

# Program Understanding as Constraint Satisfaction: Representation and Reasoning Techniques

Steven Woods

Qiang Yang

Department of Computer Science  
University of Waterloo  
Waterloo, ONTARIO N2L 3G1  
{sgwoods, qyang}@logos.uwaterloo.ca

## Abstract

The process of understanding a source code in a high-level programming language involves complex computation. Given a piece of legacy code and a library of program plan templates, understanding the code corresponds to building mappings from parts of the source code to particular program plans. These mappings could be used to assist an expert in **reverse engineering** legacy code, to facilitate **software reuse**, or to assist in the *translation* of the source into another programming language. In this paper we present a model of program understanding using constraint satisfaction. Within this model we intelligently compose a partial global picture of the source program code by transforming knowledge about the problem domain and the program itself into sets of constraints. We then systematically study different search algorithms and empirically evaluate their performance. One advantage of the constraint satisfaction model is its generality; many previous attempts in program understanding could now be cast under the same spectrum of heuristics, and thus be readily compared. Another advantage is the improvement in search efficiency using various heuristic techniques in constraint satisfaction.

## 1 Introduction

Humans are particularly adept at successfully interpreting explicit representations of knowledge created by other intelligent agents. A shared understanding of the terms of reference and subject material provides a basis for this interpretation. In software engineering, experts often apply such skill to the task of *program understanding*. As shown in Figure 1, it is possible to conceptualize an expert's understanding of a given source program as a successful construction of a *mapping* between the expert's store of relevant knowledge and the

structures and components inherent in the source code. The expert or agent can use this mapping to infer the source program's high-level goals. This mapping essentially raises the level of abstraction of the understanding of the source from the level of actual code to the more abstract level of the existing representation (or language of expression) of the domain knowledge. This abstract understanding may be exploited as part of the process of: (1) translating the program into the source code of another programming language, (2) recognizing errors in the legacy code and assisting in debugging the code at the more abstract level, and (3) replacing understood code portions with generic application code or calls to other code libraries. We know that in many real-world circumstances, a reduction in the size of an existing source code library by only a small percentage can result in a substantial reduction of the maintenance cost, and consequently creating a mapping (even a very partial one) between existing domain knowledge and a particular legacy source offers many possible levers to the experts charged with dealing with this source.

In Artificial Intelligence research, the problem of program understanding has been approached indirectly from the perspective of plan recognition [7, 1, 2, 25]. In many of these works, existing human knowledge in a particular domain is represented as hierarchies of plans that describe relevant actions and goals. Given such a hierarchy, and an observation of another agent's plan, a plan-recognizer would typically construct a mapping from input plan fragments to the leaf nodes of the knowledge-base and infer upwards toward a goal. To disambiguate among alternative goals, the mapping processes may employ knowledge about the temporal relations between parts of the plan. These plan recognition programs have been applied mostly to *toy domains* (such as the cooking domain), involving small knowledge bases and a small amount of search.

Recently, researchers have adopted a more direct approach to program understanding. In this direction, an explicit library of programming plan templates and concepts is constructed, and various top-down and bottom-up search strategies are utilized to implement the mapping process. Notable examples are Quilici[18], Kozaczynski and Ning[8], Rich and Waters[23] and Wills[27, 28]. To some extent, all are aimed at improving the effectiveness of the mapping process through heuristic knowledge.

In Figure 2 a subset of expert knowledge about a particular application domain is represented in a fragment of a hierarchical library of program templates. One possible mapping is shown between a plan template from the library and a specific legacy source fragment, in this case a single source statement. The existence of such a mapping essentially *explains* the presence of the low-level source statement at a higher level of abstraction, in this case as an instance of the plan template **copy-character** specified in the library.

Much of the previous program understanding work has failed to demonstrate heuristic adequacy in even partially generating "understanding" of large problems. Specifically, many recognition algorithms presented may be viewed as collections of heuristic tricks. This construction makes it difficult for one to perform a systematic analysis of different search methods within a particular approach, or to understand how the addition or deletion of certain types of domain-specific knowledge may affect performance. We are unaware of concrete

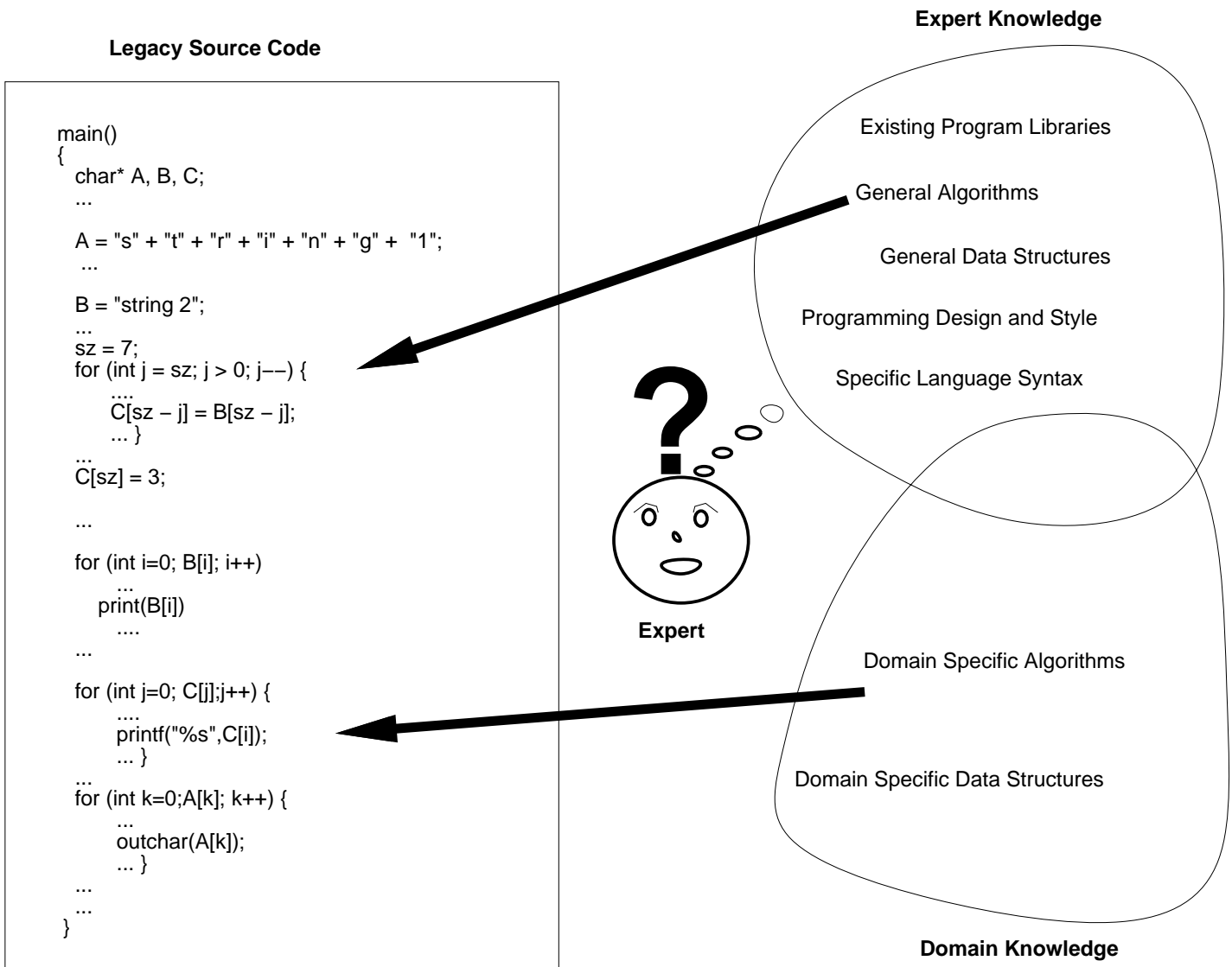


Figure 1: Conceptualizing source with expert knowledge.

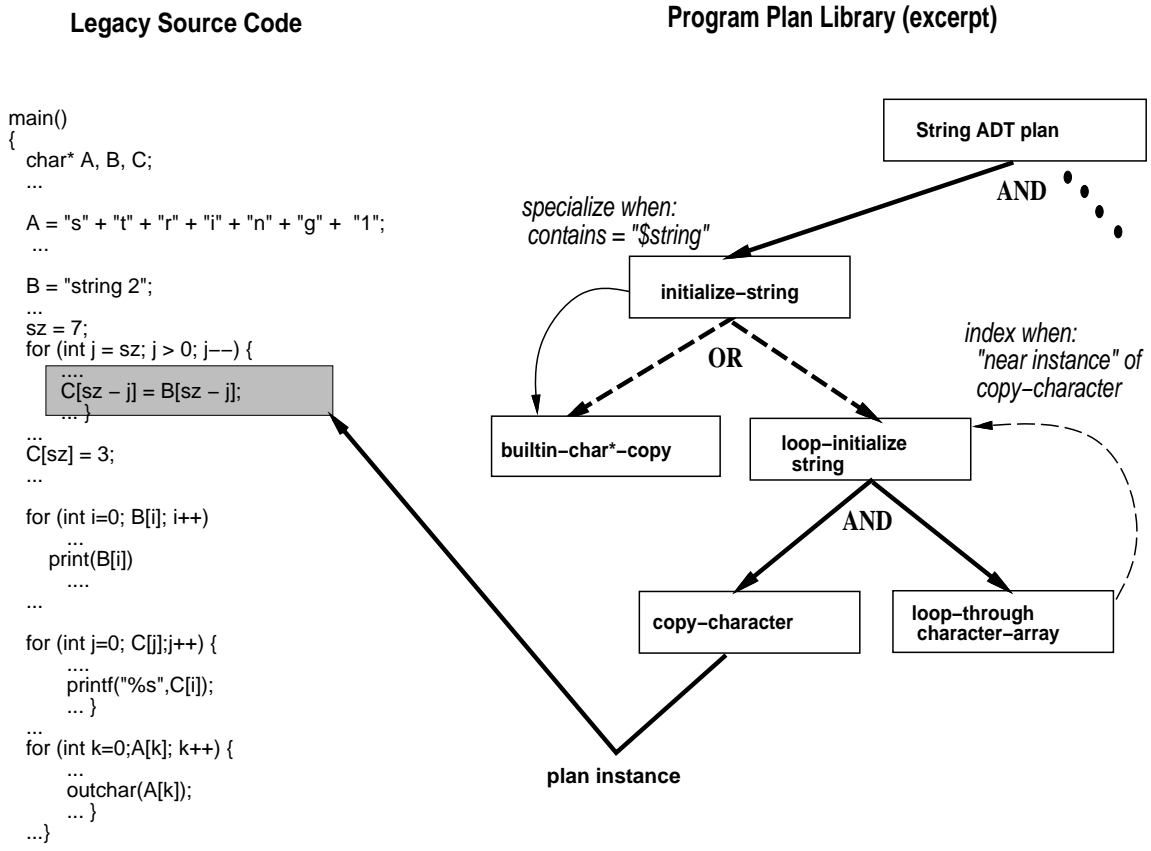


Figure 2: Conceptualizing source with a plan library.

examples or experiments which might suggest that these approaches might scale up for specific uses in large sources. One exception might be Wills[28] who presents empirical results promising in identifying partial mappings of reasonably sized legacy sources to a library of program plans.

The work presented in this paper is part of the initial phase of work focused on demonstrating that an effective approach to partial program understanding is possible with large legacy code examples. Specifically, we intend to clearly categorize the circumstance in which this use is possible, and the preconditions which must first be met in terms of representation and application of domain knowledge. We present a generalized representation of program understanding as a *Constraint Satisfaction Problem (CSP)*[10]. For a given legacy source code, the program components (explained later) are variables in the CSP. The domain values are the known program plans that may *explain* each component. The CSP constraints are either *knowledge constraints* which describe how program plans may fit together to form larger plans, or *structural constraints* which describe how program components are structurally related. We refer to the program understanding CSP as PU-CSP.

In addition, we present and empirically evaluate a mapping algorithm (as part of the PU-CSP), also formulated as a CSP, which provides the ability to locate all instances of a specific general programming plan template, and to map the plan's structure to actual source program components. We refer to this mapping CSP as MAP-CSP. Some earlier works also attempt to define and recognize abstract concepts as part of program understanding[8, 28]. For a given program plan template (explained later), the different parts of the template are the variables in the MAP-CSP. The various syntactically known pieces of the source code correspond to domain values for each variable. The constraints among the different parts of the program plan are constraints in the MAP-CSP.

There are at least two advantages in our constraint-based approach. The first is its **generality**; most of the previous recognition methods and heuristics can now be unified under the constraint-based view. Another advantage is an increased ability to address **heuristic adequacy**, or **scalability**; by casting program understanding as a CSP, the previously known constraint propagation and search algorithms could be easily adapted. We may now perform a systematic study of different search heuristics, including both top-down and bottom-up as well as many other hybrids, in order to discover their applicability to a particular source code.

## 2 The Program Understanding Problem

### 2.1 An Illustrative Example

Consider the C program outlined on the left hand side of Figure 3. This example program contains declarations, initializations and an embedded print loop for *each* of three strings. As an illustration, strings are treated as a primitive data type by the programmer, with no

shared functionality for printing.

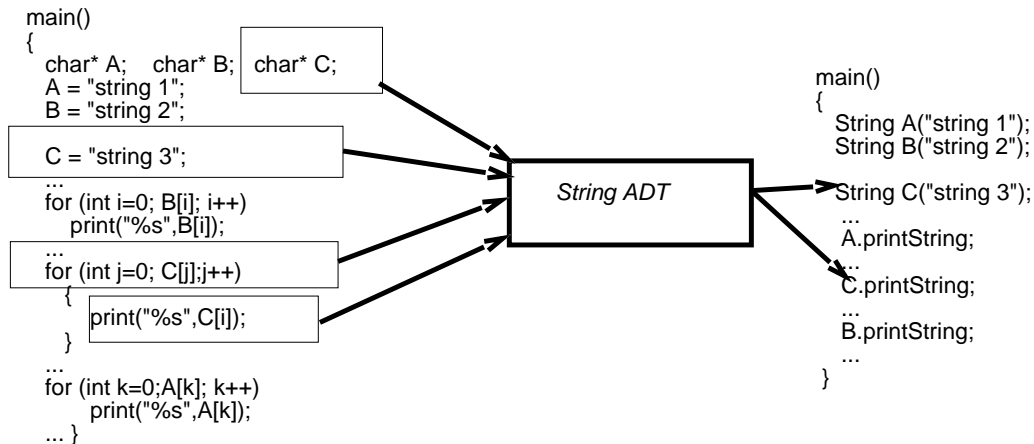


Figure 3: C legacy code mapped as String ADT instance to C++ code.

To understand this program, one might use as a basis a library of program plans as shown in Figure 4 which represents previously compiled knowledge about program composition within a particular domain. Figure 1 shows a program plan for the Abstract Data Type (ADT) or class **String** which is part of this library of plans. Once a mapping is constructed between the source and compiled knowledge, one could translate the redundant source code to one with a single inclusion of the ADT, as shown in the C++ code on the right hand side of Figure 3.

```

Class String {
  char localStr [MAXSIZE];

  String( char* inStr )
  {
    for (int j=0; inStr[j]; j++)
      localStr[j] = inStr[j]; }

  printString()
  {
    for (int j=0; localStr[j]; j++)
      printf("%s",localStr[j]); } }

```

Table 1: Example abstract data type.

Given the legacy source code on the left side of Figure 3, we would like to *understand* or *explain* some portions of the source program within the known context of the program plans such as represented by the **String** ADT. Successful identification could result in the replacement of much redundant source code with a single inclusion of the ADT. The C++

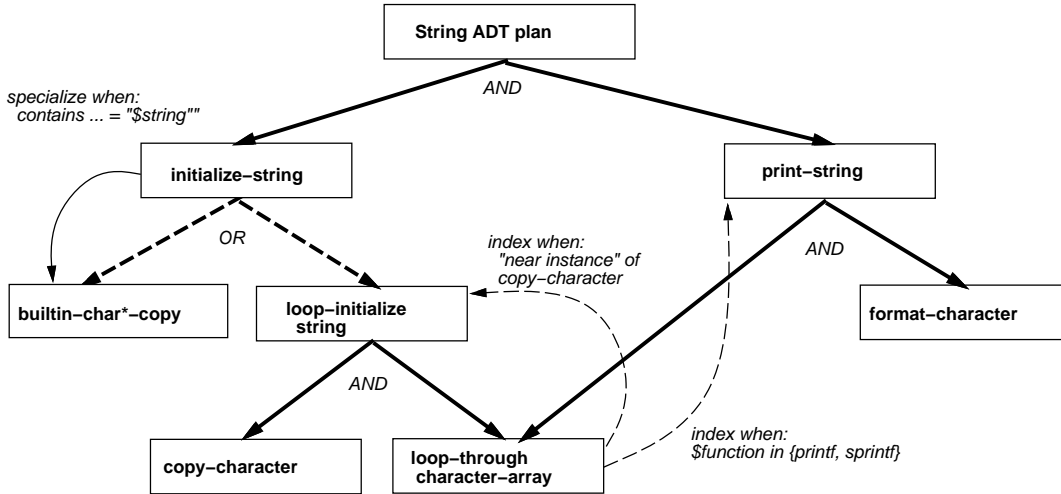


Figure 4: **String** ADT within a hierarchical program plan library.

code shown at the right of Figure 3 is obtained with replacement of C source with references to **String** ADT functionality. This understanding process might be executed in two steps. First, one identifies all instances of a particular abstract program plan in a source code. We refer to this problem as the *MAP-CSP* problem. Second, one relates some set of identified plan blocks (or program slices) to conform to the hierarchical structure in a given program-plan knowledge base. The latter we refer to as the *PU-CSP* problem.

We identify two important benefits of locating mappings between a programming plan library and an existing source or legacy code. First, the resulting replacement of legacy code with ADT instances can result in substantial reduction in code. This size savings can reduce the amount of effort required for subsequent code understanding or maintenance by programmers. Second, the mapping between source and library plan can be used as a building block in attempting to understand and translate the legacy code. The intent of this work is twofold. We describe how various types of individual mappings can be identified efficiently, and we outline how this mapping process may be integrated into the larger task of program understanding.

## 2.2 Quilici’s Memory-based Method

Quilici’s method is representative of other earlier work in this area, including work by Kozaczynski and Ning[8]. This approach [18, 20, 19, 21] is based on a construction of an explicit library of programming plan templates, complete with an indexing ability, which can quickly associate a particular recognized source code with program plan templates in the knowledge base. Furthermore, a combination of top-down and bottom-up search strategies is utilized to implement the matching process. With this system Quilici demonstrated how simple C programs could be translated to C++ programs.

Program plans (such as embedded in ADTs) are organized hierarchically in a library as shown in Figure 4. Legacy source code in the form of an abstract syntax tree is mapped to the plan library through the use of indices, which are pointers from the source code to parts of the plan library. Index tests indicate when to *specialize* or to *infer* the existence of other plans according to a set of conditions. As an example of specialization, consider Figure 4 in which the program plan **initialize-string** is specialized to **builtin-char\*-copy** when a direct string assignment is observed in the source code. An example of an inference test is also shown in Figure 4, where the existence of **loop-initialize-string** is inferred when an instance of **loop-through-character-array** is “near” a related instance of **copy-character** in the source code.

Given a source code and a program plan, Quilici describes an approach to understanding the legacy source based on a search in the plan library. Search behaves *bottom-up* when existing index tests indicate possible higher-level explanation plans for a particular lower-level component in the library. Quilici observes that people only make bottom-up inferences in particular “well-known” circumstances, and consequently limits the number of upward explanations by inferring only those specified by explicit indexes. On the other hand, search behaves *top-down* when low-level components are indexed and subsequently matched based on some hypothesized high-level plans. Quilici’s algorithm attempts to specialize any matched plan as much as possible according to predefined specialization tests, and directs search for low-level plans based on high-level hypothesized plans. This approach marks one of the first cognitively motivated attempts to program understanding using a hierarchical library of program plans. There are, however, a number of shortcomings. First, the lack of a general mathematical model of the indexing and search process makes it unclear as to how one should coordinate the top-down and bottom-up search. Second, Quilici’s algorithm depends on a number of heuristics, such as specializing a plan as much as possible. It is not clear how these heuristics integrate or how they scale-up when the problem size increases. Finally, Quilici makes a substantial effort in capturing actual programmer’s methodologies as heuristic enhancements to search control, but presents no empirical results.

While studying this work, it occurred to us that the program understanding problem could be broken down into a number of choice points. Examples of these choices include: (1) choosing between candidate unexplained components, (2) choosing between multiple initial plan assignments for a component, (3) choosing between several plans whose existence is implied top-down, and (4) choosing a particular index or specialization test from a candidate set. The existence and interactions of these decisions are buried in Quilici’s presentation, but are very important in addressing the efficiency of the search problem. In the next section, we explore how to represent and exploit these choice points using a simple and elegant mathematical model known as *constraint satisfaction*.



## 2.3 Wills' Graph Parsing Method

Wills[23, 27, 28] outlined an approach to recognition in which stereotypical program or data structures known as *clichés* are represented as a type of graph grammar. A source program is translated into an intermediate representation as a flow graph. These flow graphs are parsed so as to identify all possible derivations of the flow graph based on the known *clichés*. These derivations each represent a possible *partial* interpretation of the source program or mapping to the library of clichés. Wills notes that although the parsing problem is NP-complete in general, experience suggests that attribute constraint checking significantly prunes the search space. Wills evaluates the effectiveness of such an approach empirically for two medium-size source code examples.

Wills' work differs from our approach in at least 3 important ways: (1) cliché and program representation, (2) library knowledge representation and exploitation during search, and (3) method of integrating cliché instances in the larger understanding problem.

## 2.4 Other Related Work

Kozaczynski and Ning[8] describe a method of automatically recognizing abstract concepts in source code given a library of concepts and rules for how to recognize the higher-level concepts in lower-level language concepts, essentially controlling the concept search in a top-down fashion. Muller and others[16, 15, 14] are involved in the construction of Rigi, a system for analyzing software systems which includes visual representations of data and control flow structures in code towards the identification of subsystems and hierarchies of structure in code. Rich and Waters[22, 23] headed the Programmer's Apprentice project which focused on the development of a demonstration system (Knowledge-Based Editor in Emacs or KBEmacs) with the ability to assist a programmer in analyzing, creating, changing, specifying and verifying software systems. In addition, Rich and Waters[23][pp.171-188] describe a cliché recognizer Recognize based in KBEmacs. Rugaber, Stirewalt, Wills and others are part of an effort in reverse engineering being conducted at the Georgia Institute of Technology. Recent work[24] describes one major research area in program understanding known as interleaving in which program plans intertwine.

# 3 An Introduction to Constraint Satisfaction

Constraint satisfaction problems (CSPs) provide a simple and yet powerful framework for solving a large variety of AI problems. The technique has been successfully applied to machine vision, belief maintenance, scheduling, and planning, as well as many design tasks. For a successful of this technique to knowledge-based planning, see [34].

A constraint satisfaction problem (CSP) can be formulated abstractly as three components:

1. a set of **variables**,  $X_i, i = 1, 2 \dots n$ ,
2. for each variable  $X_i$  a set of values  $\{v_{i1}, v_{i2}, \dots v_{ik}\}$ . Each set is called a **domain** for the corresponding variable, denoted as  $\text{domain}(X_i)$ ,
3. a collection of **constraints** that defines the permissible subsets of values to variables.

The goal of a CSP is to find one (or all) assignment of values to the variables such that no constraints are violated. Each assignment,  $\{x_i = v_{ij}, i = 1, 2, \dots, n\}$ , is called a **solution** to the CSP.

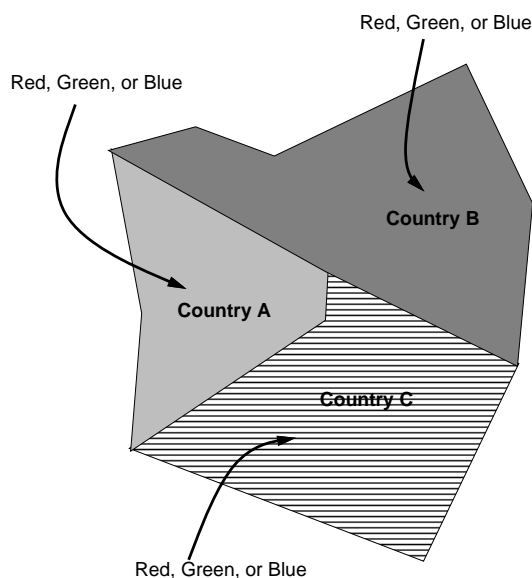


Figure 5: A Map Coloring Problem.

As an example of a CSP, consider a map-coloring problem, where the variables are regions  $R_i, i = 1, 2, \dots, n$  that are to be colored (see Figure 5). In a final solution every region must be assigned a color such that no two adjacent regions share the same color. A domain for a variable is the set of alternative colors that a region can be painted with. For example, a domain for  $A$  might be  $\{\text{Green, Red, Blue}\}$ . A constraint exists between every pair of adjacent variables, which states that the pair cannot be assigned the same color. Between adjacent regions  $A$  and  $B$ , for example, there is a constraint  $A \neq B$ . A solution to the problem is a set of colors, one for each region, that satisfies the constraints.

Let  $\text{Vars} = \{X, Y, \dots Z\}$  be a set of variables. A constraint on  $\text{Vars}$  is essentially a *relation* on the domains of the variables in  $\text{Vars}$ . If a constraint relate only two variables then it is called a **binary constraint**. A CSP is binary if all constraints are binary. For any two variables  $X$  and  $Y$ , we say  $X = u$  and  $Y = v$  is **consistent** if all binary constraints between  $X$  and  $Y$  are satisfied by this assignment.

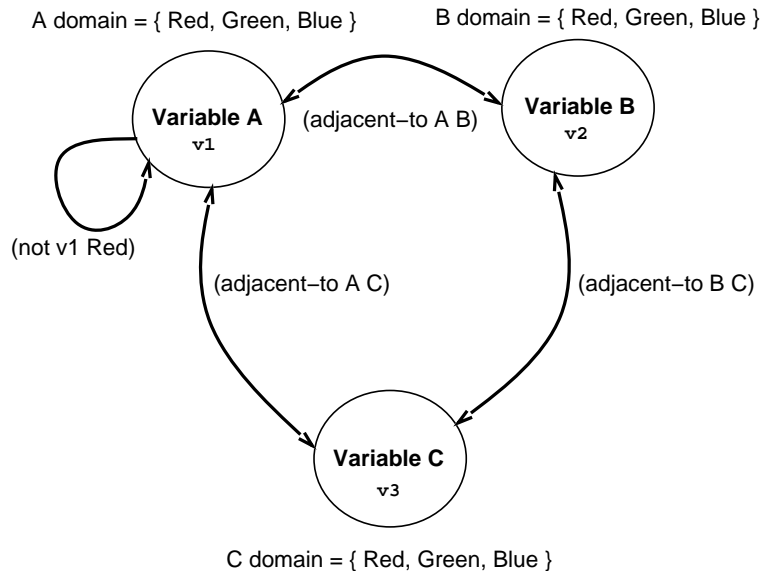


Figure 6: Map-Coloring CSP.

A variety of techniques have been developed for solving CSPs. They can be classified as local *consistency-based methods*, global *backtrack-based methods* or *local-search methods*. Local-search methods [12] is a kind of greedy algorithm which is gaining popularity. We do not review this method here, but we do intend for our CSP modeling to be general enough to include local-search as a reasoning method.

### 3.1 Local Consistency Methods

Local consistency methods follow the theme of *preprocessing*. That is, before a more costly method is used, a consistency-based method could be applied to simplify a CSP and remove any obviously incompatible values. Often these methods yield tremendous headway toward eventually solving the problem.

Let  $X$  and  $Y$  be two variables. If a domain value  $A$  of  $X$  is **inconsistent** with all values of  $Y$ , then  $A$  cannot be part of a final solution to the CSP. This is because in any final solution  $S$ , any assignment to  $X$  must satisfy all constraints in the CSP. Since  $X = A$  violates at least one constraint in all possible solutions,  $A$  can be removed from the domain of  $X$  without affecting any solution.

If for a pair of variables  $(X, Y)$ , for every value of  $X$  there is a corresponding **consistent** value of  $Y$ , then we say  $(X, Y)$  is arc-consistent. By the above argument, enforcing arc-consistency by removing values from variable domains does not affect the final solution. The process of making every pair of variables arc-consistent is called *arc-consistency*.

### 3.2 Backtrack-based Algorithms

Arc-consistency algorithms only work on pairs of variables, and as such can only handle binary constraints and cannot always guarantee a final solution to a CSP. A more thorough method for solving a CSP is backtracking, where a depth-first search is performed on a search tree formed by the variables in the CSP. A thorough examination of these techniques can be found in [17] and [9]. During a backtracking search, each variable instantiation is interpreted as extending the current understanding of a legacy program one step further.

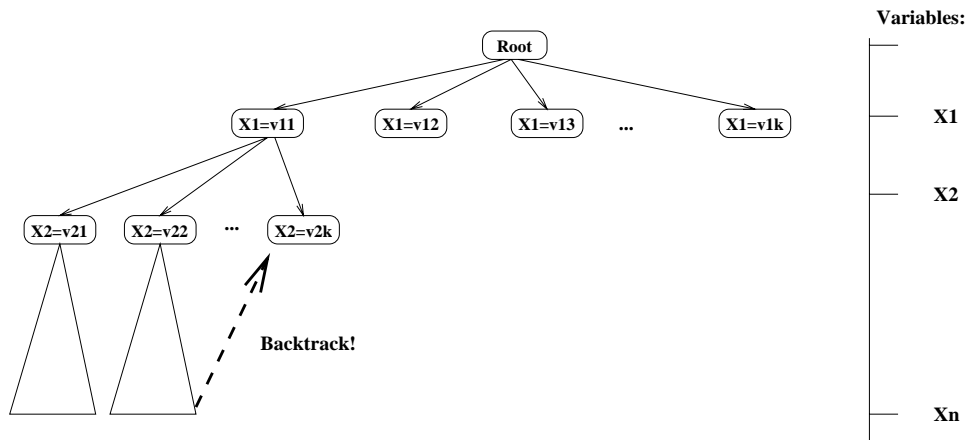


Figure 7: A search tree for a backtrack-based algorithm.

A backtracking algorithm instantiates the variables one at a time in a depth-first manner. It backtracks when the constraints accumulated so far signal inconsistency. In Figure 7 we show this process. First, variables are ordered in a certain sequence. Different orders of variables might entail different search efficiency, and heuristics for good ordering of variables are called **variable-ordering heuristics**. Similarly, for each variable, the values are tried out one at a time, and the heuristics for a good ordering of values are called **value-ordering heuristics**.

Using the CSP representation, we can also consider a more systematic study of different search algorithms. Figure 2 provides a general backtracking algorithm for solving a CSP. In this algorithm, we have a number of hooks where we could place different search heuristics. They correspond to heuristics for ordering variables and constraints, as well as heuristics for deciding the amount of constraint propagation.

There are several choice points which both individually and in combination affect the resulting search performance. These choice points are explained as follows:

1. *Initialization* and *Initial Constraint Propagation* are the determination of variables and domain values before the search starts. It can be viewed as a special type of localized constraint propagation algorithm, but one that is directed according to pre-defined domain knowledge. The determination of the set  $V$  and of  $Dom(X)$  controls how

## Generic CSP Search

$V$ : variables in a CSP,  $Dom(X)$ : the domain values of  $X$ .

1. [**Initialization**] for each variable  $X_i \in V$ , find the set of domain values for  $X_i$ ;
2. [**Initial Constraint Propagation**] Reduce  $Dom(X)$  by constraint propagation.
3. Solution = NULL
4. [**Variable Selection**] Select and remove a variable  $X$  from  $V$
5. [**Value Selection**] Select and remove a value of  $X$  from  $Dom(X)$ .  
The value must be consistent with all assignments in Solution.
6. [**In-search Propagation**] Apply a subset of constraints to  $V$ .
7. [**Backtrack Point Selection**] Backtrack if any  $Dom(X)$  in  $V$  becomes empty.
8. [**Solution Evaluation**] If  $V$  is empty, exit with Solution (if all-solution, continue); else, goto Step 4.

Table 2: Generic CSP Search Algorithm.

much work is done in advance. This reduction could also be performed as an in-search propagation at Step 6 of the Generic CSP algorithm.

2. *Constraint Propagation* is the reduction of domains locally or globally within the CSP problem graph. Existing algorithms include AC-3[10], AC-4[13], AC-5[26], and other variations[4, 3].
3. *Variable Selection* is the determination of which component variable should be chosen next for instantiation during search. The decision may be based on domain independent measures, such as the size of a variable's domain; on information specific to the instance and domain plan library, such as frequency of occurrence of particular plan templates in the variable domain set, or on some combination of these types of information.
4. *Domain Value Selection* is the determination of a particular plan explanation, taken from the plan library, to assign to the component variable. Typically this selection should be made so as to most effectively limit the remaining variable ranges, that is, to be the most context limiting. In terms of our plan library this means a plan that is as *specific* as possible.
5. *In-search Propagation* is the reduction (as for Step 2) of the remaining uninstantiated variable domains according to some constraint propagation algorithm. Problem characteristics such as variable domains that exceed some average or absolute bounds are potential signals that constraint propagation may be useful before continuing search. In [17] the advantages of exploiting various algorithms for achieving a limited degree of partial consistency amongst variable sets are examined.

6. *BackTrack point selection* is the determination, after it has become evident that no possible solution exists along a particular variable-instantiation path, of which instantiation to retract. Intelligent backtracking approaches such as BackJumping and BackMarking<sup>1</sup> attempt to determine the origin of the conflict that caused the failure, and to BackTrack as far up the search tree as possible to avoid a repeated failure of the same condition.
7. *Solution Evaluation* determines whether or not a particular solution is satisfactory. In a cooperative interactive approach to program understanding, it is at this point that an expert might interact and evaluate a particular partial solution for adequacy. Similarly, if there exists particular measures of adequacy or *soft, preferential* constraints that may have been relaxed during search, such measures may be applied here.

There are in addition several other ways to improve the search efficiency. One method is to employ the particular hierarchical structure of the plan library, and using a *hierarchical constraint satisfaction algorithm*[11]. In this approach, the plan library represents plans at varying levels of abstraction. A set of low-level program components which have been mapped to the program library may be grouped according to their functional relationships and form a higher-level component. This component (or variable) may now be explained by a more abstract plan (or domain value) according to both the structural constraints imposed in source structure and the knowledge constraints present in the program plan library. We pursue this type of constraint application more completely in future work.

In the generic search algorithm a set of choice points are presented in the new context of CSP solving. In the next section of this paper we discuss and evaluate several selection variations for recognition of one particular template in sets of generated source code examples. We examine variations that include applying AC-3 as Step 1 combined with BackTracking and also another more intelligent search algorithm known as Forward Checking[5], which performs a limited amount of in-search propagation at Step 6. In addition, the intelligent search algorithm dynamically rearranges the order of variables during search according to the size of the variable domains, selecting the shortest first.

The order in which constraints are applied can also dramatically affect search. Constraint ordering or selection would occur at Step 6. In particular, it is advantageous to apply constraints that are inexpensive computationally and that (potentially) prune a large number of domain values. In a particular domain it may be possible to determine or estimate such relative benefits either from past empirical results or through analysis of the domain structure itself. For instance, the property that program template features tend to be found *spatially* near each other can be exploited through heuristics that limit the range of search for related components. The effectiveness of such abstraction heuristics has been reported elsewhere[6, 29].

---

<sup>1</sup>These and other intelligent backtracking algorithms are described in detail by Nadel in[17].

### 3.3 Program Understanding as CSP

We view the entire program understanding problem as a constraint satisfaction problem. In this model, a long program code is first divided into blocks, where each block is a set of closely related source code. The program understanding problem is to identify the top-level function of each of these program blocks, so that not only the inter-relationships between the blocks are explained, but also the constraints specified by a program library on the program plans are respected. A key problem, then, is to assign one plan component to each block, subject to a set of constraints. This problem we call the **program-understanding CSP**, or PU-CSP.

The number of program plan components that one could assign to each block could be enormous. To be practical, it is crucial to first reduce the number of explanations for each block as much as possible. This process could be helped by a related constraint satisfaction problem, one that we will explain in detail in Section 5: the problem of finding all instances of a given program plan or pattern in the entire source code. This problem we call the MAP-CSP problem.

Below, we explain both problems in detail.

## 4 Program Understanding CSP: PU-CSP

PU-CSP is formed in the following way. Suppose that an initial decomposition or slicing of the source code is given. Each block of source code corresponds to a *variable* in the PU-CSP. The *Variable domains* correspond to all possible explanations of an individual source code block. As an example, consider the legacy code program statements of Figure 3 as the blocks. We take each block as a PU-CSP variable which ranges over all possible program plans of corresponding statement type, such as “declaration”, “assignment”, “print”, etc, in the plan library of Figure 4.

### 4.1 The Modeling Process

A Program Understanding CSP (PU-CSP) is formulated via four distinct steps shown in Figure 8. First, the legacy source is pre-processed creating a set of artifacts that describe some precise interrelationships in the source regarding data flow relationships between functional blocks, control flow among the functional blocks, and the creation of an abstract syntax tree in an intermediate abstract language via parsing of the source. Second, the source code is partitioned according to existing program slicing methodologies into spatially localized blocks of code which are known to exhibit functional relationships among one another, and cohesive properties within one’s boundaries. Third, a skeleton CSP is formulated consisting of one variable for each identified source block, and constraints between these variables are derived from the intermediate representation level artifacts. Each variable ‘typed’ via the addition of reflexive constraints on the variable which describe properties of the block such

as *kinds* of input or output. Finally, each CSP variable is compared against the templates in the program plan library, with any templates which potentially match a variable with regards to input and output typing are composed as the domains of that variable.

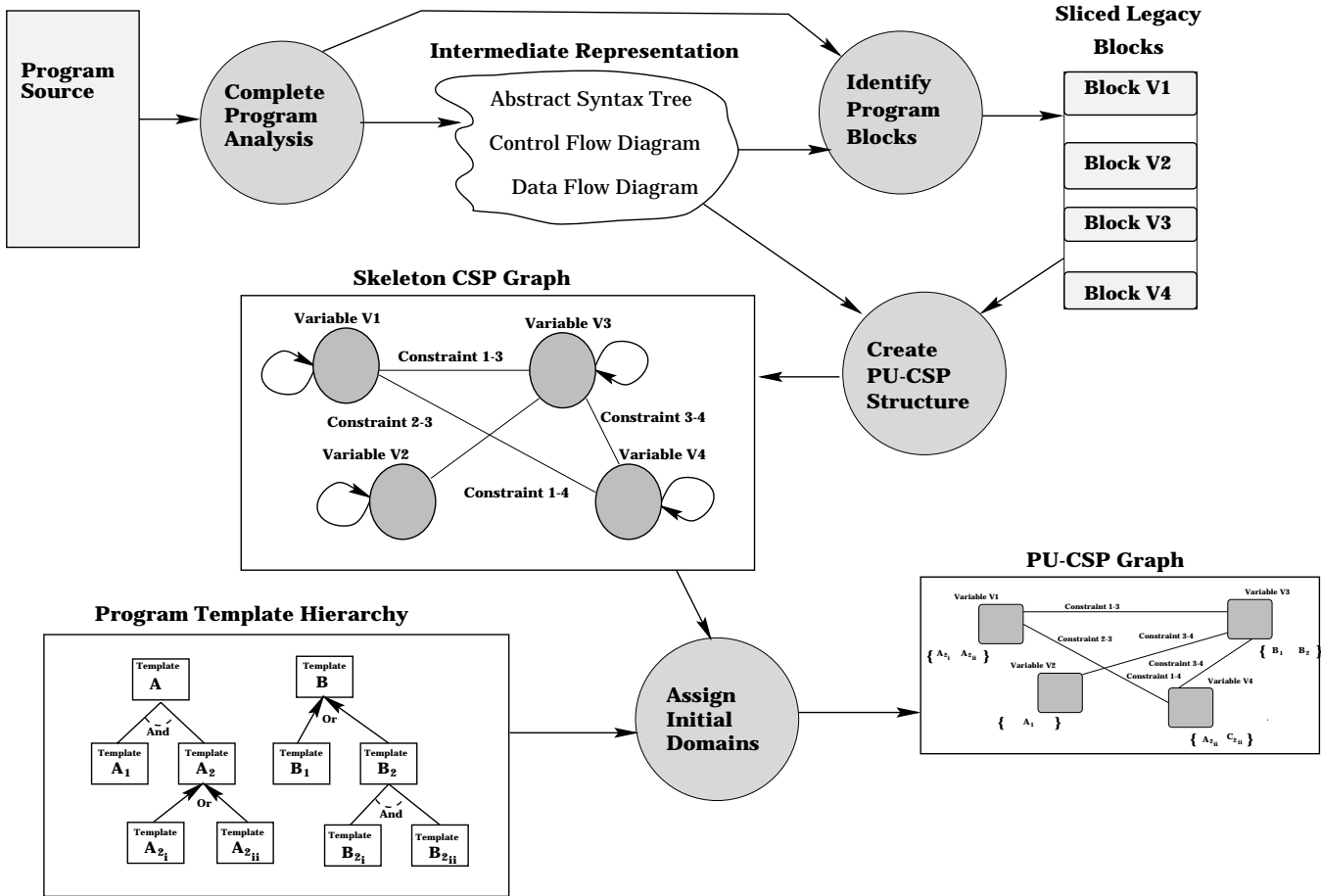


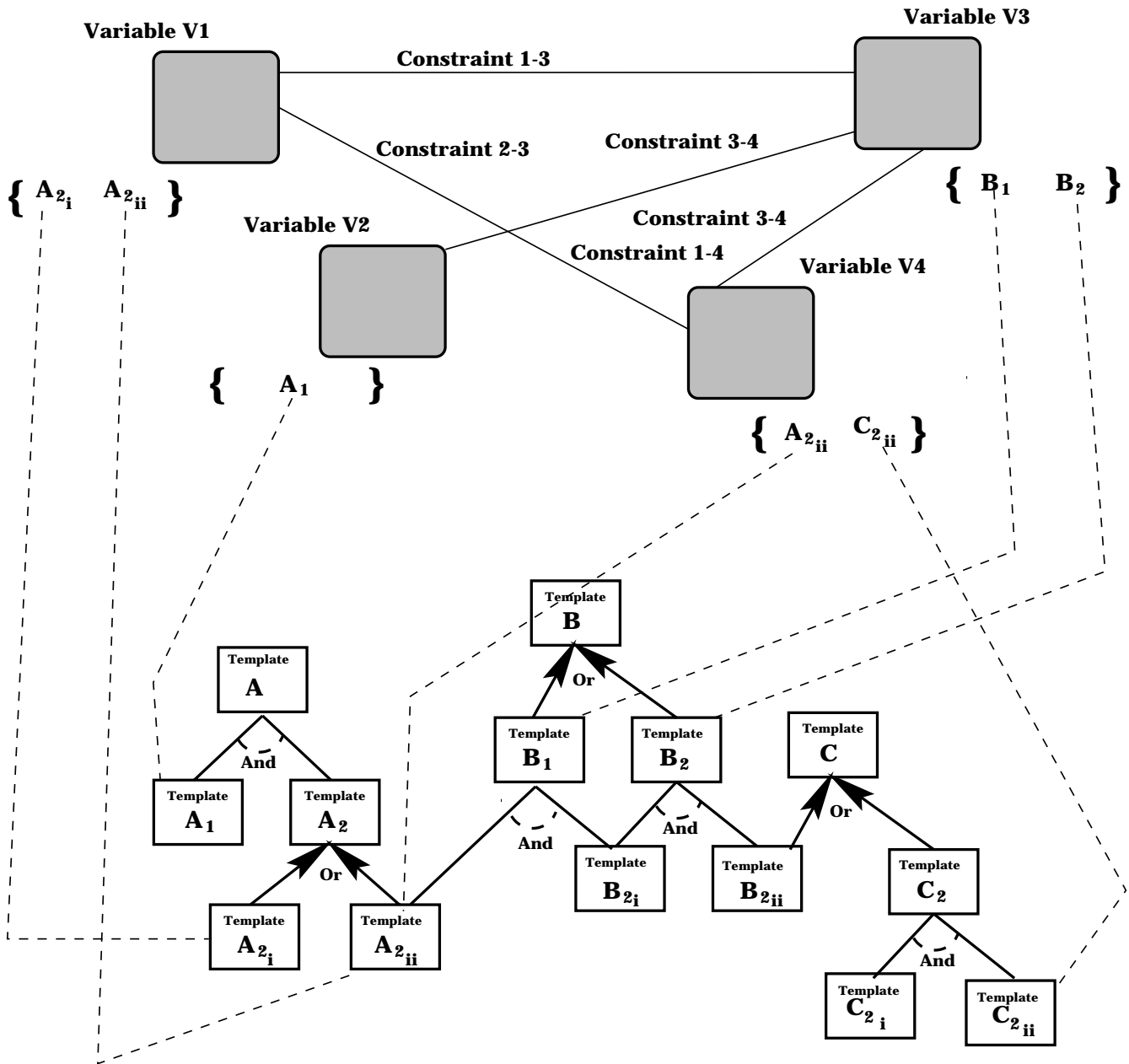
Figure 8: PUCSP Formulation; CSP Graph exploded in Figure 9.

Figure 9 shows an example formulated PU-CSP in which the domains of each variable are shown as instances identified in the program template hierarchy. During discussion of the PU-CSP we will heretofore discuss two distinct types of constraints: *structural constraints* depicted in Figure 9 as the inter-variable constraints, which are exactly those constraints derived from the intermediate source representation and which describe how program components are structurally related, and *knowledge constraints*, depicted in the figure as the compositional and specialization constraints in the program template hierarchy, which describe how program plans or templates may fit together to form larger (more abstract) plans in this domain.

The program template hierarchy is composed of hierarchically related plan templates (for a formalization of hierarchical planning knowledge base, see [33]). A template plan may be



## PU-CSP Graph (node consistent)



## Program Template Hierarchy

Figure 9: PUCSP Graph.

broken down into several sub-plans, in which case this is recorded as an **And** relationship between the sub-plans and the parent plan. Further, any required structure between the sub-plans such as necessary ordering, data flows between the sub-plans or control-flow between the sub-plans is recorded with the **And** relationship. Similarly, a template plan may be a specialization of another plan (or one of many such specializations), and in this case the constraints that constitute the specialization such as restriction of variable type or a particular restriction of data or control flow is recorded with the **Or** relationship. Figure 10 shows a simple **And** example in which **Template A** is composed of two subplans  $A_1$  and  $A_2$  where  $A_1$  provides the data flow  $r$  which  $A_2$  requires, and a simple **Or** example in which **Template A** may be specialized by either of the plans  $B_1$ , which also exports  $n$  in addition to the primary exports of  $B$  or  $B_2$ , which exports  $p$ .

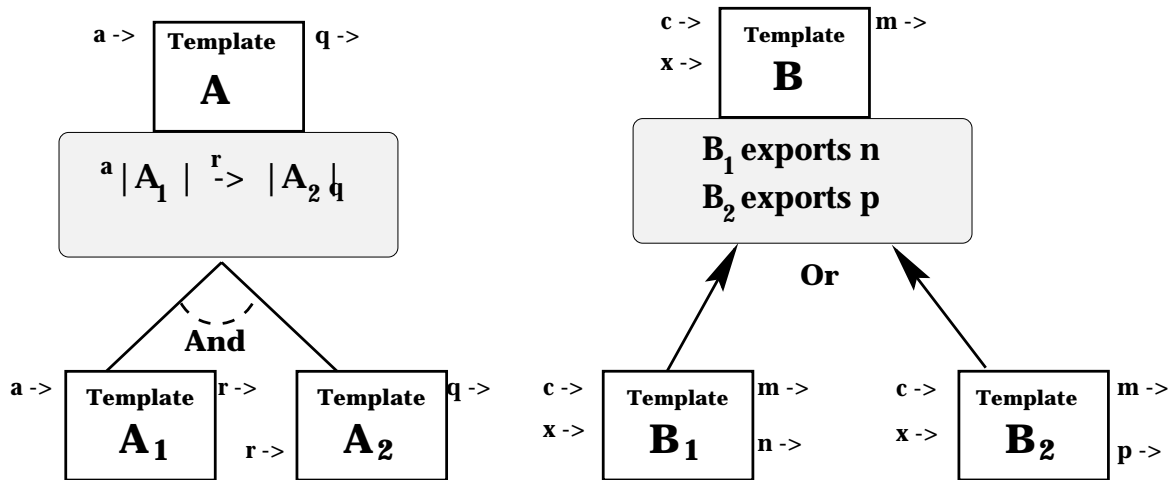


Figure 10: Library knowledge constraints.

## 4.2 More on Constraints

In a PU-CSP, the constraints among variables are of two types:

- *Structural* constraints are determined from the legacy code. They include such things as scope or called/calling relations, precedence relations, or shared information relations between component blocks. For instance, in the legacy source in Figure 3, the **print** statements appear within the scope of **for** statements, **declarations** precede their initial **assignment**, and print statements act upon array positions indexed by corresponding **for** statement indexes.
- *Knowledge* constraints are independent of the legacy code. They are program plans restricted in their relationship by the AND/OR structure given in the plan library. AND constraints are for composing program plans into higher level plans, and OR's

are for specializing an abstract plan in one of several ways. Assigning one program plan as an explanation of a particular PU-CSP variable thus constrains consistent assignments of other component variables.

As an example of a knowledge constraint mandated from the library structure, if a variable corresponding to program component **A** = “**string 1**” in Figure 3 were instantiated to program plan **builtin-char\*-copy** as shown in Figure 4, then it is consistent to assign the last **for-loop** variable an explanation of **print-string**, where the strings are the same.

A solution to the PU-CSP is an assignment to each variable by one program plan component in the plan library, such that no structural constraint from the source code, or knowledge constraint from the plan library is violated.

Representing program understanding as PU-CSP provides a convenient framework for interpreting Quilici’s index tests as constraint applications as part of search strategies typically used for solving CSPs. Specialization tests are specific instances of knowledge constraints that may be used to systematically reduce the range of domain variables in a hierarchical CSP. Inference tests identify “related” program plan templates according to earlier component instantiation, and can be interpreted as a special kind of variable-ordering heuristic.

## 5 Program Template Matching as CSP: MAP-CSP

We have seen how PU-CSP resolves integration of “local” explanations of source code blocks. We represent the process of matching particular abstract program plans to our legacy source as the MAP-CSP. We view MAP-CSP as an integral part of the more ambitious understanding task. Successful matches “locally explain” certain program blocks, and these local solutions can then be exploited to restrict the larger PU-CSP.

A MAP-CSP or program template matching problem can be stated as follows: given a plan template with a number of elements and constraints among the elements, find all instances of the template in a source code. As an example, consider finding all instances of an abstract data type in a C program. Figure 11 is a **String** ADT plan template taken from a plan library. The ADT is described in terms of 5 features describing various key components of a string class. In addition, there are constraints among the different parts as well, such as the one that requires one component to go before another.

We could model this problem as a CSP. For the given plan template (or ADT), each feature is a variable in our MAP-CSP. The *domain range* consists of all possible source program statements. Variables here can have attributes such as (**print,for**) that may be seen as *constraints* on allowable assignment of program statements (values) to template features (variables). Other *constraints* are on the sharing of information among variables, and on the order in which template features or variable are expected to appear in legacy source.

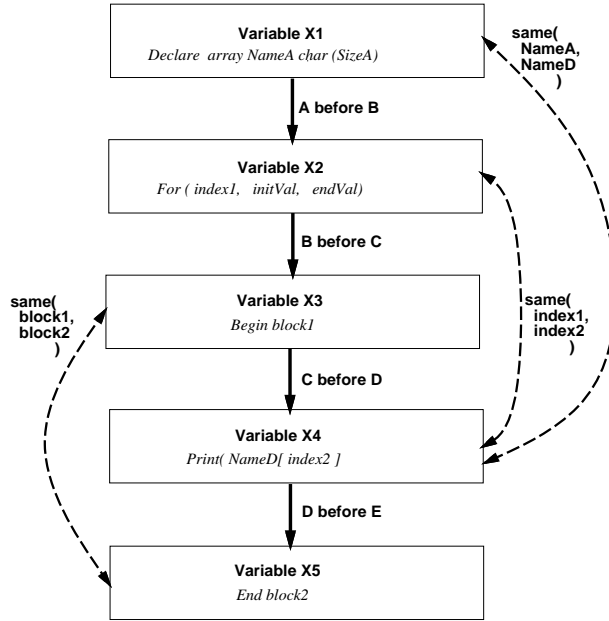


Figure 11: The **String** ADT in MAP-CSP.

A solution to the MAP-CSP consists of the set of all assignments of plan template features by source code statements, where each assignment must satisfy all constraints. As an example, consider the ADT of Table 1. When represented as a plan template as in Figure 11, the variables of the MAP-CSP are:  $X_i, i = 1, \dots, 5$ . Initially the domain for each variable ranges through all source statements in Figure 3. The constraints are as shown in the figure. The solution to this problem corresponds to the three alternative consistent assignments to the variables, one for each character string  $A$ ,  $B$  and  $C$ , respectively. Thus, the solution to a MAP-CSP provides a mapping that *explains* the matched source statements as parts of an instance of the abstract program plan or ADT.

## 6 Empirical Results of MAP-CSP

In this section we present and discuss experiments which are intended to show the feasibility of the MAP-CSP representation and related algorithms.

In Figure 11 a CSP is described for the **String** ADT. This CSP contains 5 variables each corresponding to a part of the program plans contained in the ADT in Table 1. The domain values are made up of source program statement blocks. For this test problem, there are 4 precedence constraints amongst the template variables, along with one additional constraint that the **begin** block corresponding to the **for** statement be within 15 program lines (the number 15 is arbitrarily chosen).

A test case is produced by instantiating 3 instances of the program template in a sample

source code, and by adding some variable amount of additional program statements as “noise” around the template instances<sup>2</sup>. For example, the legacy source example shown in Figure 3 contains three separate instances of the **String** ADT of Figure 11.

Our experiments are undertaken using a specific version of the generic CSP search algorithm. This algorithm has been implemented in Common Lisp on a SPARC server 470 workstation with four different heuristic configurations:

- Standard BackTracking
- Forward Checking with Dynamic Rearrangement
- AC-3 in advance of Standard BackTracking
- AC-3 in advance of Forward Checking

Each configuration has been applied to legacy sources ranging in size from 50 source lines to 1000 source lines, in 50 source line increments. Each increment was tested with 10 different random problem instances.

The MAP-CSP experiments are detailed in Figure 12<sup>3</sup> for *Standard BackTracking*, in Figure 13 for Forward Checking with Dynamic Rearrangement, in Figure 14 for AC-3 constraint propagation in advance of Standard BackTracking, and in Figure 15 for AC-3 in advance of Forward Checking. Each figure shows the number of CPU seconds required to find all template instances, with a 95% confidence interval charted.

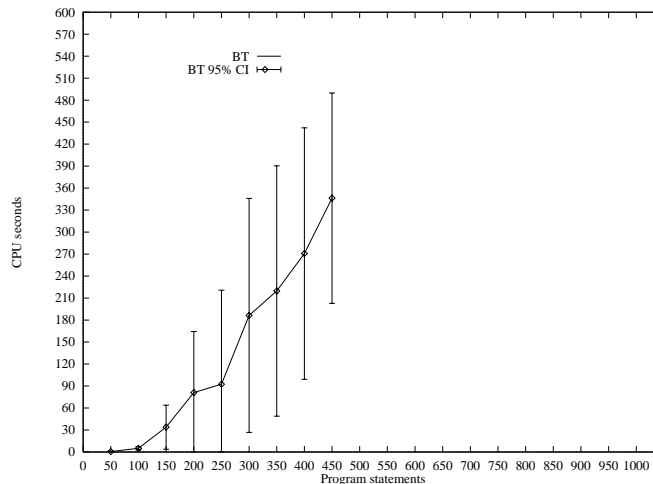


Figure 12: Standard BackTrack (95% conf. interval).

<sup>2</sup>The experiments presented in this paper are undertaken using artificially generated, potentially nonsensical, legacy code. Our ongoing research effort is focused on testing the integrated PU-CSP and MAP-CSP techniques to large existing commercial source libraries.

<sup>3</sup>If more than 2/10 of the experiments for a particular level of program statements did not complete in 600 CPU seconds, those results are not charted. This occurs in the results shown in Figures 12 and 14.

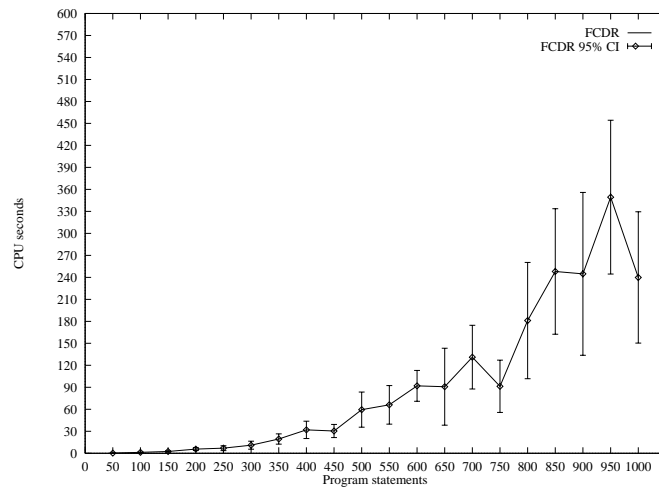


Figure 13: Forward Checking, DR (95% conf. interval).

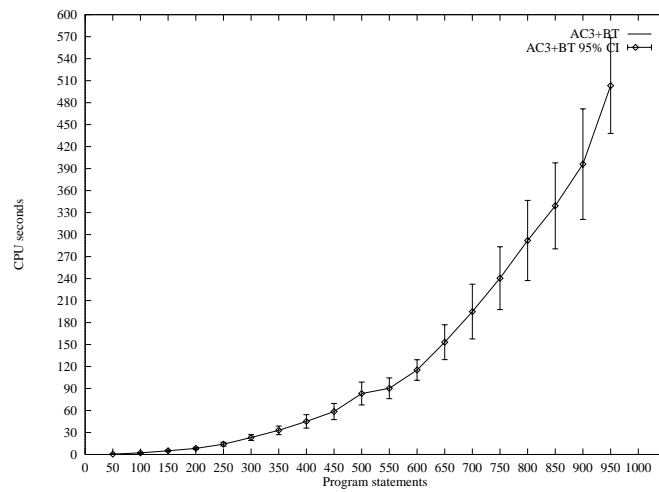


Figure 14: AC-3 with BT (95% conf. interval).

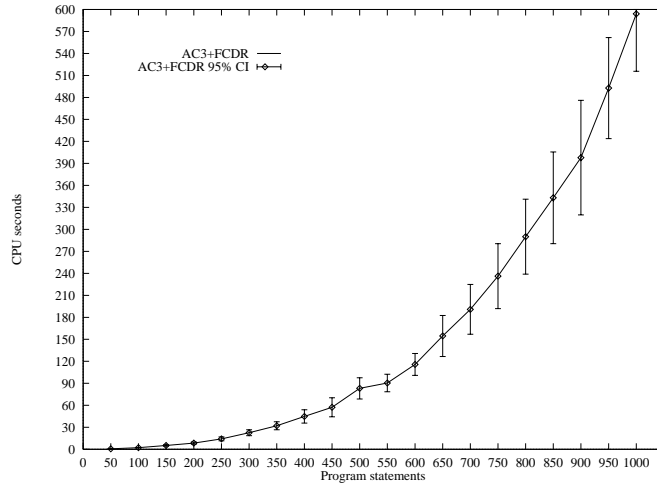


Figure 15: AC-3 with FC, DR (95% conf. interval).

We wish to demonstrate that the MAP-CSP representation and algorithm is capable of providing all-instance results in moderately sized program slices. An efficient MAP-CSP algorithm could make the execution of the larger PU-CSP algorithm more feasible. In addition, the MAP-CSP algorithm for template matching could potentially be stand-alone as a tool for assisting in the identification of legacy source portions that may be replaced with existing source library objects.

Several observations can be made from our test results:

- Standard Backtracking exhibited very unstable performance in examples of the same size. As hoped, more intelligent strategies behaved in a more stable manner. Forward Checking was considerably more stable, and the applications using AC-3 in advance of search exhibited very small variance across test cases of similar size. Stability is an important factor in any application that may be used as part of an online or interactive tool. In addition, Standard Backtracking was unable to complete in less than 600 CPU seconds for source instances exceeding 500 program statements.
- For legacy source examples of up to approximately 500 lines of code, the intelligent strategies located all instances of the ADT in approximately one minute of CPU time. In examples of up to 300 lines of code, all instances were identified in approximately 30 seconds. In such near-real-time circumstances it would appear that a tool could be fashioned that could be called up to run as a background process supporting an expert working with some legacy code.
- In experiments where the number of source lines exceeded 200, the appearance of *false solutions*<sup>4</sup> started to become apparent. These solutions arise through combinations of

---

<sup>4</sup>A “false solution” is a satisfying assignment of template variables to program parts such that the template

actual template instance components and nearby program statements that meet all of the constraints of the ADT. However, the number of false solutions never exceeded 10 in programs of 900 or less lines, and rarely exceeded 5 in smaller sources. These results suggest that *either* our template specifications need to be tightened somewhat so as to exclude these false solutions, or the system should be capable of interacting with an expert who may verify solutions before they are adopted. It is important to note that in the solution of the larger PU-CSP it is expected that these false solutions will be identified and discarded, primarily on the strength of *knowledge* constraint restrictions.

## 7 Conclusions

In this paper we have constructed a general representation of the program understanding task as a constraint satisfaction problem. Two versions of the task are identified: one is to find all instances of a given program plan template in a source code, and the other is to construct or verify an explanation of the source code in terms of a program plan library. In addition, we have modeled various search heuristics for program understanding as instances of a generic CSP search algorithm. We believe that the algorithm subsumes the previously proposed methods for the same problem, and can be systematically studied on a spectrum of heuristics.

We have also implemented the all-instances template matching problem, MAP-CSP and demonstrated that MAP-CSP can be solved for problems of non-trivial size using intelligent backtracking and constraint propagation within a reasonably stable and reasonably short time period. MAP-CSP has potential application both as a stand-alone tool for legacy code reduction and as a key component within the program understanding task.

We summarize some of the advantages of our approach below.

**Scalability** Our empirical results demonstrated that the MAP-CSP problem can be scaled up for legacy code of useful sizes. This efficiency gain is achieved by viewing the recognition problem as constraint satisfaction, and applying known constraint satisfaction algorithms. In our experiment, we haven't utilized the full range of constraints inherent in a program source code, such as those derived from program parsing, a technique employed by Kozaczynski & Ning[8] and Wills[28]. More extensive consideration is given to the specific use of these constraints in [30]. We expect the empirical results to improve further with use of these constraints.

**Usability** We envision our system as one part of a programmer's assistant toolset. For the MAP-CSP problem, a programmer could use the system to identify abstract program

---

constraints are satisfied, however, the found mapping is in fact not an actual instance of the program template. They arise as a result of overly abstracted template specifications.



plans in legacy programs up to around 500 lines of code in almost real-time, and can apply the system in batch-mode to much larger programs.

We are currently engaging in cooperation with a main telecommunications provider to investigate the applicability of this approach to extremely large source code in the telephony domain. Achieving partial automatic recognition of even 5% of the code would greatly benefit software maintainers.

We are currently implementing the search algorithm for PU-CSP. We expect to see similar effective results from constraining the search with hierarchical plan knowledge, particularly when this algorithm is fully integrated with the MAP-CSP solutions.

## Acknowledgments

We thank Alex Quilici and Jim Ning for their insight and comments and Grant Weddell for many helpful discussions. This research has been carried out with the support of the Natural Sciences and Engineering Research Council of Canada and the Institute for Computer Research (ICR).

## References

- [1] Sandra Carberry. Modeling the user's plans and goals. *Computational Linguistics*, 14(3):23–37, 1988.
- [2] Sandra Carberry. Incorporating default inferences into plan recognition. *Proceedings of the 8th AAAI*, 1:471–478, 1990.
- [3] Martin C. Cooper. An optimal k-consistency algorithm. *Artificial Intelligence*, 41:89–95, 1989.
- [4] E.C. Freuder. A sufficient condition of backtrack-free search. *Journal of the ACM*, 29(1):23–32, 1982.
- [5] R.M. Haralick and G.L Elliott. Increasing tree-search efficiency for constraint satisfaction problems. *Artificial Intelligence*, 14:263–313, 1980.
- [6] R. Holte, T. Mkadmi, R. Zimmer, and A. MacDonald. Speeding up problem-solving by abstraction: A graph-oriented approach. Technical report TR-95-07, University of Ottawa, March 1995.
- [7] Henry Kautz and James Allen. Generalized plan recognition. In *Proceedings of the Fifth National Conference on Artificial Intelligence*, pages 32–37, Philadelphia, Pennsylvania, 1986.

- [8] Wojtek Kozaczynski and Jim Q. Ning. Automated program understanding by concept recognition. *Automated Software Engineering*, 1:61–78, 1994.
- [9] Vipin Kumar. Algorithms for constraint-satisfaction problems. *AI Magazine*, pages 32–44, Spring 1992.
- [10] A.K. Mackworth. Consistency in networks of relations. *Artificial Intelligence*, 8:99–118, 1977.
- [11] Alan Mackworth, Jan Mulder, and William Havens. Hierarchical arc consistency: exploiting structured domains in constraint satisfaction problems. *Computational Intelligence*, 1:188–126, 1985.
- [12] Steve Minton. Quantitative results concerning the utility of explanation-based learning. *Artificial Intelligence*, 42:363–391, 1990.
- [13] R. Mohr and T.C. Henderson. Arc and path consistency revisited. *Artificial Intelligence*, 28:225–233, 1986.
- [14] H. Muller, K. Wong, and S.R. Tilley. Understanding software systems using reverse engineering technology. In *Proceedings of the Colloquium on Object Orientation in Databases and Software Engineering*, pages 88–98. The 62nd Congress of the "L'Association Canadienne Francaise pour l'Avancement des Sciences (ACRAS)", December 1994. May 16-17, 1995 Montreal, Quebec, Canada.
- [15] Hausi Muller, M.A. Orgun, S.R. Tilley, and J.S. Uhl. A reverse engineering approach to subsystem structure identification. *Journal of Software Maintenance*, 1993.
- [16] Hausi Muller, M. Tilley, M.A. Orgun, B. Corrie, and N. Madhavji. A reverse engineering environment based on spatial and visual software interconnection modules. In *Proceedings of the Fifth ACM SIGSOFT Symposium on Software Development Environments (SIGSOFT '92)*, ACM Software Engineering Notes, pages 88–98, December 1992.
- [17] Bernard A. Nadel. Constraint satisfaction algorithms. *Computational Intelligence*, 5:188–224, 1989.
- [18] Alex Quilici. A memory-based approach to recognizing programming plans. *Communications of the ACM*, 37(5):84–93, May 1994.
- [19] Alex Quilici. Toward practical automated program understanding. *Proceedings of the 1995 IJCAI Workshop on AI and Software Engineering (AISE-95)*, August 1995.
- [20] Alex Quilici and David Chin. A cooperative program understanding environment. In *Proceedings of the Ninth Knowledge-Based Software Engineering Conference*, pages 125–132, Monterey, CA, 1994.

- [21] Alex Quilici and David Chin. DECODE: A cooperative environment for reverse-engineering legacy software. In *Proceedings of the Second Working Conference on Reverse-Engineering*, pages 156–165, 10662 Los Vaqueros Circle, Los Alamitos CA 90720-1264, July 1995. IEEE Computer Society Press.
- [22] C. Rich and R.C. Waters. The Programmer’s Apprentice: A research overview. *IEEE Comput.*, 21(11):10–25, 1988.
- [23] C. Rich and R.C. Waters. *The programmer’s apprentice*. Addison-Wesley, Reading, Mass., 1990.
- [24] Spencer Rugaber, Kurt Stirewalt, and Linda Wills. The interleaving problem in program understanding. In *Proceedings of the Second Working Conference on Reverse-Engineering*, pages 166–175, 10662 Los Vaqueros Circle, Los Alamitos CA 90720-1264, July 1995. IEEE Computer Society Press.
- [25] Peter van Beek, Robin Cohen, and Ken Schmidt. From plan critiquing to clarification dialogue for cooperative response generation. *Computational Intelligence*, 9(3), 1993.
- [26] P. Van Hentenryck, Y. Deville, and C-M. Teng. A generic arc-consistency algorithm and its specializations. *Artificial Intelligence*, 57:291–321, 1992.
- [27] L. M. Wills. Automated program recognition: A feasibility demonstration. *Artificial Intelligence*, 45(2):113–172, February 1990.
- [28] L. M. Wills. *Automated program recognition by Graph Parsing*. PhD thesis, MIT, July 1992.
- [29] Steven Woods. A method of interactive recognition of spatially defined model deployment templates using abstraction. In *Proceedings of the Knowledge Based Systems and Robotics Workshop*, pages 665–675. Government of Canada, November 1993.
- [30] Steven Woods. A constraint-based approach to program plan recognition in software reverse engineering. Ph.D. Thesis Proposal, University of Waterloo, February 1995.
- [31] Steven Woods and Qiang Yang. Constraint-based plan recognition in legacy code. *Proceedings of the 1995 IJCAI Workshop on AI and Software Engineering (AISE-95)*, August 1995.
- [32] Steven Woods and Qiang Yang. Program understanding as constraint satisfaction. In *Proceedings of the IEEE Seventh International Workshop on Computer-Aided Software Engineering (CASE-95)*, pages 318–327. IEEE Computer Society Press, July 1995. Also appears in the *Proceedings of the 2nd Working Conference on Reverse Engineering (WCRE)*, July 1995.

- [33] Qiang Yang. Formalizing planning knowledge for hierarchical planning. *Computational Intelligence*, 6, 1990.
- [34] Qiang Yang. A theory of conflict resolution in planning. *Artificial Intelligence*, 58(1-3):361–392, 1992. Special Issue on Constraint-directed Reasoning.