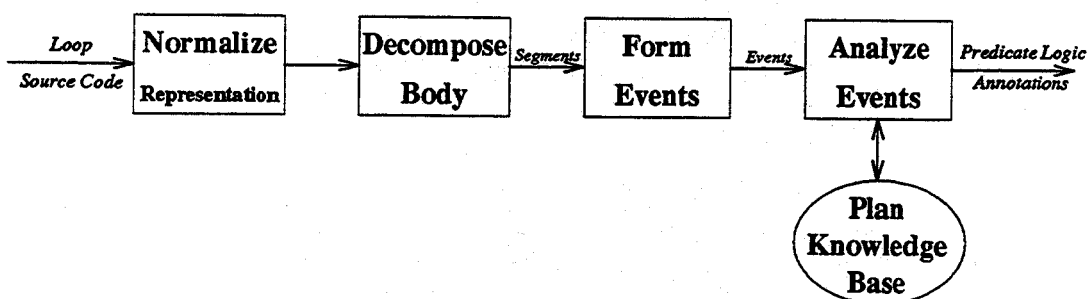


Figure 2: Analysis of flat loops

As depicted in Figure 2, the analysis of flat loops is performed in a step by step process divided into four main phases. Descriptions of these phases and their application to the example shown in Figure 3 are given in the remainder of this section [3]. In this example, a simple loop with a noncomposite condition scans a segment of the array *capacity* searching for its minimum.

```

j, index, min, num_of_rooms : integer;
capacity : array[1..max_rooms] of integer;
:
while j < num_of_rooms + 1 do begin
  if capacity[j] < min then begin
    index := j;
    min := capacity[j];
  end;
  j := j + 1
end;

```

Figure 3: Example of a simple loop

4.1 Normalization of the loop representation

The purpose of this phase is to make the loop representation independent of the programming language and the implementation specific details.

Normalization of the loop condition: The loop condition is converted into an arbitrarily chosen normal form, which is the *conjunctive normal form* [33]. This normal form converts a well-formed formula (wff) in predicate logic into a conjunction of clauses where a clause is defined to be a wff in conjunctive normal form but with no instances of the *and* connector. For example, the loop condition $x < a$ or $(y < b$ and $z < c)$ is transformed to the conjunction of two clauses $(x < a$ or $y < b)$ and $(x < a$ or $z < c)$.

Normalization of the loop body: A single unwinding of the loop body is performed by symbolic execution [1, 5] which gives the net modification performed on each variable in one iteration of the loop, if any [7]. We use the conditional assignment notation to represent the result of this symbolic execution.

For the loop given in Figure 3, the condition is already in conjunctive normal form containing the one clause $j < num_of_rooms + 1$. The symbolic execution does not affect the body of the loop. However, the net modification performed on each variable is given in the form of a conditional assignment as follows:

| Name | Conditional assignment |
|-------|--|
| C_1 | $capacity[j] < min \implies index := j,$ |
| C_2 | $capacity[j] < min \implies min := capacity[j],$ |
| C_3 | $true \implies j := j + 1$ |

4.2 Decomposition of the loop body

To facilitate the mechanical generation of loop annotations, the symbolic execution result is uniquely decomposed into *segments* of code which can be analyzed separately. Each segment encapsulates the statements which are interdependent with respect to data flow. The *loop segments* are partitions of the loop body symbolic execution result. Each segment consists of a maximal set of conditional assignments such that any two conditional assignments in the set are circularly dependent.

To obtain the loop segments, we assume that the conditional assignments of the symbolic execution result correspond to the nodes of a directed graph. An edge from node C_i to node C_k exists if and only if there is data flowing out of C_i into C_k and C_i and C_k are distinct. The strongly connected components of this graph correspond to the loop segments [10].

For the loop shown in Figure 3, the three conditional assignments of the symbolic execution result form a directed graph G with three edges: two from C_3 to C_1 and C_2 , and one from C_2 to C_1 . Since there are no cycles in G , its strongly connected components correspond to its nodes. Thus, the loop segments correspond to C_1 , C_2 , and C_3 .

Because the analysis of a segment might be dependent on the analysis results of other segments, a segment analysis result should be obtained before analyzing the segments dependent on it. That is why we need to order the segments according to their data flow dependencies [41]. Assuming the S is the set of segments in the loop body, the order of each segment is determined by the following algorithm:

1. Set m to 1.
2. While the number of segments in S is ≥ 1 do
 - (a) Identify the maximal subset of S such that every segment does not have data flowing out of it and into other segments of S .
 - (b) Let the order of the identified segments be m .
 - (c) Remove the identified segments from S .
 - (d) increment m .
3. Let the final order of each segment be (m - old order). \square

Step 2 of the above algorithm assigns unique orders to the segments such that order of $S_i >$ order of S_k if and only if there is data flowing, either directly or indirectly, from segment S_i to segment S_k . Step 3 produces an irreflexive partial order of the segments. The resulting ordering relation 'analyzed before', is denoted by ' \rightarrow '. It is irreflexive because it is meaningless for a segment to be analyzed before itself. It satisfies the antisymmetric property because any two distinct segments, by definition, have no circular dependencies. The design of the above algorithm ensures the satisfaction of the transitive property. Moreover, it is possible for two segments to be unrelated (i.e., they can have the same order).

In the example given in Figure 3, let the segments of the loop be S_1 , S_2 and S_3 which correspond to C_1 , C_2 , and C_3 , respectively. The orders assigned to these segments using the above algorithm are:

| Order | Name | Segment |
|-------|-------|--|
| 1 | S_3 | $j := j + 1$ |
| 2 | S_2 | $capacity[j] < min \Rightarrow min := capacity[j]$ |
| 3 | S_1 | $capacity[j] < min \Rightarrow index := j$ |

Notice that the segment that defines j , S_3 , has the lowest order because the other two segments, S_1 and S_2 , reference j (i.e., $S_3 \rightarrow S_1$ and $S_3 \rightarrow S_2$). Similarly, $S_2 \rightarrow S_1$ because min is defined in S_2 and referenced in S_1 . Since the premise of the conditional assignment that modifies j is *true*, it is removed.

4.3 Formation of the loop events

To represent the abstract concepts in a loop, we use the loop body segments and the clauses of the loop condition to form the loop *events*. We define two categories of loop events: *basic events* and *augmentation events*.

Basic Events (BE's) are the fragments that constitute the control computation of the loop. A BE consists of three parts: the *condition*, the *enumeration*, and the *initialization*. The *condition* is one clause of the loop condition. The *enumeration* is a segment responsible for the data flow into the *condition*. The *initialization* is the initialization of the variables defined in the *enumeration*.

To form BE's, each clause of the loop condition is combined with the highest order segment(s) having data flow into it. If a clause has no segment responsible for the data flow into it, this means that this clause is redundant and should be removed from the loop condition. If a segment is responsible for the data flow into the condition but remains with no clause associated with it, its condition is set to *true*. The initializations of the control variables defined in a BE are included in the initialization part.

The BE of the loop given in Figure 3 is formed by combining the unique condition clause, ($j < num_of_rooms + 1$), with the only segment which is responsible for the data flow into it, S_3 . Since the loop under consideration has no initializations, we use the notation $j?$ to denote the initial value of a variable j . As a result, the BE has the following form:

condition: $j < num_of_rooms + 1$
enumeration: $j := j + 1$
initialization: $j := j?$

An *Augmentation Event (AE)* consists of two parts: the *body* and the *initialization*. The *body* is one segment of the loop body which is not responsible for the data flow into the loop condition. The *initialization* is the initialization of the variables defined in the *body*.

After identifying the BE's, The AE's bodies are the remaining segments of the loop. To complete the formation of the AE's, the initialization of each variable defined in an event is included in it.

For the loop shown in Figure 3, the remaining segments S_2 and S_1 constitute the bodies of two AE's given below. The notation $min?$ and $index?$ are used to denote the initial values of the variables min and $index$.

| | |
|--|--|
| <p>AE 1</p> <p>body: $capacity[j] < min \implies min := capacity[j]$ initialization: $min := min?$</p> | <p>AE 2</p> <p>body: $capacity[j] < min \implies index := j$ initialization: $index := index?$</p> |
|--|--|

Finally, we give each event (basic or augmentation) the same order as the segment it utilizes. This enforces the condition that the variables referenced in an event are either defined in a lower order event or not modified within the loop at all. As mentioned in the previous subsection, this makes it possible to propagate the results of analyzing an event to the analysis of other events dependent on it.

The three events of the loop shown in Figure 3 are thus ordered as follows:

- | | |
|---|---|
| <p>1. BE (order 1)</p> <p>condition: $j < num_of_rooms + 1$ enumeration: $j := j + 1$ initialization: $j := j?$</p> | <p>2. AE 1 (order 2)</p> <p>body: $capacity[j] < min \implies min := capacity[j]$ initialization: $min := min?$</p> |
|---|---|

3. AE 2 (order 3)

body: $capacity[j] < min \implies index := j$

initialization: $index := index?$

4.4 A knowledge base of plans

To analyze the loop events, we utilize plans stored in a knowledge base. The term 'plan' refers to a unit of knowledge required to identify an abstract concept in a program. Our plans are used as inference rules [15, 16]. Their basic structure is divided into two parts: the antecedent and the consequent. When a loop event matches a plan antecedent, the plan is fired. The instantiation of the information in the consequent represents the contribution of this plan to the loop specifications.

The Knowledge base is designed so that any two plans do not have similar antecedents. Thus, a loop event can only match the antecedent of a unique plan. It should also be noted that the possibility of designing as many plans as the number of loop events in a specific program is reduced because the loop events encapsulate abstract concepts which can occur in different loops. Section 7 will examine this issue of the knowledge base size in more detail.

Corresponding to the two event categories, we have two plan categories: *Basic Plans (BP's)* and *Augmentation Plans (AP's)*. BP's analyze BE's and AP's analyze AE's. Plans are further classified according to the kind of loops they analyze.

In case of simple loops, the sequences of values scanned by the control variable during and after the execution of a simple loop can be easily written because the control computation is isolated from the rest of the loop. The loop condition, the control variable's initial value, and the net modification performed on the control variable in one loop iteration, if any, provide sufficient information for writing these sequences. This specific information about the control computation of the loop can be used to produce equally specific loop specifications. The plans that analyze simple loops can include these sequences and utilize them in writing the loop specifications.

The analysis of *general loops* is not as straightforward as that of simple ones. In many cases, it might not be easy, or even possible, to obtain such specific knowledge because the control computation of the loop is not as determinate and isolated as in the case of simple loops. The sequences of values scanned by the control variable(s) and the program state at the end of the loop are usually dependent on the combined indeterminate effects of several events and the values of some program variables. As a result, the plans that analyze general loops neither include the aforementioned sequences nor utilize them in writing the loop specifications. The loop postcondition can only be deduced after the synthesis of the loop invariant. The postcondition is formed by taking the conjunction of the loop invariant with the negation of the loop condition [19]. Using this condition to obtain the loop postcondition yields predicates which might not be as informative and concise as those of simple loops. As a result, additional simplifications might be needed to reduce the complexity and improve the readability of general loops postconditions.

For instance, consider the simple loop shown in Figure 3. The sequence scanned by the control variable at any point during the loop execution is $j?$ to $j-1$. This sequence is needed to write the part of the invariant: $min = MIN(\{min?\} \cup \{capacity[j?..j-1]\})$, where $MIN(s)$ is the minimum of the set s and \cup is the set union operator. The final sequence of values scanned by the control variable in this loop is $j?$ to num_of_rooms . This sequence is needed to write the part of the postcondition: $min = MIN(\{min?\} \cup \{capacity[j?..num_of_rooms]\})$.

```

while (j < num_of_rooms + 1) and (flag = false) do begin
  if capacity[j] < limit then begin
    index := j;
    flag := true;
  end;
  j := j + 1
end

```

Figure 4: Example of a general loop

In the general loop given in Figure 4, however, there is no guarantee that the final sequence scanned by the control variable *j* will be *j?* to *num_of_rooms*. The value of the final sequence is dependent on the interaction of the two events which modify *flag* and *j*, and the contents of the variables *capacity* and *limit*. As a result of this generality of the control computation, the sequences of values scanned by the control variable(s) and, consequently, the postcondition parts of the individual events cannot be written.

To accommodate the differences between simple and general loops, we have two categories of BP's. Determinate BP's (DBP's) contain in their consequents information regarding the postcondition and the sequences of values scanned by the control variable. Indeterminate BP's (IBP's), on the other hand, do not contain such information. We also have two categories of AP's. Simple AP's (SAP's) utilize the above sequences in writing the loop specifications, including its postcondition. General AP's (GAP's) do not include the loop postcondition part or utilize the above sequences. These plan categories are shown in Figure 5. It should also be noticed that if we neglect the information regarding the control sequences and the postcondition, DBP's can be used in analyzing general loops. However, the reverse is not true because DBP's are more specific than IBP's.

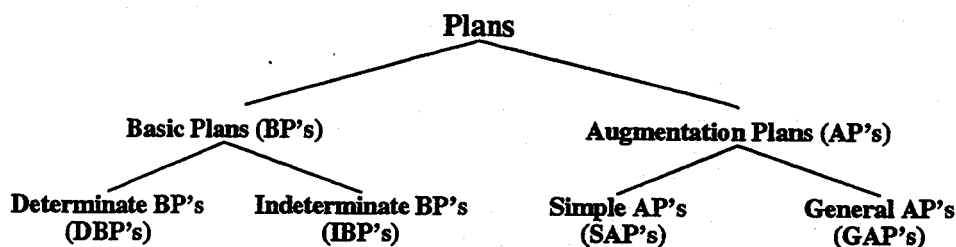


Figure 5: Plan categories

In general, The information included in a plans antecedent and consequent are described below. In this description, the words printed in **bold** correspond to fields in the plans (see Figures 6 and 7).

An antecedent contains the following information:

1. An individual listing of the control variables, in the **control-variables** part, which serves to underscore their importance and to facilitate the design, readability, and comprehension of the plan.
2. Abstractions of BE's and AE's which are used to recognize stereotyped loop events.

| | |
|--------------------------|--|
| plan-name | DBP ₁ (ascending enumeration) |
| antecedent | |
| control-variables | <i>var#</i> |
| condition | <i>var# R# exp#</i> |
| enumeration | <i>var# := SUCC(var#)</i> |
| initialization | <i>var# := var?#</i> |
| firing-condition | (<i>R#</i> equals \leq or $<$) and (<i>var#</i> is of a discrete ordinal type) and (<i>B</i> noncomposite or general loop) |
| consequent | |
| precondition | <i>PRED(var?#) R# exp#</i> |
| invariant | <i>var?# \leq var# R# SUCC(exp#)</i> |
| postcondition | <i>var# = SUCC(SHIFT(exp#))</i> |
| sequence | <i>var?# .. PRED(var#)</i> |
| final-sequence | <i>var?# .. SHIFT(exp#)</i> |
| where, | |
| <i>i..j</i> | Sequence of integers from <i>i</i> up to <i>j</i> inclusive. |
| <i>SUCC(x)</i> | The successor of <i>x</i> . |
| <i>PRED(x)</i> | The predecessor of <i>x</i> . |
| <i>SHIFT</i> | The identity function if <i>R#</i> equals \leq . Equals PRED otherwise. |

Figure 6: A Determinate Basic Plan

3. Knowledge needed for the correct identification of the plans such as data types' information and the previous analysis knowledge of a variable. This knowledge is given in the **firing-condition**

A consequent includes the following information:

1. Knowledge necessary for the annotation of loops with their Hoare-style [19] specifications. The **precondition** and **invariant** have the usual meaning [19]. The **postcondition** part gives information, in case of simple loops, about the variables' values after the loop execution ends. It is correct provided that the loop executes at least once. If the loop does not execute, no variable gets modified.
2. In case of DBP's, knowledge about the sequence of values scanned by the control variables at any point during and after the loop execution is captured in **sequence** and **final-sequence**, respectively.

Figures 6 and 7 show two example plans of the categories DBP and SAP, respectively. To convey the basic analysis ideas within a reasonable space limit, we only show simplified versions of the plans. The suffix '# ' is used to indicate terms in the antecedent (or consequent) that must be matched (or instantiated) with actual values in the loop events.

| | |
|--------------------------|--|
| plan-name | SAP ₅ (find minimum) |
| antecedent | |
| control-variables | $v\#$ |
| body | $a\#[exp\#] R\# lhs\# \implies lhs\# := a\#[exp\#]$ |
| initialization | $lhs := lhs?\#$ |
| firing-condition | $(R\# \text{ equals } \leq \text{ or } <)$ |
| consequent | |
| precondition | <i>true</i> |
| invariant | $lhs\# = MIN(\{lhs?\#\} \cup \{a\#[exp\#]_{sequence}^{v\#}\})$ |
| postcondition | $lhs\# = MIN(\{lhs?\#\} \cup \{a\#[exp\#]_{final-sequence}^{v\#}\})$ |
| where, | |
| $MIN(s)$ | The minimum of the set s . |

Figure 7: A Simple Augmentation Plan

The plan DBP₁ (Figure 6) represents an enumeration construct that goes over a sequence of values of a discrete ordinal type in an ascending order with a unit step. In the case where the loop has a composite condition, the **sequence**, **final-sequence** and **postcondition** of this plan are written in a more general form that enables deducing the corresponding **sequence**, **final-sequence** and **postcondition** of the loop from the multiple BE's it contains. The plan SAP₅ (Figure 7) searches for the minimum of a segment of the array $a\#$ and stores it in the variable $lhs\#$.

The knowledge base in a specific application domain should be created by an expert in both formal specifications and this domain. The expert should analyze the commonly used events in this domain and create new plans or improve on already existing ones. In creating this knowledge base, its size should be controlled by increasing the utilization of the designed plans. The loop decomposition method was designed for this purpose; to reveal the common algorithmic constructs that can be incorporated in many different loops. The hypothesis is that this decomposition can have a positive effect on plan utilization and, hence, on the size of the knowledge base. Improvements on the structure and/or the knowledge represented in the plans can also make the plans applicable to a larger set of events.

Knowledge representation improvements, called *abstractions*, involve replacing some of the terms in a plan with more abstract ones that make the plan capable of analyzing more cases. For example, replacing the addition operator, $+$, in a plan which analyzes an accumulation by summation event by a more abstract one which denotes either addition or multiplication represents an abstraction of this plan. The new plan can analyze both accumulation by summation and accumulation by multiplication events.

Structural improvements to a plan modify the basic structure into a tree structure which allows the inclusion of several similar plans in one tree-structured plan. The root of the tree corresponds to an antecedent part that should match loop events. The edges of the tree correspond to **firing-conditions** which control the selection of the appropriate loop annotations given in the remaining tree nodes. In other words, a tree-structured plan consists of a single antecedent and several

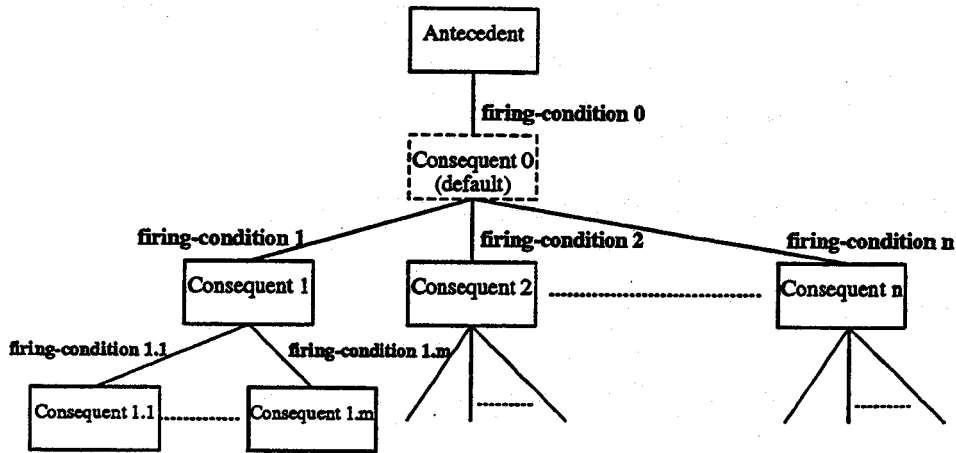


Figure 8: The tree structure of a plan

consequents organized into one or more tree structures as shown in Figure 8. The consequents are organized into one tree if the default consequent exists. Otherwise, they are organized into more than one tree (forest). In order to select a specific tree-structured plan, a match with the antecedent should occur first. Then, **firing-condition 0** must be satisfied. Within the plan, local **firing-conditions** of the consequents guide the search for the suitable consequent. The more general the consequent, the closer it is to the root of its tree (e.g., consequent 1 of Figure 8 is more general than consequent 1.1). Consequents located at the same level have mutually exclusive **firing-conditions**. This means that only forward search is needed and no backtracking is required. When the event matches the antecedent and **firing-condition 0** of the tree-structured plan is satisfied, the search for the appropriate consequent starts at the appropriate root going down in the tree as far as possible. The path between a parent and a child can only be taken if the local **firing-condition** associated with the child consequent is satisfied.

Tree-structured plans can be used to detect special cases and output loop specifications that are simple and concise. They can also be used to analyze similar events whose specifications vary depending on their environment (e.g., data types, control computation of the loop, ..., etc).

For instance, the plan SAP_5 (Figure 7) can be structurally improved as shown in Figure 9. The antecedent is similar to that shown in Figure 7 except for the firing condition. **Firing-condition 0** allows $R\#$ to be matched with more relational operators. Three local **firing-conditions** and their consequents cover three different variations. Consequent 1, which is similar to the consequent of the basic plan in Figure 7, is for finding the minimum. Consequent 1.1 further simplifies the resulting annotations based on special values of $lhs\#$ and the analysis information of the control variable $v\#$. Consequent 2 is for finding the maximum.

Using the tree-structured plans can lead to a reduction in the size of the knowledge base since several plans can be combined together into a larger one having a unique antecedent. However, the instantiation of the proper consequent becomes more complicated.

The construction of the tree-structured plans, especially in case of large knowledge bases, can be facilitated by the design of automated techniques which assist in their acquisition and development. For instance, automatically identifying similar plans and combining them into more sophisticated

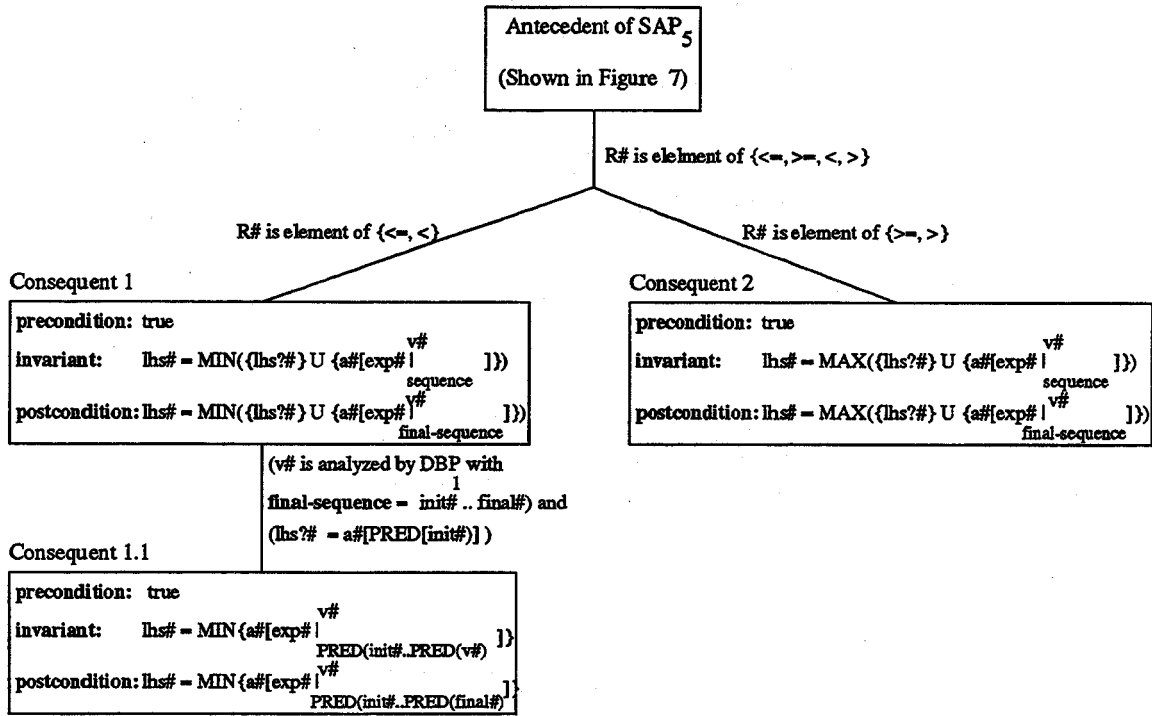


Figure 9: Structural improvement to the plan SAP₅

tree structures is an interesting topic for future study.

4.5 Analysis of the events

The events are analyzed by trying to match them with the antecedents of the knowledge base plans. When an event matches the antecedent of a plan and **firing-condition** 0 is satisfied, the appropriate consequent of the matched plan is instantiated giving the contribution of the event to the loop specification. The precondition, invariant and postcondition of the loop are formed by taking the conjunction of the corresponding parts of the event analysis results. When some event(s) do not match any library plans, we get partial specifications of the loop.

To represent the results of matching the loop events with the plans antecedents, we define the *Analysis Knowledge* notation. The *Analysis Knowledge*, $AK(v)$, of a variable v modified by a certain loop event consists of an n -tuple where n is dependent on the specific matched plan. The first term of the tuple is the name of the matched plan. The remaining $(n-1)$ terms are the results of matching the $\#$ terms with the actual values in the event.

The resulting AK tuples and instantiations for the events of the loop given in Figure 3 are shown below. The event and plan responsible for the production of each predicate are shown to its left. The first two events are matched with the plans DBP_1 (Figure 6) and SAP_5 (Figure 7), respectively. The plan, SAP_{n1} , which analyzes the third event is not shown here because it is similar to the plan SAP_5 . It searches for the location of the minimum instead of the minimum.

$$AK(j) = (DBP_1, var\#: j, var?\#: j?, R\#: <, exp\#: num_of_rooms + 1)$$

$AK(min) = (SAP_5, v\# : j, a\# : capacity, exp\# : j, lhs\# : min, lhs?\# : min?)$

$AK(index) = (SAP_{n1}, v\# : j, a\# : capacity, exp\# : j, rhs\# : min, rhs?\# : min?, lhs\# : index, lhs?\# : index?)$

Precondition:

| Event | Plan | Predicate |
|-------|-------------------|-------------------------------|
| 1 | DBP ₁ | $j? - 1 < num_of_rooms + 1$ |
| 2 | SAP ₅ | <i>true</i> |
| 3 | SAP _{n1} | $capacity[index?] = min?$ |

Invariant:

| Event | Plan | Predicate |
|-------|-------------------|--|
| 1 | DBP ₁ | $j? \leq j < num_of_rooms + 2$ |
| 2 | SAP ₅ | $min = MIN(\{min?\} \cup \{capacity[j?..j - 1]\})$ |
| 3 | SAP _{n1} | $capacity[index] = min$ |

Postcondition:

| Event | Plan | Predicate |
|-------|-------------------|---|
| 1 | DBP ₁ | $j = num_of_rooms + 1$ |
| 2 | SAP ₅ | $min = MIN(\{min?\} \cup \{capacity[j?..num_of_rooms]\})$ |
| 3 | SAP _{n1} | $capacity[index] = min$ |

The synthesized specifications of the inner loop are:

Precondition:

$(j? - 1 < num_of_rooms + 1)$ and
 $(capacity[index?] = min?)$

Invariant:

$(j? \leq j < num_of_rooms + 2)$ and
 $(min = MIN(\{min?\} \cup \{capacity[j?..j - 1]\}))$ and
 $(capacity[index] = min)$

Postcondition:

$(j = num_of_rooms + 1)$ and
 $(min = MIN(\{min?\} \cup \{capacity[j?..num_of_rooms]\}))$ and
 $(capacity[index] = min)$

5 Analysis of nested loops

To rigorously analyze nested loops, the following problems need to be solved:

1. **How to represent and utilize the analysis results of inner loops?** A technique for analyzing flat loops has been described in Section 4. Can the same basic technique be used for outer loops (loops containing other loops)? What modifications, if any, need to be performed on the basic analysis technique to analyze outer loops?

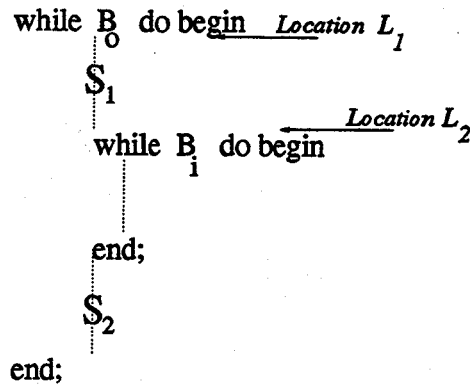


Figure 10: A nested structure of while loops

2. How to modify the resulting specifications to facilitate Hoare-style verification [19]? This problem can be further divided into two subproblems, which are explained using the nested construct shown in Figure 10. In this nested construct, let I_i and I_o be the invariants of the inner and outer loops, respectively.

- (a) Can the above invariants be used to satisfy Hoare verification conditions which connect the specifications of inner and outer loops in the nested construct? In other words, is it possible to prove the following rules:

$$(I_i \text{ and } \neg B_i) \quad \{S_2\} \quad I_o \quad (1)$$

$$(I_o \text{ and } B_o) \quad \{S_1\} \quad I_i \quad (2)$$

In these rules, the notation $P\{S\}Q$ means that if the predicate P is true before executing the first statement of the program part S , and if S terminates, then the predicate Q will be true after execution of S is complete.

- (b) If the above invariants use the notation $var?$ to denote the initial value of a variable var , does this notation consistently refer to the value of var before the start of the outermost loop in the nested construct? If not, how can this inconsistency be removed?

To solve these problems, the analysis of nested loops is performed by recursively analyzing the innermost loops and replacing them with sequential constructs which represent their functional abstraction. The functional abstraction of an outer loop depends on the functional abstraction of the inner ones and not on the details of their implementation or structure.

Since this recursive analysis approach is performed bottom-up, complete knowledge of the inner loops functions is available during the analysis of an outer loop. Thus, the invariant of an outer loop can be directly designed to satisfy the verification rules which are similar to rule (2) listed above. However, inner loops are analyzed in isolation of the outer ones enclosing them. As a result, their invariants and, consequently, postconditions might not be strong enough to satisfy the verification rules which are similar to rule (1). Some predicates might need to be added to the inner loops invariants and postconditions to enable the verification of such rules. The *context adaptation* phase derives these predicates and adds them to the inner loops specifications. Moreover, the

```

i, j, index, min, num_of_rooms : integer;
capacity : array[1..max_rooms] of integer;
:
i := 1;
while i ≤ num_of_rooms - 1 do begin
    index := i;
    min := capacity[i];
    i := i + 1;
    j := i;
    while j < num_of_rooms + 1 do begin
        if capacity[j] < min then begin
            index := j;
            min := capacity[j];
        end;
        j := j + 1
    end;
    capacity[index] := capacity[i - 1];
    capacity[i - 1] := min
end;

```

Figure 11: Example of a nested loop

consistency of using the notation *var*? to denote the initial value of a variable *var* is ensured using the *initialization adaptation* phase.

We start in Section 5.1 with some definitions which explain how we extract the initialization of a loop in a nested construct, whether it is the outermost loop or an inner one. Sections 5.2-5.4 present solutions to the two research problems mentioned above. Sections 5.2 and 5.3 offer a solution to the first research problem. Section 5.4 presents a partial solution to the second problem. In these sections, the descriptions of the analysis steps are interspersed with their application on the selection sorting example given in Figure 11. In this example, a simple nested loop repeatedly scans an array segment searching for its minimum. It interchanges the minimum with the first element in the segment. It stops after the array *capacity*[1..*num_of_rooms*] has been sorted in ascending order. The inner loop of this example is the same one given in Figure 3.

5.1 Definitions

In the following definitions, we limit the initialization of a loop to assignment statements which serve to make the resulting loop specifications both informative and readable. Conditional statements are not considered as initializations to limit the complexity of the resulting loop specifications and make them representative of the loop functionality without composing them with specifications of external compound statements.

The *initialization of a loop that is not enclosed by another loop* is assumed to be a set of assignment statements of the form *identifier* := *expression*, which are immediately placed before its start. These statements give initial values for identifiers that get modified within the loop body.

If this assumption cannot be satisfied or, equivalently, the loop initialization is unavailable, the notation $v?$ is used to denote the initial value of a variable v just before the start of the loop.

If we have two nested while loops, the *adaptation path* of the inner loop is a sequence of statements extracted from their control-flow graph representation. This sequence contains all the statements, simple or compound, which are completely located on the paths starting from the control predicate node of the outer loop and ending at the control predicate node of the inner loop exclusive. In addition, the relative order of the statements is kept unchanged.

The *initialization of an inner loop* in a nested construct is obtained by, first, symbolically executing its adaptation path to produce the net modification performed on each variable, if possible. Statements of the form *identifier := expression* are, then, extracted from the symbolic execution result. Statements are extracted if they satisfy the following two conditions:

1. The *identifier* is one of the variables modified within the inner loop body.
2. The *expression* does not reference any of the variables modified along the adaptation path.

If the initialization of a variable v that gets modified within the loop body is not given by the extracted statements, the notation $v?$ is used to denote its initial value just before the start of the loop.

The first condition, in the above definition, ensures that the initialization statements are utilized by the inner loop events. The second condition ensures that the values of *identifier* and *expression*, just before the start of the inner loop, are equal. For example, if the adaptation path is $i := i+1; j := i$, then its symbolic execution gives the concurrent assignment $i, j := i+1, i+1$. Taking $j := i+1$ as an initialization statement is not allowed because the values of j and $i+1$, just before the start of the loop, are not equal (the values of j and i are equal). The second condition also prevents using statements of the form, say, $x := x + 1$ as initializations.

To extract the initialization of the inner loop given in Figure 11, we use the above definitions. The adaptation path is:

$index := i; min := capacity[i]; i := i + 1; j := i.$

The symbolic execution of the adaptation path yields the following concurrent assignment:

$index, min, i, j := i, capacity[i], i + 1, i + 1.$

According to the definition of the initialization of an inner loop in a nested construct, all the resulting assignments are unacceptable initializations because the second condition is not satisfied. In addition, the assignment to i does not satisfy the first condition. As a result, the initialization is written by using the notation $v?$ to denote the initial value of a variable v .

5.2 Analysis of inner loops and representation of their analysis results

The analysis of inner loops is performed using the same four phases described, in Section 4, for flat loops. To analyze an outer loop in a nested construct, the analysis results of its inner loops must be represented in a way that reveals the functionality of the inner loops and the flow of data into and out of the inner loops. The data flow information is needed to perform the decomposition of the outer loop body.

Though the resulting AK tuples or predicate logic annotations can be used to represent the inner loop analysis results, they either include too much detail or the deduction of the required information is difficult, respectively. Hence, the solution is to use a formalism that is similar to function calls; the name encapsulates the functionality while the arguments indicate the data flow information. The formalism used for this purpose is called an *Abstraction Class (AC)*.

An AC is a knowledge base object which transforms the detailed analysis results of an inner loop to a more abstract representation that facilitates the analysis of outer loops. It groups AK tuples based on some common functionality and ignores the unnecessary implementation specific details. The common functionality is documented to explain the purpose of designing the AC and to enhance its modifiability. Furthermore, the definition of an AC offers an abstract representation of its elements that specifies the data flow information to facilitate their mechanical manipulation. An *Abstraction Class (AC)* consists of three parts:

1. The **elements** part consists of generic AK tuples which are separated by the symbol '|’.
2. The **Common-function** describes the functionality that the elements of this class share by using common instantiated **final-sequence**, **postcondition**, or **invariant** parts of the matched plans.
3. The **representation** is a unique abstract representation that gives the class name, followed by the following arguments (separated by semicolons and enclosed between two parentheses): the list of expressions responsible for the data flow into this AC, the list of variables defined by the AC, the control variables of the loop under consideration, and a unique number identifying the loop being analyzed. □

The representation part contains the class name which is an arbitrary and unique name. It also contains the arguments responsible for the data flow into and out of the AC so that they can be used during the data flow analysis. The control variables and unique number of the loop are used in the design of some plans’ consequents. To simplify the presentation, the last 2 arguments are only listed when needed.

The AK of some variable belongs to a specific AC if it matches an AK tuple existing in the **elements** part. The symbol ‘*’ is used to denote irrelevant information. An expression, *exp*, enclosed between two brackets in the **elements** part implies that the expression should be matched with the corresponding instantiated element of the actual AK to deduce the value of the variables defined in it. Some of these variables are utilized in forming the AC arguments.

The AK of the variable *j* analyzed in the inner loop of Figure 11 has the following form:

$$AK(j) = (DBP_1, var\# : j, var?\# : j?, R\# : <, exp\# : num_of_rooms + 1).$$

This AK belongs to the AC in Figure 12 because it matches the first AK tuple of the **elements** part. If we had implemented this loop with the condition $j \leq num_of_rooms$ instead of $j < num_of_rooms + 1$, it would have belonged to the same AC. This is because it matches the second AK tuple of the **elements** part. These two different implementations belong to the same AC because they have the common function of going over a sequence of values of a discrete ordinal type in an ascending order.

Using similar analysis, the AK of the variable *min* is found to belong to AC_{SL_5} (Figure 13). The AK of the variable *index* belongs to $AC_{SL_{n_1}}$. Because $AC_{SL_{n_1}}$ is similar to AC_{SL_5} , it is not shown

| | |
|------------------------|--|
| elements | $(DBP_1, var\# : [v], var?\# : *, R\# : \leq, exp\# : [final])$ $(DBP_1, var\# : [v], var?\# : *, R\# : <, exp\# : [SUCC(final)])$ |
| common-function | The instantiated final-sequence of the plan is: $v? .. final$ |
| representation | $AC_{SB_1}(v, final; v)$ |

Figure 12: An abstraction class for ascending enumeration

| | |
|------------------------|--|
| elements | $(SAP_5, v\# : [v], a\# : [a], exp\# : [PRED(v)], lhs\# : [lhs], lhs?\# : *)$, where $AK(v) \in AC_{SB_2}([SUCC(final)], [SUCC(init)])$ $(SAP_5, v\# : [v], a\# : [a], exp\# : [v], lhs\# : [lhs], lhs?\# : *)$, where $AK(v) \in AC_{SB_2}([final], [init])$ or $AK(v) \in AC_{SB_1}([init], [final])$ $(SAP_5, v\# : [v], a\# : [a], exp\# : [SUCC(v)], lhs\# : [lhs], lhs?\# : *)$, where $AK(v) \in AC_{SB_1}([PRED(init)], [PRED(final)])$ |
| common-function | The instantiated postcondition of the plan is: $lhs = MIN(\{lhs?\} \cup \{a[init..final]\})$ |
| representation | $AC_{SL_5}(a, init, final, lhs; lhs)$ |

Figure 13: An abstraction class for finding the minimum

here. AC_{SL_5} includes the AK tuples which have the common function of finding the minimum of an array segment irrespective of the enumeration direction (ascending or descending) and the index of the array element being checked (v , $PRED(v)$, or $SUCC(v)$). It should be mentioned that AC_{SB_2} is similar to AC_{SB_1} but for descending enumeration. The AC's of the variables modified in the inner loop of Figure 11 are, thus, as follows:

- $AK(j) \in AC_{SB_1}(j, num_of_rooms; j)$
- $AK(min) \in AC_{SL_5}(capacity, j, num_of_rooms, min; min)$
- $AK(index) \in AC_{SL_{n1}}(capacity, j, num_of_rooms, min, index; index)$

After analyzing an inner loop, it is replaced with the concurrent assignment that assigns to the list of variables modified by it the result of their analysis. If the AK of a variable belongs to a predefined AC, its abstract representation, as deduced from the identified AC, is assigned to it. If the AK of the variable, var , does not belong to a predefined AC, we assign the form $UAC(ak-$

list; var) to it, where *UAC* stands for Unknown AC and *ak-list* is a list representing the AK data. The *ak-list* and *var* are used, during automatic analysis, to provide information on the unanalyzed parts of the loop.

Conceptually, the described replacement is equivalent to replacing the inner loop with a set of function calls that assign to each variable changed in the inner loop the desired value. This replacement preserves the control flow dependencies because the concurrent assignment is placed at the same relative location within the outer loop body. It also preserves the data flow dependencies between the variables because the AC's clearly state what variables are responsible for the data flow into and out of it.

Replacing the inner loops given in Figure 11 with the described concurrent assignment gives the following modified outer loop:

```

i := 1;
while i < num_of_rooms - 1 do begin
  index := i;
  min := capacity[i];
  i := i + 1;
  j := i;
  j, min, index := ACSB1(j, num_of_rooms; j),
                  ACSL5(capacity, j, num_of_rooms, min; min),
                  ACSLn1(capacity, j, num_of_rooms, min, index; index);
  capacity[index] := capacity[i - 1];
  capacity[i - 1] := min
end;

```

5.3 Analysis of outer loops

After modifying an outer loop body, we analyze it using the previously described method for analyzing flat loops (Section 4), as if it does not contain any other loops inside it. This can be done since the inner loop(s) have been replaced by ordinary sequential constructs. The only difference, in this case, is that high-level plans are used in addition to the usual (low-level) ones. High-level plans are those which utilize AC's.

Adding another classification level, based on whether the plan is low-level or high-level, to the four plan categories shown in Figure 5, we get 8 plan categories. These new plan categories are shown in Figure 14.

The strength of this approach for analyzing nested constructs is that it can scale up to handle more than two nested loops. This is because the inner loops can be recursively analyzed and replaced by sequential constructs. Any outer loop can thus be analyzed by using the high-level plans in addition to the low-level ones. An outline of the application of the analysis steps on the modified outer loop of Figure 11 is given below.

The ordered events of the modified outer loop are as follows:

1. BE (order 1)
 - condition: $i \leq \text{num_of_rooms} - 1$
 - enumeration: $i := i + 1$

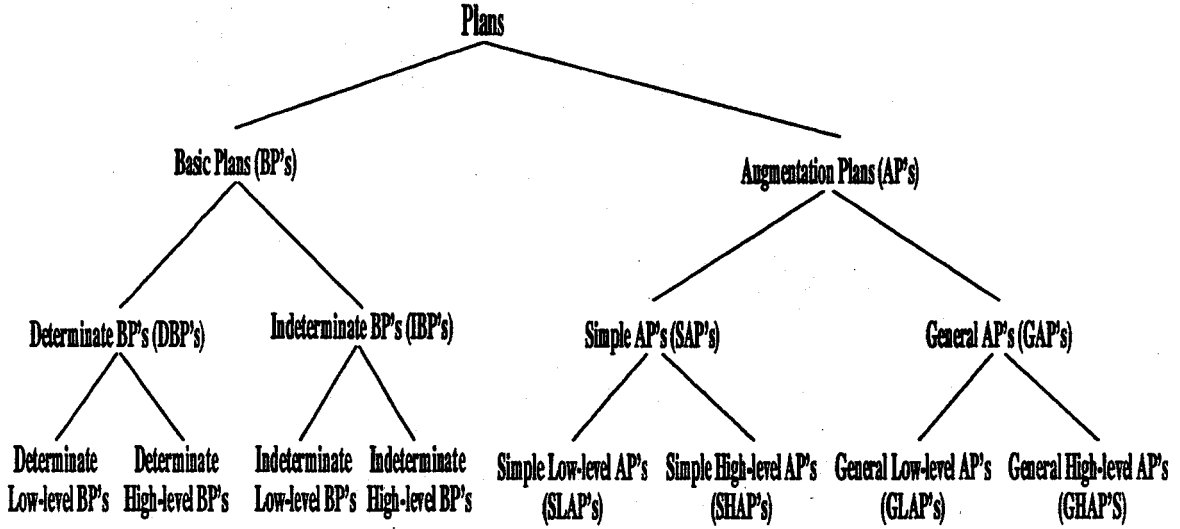


Figure 14: New plan categories

initialization: $i := 1$

2. AE (order 2)

body: $capacity[i], capacity[AC_{SL_{n_1}}(capacity, i + 1, num_of_rooms, capacity[i], i, index)]$
 $:= AC_{SL_s}(capacity, i + 1, num_of_rooms, capacity[i]; min), capacity[i]$

initialization: $capacity := capacity?$

3. AE (order 2)

body: $j := AC_{SB_1}(i + 1, num_of_rooms; j);$

initialization: $j := j?$

4. AE (order 3)

body: $min := AC_{SL_s}(capacity, i + 1, num_of_rooms, capacity[i]; min)$

initialization: $min := min?$

5. AE (order 3)

body: $index := AC_{SL_{n_1}}(capacity, i + 1, num_of_rooms, capacity[i], i, index)$

initialization: $index := index?$

The first event is matched with the antecedent of the plan DBP_1 (Figure 6). The second event is matched with a Simple High-level AP (SHAP) which represents the selection sorting concept. Because the variables j , min and $index$ do not explicitly contribute to the outer loop specifications, the last 3 events are matched with SHAP's which produce *true* predicates. These variables implicitly affect the outer loop specifications through their abstraction classes which get used by the second event. For details concerning the plans used and the event analysis results, refer to [2]. The final synthesized analysis results are given below. The first event is responsible for the production of the first conjugate of each predicate. The second event is responsible for the production of the rest of the specifications.

$(0 \leq num_of_rooms)$

Invariant:

$(1 \leq i \leq \text{num_of_rooms})$ and
 $(\text{FORALL } ind : 1 \leq ind \leq i - 1 : \text{capacity}[ind] = \text{MIN}(\{\text{capacity}[ind..\text{num_of_rooms}]\})$ and
 $\text{PERM}(\text{capacity}, \text{capacity}?)$

Postcondition:

$(i = \text{num_of_rooms} + 1)$ and
 $(\text{FORALL } ind : 1 \leq ind \leq \text{num_of_rooms} - 1 : \text{capacity}[ind] = \text{MIN}(\{\text{capacity}[ind..\text{num_of_rooms}]\})$ and
 $\text{PERM}(\text{capacity}, \text{capacity}?)$

The resulting predicate logic annotations produced for the inner and outer loops can be used to assist the understanding of the nested construct. An understanding of the sorting algorithm can be formed using the predicate

$(min = \text{MIN}(\{min?\} \cup \{\text{capacity}[j?..\text{num_of_rooms}]\}))$ and $(\text{capacity}[index] = min)$

of the inner loop postcondition and the predicate

$(\text{FORALL } ind : 1 \leq ind \leq \text{num_of_rooms} - 1 : \text{capacity}[ind] = \text{MIN}(\{\text{capacity}[ind..\text{num_of_rooms}]\})$ and
 $\text{PERM}(\text{capacity}, \text{capacity}?)$

of the outer loop postcondition. However, such specifications cannot be proved using Hoare-style

[19] axiomatic correctness. To be able to prove the outer loop invariant, the predicate:

$(1 \leq i - 1 \leq \text{num_of_rooms})$ and
 $(\text{FORALL } ind : 1 \leq ind \leq i - 2 : \text{capacity}[ind] = \text{MIN}(\{\text{capacity}[ind..\text{num_of_rooms}]\})$ and
 $\text{PERM}(\text{capacity}, \text{capacity}?)$

should be added to the invariant of the inner loop. This predicate provides information about the context of the inner loop, which is needed to prove rule (1) at the beginning of this section. In addition, $j?$, $min?$, and $index?$ in the inner loop specifications should be replaced with i , $\text{capacity}[i - 1]$, and $i - 1$, respectively.

5.4 Adaptation of inner loops specifications

To be able to prove the specifications of nested constructs, the specifications of the inner loops need to be strengthened to include information about the context of outer loops enclosing them. To ensure that the notation $var?$ is consistently used to denote the initial value of a variable var before the start of the outermost loop, variables of the form $var?$ in specifications of inner loops need to be replaced by their actual values. These tasks are performed in the context and initialization adaptation phases. The remainder of this subsection describes how to perform these adaptations. In this description, it is assumed that the adaptation path of an inner loop only includes assignment and conditional statements. The cases in which this assumption is not satisfied are discussed in Section 6.

Context adaptation: While analyzing the outer loop, we have complete knowledge of an inner loop function. Thus, this is the best time to generate the context related predicate *inner-addition*, which strengthens inner loops invariants. By studying the differences between the current outer loop invariant part and the generated inner loop invariant part, we design and add an *inner-addition* field to the consequents of the knowledge base plans. This field provides any predicates

that should be added to the invariants of inner loops to enable the verification of rules similar to: $(I_i \text{ and } \neg B_i) \{S_2\} I_o$. The predicate *inner-addition* is synthesized, by conjunction, from the *inner-additions* defined in the consequents of the outer loop plans.

For instance, assume that the plan DBP_1 (Figure 6) is used to analyze an ascending enumeration construct of an outer loop having the control variable $var\#$. While analyzing the inner loop in isolation, no knowledge exists about $var\#$ being an outer loop control variable which scans a specific sequence of values. Hence, the *inner-addition* filed of DBP_1 should provide this information in the form of the predicate: $var?\# \leq var\# R\# exp\#$.

Analyzing the BE of the outer loop given in Figure 11 using DBP_1 yields the following instantiated *inner-addition*:

$(1 \leq i \leq num_of_rooms)$.

Similarly, when the inner loop of this sorting example is analyzed in isolation, its invariant does not include any information about the sorted segment of the array *capacity*. Thus, the *inner-addition* part of the outer loop selection sorting plan should provide the following predicate:

$(FORALL\ ind : 1 \leq ind \leq i - 1 : capacity[ind] = MIN(\{capacity[ind..num_of_rooms]\}) \text{ and } PERM(capacity, capacity?))$.

By taking the conjunction of the two instantiated *inner-addition* parts, the *inner-addition* of the example given in Figure 11 is:

$(1 \leq i \leq num_of_rooms) \text{ and}$

$(FORALL\ ind : 1 \leq ind \leq i - 1 : capacity[ind] = MIN(\{capacity[ind..num_of_rooms]\}) \text{ and } PERM(capacity, capacity?))$

However, the synthesized *inner-addition* is designed to be correct at a fixed reference point which is location L_1 (see Figure 10). This is because during the design of the library plans there is no knowledge, a priori, of the statements physically located along the adaptation path. The effect of the statements along the adaptation path should be taken into consideration to get the corresponding correct predicate, *inner-addition₂*, at location L_2 .

By comparing the *inner-addition* produce for the loop given in Figure 11 to the predicate which should be added to the inner loop specifications (given at the end of Section 5.3), it is clear that they are not exactly the same. This is because the effect of the statements along the adaptation path has not been taken into consideration yet.

The context adaptation uses *inner-addition* and the adaptation path to find *inner-addition₂*. The predicate *inner-addition₂* is deduced by reversing the effect of the statements along the adaptation path on the variables in *inner-addition*. For example, if the adaptation path changes i to $i - 1$, then all the free occurrences of i in *inner-addition* are replaced by $i + 1$ to generate *inner-addition₂*.

This reversing is performed, mechanically, by introducing a set of auxiliary variables which replace all the free occurrences, in *inner-addition*, of the variables modified along the adaptation path. Conceptually, the auxiliary variables denote the state of the corresponding original ones at location L_1 .

For the example shown in Figure 11, the auxiliary variable i_1 replaces the variable i in *inner-addition*. Since the variable *capacity* is not modified along the adaptation path, no corresponding auxiliary variable is introduced for it. The modified *inner-addition*, which is called *inner-addition₁*, has the form:

$(1 \leq i_1 \leq \text{num_of_rooms})$ and
 $(\text{FORALL } ind : 1 \leq ind \leq i_1 - 1 : \text{capacity}[ind] = \text{MIN}(\{\text{capacity}[ind..\text{num_of_rooms}]\})$ and
 $\text{PERM}(\text{capacity}, \text{capacity?})$

The relation between the auxiliary variables, used at location L_1 , and the corresponding original ones, used at location L_2 , is represented in the predicate *aux-values*. The predicate *aux-values* is formed using the symbolic execution result of the adaptation path. First, the introduced auxiliary variables should replace their corresponding actual ones which are responsible for the data flow into the symbolic execution result. The predicate equivalent of the statements which modify the original variables are, then, generated and conjunctioned together.

The predicate equivalent of an assignment statement is produced by replacing the assignment sign with an equal sign. Conditional assignments can be converted into assignment statements of the form: $var := \text{choice}(\text{condition1}, \text{value1}, \text{condition2}, \text{value2}, \dots, \text{etc})$, where the right hand side is equal to *value1* if *condition1* is *true*, *value2* if *condition2* is *true*, and so on. The resulting assignment statement is converted into a predicate as described before.

In the example shown in Figure 11, the symbolic execution result of the adaptation path is:

$index, min, i, j := i, \text{capacity}[i], i + 1, i + 1.$

The context adaptation replaces i by i_1 in the right hand side to produce:

$index, min, i, j := i_1, \text{capacity}[i_1], i_1 + 1, i_1 + 1.$

The statement which modifies the original variable i is $i := i_1 + 1$. The predicate equivalent of this statement, $i = i_1 + 1$, is the predicate *aux-values*.

The required correct predicate *inner-addition₂* is the conjunction of *aux-values* and *inner-addition₁*. The predicate *inner-addition₂*, which is actually added to the inner loop invariant, has the form:

$(1 \leq i_1 \leq \text{num_of_rooms})$ and
 $(\text{FORALL } ind : 1 \leq ind \leq i_1 - 1 : \text{capacity}[ind] = \text{MIN}(\{\text{capacity}[ind..\text{num_of_rooms}]\})$ and
 $\text{PERM}(\text{capacity}, \text{capacity?})$ and
 $i = i_1 + 1$

Initialization adaptation: The initialization adaptation replaces each variable of the form *var?*, in an inner loop specification, with its value as deduced from its adaptation path and the invariant of the enclosing loop. After this replacement, the notation *var?* is reserved for referring to the state of a variable *var* before the start of the outermost loop. The notation var_{outer} is used to refer the value of *var* as deduced from the invariant of the loop enclosing it.

The initial value of a variable *var* is extracted from the symbolic execution result of the adaptation path. If the symbolic execution result assigns the value var_{adapt} to *var*, then var_{adapt} is the needed initial value. However, var_{adapt} needs to be modified so that it is expressed in terms of the program state at location L_2 and not location L_1 . This modification is performed in the same way we modified *inner-addition*. That is, var_{adapt} is modified to $(var_{adapt})_{\text{auxiliary variables}}^{\text{original variables}}$. However, if *var* itself occurs in var_{adapt} , it should, first, be replaced by var_{outer} to avoid a circular definition of the initial value of *var*. In short, every *var?* in the inner loop specification is replaced by $((var_{adapt})_{var_{outer}}^{\text{original variables}})_{\text{auxiliary variables}}$.

For instance, the symbolic execution result of the adaptation path of the example shown in

Figure 11 is:

$index, min, i, j := i, capacity[i], i + 1, i + 1.$

The variable $j?$ in the inner loop specification is replaced by $((i + 1)_{j_{outer}}^j)_{i_1}^i$, where $i = i_1 + 1$. So, $j?$ is effectively replaced by i . Similar analysis shows that $min?$ and $index?$ should be replaced with $capacity[i - 1]$ and $i - 1$, respectively.

In summary, the specification of the inner loop shown in Figure 11 is adapted by adding the predicate $inner-addition_2$ which is simplified to:

$(1 \leq i - 1 \leq num_of_rooms)$ and
 $(FORALL\ ind : 1 \leq ind \leq i - 2 : capacity[ind] = MIN(\{capacity[ind..num_of_rooms]\})$ and
 $PERM(capacity, capacity?)$.

The initial variables $j?$, $min?$, and $index?$ are replaced with i , $capacity[i - 1]$, and $i - 1$, respectively. These adaptation results are exactly the ones described at the end of Section 5.3.

6 Discussion

In this paper, a knowledge-based program understanding approach has been described. The resulting predicate logic annotations are unambiguous and have a sound mathematical basis which allows correctness conditions to be stated and verified, if desired. The analysis approach does not rely on user-supplied information and can analyze non-adjacent loop parts.

However, there are limitations to this approach. These limitations are:

- Practical limitations which are related to the effort and ingenuity needed for the design of the plans.
- Theoretical limitation which is related to the ability to produce concise postconditions for general loops.
- Theoretical limitation which is related to the ability to perform the context and initialization adaptations for some nested loops.

The practical limits are based on the plan designers ability to formally analyze complicated loops and find their invariants. The resulting specifications are as accurate, readable, and correct as the plans are. That is why the tasks of designing plans and managing the knowledge base, for a specific application domain of interest, should be performed by an expert in both the desired domain and formal specifications.

To develop the knowledge base, the desired application domain should be analyzed to design an initial set of plans which is believed to cover a considerable number of loop constructs that might occur in it. After adding this initial set, the knowledge base should evolve over time. It should undergo a process of controlled usage where the knowledge base manager closely monitors its utilization.

The first theoretical limit was discussed in Section 4.4. In case of general loops, we cannot produce loop postconditions as intelligently and concisely as for simple loops because it was not possible to include postcondition parts in the plans designed for analyzing individual events of

general loops. This, in turn, can require additional simplifications of the postconditions which transforms them into more readable ones. The second theoretical limit occurs in nested structures having the following characteristic: the adaptation path of an inner loop contains statements other than assignment and conditional statements (e.g., loops or procedure calls). The context and initialization adaptation cannot, in general, be performed for such cases. These limits only affect the ability to prove the resulting specifications. They do not affect the ability to assist the understanding of nested loops. This is because the approach still produces meaningful specifications of the whole construct. For instance, it has been shown that an understanding of the sorting algorithm in our example was possible before performing the adaptation steps.

The reason for the second limitation is that the context and initialization adaptations are based on the fact that assignment statements and, to a lesser extent, conditional statements can be easily inverted in a mechanical way [14]. However, if there are loops, procedure calls, or function calls, this inversion cannot be performed mechanically. Performing such an inversion is equivalent to finding the specifications of arbitrary program fragments containing nonsequential constructs and representing their analysis results in terms of equational specifications that can be easily inverted. The presented approach can perform symbolic execution of sequential constructs and can produce first order predicate logic specifications of loops. However, these two different capabilities have not been integrated to produce invertible equational specifications of arbitrary program fragments. The context and initialization adaptation can be performed in some special cases. One special case occurs when the variables used in the *inner-addition* do not get modified along the adaptation path. Another special case happens when variables, whose initial values need to be replaced, do not get modified along the adaptation path. In the first case, the context adaptation does not need to modify the predicate *inner-addition*. In the second case, the initialization adaptation directly replaces *var?*, if any, with its value as deduced from the outer loop invariant. A third special case occurs when the loops located on the adaptation path are simple ones. In this case, the adaptation of an inner loop specification can be performed using postcondition parts of its preceding loop, which are in equational form, instead of its outer loop invariant. It should be noted that the first theoretical limit partly affects the second one. If we were able to include equational postconditions in the plans that analyze general loops, they could have been used in the adaptation steps

7 Case Study

The program chosen as a case study of our loop analysis process deals with scheduling a set of university courses. It has about 1400 lines of executable Pascal source code. There are a total of 39 modules (functions and procedures). A complete listing of the requirements, specifications, design, and code documents is given elsewhere [21]. In this program, there are 77 loops that cover all the classes in our taxonomy.

7.1 Objectives

The main objective of this case study was to test our analysis approach and to assess its effectiveness when applied to a fixed set of loops in a real and pre-existing program of some practical value. To this effect, we collected the data needed for performing the following validations and characterizations:

- Test the hypothesis that a loop classification class is an indicator of its amenability to analysis.
- Test the hypothesis that the loop decomposition and plan design methods of our approach can make the plans applicable in many different loops and, hence, increase their utilization.
- Characterize the practical limits of the analysis approach.

7.2 Method

The case study was performed, manually, prior to the implementation of the prototype tool. Case study results are, thus, not affected by the limits of the implementation. The set of 77 loops in the described program were extracted along with their initializations. This set included 25 *for* loops, which were transformed to their equivalent while loops. The loops analyzed had the usual programming language features such as pointers, procedure and function calls, and nested loops.

During the study, every loop under consideration was first decomposed into its basic and augmentation events. Then, every event was analyzed in order to design a plan suitable for it. If no plan was available in the knowledge base to match the event under consideration, or a similar event, a new plan was developed with designer defined, candidate specifications. The plan was then modified and tailored to give correct specifications by trying to prove the loop invariant using Hoare techniques [19]. If a plan that matched a similar event, but not the exact one under consideration, existed in the knowledge base, improvements on the structure and/or the knowledge represented in the existing plan were considered.

As the number of analyzed loops increased, the experience gained led to the evolution of the knowledge base. The controlled utilization of the knowledge base served to adapt some of the plans in terms of their structure, knowledge content, number, and naming conventions. As a result, the knowledge base was more suitable for the domain under consideration.

The designed plans were not only limited to those which provided functional specifications but also included plans which discarded unnecessary detail about temporary variables and plans which provided warning and error messages. It should also be mentioned that the resulting specifications were not formulated in terms of concepts specific to the application domain. Even though such specifications can increase the chance of reusing the plans, they sometimes have the disadvantage of being more difficult to read.

We decided not to specifically design plans for the analysis of 12 loops (15.6%) in the case study. The major reason for this was the complexity and specificity of these loops that made the reuse of their plans unlikely. These loops are arbitrarily numbered from p1 through p12. They were analyzed using the available set of plans to determine whether or not useful partial specifications could be obtained.

7.3 Results and analysis

Tables 2 and 3 give the data that was collected to examine the relationship between a loop classification class and its amenability to analysis. Table 2 gives the number of loops completely analyzed in each class defined by our taxonomy. Along the first dimension, the available and analyzed numbers of Simple (S) and General (G) loops are given. In the second dimension, the

| Analysis statistics | Dimension | | | | | |
|---------------------|---------------|-------------|---------------|---------------|---------------|---------------|
| | 1 | | 2 | | 3 | |
| | S | G | N | C | F | N |
| Available loops | 52 | 25 | 46 | 31 | 53 | 24 |
| Loops analyzed | 48 (92.3%) | 17 (68%) | 42 (91.3%) | 23 (74.2%) | 52 (98.1%) | 13 (54.2%) |

Table 2: Number of completely analyzed loops along the three dimensions

| Analysis statistics | Equivalence class | | | | | | | |
|---------------------|-------------------|-------------|---------------|-----|-----|---------------|-----|--------------|
| | SNF | SCF | SNN | SCN | GNF | GCF | GNN | GCN |
| Available loops | 31 | 6 | 15 | 0 | 0 | 16 | 0 | 9 |
| Loops analyzed | 31 (100%) | 6 (100%) | 11 (73.3%) | - | - | 15 (93.8%) | - | 2 (22.2%) |
| Number of events | 75 | 18 | 61 | - | - | 51 | - | 8 |
| Average events/loop | 2.42 | 3 | 5.55 | - | - | 3.4 | - | 4 |

Table 3: Number of completely analyzed loops in the available classes

available and analyzed numbers of loops with Noncomposite (N) and Composite (C) conditions are given. Finally, the available and analyzed numbers of Flat (F) and Nested (N) loops are given along the third dimension.

Using the three classification dimensions, any loop must belong to one of the 8 (2^3) equivalence classes given in Table 3. In this table, the available and analyzed numbers of loops in each of these equivalence classes are shown. The table also gives the total numbers of events and their averages for the analyzed loops in each class.

The results given in Tables 2 and 3 support the hypothesis that the classification taxonomy helps in predicting a loop amenability to analysis. Table 2 shows that the presumably more complex classes always have lower percentages of completely analyzed loops than the presumably less complex ones. For example, the percentages of completely analyzed flat and nested loops are 98.1 and 54.2, respectively. All flat loops were completely analyzed except for one loop (loop 10p) which contained a call to a procedure with a partially analyzed nested loop (loop 9p). This percentage variation is even more notable when further investigated along the five available equivalence classes of Table 3. Percentages range from 100% for SNF and SCF to 22.2% for GCN. The numbers of events in the analyzed loops further support the interpretation that the classification of a loop is an indicator of its complexity and, correspondingly, its amenability to analysis. For example, while SNF loops (Flat) have an average of 2.42 events/loop, SNN loops (Nested) have an average of 5.55 events/loop.

Table 4 summarizes the data that was collected to examine the plan utilization issue. It shows the number of events analyzed by each of the designed plans. It also shows the total utilization of the plans in each of the six available categories. Since only one high-level basic plan (IBP₆) was designed, we do not differentiate between low and high-level BP's. The * or + superscript is used to denote plans which underwent structural or knowledge representation improvements, respectively, during the iterative process of their design. For example, plan DBP₁ was used 45 times and had a tree-structured design.

| Plan name (subscriber) | Plan category | | | | | |
|---------------------------|---------------|-----|------|------|------|------|
| | DBP | IBP | SLAP | GLAP | SHAP | GHAP |
| 1 | 45* | 4 | 23** | 4 | 3+ | 3 |
| 2 | 6* | 15* | 19+ | 13** | 13** | 2 |
| 3 | 8 | 1 | 3* | 1 | 1 | — |
| 4 | 9* | 2 | 1 | 2 | 1 | — |
| 5 | 1 | 2 | 1 | 1 | 1 | — |
| 6 | — | 2 | 1 | 1 | 1 | — |
| 7 | — | — | 1 | 1 | 2 | — |
| 8 | — | — | 20 | 1 | 2 | — |
| 9 | — | — | 3 | — | 2 | — |
| 10 | — | — | — | — | 1 | — |
| 11 | — | — | — | — | 2 | — |
| 12 | — | — | — | — | 1 | — |
| 13 | — | — | — | — | 1 | — |
| 14 | — | — | — | — | 1 | — |
| 15 | — | — | — | — | 2 | — |
| 16 | — | — | — | — | 1 | — |
| 17 | — | — | — | — | 3 | — |
| 18 | — | — | — | — | 1 | — |
| Total | 69 | 26 | 72 | 24 | 39 | 5 |

Table 4: Utilization of the designed plans

The 48 plans designed were utilized in analyzing a total 235 events. A closer examination of the results in Table 4 shows that a set of 27 plans (56.3%) analyzed 214 events (91.1%). The remaining 21 plans were only used once. These results indicate that if we focus on a specific application domain, there is bound to be a kernel of events which can be captured by a relatively reasonable number of plans. On the other hand, there will also be plans which, as in our study, may be used just once. The emphasis should be on the design of the plans that cover the kernel.

The 10 plans which underwent improvements to their structure and knowledge representation (20.8%) analyzed 149 events (63.4%). The average number of utilizations of the plans vary from 4.9 (with standard deviation of 7.97) for all 48 plans to 14.9 (with standard deviation of 11.8) for the 10 improved plans which are marked with the * and + superscripts. These numbers support the argument that commonly used plans get more chances to be revised and adapted and this, in turn, leads to their higher utilization.

We also notice, from Table 4, that even though 9 SLAP's analyzed 71 events, double the number of SHAP's (18) only analyzed 39 events. This indicates that simple 'low-level' blocks of code are more frequently utilized than the more complex 'high-level' ones.

In general, the results in Table 4 show that the events/plan ratio is high (4.9), especially in case of the plans that underwent structural and knowledge representation improvements (14.9). This indicates that the decomposition and plan design methods tend to have a positive effect on plan utilization and, consequently, on the size of the knowledge base. However, since our main objective was to validate and evaluate the analysis approach, we designed many plans (21) which

| Loop # | Characteristics | | | | | | | | |
|--------|-----------------|--------|-----------------|--------------------|-------------|-------------------|-----------------|----------------|-------------|
| | Class | Events | Executable SLOC | Modified variables | | Pointer variables | Procedure calls | Function calls | Inner loops |
| | | | | control | non-control | | | | |
| p1 | GCN | 13 | 48 | 3 | 10 | 0 | 5 | 4 | 1 |
| p2 | GCN | 9 | 30 | 3 | 6 | 0 | 2 | 2 | 1 |
| p3 | GCN | 13 | 46 | 3 | 10 | 0 | 3 | 2 | 2 |
| p4 | GCN | 9 | 32 | 3 | 6 | 0 | 2 | 2 | 1 |
| p5 | GCN | 13 | 49 | 3 | 10 | 0 | 4 | 2 | 2 |
| p6 | SNN | 17 | 53 | 1 | 16 | 2 | 7 | 2 | 1 |
| p7 | SNN | 20 | 53 | 1 | 19 | 4 | 0 | 1 | 1 |
| p8 | GCN | 8 | 36 | 2 | 7 | 2 | 4 | 0 | 1 |
| p9 | SNN | 5 | 29 | 1 | 4 | 3 | 0 | 0 | 2 |
| p10 | GCF | 5 | 13 | 2 | 4 | 0 | 1 | 2 | 0 |
| p11 | GCN | 12 | 52 | 3 | 11 | 4 | 1 | 4 | 2 |
| p12 | SNN | 19 | 77 | 1 | 20 | 4 | 1 | 4 | 3 |

Table 5: Characteristics of the 12 partially analyzed loops

| Characteristics (in terms of average numbers) | Completely analyzed loops | Partially analyzed loops |
|--|---------------------------|--------------------------|
| Events | 3.28 (SD = 2.05) | 11.92 (SD = 4.77) |
| Executable SLOC | 10.45 (SD = 8.29) | 43.2 (SD = 15.7) |
| Modified variables | 3.42 (SD = 2.45) | 12.4 (SD = 4.9) |

Table 6: Comparison between the completely and partially analyzed loops

were only used once. These plans helped us in evaluating the analysis approach in loops with, say, high nesting level or a large number of procedure calls. Since these plans were designed to handle single specific events, they are probably not fully developed. The analysis of more loops in the same application domain should either eliminate or improve them.

Tables 5 and 6 summarize the data that was collected to determine which kinds of loops are more appropriately analyzed by the approach. Table 5 provides some insight into the practical limits of the approach. It gives different characteristics of the partially analyzed loops. Table 6 compares some of these characteristics to the corresponding ones of the completely analyzed loops. Since it was difficult to rigorously decompose some of the partially analyzed loops into events, the numbers of events were estimated by thoroughly reading the code and inspecting the data flow dependencies between the different statements.

The theoretical limitation, described in Section 6, only occurred in loop p9. That is, the partial analysis of the 12 loops in this case study was mainly because of practical limitations. Analyzing loops p1-p6 and p8-p9 using the current set of plans yielded no partial results. Loop p10, whose characteristics are compatible with those of the completely analyzed loops, was almost completely analyzed; four out of five events were analyzed. The fifth event was not analyzed because it

contained a call to an unanalyzed procedure. Loops p7 and p12 yielded some minor partial analysis results. Loop p11 gave considerable partial analysis results.

It is clear from Table 5 that almost all of the partially analyzed loops are nested (11 out of 12) and contain procedure calls (10 out of 12). They have an average size of 43.2 executable source lines of code and an average of 12.4 modified variables. Table 6 shows that some of these characteristics are considerably different from the corresponding ones for the completely analyzed loops. For example, the completely analyzed loops have an average size of 10.45 executable source lines of code and an average of 3.42 modified variables. While the average number of events in the completely analyzed loops is 3.28, the partially analyzed loops have 11.9 events on the average. This case study has given us the impression that loops of up to 5 events were more easily analyzed than others.

However, we noticed in some loops (p7, p8, p11 and p12) that some events closely match some of the designed plans. A larger domain of study could have improved those plans or resulted in designing similar ones that can contribute more to the specifications of such loops.

Even though the results of the case study are encouraging, further experimentation is, in our opinion, needed to investigate the generality and efficiency of the presented approach with respect to various application domains. This experimentation can serve to characterize the cases in which this approach can work best.

8 Conclusion

In this paper, a knowledge-based loop analysis approach has been described. This approach mechanically generates rigorous unambiguous predicate logic annotations of computer programs. It is a bottom-up analysis approach which does not rely on user-supplied information that might not be available at all times (e.g., the goals a program is supposed to achieve). In addition, it enables partial recognition and analysis of stereotyped, non-adjacent program parts.

A case study was performed on a real and pre-existing program of some practical value. This case study served to partially validate the analysis approach and to characterize its practical limits. To demonstrate the feasibility of automating our knowledge-based analysis approach, a prototype tool, which annotates loops with predicate logic annotations, has been designed and implemented [4].

The approach can assist in the maintenance and reuse activities by producing semantically sound and expressive predicate logic annotations of programs. Since many programs are undocumented, underdocumented, or misdocumented, a major part of the maintenance task is spent in recognizing and understanding abstract programming concepts [6, 26]. Automation of program analysis and understanding can, thus, contribute to maintenance tools and methods and provide support for various maintenance activities. Program analysis and understanding is also crucial for code reuse since the reuser must be aware of what a code component does. Understanding reusable code components can be achieved by augmenting them with a precise and clear description of their functionality [7]. If these descriptions are in the form of formal specifications, they can be further used in generating test cases and assessing the correctness of the implementation. Automation of program understanding is needed to facilitate the quick and efficient population of a reuse repository with well documented components [5, 11].

However, when annotating complicated and large program parts, these formal specifications can become hard to read. The readability of such specifications can be enhanced if they are further abstracted. This abstraction can be performed by replacing a formal statement with another one that is formulated in terms of a more widely known and understood concept [13]. Domain abstractions can further abstract the formal specifications with concepts specific to the application domain. The domain specific replacements can be explicitly performed by producing the abstract and then the domain specific ones. Otherwise, they can be implicitly performed by designing the plans such that their consequents are directly written in terms of the domain specific terms. In the former case, the knowledge base plans are more general and can be used in several different domains. The last stage which performs the higher level abstractions can be tailored to the needs of different domains and thus enhances the portability of the system. The latter approach, however, is easier to implement mechanically but reduces the generality of the plans.

With respect to software development, predicate logic plays an important role in development of software using such languages as VDM and Z [24, 38, 44]. Since our loop analysis technique produces predicate logic annotations, it can assist such formal development methods [4]. Our reverse engineering approach can provide assistance in the last development stage which moves from operation specifications to imperative programming language implementations. That is, the presented loop analysis technique can help in showing that the proof obligations generated during the operation refinement process are satisfied. It should be noted, however, that the mathematical notations used in VDM, Z, and our plans are not the same. To transform one mathematical notation to another, simple syntactic variations need to be performed.

There are some practical and theoretical limits to the presented approach. The practical limits are due to the difficulty of designing the knowledge base plans. The theoretical limits occur in nested structures with adaptation paths that contain statements other than assignment and conditional statements. They also occur while deriving the postconditions of general loops.

Future work includes extensions and improvements of the analysis approach, experimenting with the techniques in various application domains, and improvements on the prototype tool. The analysis approach needs to be expanded to perform an intelligent analysis of complete program modules which include non-algorithmic constructs such as stacks and queues. We need to investigate the utilization of additional information and knowledge in the source code (e.g., comments, variable names) to assist in plan recognition. Performing empirical studies in various application domains can serve to address and investigate several issues related to the acquisition and development of plans and the generality and efficiency of the presented approach with respect to different application domains. Finally, the developed tool served to demonstrate that the analysis techniques can be automated [4]. For practical utilization of such a tool, it needs to be enhanced to support additional programming language features and improve the user interface.

Acknowledgement

We would like to thank Lionel Briand, Gianluigi Caldiera, Walcelio Melo, Carolyn Seaman and Barbara Swain for their helpful contributions to a variety of aspects presented in this paper. This research was supported in part by the ONR grant N00014-87-k-0307 to the University of Maryland.

References

- [1] S. K. Abd-El-Hafiz. A Tool for Understanding Programs Using Functional Specification Abstraction. Master's thesis, University of Maryland, College Park, MD 20742, December 1990.
- [2] S. K. Abd-El-Hafiz. *A Knowledge-Based Approach to Program Understanding*. PhD thesis, University of Maryland, College Park, MD 20742, May 1994.
- [3] S. K. Abd-El-Hafiz and V. R. Basili. Documenting Programs Using a Library of Tree Structured Plans. In *Proceedings of the Conference on Software Maintenance*, pages 152–161, Montréal, Canada, September 27-30 1993. IEEE Computer Society Press.
- [4] S. K. Abd-El-Hafiz and V. R. Basili. A Tool for Assisting the Understanding and Formal Development of Software. In *Proceedings of the Sixth International Conference on Software Engineering and Knowledge Engineering*, Jurmala, Latvia, June 21-23 1994.
- [5] S. K. Abd-El-Hafiz, V. R. Basili, and G. Caldiera. Towards Automated Support for Extraction of Reusable Components. In *Proceedings of the Conference on Software Maintenance*, pages 212–219, Sorrento, Italy, October 15-17 1991. IEEE Computer Society Press.
- [6] A. Abran and H. Nguyenkim. Analysis of Maintenance Work Categories Through Measurement. In *Proceedings of the Conference on Software Maintenance*, pages 104–113, Sorrento, Italy, October 1991. IEEE Computer Society Press.
- [7] V. R. Basili and S. K. Abd-El-Hafiz. Packaging Reusable Components: The Specification of Programs. Technical Report CS-TR-2957, Department of Computer Science, University of Maryland, College Park, MD 20742, September 1992.
- [8] S. K. Basu and J. Misra. Proving Loop Programs. *IEEE Trans. on Software Engineering*, SE-1(1):76–86, March 1975.
- [9] K. Bertels, P. Vanneste, and C. De Backer. A Cognitive Approach to Program Understanding. In *Proceedings of the Working Conference on Reverse Engineering*, pages 1–7, Baltimore, Maryland, May 21-23 1993. IEEE Computer Society Press.
- [10] G. Brassard and P. Bratley. *Algorithmics: Theory & Practice*. Prentice Hall, 1988.
- [11] G. Caldiera and V. R. Basili. Identifying and Qualifying Reusable Software Components. *IEEE Computer*, 24(2):61–70, February 1991.
- [12] D. D. Dunlop and V. R. Basili. A Heuristic for Deriving Loop Functions. *IEEE Trans. on Software Engineering*, SE-10(3):275–285, May 1984.
- [13] R. B. France and V. R. Basili. A Pattern-Driven Approach to Code Analysis for Reuse. Technical Report CS-TR-2802, Department of Computer Science, University of Maryland, College Park, MD 20742, November 1991.
- [14] D. Gries. *The Science of Programming*. Springer-Verlag, 1981.

- [15] M. T. Harandi and J. Q. Ning. PAT: A Knowledge-Based Program Analysis Tool. In *Proceedings of the Conference on Software Maintenance*, pages 312–318, Phoenix, Arizona, October 1988. IEEE Computer Society Press.
- [16] M. T. Harandi and J. Q. Ning. Knowledge-Based Program Analysis. *IEEE Software*, 7(1):74–81, January 1990.
- [17] J. Hartman. Understanding Natural Programs Using Proper Decomposition. In *Proceedings of the 13th International Conference on Software Engineering*, pages 62–73, Austin, Texas, May 13-16 1991. IEEE Computer Society Press.
- [18] P. A. Hausler, M. G. Pleszkoch, R. C. Linger, and A. R. Hevner. Using Function Abstraction to Understand Program Behavior. *IEEE Software*, 7(1):55–63, January 1990.
- [19] C. A. R. Hoare. An Axiomatic Basis for Computer Programming. *Communications of the ACM*, 12(10):576–580, 583, October 1969.
- [20] C. Samuel Hsieh. Slice, Chunk, and Dataflow Anomaly as Datalog Rules. *The Journal of Systems and Software*, 16(3):197–203, November 1991.
- [21] P. Jalote. *An Integrated Approach to Software Engineering*. Springer-Verlag, 1991.
- [22] W. L. Johnson. *Intention-Based Diagnosis of Novice Programming Errors*. Morgan Kaufmann, 1986.
- [23] W. L. Johnson and E. Soloway. PROUST: Knowledge-Based Program Understanding. *IEEE Trans. on Software Engineering*, SE-11(3):267–275, March 1985.
- [24] C. B. Jones. *Systematic Software Development Using VDM*. Prentice Hall International, 1990.
- [25] S. Katz and Z. Manna. Logical Analysis of Programs. *Communications of the ACM*, 19(4):188–206, April 1976.
- [26] B. P. Leintz and E. B. Swanson. *Software Maintenance Management*. Addison-Wesley, 1980.
- [27] S. Letovsky. Program Understanding with the Lambda Calculus. In *Proceedings of the 10th International Joint Conference on AI*, pages 512–514, August 1987.
- [28] K. B. McKeithen, J. S. Reitman, H. H. Rueter, and S. C. Hirtle. Knowledge Organization and Skill Differences in Computer Programmers. *Cognitive Psychology*, 13:307–325, 1981.
- [29] H. D. Mills, V. R. Basili, J. D. Gannon, and R. G. Hamlet. *Principles of Computer Programming: A Mathematical Approach*. Allyn and Bacon, 1987.
- [30] W. R. Murray. *Automatic Program Debugging for Intelligent Tutoring Systems*. Morgan Kaufmann, 1988.
- [31] A. Quilici. A Hybrid Approach to Recognizing Programming Plans. In *Proceedings of the Working Conference on Reverse Engineering*, pages 126–133, Baltimore, Maryland, May 21-23 1993. IEEE Computer Society Press.

- [32] C. Rich and L. M. Wills. Recognizing a Program's Design: A Graph-Parsing Approach. *IEEE Software*, 7(1):82-89, January 1990.
- [33] E. Rich and K. Knight. *Artificial Intelligence*. McGraw-Hill, 1991.
- [34] P. G. Selfridge, R. C. Waters, and E. J. Chikofsky. Challenges to the Field of Reverse Engineering. In *Proceedings of the Working Conference on Reverse Engineering*, pages 144-150, Baltimore, Maryland, May 21-23 1993. IEEE Computer Society Press.
- [35] E. Soloway. Learning to Program = Learning to Construct Mechanisms and Explanations. *Communications of the ACM*, 29(9):850-858, September 1986.
- [36] E. Soloway, J. Bonar, and K. Ehrlich. Cognitive Strategies and Looping Constructs: An Empirical Study. *Communications of the ACM*, 26(11):853-860, November 1983.
- [37] E. Soloway and K. Ehrlich. Empirical Studies of Programming Knowledge. *IEEE Trans. on Software Engineering*, SE-10(5), September 1984.
- [38] J. M. Spivey. An Introduction to Z and Formal Specifications. *Software Engineering Journal*, pages 40-50, January 1989.
- [39] J. M. Spivey. *The Z notation: A reference Manual*. Prentice Hall International, 1992.
- [40] M. Ward, F. W. Calliss, and M. Munro. The Maintainer's Assistant. In *Proceedings of the Conference on Software Maintenance*, pages 307-315, Miami, Florida, October 1989. IEEE Computer Society Press.
- [41] R. C. Waters. A Method for Analyzing Loop Programs. *IEEE Trans. on Software Engineering*, SE-5(3):237-247, May 1979.
- [42] B. Wegbreit. The Synthesis of Loop Predicate. *Communications of the ACM*, 17(2):102-112, February 1974.
- [43] L. M. Wills. Flexible Control for Program Recognition. In *Proceedings of the Working Conference on Reverse Engineering*, pages 134-143, Baltimore, Maryland, May 21-23 1993. IEEE Computer Society Press.
- [44] J. C. P. Woodcock. Structuring Specifications in Z. *Software Engineering Journal*, pages 51-66, January 1989.