

Refinement-Based Student Modeling and Automated Bug Library Construction

Paul Baffes
SciComp, Inc., 5806 Mesa Drive,
Suite 250, Austin, TX 78731, U.S.A.

BAFFES@SCICOMP.COM

Raymond Mooney
Department of Computer Sciences, Taylor Hall 2.124,
The University of Texas at Austin, Austin, TX 78712, U.S.A.

MOONEY@CS.UTEXAS.EDU

Abstract

A critical component of model-based intelligent tutoring systems is a mechanism for capturing the conceptual state of the student, which enables the system to tailor instruction to suit individual strengths and weaknesses. To be useful such a modeling technique must be practical in the sense that models are easy to construct and use, and that using the model actually impacts student learning. This research presents a modeling technique which can automatically capture novel student errors using common domain knowledge, and can automatically compile trends across multiple students. This approach has been implemented as a computer program, a machine learning technique called refinement, which is a method for automatically revising a knowledge base to be consistent with a set of examples. Using a knowledge base that directly defines a domain and examples of a student's behavior in that domain, models student errors by collecting any refinements to the correct knowledge that are necessary to account for the student's behavior. The efficacy of the approach is demonstrated by evaluating 100 students tested on a classification task covering concepts from an introductory computer programming language. Students who received feedback based on the models automatically generated by A significantly better on a post test than students who received simple reteaching.

1 Introduction

Student modeling has a long and interesting history, dating back well into the development of intelligent tutoring systems. The best method for constructing and using a student model is still the subject of much debate. Most student modeling techniques, however, have a similar goal, which might be defined as follows:

- Given: A set of expectations regarding student behavior in some domain and, a set of observations of a specific student's behavior on one or more tasks in that domain,
- Find: A representation accounting for any discrepancies between the expectations and observations that can be used as a basis for tutoring the student.

Ideally, a unique model is built for every student who interacts with the tutor. Capturing misconceptions specific to each student pre-programmed into the tutor. Using the student model, an ITS would then modify its feedback to suit the student's strengths and weaknesses, enabling it to be truly adaptive to the individual.

Unfortunately, the difficulty of constructing and testing student models has discouraged many researchers from pursuing further investigations into the field. Despite more than two decades of research has resulted in a wide variety of student modeling techniques, the practical task of incorporating these techniques into a functional system has proved to be a major roadblock. Furthermore, neither the utility nor the necessity of student modeling as a component of an ITS is a universally accepted. To the contrary, an interview of ten well-known ITS researchers which appeared in the 1993 issue of *Artificial Intelligence Magazine* came to the conclusion that "most of the researchers no longer believe in on-line student modelling." (Sandberg & Barnard, 1993). It is difficult to conclude that "instead of becoming more integrated, the field has diverged in the last few years. It appears that scientists in the field of educational psychology no longer share a research paradigm."

Thus the current challenge for student modeling is to show that modeling can be made both practical and effective. This is precisely the contribution of this work embodied in the ASSERT algorithm. Acquiring Expert Models by Finding Theories ASSERT was designed to show that student modeling is a viable tool in an effective tutoring system. By taking advantage of some of the latest machine learning, ASSERT is able to construct student models efficiently and automatically catching both expected and novel student misconceptions. Also, it is the first system which can construct bug libraries automatically using the interactions of students, without requiring input from the author, and integrate the results into future modeling efforts. In this sense, ASSERT is a self-improving student modeler. Finally, ASSERT can be used to significantly improve student performance, as will be shown in sections which follow.

The remaining sections are organized as follows. Section 2 reviews previous student modeling as a motivation underlying the design of ASSERT. Section 3 then presents a description of ASSERT focusing on the portion of the algorithm which captures individual student errors. Next, Section 4 describes how trends across a population of student models are automatically collected into a bug library and how such a library is then incorporated into the modeling process. Finally, Section 5 presents experimental results followed by discussion and conclusions in Sections 6 and 7.

2 Previous Work

2.1 Overlay Modeling

The application of artificial intelligence to the task of student modeling has led to the progression of techniques for collecting information both about what a student does not know, and about explaining the sources of misconceptions. The earliest student models, embodied in systems such as SCHOLAR (Carbonell, 1970), WEST (Brown & Brown, 1976) and WUSOR (Carr & Goldstein, 1977), used a form of modeling now generally referred to as overlay modeling. An overlay model relies on the assumption that a student's knowledge is always a subset of the correct domain knowledge. As a student performs actions which illustrate that he or she understands particular domain knowledge, these are marked in the overlay model. More sophisticated models can express a range of values indicating the extent to which the system understands a given topic using some form of truth-maintenance scheme (Finlayson, 1991). However the marking is achieved, typically the unmarked element

are used to focus tutoring on new problem areas for the student, or to ensure the domain.

The advantage of the overlay is its simplicity; the elements of the model directly on to the knowledge used to engineer the system. The disadvantage, a restriction placed on the model—only missing elements of the core can be modeled. Alternative notions which a student might have cannot be captured, this means that concepts cannot be modeled. Thus, overlay models can capture the notion of a student's lack of knowledge, but they cannot be used to model who knows of a topic but misunderstands it.

2.2 Bug Libraries

To address the limitation of overlay models, other researchers focused on databases of student misconceptions typically represented as classic bug-library work was done by Brown, Burton and VanLehn (Brown & Burton, 1978; Burton, 1982; VanLehn, 1980), Sleeman and Smith (Sleeman & Smith, 1981), and Young and O'Shea (Young & O'Shea, 1981), but a number of other systems can be said to incorporate forms of stored misconceptions (Rich, 1989; Liang & Taotao, 1991; Miller, 1977; Quilici, 1989; Soloway & Johnson, 1984). With a bug library, models match student behavior against a catalog of expected bugs which are handled through an analysis of student errors.

The idea is a very powerful one, especially if specific responses can be encoded. However, two important problems remain with the bug-library approach. First, the construction of such catalogs by hand is a consuming task which must be repeated for every new domain. Second, even if taken, the resulting library may still fail to cover a wide enough range of misconceptions successfully. That is, a student may exhibit a misconception which was not in the library. As with overlay models, the static nature of them is incapable of modeling unanticipated student behaviors.

2.3 Dynamic Modeling

To capture novel student misconceptions, one must turn to some kind of search space of possible bugs. Two methods have been tried to date: one attempts to build a library and the other attempts to infer a model of the student from scratch using machine-learning techniques. In both cases, novel errors are modeled by capturing buggy information dynamically, using data from a student's behavior to bound the search space.

2.3.1 Extending a bug library

Sleeman et al. (Sleeman et al., 1990) describe two extensions to their PIXIE INFER* and MALGEN, both of which can be used to extend a bug library. PIXIE is a tutoring system designed for the domain of high-school algebra whose goal is to provide appropriate feedback to improve student performance. PIXIE's underlying representation is a state-space paradigm, where the domain theory is a set of operators implemented as mal-rules. Misconceptions which comprise the bug library are encoded as mal-faulty rule rules. Both INFER* and MALGEN attempt to generate new mal-rules when the student exhibits a problem that cannot be modeled using the mal-rules already in the library. The difference between the two extensions is that INFER* attempts to patch specific student solutions, whereas MALGEN generates and tests new mal-rules by altering the domain theory.

domain theory. INFER* uses the rules it has to work forward from the problem backward from the student's solution as far as gaps are filled by maintaining a new mal-rule. In MALGEN, formalized perturbation operators are used to change the domain theory to generate new mal-rules for the bug library.

The disadvantage of both systems is their reliance upon a user to decide which rules are appropriate extensions for the bug library. To their credit, the systems allow the user to issue and discuss potential filters that might be used to cut down on the number of rules presented to the user. Unfortunately, to this point no general-purpose system which might be usable across domains has been found. In the end, the user is presented with a number of proposed mal-rules and must decide which ones are the "keepers," since both systems were developed to address the problem of extending the bug library, they remain strongly tied to the arduous task of preconstruction by hand.

2.3.2 Modeling by Induction

In an effort to avoid the cost associated with hand-constructed bug libraries, researchers have turned to machine learning. Their inspired work on the ACM system (Langley 1984; Langley et al., 1990) was the first effort to harness machine learning for the diagnosis of misconceptions through induction. The underlying idea is to invent a student model on the fly by automating the process of analysis (Newell & Simon, 1972) which is used by human instructors as a means of unearthing student misconceptions.

ACM uses a domain-independent induction algorithm to induce control knowledge representing how students apply operators in a given domain. Starting with a set of general operators, the goal is to induce a set of control rules that will generate a sequence that produces the same solutions as the student. To model a particular student, ACM starts with only general knowledge of how the operators can be applied for a path of operator applications connecting the problem specification to the solution. Given this "solution path," induction is performed by noting whether each operator application is used by an operator application lies on or off the path. The output of induction is a set of conditions which predict when an operator will produce a state which lies on the solution path. The conditions found by induction are then used to specialize the general operators. The result is a procedure that models the student's unique problem solving behavior.

By using induction, ACM can operate automatically, constructing models of both correct and buggy knowledge. However, because the operators must initiate a search process enough to model many kinds of behavior, both correct and incorrect, the problem space is huge. Langley et al. note this, and suggest various "psychological constraints" which can be applied to the operators to limit the search. They also suggest additional underlying first principles, representing a deeper level of knowledge which provide the basis for automated addition of such psychologically plausible constraints. However, the system is still fundamentally limited by the complexity of having to build a model completely from scratch. This can only be remedied by collecting large amounts of data on each student or by imposing further constraints on the search space. The latter require finding such constraints by using the very human-intensive methods that the system is designed to avoid.

2.4 Tracing Techniques

One final style of student modeler bears mentioning because it represents a

ious techniques described to this point. In what might be termed the underlying philosophy is to follow along with the student during his or her problem solving, stopping whenever the student deviates from the correct procedure. A tracing system must have knowledge of both correct and incorrect actions like a modeler and must also have a mechanism for reproducing the steps followed by the student's solution paths. By focusing on the correct path as a bias, tracing systems can be very efficient.

The pioneering efforts in this area are those of Anderson (Anderson, 1983; Anderson et al., 1985; Reiser et al., 1985), which follow student behavior during interaction with the tutor to occur through menu selection. However, building a model of student is not the primary focus of model tracing. Instead, buggy information system acts mostly to support the tracking process, with the goal being to keep the student from straying too far off the correct path. Other tracing systems utilize a logic-based approach (Costa et al., 1988; Ikeda & Misoguchi, 1993; Hoppe, 1994) where the modeler uses an analytical approach, such as deduction or resolution, to search through a space where a misconception lies. Essentially, whenever the rule-based modeler fails to produce a "proof" which mimics the student's actions, the points where the proof fails are used as dates for querying the user about his or her beliefs. Again the emphasis is on detecting deviations from an expected path of correct student behavior.

Unlike the previous methods, tracing techniques do not dynamically construct models. Instead, they rely upon either the assumption that the student can follow along the correct path or querying an oracle whenever a deviation is detected. They lack the ability to handle novel student misconceptions independently.

3 Refinement-Based Modeling

This previous work on student modeling has resulted in three important ideas for research presented here. First, modeling systems can increase their coverage of student behavior by incorporating knowledge from a library of expected misconceptions. Second, to be truly adaptive and to avoid the costs of bug library construction, one needs a sort of dynamic modeling or learning algorithm. And third, tracing student behavior in comparison to expected correct behavior can be an effective tool for detecting deviations without requiring a great deal of searching. These ideas by using a relative to the new machine learning technique of theory refinement (Ginsberg, 1990; Ourston and Mooney, 1994; Craw and Sleeman, 1991; Towell and Shavlik, 1991). Theory refinement is a method for automatically revising a knowledge base to be consistent with a set of examples. Typically, the knowledge base is considered incorrect or incomplete, and the modeler presents correct behavior which the knowledge base should be able to emulate. The refinement procedure itself is blind to whether or not the input knowledge is correct in any absolute sense; the theory-refinement process merely modifies the knowledge base until it is consistent with the examples. Thus, one can also use theory refinement to refine a correct knowledge base and examples of behavior, and theory refinement will introduce whatever modifications are necessary to cause the knowledge base to be consistent with the incorrect examples.

Theory refinement, then, provides a basis for the development of a modeler. Starting with a representation of the correct knowledge of the domain and a set of examples of erroneous student behavior, theory refinement will revise the knowledge base to make it consistent with the student, i.e., introduce "faulty" knowledge to

student's mistakes. The refinements made to the knowledge base then represent a model student, and can be used directly to guide tutorial feedback by comparing with whatever elements of the correct knowledge base they replaced.

Using theory refinement, ASSERT combines the methods used in previous modeling systems. A theory-refinement learner combines the power of both analytic (as in INFER) and empirical (as in ACM) learning techniques in an integrated, domain-specific way. ASSERT can model any misconception consistent within the primitive domain. And finally, ASSERT provides an extension to theory refinement that combines the results of different student models from different students. This mechanism allows ASSERT to construct a bug library automatically, without the necessity on the part of the author. Section 4 describes this algorithm in detail. Finally, we draw our attention to the mechanism of theory refinement and its role in the design of ASSERT.

3.1 Outline of Theory Refinement

Having outlined the philosophy behind ASSERT, we can now turn our attention to the theory refinement algorithm around which ASSERT is constructed. It is important to point out from the start that basic design of ASSERT is tied to a particular theory refinement algorithm. Other theory refinement systems which differ from the one presented here can provide ASSERT with different or enhanced capabilities.

ASSERT uses the HETHER algorithm (Baffes, 1994; Baffes & Mooney, 1993) which is based on the HETHER theory-refinement system (Ourston & Mooney, 1990; Ourston, 1991). HETHER was chosen as a starting point because it was the most complete symbolic theory refinement system available. HETHER is designed to work with a propositional Horn-clause knowledge representation. It takes two inputs, a propositional Horn-clause theory which is repaired using a set of examples. The examples are assumed to be lists of feature-value pairs chosen from a domain of features. Each example has an associated category which should be provable using the theory with the feature values in the example. HETHER can generalize or specialize a theory, without user intervention, and is guaranteed to produce a set of refinements which are consistent with the examples.

Figure 1 shows an example theory and four input examples. The top of the figure shows part of a rule-based theory built for teaching a subset of C (for a complete listing of the rules, see Appendix A). The rules, numbered R1 through R7, consist of a consequent which is considered true for an example only when the conditions to its right are provable from the feature values of that example. Propositions represent either intermediate concepts or are shorthand for binary feature values as a value. This simplified theory has only one category, "compile-error," and is used to classify examples. The input examples, shown in the table below the rules, are classified as compile-errors only if they can satisfy rules R1 or R2 or R3 or R4 or R5 or R6 or R7. Under a closed-world assumption is used to classify the example as non compile-error, example 1 is correctly classified as a compile-error because it satisfies either R6 or R7. Likewise, example 2 fails to satisfy any of the rules and is correctly classified as non compile-error.

However, examples 3 and 4 are misclassified by the theory in its current state.

1. Keep in mind that the language used here is highly subjective in nature. One need not take any actions are "mistakes." The central point is that theory refinement can be used to detect actions which are inconsistent with its given knowledge base.

- R1: ~~compile-error~~constant-not-init
- R2: ~~compile-error~~constant-assigned
- R3: constant-not(~~pointer~~ constant(~~pointer~~ pointer-init false))
- R4: constant-not(~~integer~~ constant(~~integer~~ integer-init false))
- R5: constant-assigned(~~integer~~ constant(~~integer~~ integer-init integer-set yes))
- R6: constant-assigned(~~integer~~ constant(~~integer~~ integer-init integer-set through-po))
- R7: constant-assigned(~~pointer~~ constant(~~pointer~~ pointer-init pointer-set))

	Example 1	Example 2	Example 3	Example 4
compile-error	true	false	true	true
pointer	constant	non-constant	non-constant	non-constant
pointer-init	true	false	true	false
pointer-set	true	true	true	true
integer	constant	non-constant	non-constant	non-constant
integer-init	true	true	true	true
integer-set	through- pointer	yes	no	no

Figure 1: A Theory and Examples. The desired classification is shown in italics (thus, Examples 3 and 4 are misclassified).

rule base is "incorrect" since it does not produce the desired classification for the examples. Propositional Horn-clause theories can have four types of errors (Figure 2). An overly-general theory is one that causes an example to be classified other than its own, i.e., a false specialization. Existing antecedents, add new antecedents, and deletes rules to fix such problems. An overly-specific theory causes an example not to be classified in its own category, i.e., a false generalization. Existing antecedents and learns new rules to fix these problems. Theory-refinement systems that are subject to local minima are designed to fix any arbitrarily incorrect propositional Horn-clause theory (Ourston, 1991).

Symbolic theory-refinement systems like EITHER use a combination of three computations to determine how to modify a theory. The first step is to analyze a single failing example by analyzing the rule base to determine what needs to be changed to fix the theory for the example. For a failing positive example, a set of antecedents is found which, if deleted, will fix the theory for that example. For a failing negative example, a set of rules is computed which, if deleted, will fix the theory for that example.

The second step is to repair for a single example against all the other examples to see if the repair will cause new misclassifications to occur. If not, the repair can be applied to the theory directly. Otherwise, a third step is taken using a set of additional conditions which will separate the examples fixed by the repair from the examples for which the repair causes new misclassifications. These additional conditions are then used to modify the repair.

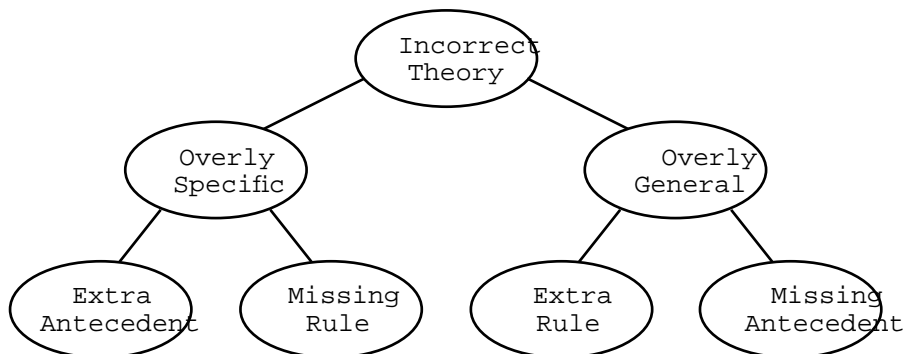


Figure 2: Theory error taxonomy for propositional Horn-

As an example, in Figure 1, notice that both example 3 and example 4 are examples since neither is classified as compile-error. This indicates that t specific and must be generalized. One way to repair the theory for example delete the "(pointer constant)" condition from rule R7. This allows rule R7 the example, without hindering the classification of example 1, and without to become so general that it would be satisfied for example 2. Testing th. examples 1, 2 and 4 yields no new misclassifications, and it can be applie theory.

Finding a repair for example 4, however, yields a different result. The to delete the "(pointer constant)" condition from rule R3. However, when th against examples 1, 2 and 3, example 2 is erroneously classified as compile way to fix the repair is to remove the "(pointer constant)" condition fron example 4, plus add a new condition to the rule which keeps example 2 from rule. Passing examples 2 and 4 to an induction algorithm would return "(in the condition which can discriminate between the examples. The final revis which correctly classifies all four examples, is shown in Figure 3.

Notice that the repairs chosen for examples 3 and 4 in Figure 3 are not repairs for these examples. For instance, example 3 could have been classi error by removing the conditions "(integer constant)" and "(integer-set t from rule R6, or by removing the conditions "(integer constant)" and "(i from rule R5. For that matter, removing all of the antecedents from any on through R4 would also have repaired the theory for example 3, by making ei pile-error" or "constant-not-init" concepts provable by default. In fact, c ble repairs for an example in the general case is exponential in the si Consequently, the way in which repairs are calculated, as well as when a re the theory in relation to computing repairs for other examples, can have a on the accuracy and performance of the theory refinement algorithm.

NEITHER and I~~HER~~ differ in the way the repair, testing, and induction comp theory refinement are combined. While a detailed comparison between the two scope of this article (see Baffes, 1994, chapter 3; Baffes & Mooney, 1993), NEITHER algorithm is provided in E~~Figure 4~~ For speed in computing repairs, focusing on quickly finding one good repair for each example. The goal is to est repair in the deepest possible part of the theory. After converting the


```

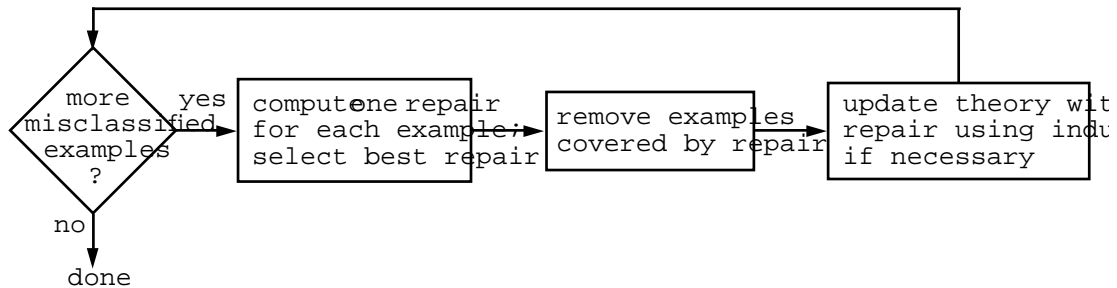
R1: compile-errorconstant-not-init
R2: compile-errorconstant-assigned
R3: constant-not-(pointer constant)pointer-init false)
R4: constant-not-initinteger constantinteger-init false)
R5: constant-assignedinteger constantinteger-initinteger-set yes)
R6: constant-assignedinteger constantinteger-initinteger-set through-pc
R7: constant-assignedpointer constantpointer-initpointer-set
-----
R1: compile-errorconstant-not-init
R2: compile-errorconstant-assigned
R3: constant-not-(pointer constant)pointer-init false)integer-set no)
R4: constant-not-initinteger constantinteger-init false)
R5: constant-assignedinteger constantinteger-initinteger-set yes)
R6: constant-assignedinteger constantinteger-initinteger-set through-pc
R7: constant-assignedpointer constantpointer-initpointer-set

```

Figure 3: Example of refinement. Above the dashed line : base before refinement; below are the rules after refinement shown in boldface.

NEITHER uses a recursive routine which starts at the leaf rules of the theories. The repairs are collected at each rule and passed up to parent rules. NEITHER always chooses the smaller repair, picking randomly to break ties. Each repair is applied only once, making the repair computation linear in the size of the theory.

Once a repair has been calculated for each example, NEITHER selects the best repair from among the set to apply to the theory. A selection is made by temporarily modifying the theory with each repair, calculating how many examples it fixes and how many new misclassifications it causes. These results are combined with the size of the repair. The repair which fixes the most examples with the fewest new misclassifications is selected. This repair is then tested against the rest of the examples, and induction is performed.



Criterion for computing a repair for ONE example
 Find the deepest, shortest, repair which causes the fewest new misclassifications

Criterion for selecting best repair from AMONG examples
 Select the shortest repair fixing the most examples with the fewest new misclassifications

Figure 4: NEITHER main loop.

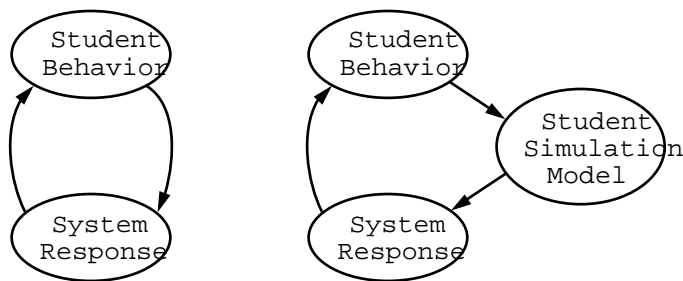


Figure 5: Abstract view of student-tutor interaction

to modify the repair to avoid new misclassifications. The whole process is a loop which continues until all misclassified examples have been accounted for. EITHER runs more than an order of magnitude faster than using accuracy (Baffes & Mooney, 1993), giving response times that are on the order of seconds. This is critical to an interactive tutoring environment where feedback must be given to the student in a timely fashion.

3.2 Overview of ASSERT

Having reviewed the basics of theory refinement, we can now turn to the details of ASSERT. ASSERT views tutoring as a process of communicating knowledge to a student. The contribution of the modeling subsystem is to pinpoint elements of the internal knowledge base to be communicated. At its most abstract level, such a tutorial can be viewed as an interaction between the student and the system as shown in Figure 5. Many details are omitted from the diagram and different researchers have chosen to emphasize different parts of the interaction. A first design decision then is one of emphasis: on the question of how to construct a useful interpretation of the student's actions. This decision is dependent on the component inserted into the diagram as shown in the right half of Figure 5. The student simulation model component models that the system contains a knowledge base that can be used to solve problems in the same context as the student. The student's knowledge base can be modified to replicate the solutions furnished by the system.

3.2.1 The Student as a Classifier

Figure 6 depicts how ASSERT views student behavior. It is assumed that all actions a student can be broken down into classifications. That is, given a set of inputs, called problems, the student will produce labeled examples which classify each of the problems into a category. Each problem consists of a feature vector describing some aspect of the problem. The task of the student is to produce a label for each feature vector, selected from among some predetermined set of legal labels. The resulting set of labeled examples pairs each feature vector with a label selected by the student.

In its simplest form, a problem consists of a single feature vector presented in a multiple-choice format, where the answers available to the student are a list of possible categories. Thus, for example, the classic diagnosis procedure



Figure 6: Student behavior diagram.

a patient's symptoms (the feature vector) and ask the student to select a d of diseases (the label). This allows us used in concept learning domains, whi are common applications for automated training systems. It also means that will translate directly into a form usable by theory refinement, which requi ples as one of its inputs.

3.2.2 Modeling by Theory Refinement

Once collected, the labeled examples generated by the student are passed refinement component depicted in Figure 7. As discussed previously, refinement will take an incoming knowledge base, plus an incoming set of e refine the knowledge base until it is consistent with the examples. The N refinement system is used to add or remove rules or parts of rules until th duces the same answers as the student, i.e., will classify each feature ve category label as the student. The resulting refined rule base is thus able dent's behavior.

The use of a propositional theory-refinement algorithm for modeling carr. assumption that the author of the tutoring system will be able to provide a tion of the domain using propositional Horn-clauses. This places a good dea how the correct rule base is constructed since it becomes the language ASSERT interprets the student's actions. If the correct rules are expressed low a level of detail, the ability of the system to form accurate models wi. course, this type of knowledge representation problem exists for all tutor. ever, ASSERT gains an advantage by purposely isolating the correct domain kno separate component: the author can easily change the focus of the tutor by rect rule base. Moreover, if students possessing different levels of unders tutor, multiple rule bases can be written to give the system more flexibilit

3.2.3 Refinement-Based Remediation

The last component of the system response, is outlined in Figure 8. Us

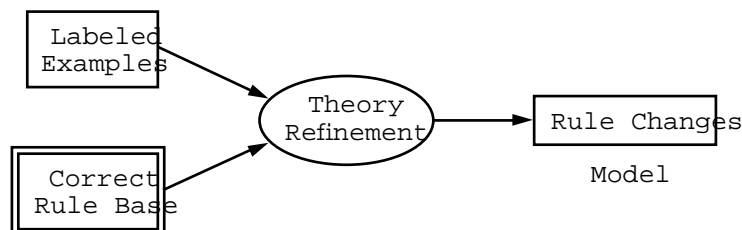


Figure 7: The student simulation model.

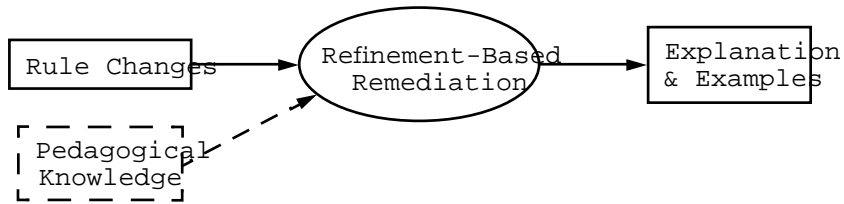


Figure 8: System response diagram.

refinements produced by N ASERT generates explanations and examples to reinforce the correct form of the rule or rules modified. The underlying approach, called refinement-based remediation, is based on fundamental units of explanation called *refinements*. Rather than implementing any particular remediation strategy, the system provides the most elementary information required to explain a refinement with one or more examples. For each refinement detected by NER, ASERT provides two functions: the ability to explain a correction to the rule which was changed, and the ability to generate an example which uses the designer of a tutoring system using the option to generate multiple explanations or examples, to determine the circumstances when such feedback is given, whether the system or the student controls which explanations and examples are generated. By providing such explanation-example pairs, the system supplies the raw materials for a variety of remediation techniques.

The specifics of generating explanations and examples for each refinement are described in detail in (Baffes, 1994), but the underlying idea is straightforward. Explaining a refinement involves describing how the correct form of the rule (not the revised version) fits into the current rule base. Each rule has an associated piece of stored text, describing the rule's role in the rule base. A full explanation is generated by chaining together the stored text for each refinement lying on the proof path for the correct label (not the student's label) for the given feature vector. The label which is produced by the correct rule base for the given feature vector is the label for the refinement.

Generating examples is a bit more complicated since they are constructed rather than being drawn from storage. Recall that each refinement made by NER involves the addition or deletion of propositional literals from a rule in the theory. In addition, the literals to be added or deleted as well, but this is the same as adding or removing literals from a rule. Using normal deductive methods, the added and removed literals can be traced back to the original feature vector. The result is a set of conditions in the feature vector which are necessary for the refinement, ignoring or a set of extra conditions not present in the feature vector which are necessary for the refinement. ASERT can thus generate an example which is correct in every way except for the added or missing conditions in the refinement. The result is then be counter example to the student, and the various added or missing conditions are explained. Note that this corresponds very closely to tutorial methods outlined for counterexamples (Tennyson, 1971).

Figure 9 shows an explanation and example pair, corresponding to one of the refinements depicted previously in Figure 3. Recall that the last rule of Figure 3 was one of its antecedents. Figure 9 is the feedback generated for this missing antecedent, illustrating how the condition represented by the antecedent is essential to drawing a conclusion. The top half of Figure 9 shows the text which explains how the refinement fits into the overall rule base. As part of the explanation, the three conditions of the refinement and its three correct antecedents, are itemized at the end of the explanation.

EXPLANATION

One way to detect a compilation error is to look for an identifier declared constant and initialized, then later assigned a new value.

A constant identifier is erroneously assigned when it is declared a pointer to an integer, initialized to the address of some integer, and then later set to the address of another integer. It does not matter if the identifier is a pointer declared to point to a constant integer or a non-constant integer. Once a constant pointer is initialized it cannot be reset to point to another integer.

Specifically, note the following which contribute to this type of error:

- * There must be a pointer declared to be constant (but not necessarily pointing to a constant object).
- * A pointer declared to be constant must be initialized.
- * A pointer declared a constant and initialized must be set after initialization.

Here is an example to illustrate these points:

Example

Here is an example which might appear to be a compile error but is actually CORRECT:

```
void main()
{
  const int x = 5, y, w, *z = &x;
  z = &w;
  cin >> w >> y;

  cout << ((y *= x) || (y > w));  cout << (w -= x);
}
```

This example is NOT a compile error because:

- * The pointer 'z' is declared as a NON-CONSTANT pointer to a constant integer, so it does not have to be initialized and it can be assigned a new value.

Figure 9: Example remediation given to a student.

generated which highlights how the "(pointer constant)" condition bears answer to the example, showing how the truth or falsity of the condition leads to the conclusion.

Finally, Figure 10 combines the components from Figures 6, 7 and 8, showing dialog flows between the student and the system. Problems given to the student are represented by labeled examples, which are passed to the system. The system uses these to refine a rule base representing correct knowledge of the domain to produce a modified rule base for the student. The refinements are then used to generate explanations and examples which gets passed back to the student.

4 Extending SMART's Modeler

The previous three sections have described the algorithm, showing how

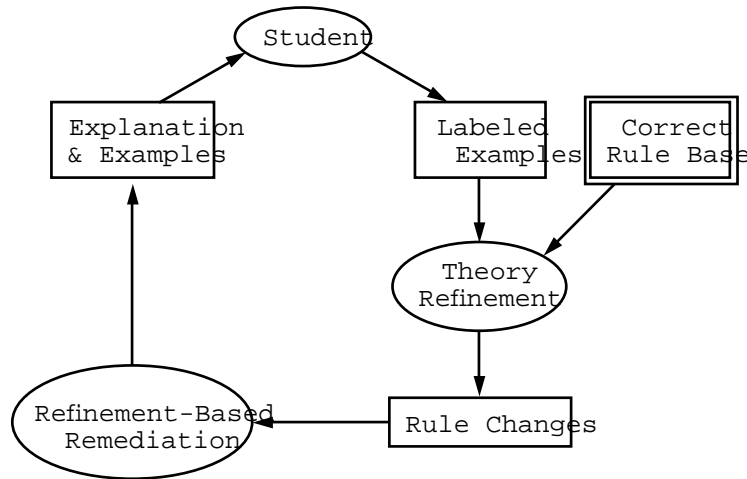


Figure 10: Basic design of the algorithm.

the flow of information between student and system can be implemented as refinements that highlight the differences between how the system and the student solve the same set of problems. As such this system constructs a student model by tracing the student's behavior against a known standard rule base. Nothing is said about how multiple student models are mined to construct a bug library. The bug library is incorporated back into the modeling process.

4.1 Building a Bug Library

A bug library represents a collection of data gleaned from the interactions of students with the tutoring system. The rule changes output from theory refinement for each student as the basis for constructing its bug library. This gives the system a way to communicate back to the author. First, the rule changes are closely related to the type of input generated by the tutoring system. Since a rule base must be supplied as input, expressing changes to that rule base is an effective way to communicate back to the author. Second, the rule-change format precisely captures what the student did. A bug library built of rule changes is the data which can be incorporated directly into the modeling process.

ASSERT constructs a bug library in three stages. First, it collects copies of rule changes from all the student models together, eliminating any duplicates. Second, it ranks each rule change by a measure of how frequently it occurs among the various models, called the stereotypicality of the rule change. Third, in the process of ranking the rule changes, ASSERT tests generalizations of the change to see if they result in better stereotypicality. If a generalization is found which has superior stereotypicality, the original rule change is replaced by the generalization. The final bug library contains the best generalization for each rule change with any duplicates removed.

Figures 11 through 14 illustrate how a bug library is constructed for a student model generated for illustration purposes. Figure 11 shows how stereotypicality is calculated for models. Each model is shown at the top of the diagram, each of which changes only one rule in the rule base. All the models alter the same rule, but in different ways. The first model changes the rule by deleting the set of antecedents

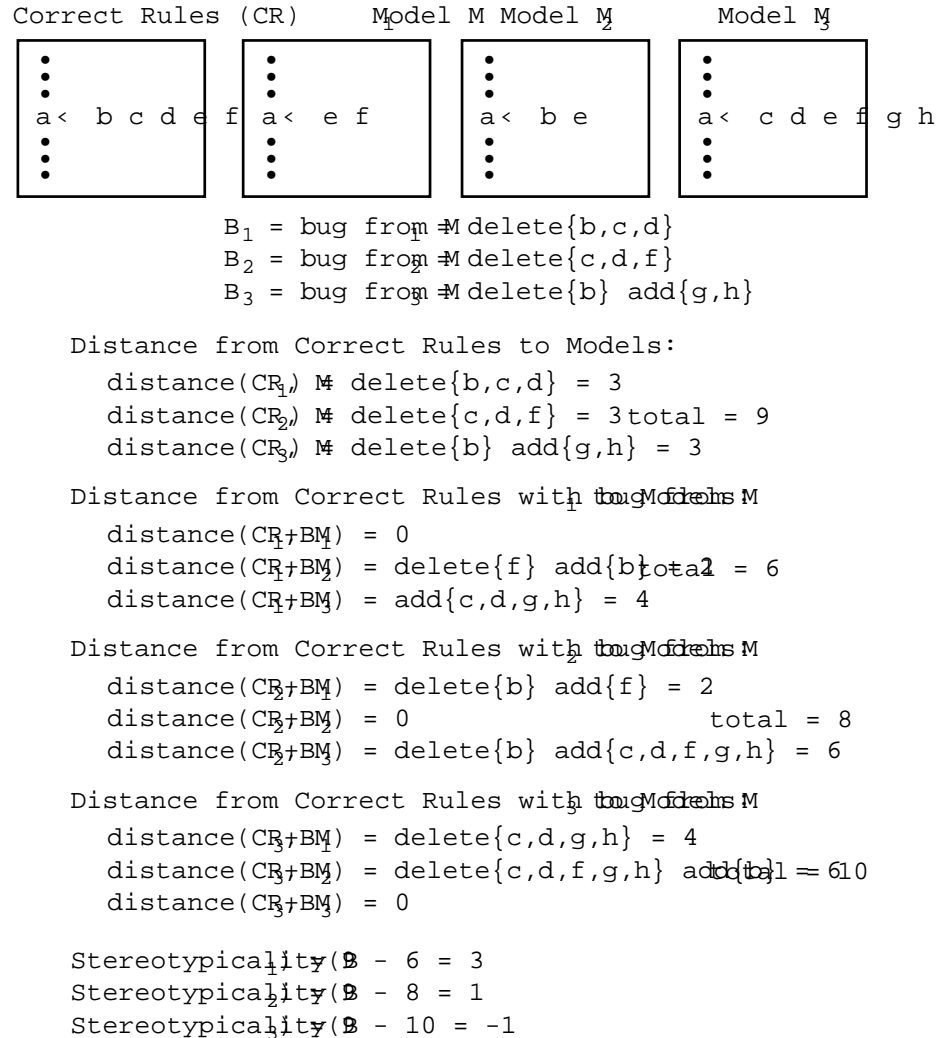


Figure 11: Stereotypicality computation.

refinement for modelable bug is the delete{b,c,d}. Below the models are the calculations for determining the stereotypicality of each of the three bugs by comparing each bug decrease distance between the models and the correct rule base. The distance between CR and the student models is shown followed by the distances to the each of the three bugs is calculated. The distance between two rule amounts to counting the number of literal changes required to convert one to the other. Changing CR+B₁ into B₂ requires changing the rule to the rule b which is done by deleting and adding. The bottom of the figure has the stereotypicality value of each of the bugs.

Note that a bug may have a negative stereotypicality. Unless a bug is present in half of the student models, it is likely to have a negative stereotypicality. A bug that overlaps between the bug and the majority of the models. Thus even a bug that is present in 30% frequency in the student population may have a negative stereotypicality.

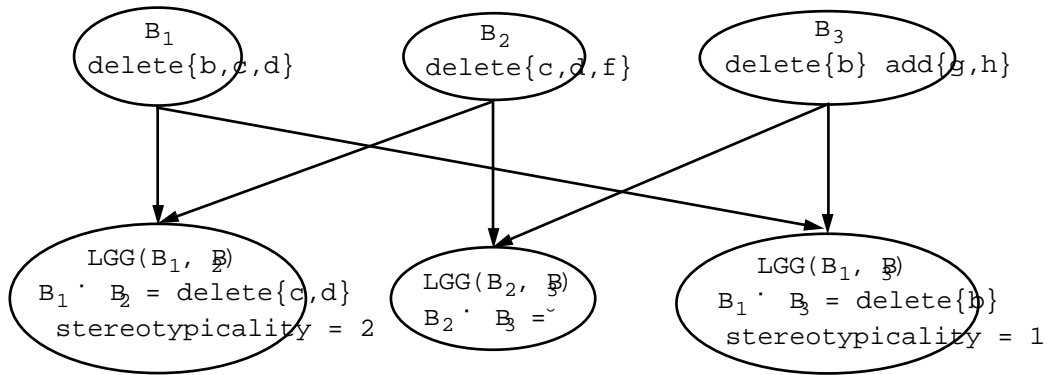


Figure 12: Bug generalization using the LGG opera

is the relative difference between stereotypicality values.

Figure 12 illustrates how these generalizations among the bugs from Figure 11. Since any refinement to a propositional theory can be expressed as a logic clause, we can compute generalizations using the least general generalization (LGG) operator (Plotkin, 1970). When two propositional logic clauses are not identical, one can form the LGG of the two by dropping literals from the clauses. Any number of literals may be dropped, but the most specific (i.e., least general) way to generalize the two clauses is to drop the literals which appear in just one of the two clauses. This is the same as the intersection between the two clauses. Since refinements consist of propositional logic literals, the LGG of two refinements is simply their intersection.

As might be expected, the LGG will often form a generalization that has better stereotypicality than a refinement from which it was taken. For example, the LGG of B_1 and B_2 is $B_1 \cdot B_2$, which has a stereotypicality of 2. Likewise, the LGG of B_1 and B_3 is $B_1 \cdot B_3$, which has a stereotypicality of 1. This will be the result whenever the LGG operation captures more of what is common among the models, and avoids more of what is uncommon. However, note that the LGG is not beneficial in all cases; the same generalization mentioned above are both worse refinements than alone, even though they both improve upon the seed.

Since the result of forming the LGG of two refinements is also a refinement, the process can be continued, forming LGG's from LGG's which can also result in better or worse refinements. Figure 13 shows the pseudocode for constructing a bug library. The fundamental operation is to perform a hill-climbing search using successive LGG operations. Starting with a seed, multiple calls are made to the LGG operator to combine the seed with refinements from the models. As long as this continues to result in a better generalization, more passes are made over all the refinements. The process halts when no generalization can be found which will improve upon the seed, which must eventually happen since continuing to refine will eventually result in a generalization that is no better than the seed.

2. As a computational note it should be pointed out that it is not necessary to compute the rule bases to calculate stereotypicality. Instead, the difference between distances can be used. For example, the distance between the correct rule base and a model is the number of literals that are not in the model. The distance between a bug and the correct rule base is the number of literals that are not in the model. By changing a single literal in the bug to a literal that is in the model, the distance between the bug and the correct rule base is reduced by one. The distance between the bug and the model is the number of literals that are not in the model. Whatever is in the bug that overlaps with the model is already in the model because it means those literal changes already exist in the model. Anything in the bug which does not overlap with the model increases the distance because those literal changes are not in the model. The computational complexity is linear in the number of literals in the refinements of the models.


```

function BuildBugLibrary (M:list of student models): bug library;
begin
  R = ~;
  for m . M do add refinements of m to R avoiding duplicates;
  for r . R do begin
    Best = r;
    S = Stereotypicality(Best);
    repeat while S continues to increase begin
      Temp = ~;
      for r . R do add Intersection(Best,r) to Temp;
      G = member of Temp with highest stereotypicality;
      if Stereotypicality(G) > S then begin
        Best = G;
        S = Stereotypicality(G);
      end;
    end;
    add Best to bug library;
  end;
  return bug library sorted by greatest stereotypicality;
end

```

Figure 13: Pseudocode for bug library constructi

between refinements will eventually produce no change or the null set. The be found starting with each refinement as the initial seed is kept and inserted in duplications in the library are eliminated and the results sorted by stereotyp

Finally, Figure 14 shows a complete example for constructing a bug libr bugs from Figure 11 plus one extra bug to highlight the hill-climbing nature. Below the four bugs are a series of boxes, each representing one iteration. Thus the first box is the iteration which computes the bug to be added to tl with B_1 as a seed, the second B_2 starts the seed, et cetera. After saving the ster icality B_1 of the inner loop is entered and an LGG is B_1 formed between three bugs. Note that there is no need to B_1 compute LGG the result B_1 is simply obviously cannot be an improvement. Once the LGG's are computed, the best this case B_1 (B_4), which has a stereotypicality of 4. This is compared with best stereotypicality, and since there is no improvement the B_1 inner loop ha the bug library. The second B_2 box also yields no improvement from generalizati resulting B_2 being added unchanged to the bug library.

The next two boxes representing the iterations B_3 and B_4 are more op for interesting B_3 . For the best LGG results from B_3 combine the resulting generali- zation delete $\{b\}$ whose stereotypicality value of 2 is an improvement of 4 poin value B_3 alone. Consequently, the inner loop repeats. A second round of LGG no further improvement, resulting in the addition of c to the generalization library. For B_4 the process is similar. A first round of LGG's produces an im which cannot be surpassed by a second iteration of the inner B_1 loop. Note th which was computed and rejected in the B_1 was also the seed, turns out to be a use ful improvement B_4 alone. The final bug library consists of the following : sorted in the follow B_1 order $\{b, B_2\}$ and delete $\{c\}$

$$\begin{aligned}
 B_1 &= \text{delete}\{b,c,d\} \quad S = 4 & B_3 &= \text{delete}\{b\} \text{ add}\{g,h\} \quad S = -2 \\
 B_2 &= \text{delete}\{c,d,f\} \quad S = 2 & B_4 &= \text{delete}\{b,c,e,f\} \quad S = 2
 \end{aligned}$$

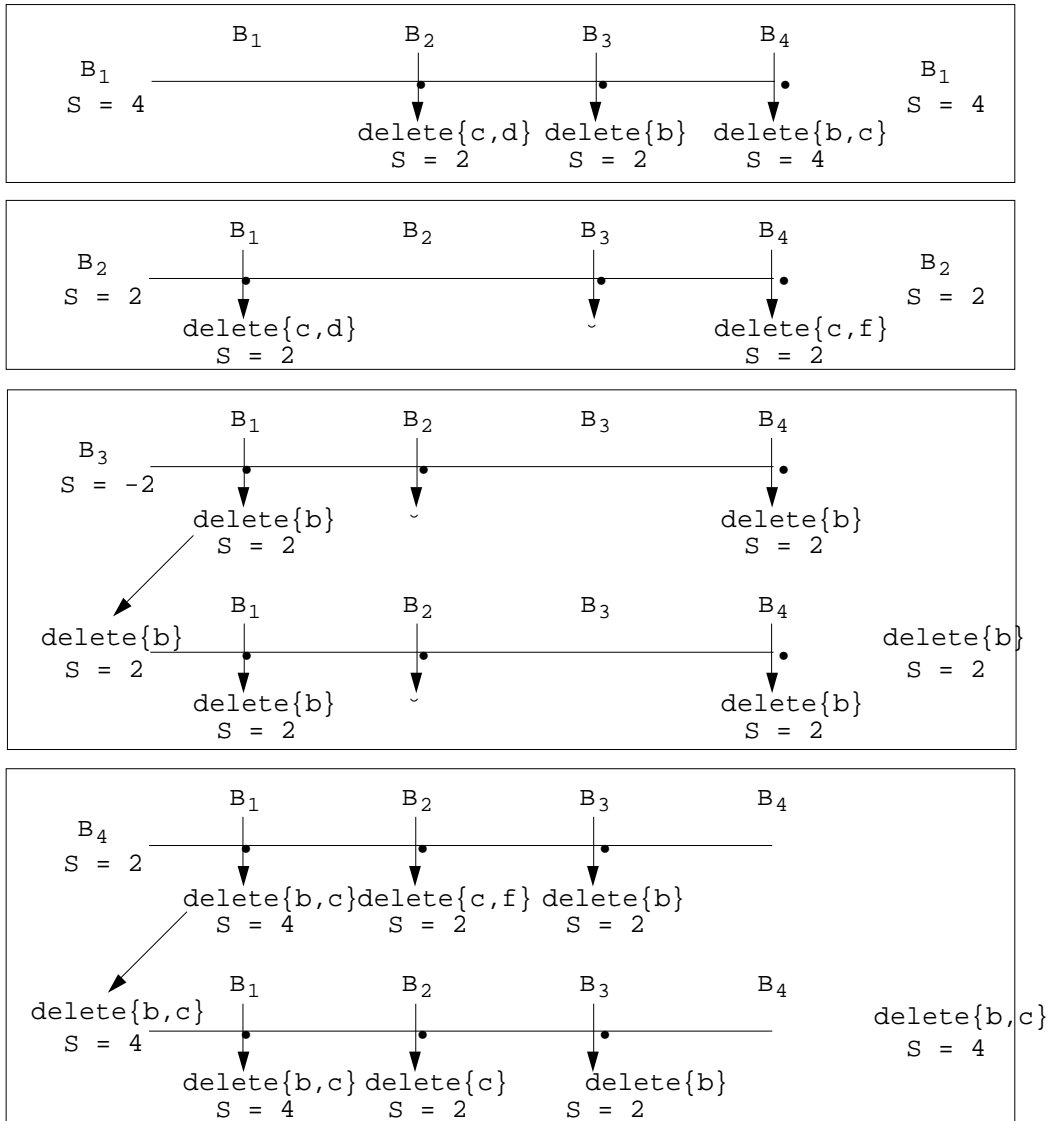


Figure 14: Bug library construction example. "S" stands :

4.2 Using the Bug Library for Modeling

Once the library is built, the question becomes how to use its information in the modeling process. Having the bugs in the library ranked by stereotypicality is an indication of how common they are. First, stereotypicality is an indication of how common a bug occurs in the student population. Bugs which occur frequently among the students are in higher stereotypicality since they will overlap more often with the models. Additionally, because stereotypicality is measured in literal add:

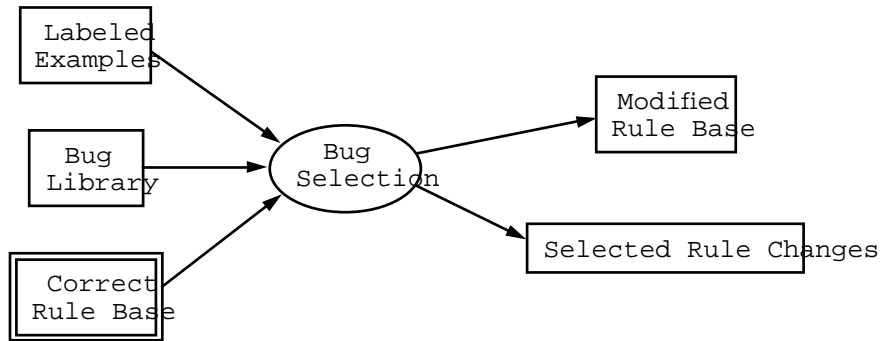


Figure 15: The bug selection module. Note that the correct rule base, when combined with the selected rule changes is equivalent to the modified rule base.

it is directly related to how the network behaves by integrating the bug with the correct rule base.

Perhaps the most obvious way to incorporate the bugs in the library is to use the refinement algorithm to use the bugs as a means for selecting repairs to fix a student. So, for example, when a repair is computed for a failing example its search might include any bug in the library which applied to any rule which could explain the failing example. The disadvantage of this approach is that it would destroy the modularity of the design. Theory refinement would no longer be an interchangeable component which could be swapped out for different refinement algorithms. A simpler approach, one used in *ISART*, is to modify one of the inputs to the refinement algorithm intact. Specifically, the correct rule base is modified before the refinement begins by incorporating elements of the bug library which are relevant to the student's behavior.

Figure 15 shows a schematic for how this is accomplished. The bug library, the correct rule base, and the student's labeled examples are input to a process which selects bugs from the library to be added to the rule base using a hill-climbing search. Upon completion, the predictive accuracy of the rule base is evaluated and the bugs which improve it are added incrementally. The process continues until a modified rule base which resembles the student's behavior more closely than the original rule base is found, which may still be incomplete. Note that the bugs which were selected are

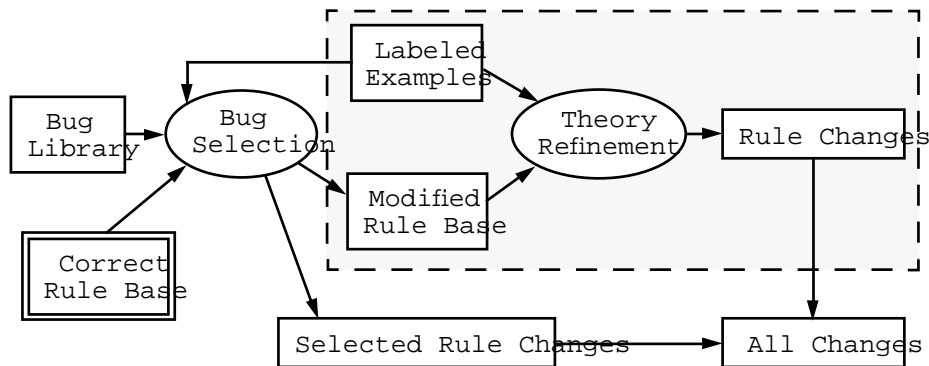


Figure 16: Extended modeling. Bug selection combined with theory refinement.

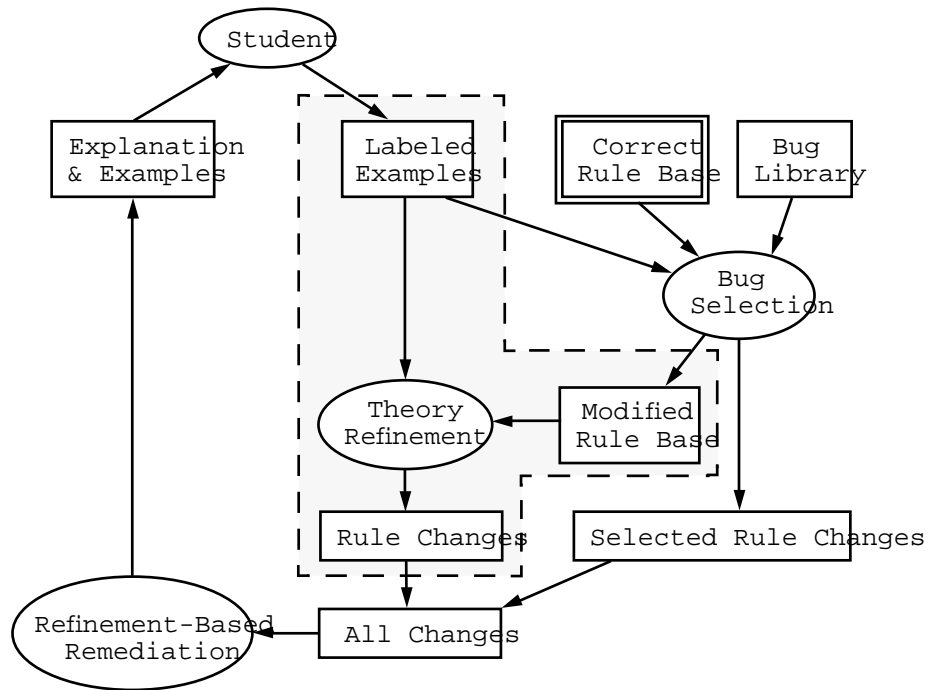


Figure 17: Overview of extended algorithm. The shaded area represents the theory refinement component.

with the modified rule base since they must be included with the final model (shown in Figure 16). Once the modified rule base is constructed, it is passed along with the labeled examples to determine any additional refinements reproducing the behavior of the student. All rule changes, whether selected or constructed by theory refinement, are returned as the final student model. the entire algorithm.

The pseudocode for modifying a rule base is shown in Figure 18. It starts with the correct rule base, and loops as long as a bug can be found which will improve the rule base on the set of labeled examples. The first step of each loop is to modify the rule base when the bug in question is added to the set of rules. All those improved accuracies are saved. Next, the bug which increases accuracy the most in the inner loop is entered to pare down the list to only those bugs whose improvement is "statistically equivalent" to the best bug, using a paired Student t-test. If there are still multiple bugs left, then the one with the greatest stereotypicality value is added to the current rule base (random selection is used as a final tie breaker). When no more bugs are found that increase the accuracy of the rule base, the routine quits and returns the modified version of the rules.

As an example of bug library selection, refer to a trace of the execution of:

3. more precisely, whose improvement in accuracy is less, but not statistically significantly
4. Since the accuracy values for all the bugs are computed using the same set of labeled examples, one can use a paired Student t-test to estimate if the difference in accuracy between two bugs is statistically significant (using the standard 0.05 level of confidence to indicate significance)


```

> (pre-model-student *student-examples* *correct-theory*)
-----iteration 1-----
Trying to beat accuracy = 80.00
bug 10, Accuracy: 85.00, Stereotypicality: -38
bug 11, Accuracy: 85.00, Stereotypicality: -38
bug 12, Accuracy: 85.00, Stereotypicality: -38
bug 20, Accuracy: 85.00, Stereotypicality: -38
bug 29, Accuracy: 85.00, Stereotypicality: -72
bug 34, Accuracy: 85.00, Stereotypicality: -128
Picked bug 20. Bug is:
  type: add antecedent to rule
  rule: compile-error <- constant-assigned
  antes: (integer-set no)
-----iteration 2-----
Trying to beat accuracy = 85.00
bug 5, Accuracy: 90.00, Stereotypicality: -32
bug 11, Accuracy: 90.00, Stereotypicality: -38
bug 12, Accuracy: 90.00, Stereotypicality: -38
bug 29, Accuracy: 90.00, Stereotypicality: -72
Picked bug 5. Bug is:
  type: delete antecedent from rule
  rule: constant-assigned <- (pointer constant) pointer-i
  antes: (pointer constant)
-----iteration 3-----
Trying to beat accuracy = 90.00
bug 29, Accuracy: 95.00, Stereotypicality: -72
Picked bug 29. Bug is:
  type: delete rule
  rule: operator-b-sets <- (operator-b auto-incr)
  antes: nil
-----iteration 4-----
Trying to beat accuracy = 95.00
done

```

Figure 19: Trace of bug selection from the bug lik

the addition of bug 20 enabled bug 5 to have its effect. This is a beautiful ordering effects inherent in selecting bugs from the library.

At the beginning of the third iteration, the updated rule base, which includes refinements from both bug 20 and bug 5, has an accuracy of 90%. Only bug 29 improve upon this accuracy so it is selected as the third bug to be added that while bug 29 continually resulted in improvements in accuracy, its relative stereotypicality prevented its addition to the rule base before this point. Final results in no further improvement.

Perhaps the most important feature of the bug library algorithm lies in its ability to model both common and unique misconceptions. As with other bug-library bas

methods, the ability to use a cache of expected errors gives the modeler a in domains where a large amount of data would otherwise be required for an nosis. But because the bug library here is used as a precursor to theory rei not restricted to using only those bugs present in the library. Any specif not in the bug library can still be captured by the theory refinement compo rithm. Thus the final theory of Figure 19 has partially accounted for the st the rest gets done by N

This completes the description of how it has been shown that a model both expected student errors as well as mistakes unique to an individual. Further fully automated scheme for bug library construction, and by integrating the its automatic modeling algorithm, it can continue to improve its modeling accura over time.

5 Experimental Results

It can be argued that the ultimate test of any tutoring system design is results in enhanced student performance. This is especially true for studer use of a model cannot significantly impact the educational experience then t son to construct one. Furthermore, this evidence must come from experime large numbers of students in a realistic setting so that the significance determined.

In this section, evidence is presented in support of the design that the be used to construct tutoring systems which significantly impact student performance. The bulk of this evidence comes from a test using 100 students who tutored with a C tutoring college-level freshmen taking an introductory course at the University of Texas at Austin. In addition to this evidence, experiments are presented from an which student responses were simulated. The advantage of this simulation d can be used to analyze the results of the C

5.1 C⁺ Tutor Tests

The C⁺ Tutor was developed in conjunction with an introductory course at the Univer- sity of Texas at Austin. The tutorial covered two concepts historically dif C⁺⁺ students: ambiguity involving statements with lazy operators and the pr and use of constants. These two concepts plus examples of correct programs categories into which example programs could be classified. A set of 27 dor developed to classify problems, using a set of 14 domain features, as beir ous, a compile error (for incorrectly declared or used constants) or cate- gory was the default category assumed for any example which could not b ambiguous or a compile error. Figure 20 shows an example question from the the complete listing of the rule base see Appendix A).

Students who used the tutorial did so on a voluntary basis and received their participation. As an added incentive, the material in the tutorial co would be present on the course final exam. This established a high level of r the students who participated in the test. Due to the large number of stud tutorial was made available over a period of four days and students were reserve time slots to use the program. In total, 100 students participated

Three major questions were the focus of the test. First, it was important establish whether a system could be an effective modeler for students in a real

```

void main()
{
  const int j = 3, *h;
  int i, k;
  h = &j;
  cin >> k >> i;

  cout << (k % j);  cout << (i %= j);
}

```

Is the above a compile error,
 (B) ambiguous (i.e., different outputs from different
 (C) neither A nor B ?

Figure 20: Example problem

ting. This was measured by testing the model produced for a student on a taken from the student which had not been given to it. Predictive accuracy of the model on such novel examples was expected to be higher than simply using the base with no modifications. Second, even with a perfect model one may not see in student performance. Though a model may be accurate a student will reach a faulty conclusion, it may not be the conclusion was reached. The only way to determine the utility of a model is to provide the student with that model and measure any change in the student's performance. Our hypothesis remediation generated using models should result in increased student performance over a control group which received no feedback. Additionally, it that students who were modeled with the benefit of a bug library would see performance increases over students who were modeled without a library. Third, as student modeling studies (Sleeman, 1987; Nicolson, 1992) we wanted to test receiving feedback based on student models would compare against students simple form of reteaching feedback. In this case, the expectation was that on modeling would result in greater post-test performance than simple reteaching.

Testing these three hypotheses was accomplished with three experiments: control, the effects of remediation, another to measure the accuracy of modeling and the utility of the bug library. In the next three sections each of these tests

5.1.1 Remediation with the tutor

For the remediation test, students who used the tutor were divided into four groups. One group received the full benefits of the tutor, the second used models formed without the benefit of a bug library, the third received reteaching and the fourth was a control group had no feedback. The expectation was that these four groups would exhibit similar performance on a post-test as the remediation ranged from full library to reteaching to no feedback.

To test whether the tutor can impact student performance, one needs to collect information for each student that has certain characteristics. To begin with, data both before and after any feedback given to the student to detect any change. Thus the tutor was constructed as a series of two tests with a remediation between. Secondly, the data from the two tests must be equally representative of the student's capability and must be collected in similar ways. The only way to de-

information from the tutoring program to the student is to have both test topics from the domain at similar degrees of difficulty.

To that end, a program was written to generate 10 example questions using format as follows. Since each question from the C was classified into one of three categories, the 10 questions were divided equally among the categories: three correctly labeled as compilation errors, four were examples of ambiguous three were questions with no errors. This process was used to generate two sets, both of which covered the same subset of the correct rule base. These two sets of questions covered the same concepts at the same level of difficulty. The two questions were identical. These two sets of questions represented the pre-test to be given to each student. One set of questions was used as the pre-test, the other as the post-test, thus the same pre-test and post-test was given to each student. To discourage cheating, the order in which the 10 questions were randomized. This meant every student answered the same two sets of questions. The difference was the feedback given between the pre-test and post-test.

Students were randomly assigned to four groups of 25, each of which received a different kind of feedback from the C. One group of 25 received no feedback, acting as a control group. This group was labeled the "No Feedback" group. The other three groups were given feedback using explanations and examples as described in Section 5.1 to ensure that the only difference between feedback groups was the type of feedback received, each group was given the same amount of feedback; specifically, four examples and four explanations for each student.

One feedback group received a form of automated reteaching. Specifying precisely what is meant by "reteaching" is extremely important, as it can have a profound effect on the results of the experiment. Furthermore, reteaching methods vary, making it difficult to specify the exact approach used. For this experiment, the essential point was that the feedback based on modeling made any difference over feedback based on no modeling. To that end, we chose an automated form of reteaching which used no information about the student, not even which answers the student got right or wrong. In such a vacuum, the option left for reteaching is to select at random from the rule base. Thus, for the "Reteaching" group, four rules at random from the rule base, and an explanation and example were generated for each rule.

The other two groups received feedback based on the models constructed from their answers to the pre-test questions. For the "Full" group (the "ASSERT-Full" group) the full ASSERT algorithm was used to build the model and for the "No Bugs" group (the "ASSERT-NoBugs" group) only ASSERT was used, i.e., no bug library information was given to the system. For both these groups, bugs were selected for remediation based on whether they were found by ASSERT.⁵ For the ASSERT-Full group, bugs from the bug library were given preference to those found by ASSERT in order of their stereotypicality value. In the ASSERT-Full and ASSERT-NoBugs groups, if fewer than four bugs were found, the remainder of the feedback was selected at random as with the Reteaching group.

Students were assigned to the four groups randomly. Since the ASSERT-Full and ASSERT-NoBugs groups required a bug library, the first 45 students to take the tutorial were randomly assigned to the ASSERT-Full, ASSERT-NoBugs, Reteaching and No Feedback groups. The models from these first 45 students were then used to construct a bug library. The remaining 55 students

5. ASSERT orders its refinements by preferring those which increase accuracy the most with the

Group	Average Pre-test Score	Average Post-test Score	Average Increase
ASSERT-Full	44.4	67.6	23.2
ASSERT-NoBugs	47.6	67.2	19.6
Reteaching	50.8	58.0	7.2
No Feedback	54.8	56.8	2.0

Table 1. Tutor remediation test. Scores indicate percent answered correctly. ANOVA analysis on average increase result between all groups except between Full and ASSERT-NoBugs and between Reteaching and No Feedback.

assigned to all four groups but at three times the rate of the other. The number of students assigned to all groups was even (25 students in each group).

Since the four groups of students each had a different average accuracy and post-test, they were compared using improvement in accuracy between pre-test and post-test. Also because each group consisted of different student between groups, significance was measured using an ANOVA test. As the only between groups was the feedback received, the significance test used was a ANOVA test at the 0.05 level of confidence using Tukey's multiple comparison (Tukey, 1953). The average improvement in performance for the four groups is shown in Table 1.

The results of the experiment confirmed most of our expectations. As predicted, performance decreased as the feedback varied from Full Library to reteaching to nothing. Moreover, both Full and ASSERT-NoBugs students improved significantly more than students in the Reteaching group. For ASSERT-NoBugs, improvement over Reteaching is more than 12 percentage points. For Full, the average improvement is even greater.

It is important to be very clear about the results in Table 1. Note that the variance among the mean pre-test scores in the four groups. Though non-significant, the variance among mean pre-test scores is a concern because the significance of the differences in average increase from pre-test to post-test is precisely why the ANOVA test was run to compare the significance. What is clear from Table 1 is that the feedback based on a model of the student significantly increase performance. There are no claims, however, as to how much improvement will get, whether the increase will always arise for every domain, how much depends upon the size of the pre-test score, or what the performance will be for other forms of modeling or reteaching. What has been illustrated is that automatic modeling and feedback performed by a tutor lead to significant performance improvements over feedback using no modeling at all.

This is the most important empirical result from this research. It illustrates that a tutor can be used to build a tutorial that significantly impacts student performance. Domain models and bug libraries are automatically constructed using only correct domain. Furthermore, it is another argument in favor of the use of student models. It shows (1) that they can have significant impact over not modeling at all and

be constructed automatically without resorting to the time-consuming task of building a library of bugs.

5.1.2 Modeling Performance using the C

The second important question to answer is whether or not there is a correlation between the ability to produce an accurate model and an improvement in student performance. This requires testing the modeling performance of the system, checking how the various features of the algorithm impact the predictive accuracy of the model. This can be accomplished using an ablation test format, in which various parts of the system are disabled and the resulting systems compared based on the accuracy of the models they produce. There are two different configurations used for modeling. The first, labeled "ASSERT-Full," uses everything available to construct the model, including means for referencing a bug library to create a modified theory which is then refined further. This method should produce the most accurate models. The second, labeled "ASSERT-NoBugs," skips the bug library and uses only the domain rules. One would expect ASSERT-Full to outperform ASSERT-NoBugs because of the additional information provided by the library.

In the C domain, only the data from the No Feedback group is useful for testing. This is because no remediation occurred between the pre-test and post-test for students in this group; thus, their 20 questions could be treated as a single training set and test set examples could be drawn. These training-test splits were chosen so as to be equivalently representative across the correct domain rules. Student performance quality is important to maintain so that any effects from modeling with theory refinement are manifested in the test set. Therefore, the 20 examples from the pre-test were grouped into 10 pairs, where each pair consisted of the two examples (one from the pre-test and one from the post-test) which covered the same domain rule. Then, training-test splits were generated by randomly dividing each pair.

The result was 10 possible training-test set splits. For each of the 25 No Feedback students, 25 training-test splits were generated, yielding 625 examples for comparison between ASSERT-Full and ASSERT-NoBugs. Each system was trained with the training set and accuracy was measured on the test set by comparing what the system predicted with what the students in the No Feedback group actually answered. The results are shown in Table 2. For these purposes, we also measured the accuracy of both an inductive learner, using training and test set splits, and the correct domain rules. The inductive learner, NEITHER with no initial theory, in which the correct domain rules are learned by induction over the training examples using a propositional version of the FOIL algorithm (Quinlan, 1980). If the correct theory no learning was performed, i.e., the correct domain rules were used without modification to predict the student's answers. Statistical significance was measured using a two-tailed Student t-test for paired difference of means at the 0.05 level.

These results illustrate that the groups with significantly better models, ASSERT-NoBugs, are precisely the groups which performed best after remediation. This is further evidence in support of the fact that more accurate student models lead directly to improved student performance via more directed remediation. This reinforces the finding of other studies (Ourston and Mooney, 1994) that indirect remediation are simply not as effective as theory refinement in terms of accuracy.

System	Average Accuracy
ASSERT-FULL	62.4
ASSERT-NoBugs	62.0
Correct Theory	55.8
Induction	49.4

Table 2: Results for modeling test. The differences between Full and ASSERT-NoBugs are not significant (all others are)

5.1.3 Bug Library Utility Test

However, note that the differences in Table 1 and Table 2 between A NoBugs are not significant. This means the use of the bug library did not significantly affect the performance of the student as expected, casting doubts as to its utility. The fact that the bug library did no harm to post-test performance, and perhaps with more data between the two groups would indeed have been significant. Thus it would be more about the conditions under which a bug library, as constructed with ASSERT, might be expected to impact the modeling process.

A series of tests designed to address this question, described in detail elsewhere, can be summarized by the results shown in Table 3. This data was generated from ablation tests like the one described in the previous section. In such a test, a theory is simulated by modifying a correct theory using six standard bugs selected plus additional random rule changes. The modified theory was then used to generate "answers" for 180 feature vectors representing a hypothetical "multiple-choice" test. The answers were then passed to a student model to see how well it could reproduce the modified theory. Once this was done for 20 students, resulting in 20 student models, the models were used to build a bug library.

Table 3 is a comparison of three libraries constructed using this technique. The first library, denoted 20-180, was formed from 20 student models built with 180 example answers per student. The second library, 20-12, used 20 students with 12 example answers per student chosen randomly for each student from among the 180. The 100-12 library was built from 100 students answering 12 random questions. Conceptually, these three libraries were designed to compare two conditions: (1) a few students answering lots of questions, a few students answering a few questions, and (2) a large number of students answering a few questions. This comparison is relevant to answering the following two-part question: (1) are bug libraries only effective when students answer a large number of questions, or (2) can one expect effective bug libraries to emerge from a large number of students answering a more reasonable number of questions. If useful bug libraries cannot be constructed from small student models, then the utility of a bug library is limited since one would still be tied to collecting large amounts of data on some students to construct the library. While more students always result in more individual models, it is important to show that a good bug library can still be built over time using less accurate models as in Part (a) of Table 3 compares the three libraries based on size and on how

Library	Total Students	Examples per Student	Common Bugs	Total Bugs Found
20-180 Library	20	180	all 6	29
20-12 Library	20	12	2	15
100-12 Library	100	12	4	48

(a)

System	Accuracy using different starting Bug Libraries		
	20-12 Library	100-12 Library	20-180 Library
ASSERT-Full	68.7	79.4	84.8
ASSERT-BugOnly	68.6	79.9	84.6
ASSERT-NoBugs	67.6	69.8	68.2
Correct Theory	63.1	63.5	62.6
Induction	25.4	26.0	23.9

(b)

Table 3: Comparison of bug libraries. Part (a) compares total number of bugs found, part (b) compares accuracy

standard bugs were found. The 20-180 library performed the best, finding a bugs. By contrast, the 20-12 library, which used the same number of student data per student, contained only two of the six common bugs. That result was due to drastically reducing the amount of information on each student, making it much less likely to be found and thus much less likely to end up in the bug library. Note that the 100-12 library contained four of the six common bugs, which is a significant improvement over the 20-12 library. Consequently, even though having smaller amount of student data reduces the chances that a common bug will be found for any given student, increasing the number of students improves the likelihood that the bug will be found for some student. Again, this is not too surprising when one compares the total number of examples used to build the student models which served as input to the three libraries. The 180 library there were 3,600 total examples, whereas the models used for the 20-12 library were built with only 240 total examples. The 100-12 library, with 1,200 examples, bridges the gap in total examples and in number of common bugs found.

However, note that the total overall size of the 100-12 bug library is larger than the other three libraries. In fact, the 100-12 library has the lowest ratio of common bugs to library size in the library. This dilution is a potential problem when the bug library is used in modeling efforts. Recall from Section 4.2 that the bug library is treated as a set of refinements which is traversed in an effort to improve the accuracy of the model. If the library is passed through a refinement process, increasing the size of the bug library widens the search space, making it potentially less likely that the common bugs will be selected, which could reduce the modeling accuracy.

This concern is addressed in Part (b) of Table 3 which shows results from tests aimed at determining if the ratio of common bugs to library size is a

of modeling accuracy. For these tests, a new crop of simulated students was used for each ablation test. In each ablation test was run with this same set of students, varying the bug student only 10 of the 180 examples were used for training, leaving the rest for testing. As the numbers in the table show, the 20-12 library results in almost the same accuracy for ~~SSART-Full and SEAT-BugOnly~~ (which used just the library without N) as opposed to ~~SEAT-NoBugs~~. By contrast, when the "better" bug libraries are used, the accuracy of ~~SSART-Full and SEAT-BugOnly~~ is significantly better, with the 20-180 library performing the best. Even the 100-12 library, with its low ratio of common bugs, resulted in a significant performance improvement over ~~SEAT-A BugOnly~~.

This implies that a bug library can be incrementally improved as more students interact with the system, even if the student interaction is moderate, resulting in more accurate modeling. And as the data from Table 1 and Table 2 shows, more accurate models result in more remediation and improved student performance. Recall, also, that in the remediation Section 5.1.1 the bug library was constructed using data from only 45 students. As more students interact with the system by taking tests that cover different subsets of the correct domain rules, it seems likely that a better bug library could be constructed which would lead to more accurate modeling and, in turn, better post-test performance.

5.2 Subjective Evaluation

Student Response to the Tutorial

Perhaps the most difficult topic to measure objectively is an evaluation of how students enjoyed using the tutorial and whether they felt the experience was valuable. A vast majority of students who used the tutorial reported that their response was positive. Many students made an unsolicited effort to express an appreciation for the opportunity to use the tutorial. Several students outside of the experimental group heard about the experiment and expressed interest in using the tutor to refresh their skills. There were even a few students who expressed disappointment that the tutorial did not cover more material.

On the negative side, students complained about their inability to back up to a previous test and post-test to change answers, which was not allowed by the interface. This was especially true during the early part of the pre-test, when students were still getting familiar with the interface. Also students expressed a boredom with redundancy in the explanations during remediation. Many students had multiple bugs detected in similar parts of the rule base. As a result, the chains of rule explanations to generate an explanation overlapped, resulting in duplications in explanations. Some of these problems could be easily fixed to make the interface more robust.

But perhaps the most important factor responsible for the positive response was that the feedback given to the student avoided negative language as much as possible. Of course, the student was told which questions he or she got wrong, but the explanation of the wrong answers did not focus on the student's mistake. Instead, the explanation focused on the correct reasoning, followed by an erroneous counter example which looked like something the student might misclassify, and explained why the counter example was wrong rather than saying something like "here's what you did wrong" the system would say "here's the right way to do something and, by the way, here's an answer which is wrong for the following reasons." This impersonal style of feedback may have made it easier for the student to accept the tutor's evaluation. And finally, by giving the student

to perform via the post-test, students were able to apply what they had learned in the average case, achieve a much better score. Such a concrete sense of improvement also contributed a great deal to the positive student response.

"Correctness" of the Bug Library

In the simulated student tests, the correctness of the contents of a bug library is measured directly by counting how many of the six common misconceptions are in the library. With the C++ Tutor domain this is not possible, since there is no a priori knowledge about what the common misconceptions might be. However, one can perform an evaluation with a domain expert to determine whether or not the bugs are correct explanations of why students made their mistakes. This was done by consulting the instructor for the course with the result that some of the bugs did, in fact, apply. For example, several bugs in the library represented missing conditions, or erroneous constant declarations, which the instructor felt students typically made. The library also contained bugs capturing the notion that students failed to use logical operators "AND" and "OR" were fully evaluated. The instructor suspected. While this is an admittedly weak evaluation, it does at least illustrate that the bugs included in the library could communicate information about trends in student modeling which made sense to the instructor.

5.3 Summary of Results

To recap, the main result presented in this chapter is that a significantly better modeler was developed with C++ Tutor. It was shown that those students for whom the modeler was able to construct significantly better models were the students whose performance improved the most. And while the use of a bug library did not significantly improve student performance, additional evidence was presented demonstrating that the contents of the library could improve over time so as to significantly improve the modeling process. This empirical evidence supports the two principal claims of the modeler: (1) that automatically constructed refinement-based models can be used to significantly improve student performance and (2) that a bug library can be constructed automatically that can enhance refinement-based modeling.

6 Discussion

Several important issues have been raised by this research that must be emphasized to place the work into a proper context. The first issue concerns the type of domain modeled. Unlike the previous modeling efforts, which focused on procedural tasks designed for use in classification domains. As an example of this difference, previous student modeling efforts have focused on the domain of writing computer programs (Goldstein, 1977; Soloway et al., 1983; Soloway and Johnson, 1984), while this research was tested using a classification task where students were asked to determine the correctness of program segments. This tie to classification domains is largely due to the fact that the most mature theory-refinement algorithms developed thus far are designed for use in classification domains and is not a limitation of the general framework. For instance, as first order logic refinement methods are enhanced and updated accordingly, enabling it to address a wider range of applications than is currently possible with PROLOG's propositional Horn-clause representation. Or, if refinement algorithms

using entirely different knowledge representations, the approach taken by ~~these approaches~~ applied to those domains as well. However, it is not immediately clear how to map ~~SAERT~~ to a procedural domain.

Shifting the focus of modeling to a concept-learning emphasis is not new. Other researchers, most notably Gilmore and Self (Gilmore & Self, 1988), have explored the potential of using machine learning for tutoring conceptual knowledge. A recent trade journal survey of applications for computer-based training indicates that commercial products are currently available for constructing computer-based tutoring systems (IDS, 1990) in which concept learning is a primary task. Concept learning is fairly well explored pedagogy (Dick & Carey, 1990) and effective techniques to correct incorrect student classifications are already extant in the literature (Laird, 1980). Thus a general technique for modeling in concept domains has a wide potential for commercial impact, and can be coupled with instructional techniques that can be effective in the presentation of conceptual material (Tennyson, 1971).

The second issue of importance is the comparison of the present experiments to previous studies performed on the utility of student modeling. Much of the controversy surrounding student modeling stems from a popular belief that detailed student models are difficult to build and yet result in little or no practical value. In truth, there are few studies in number and have reached disparate conclusions: one shows modeling to be ineffective (Sleeman, 1987), another shows that modeling can have a positive effect (Nicolson, 1992). Both of these studies used a bug library approach. An extensive library was built by hand rather than automatically. Furthermore, in both previous studies compared reteaching and modeling using the questions answered incorrectly, whereas in the present study the reteaching method was divorced from all student input to isolate the effects of the automatic modeling. Thus the present study can be seen as a direct comparison to previous experiments. In fact, the present study provides evidence that effective student models can be constructed automatically which will positively impact student performance.

Which leads to the third important issue; namely, the kind of modeling domain in which it is used. The empirical results here show simply that automatic modeling has a significant impact on student performance. This says nothing, however, about the models one ought to build nor about the best way to use them. For example, equally significant results could be achieved by using a far simpler modeler that far better results could be achieved if the student were allowed more control over the number or type of counter examples presented as feedback. It is a significant feature of the present study that it is a general-purpose method, that it works automatically, and that it can enhance student performance.

7 Conclusions

In conclusion, ~~SAERT~~ is a general-purpose method for constructing student models which operate in concept learning domains. It is able to construct student models automatically, catching both expected and novel student behavior. It is a system which can construct bug libraries automatically using the interaction of students, without requiring input from the author, and integrate the results of its efforts. Finally, the empirical evidence presented supports the two principles of research: (1) that automatically constructed refinement-based models can be used to significantly increase student performance and (2) that a bug library can also be

matically using multiple student models as input.

Acknowledgments

This research was supported by the NASA Graduate Student Researchers Program number NGT-50732.

Appendix A. [†] Tutor Domain

Domain Features:

```

pointer:          (constant non-constant absent)
integer:         (constant non-constant)
pointer-init:    (true false)
integer-init:    (true false)
pointer-set:     (true false)
integer-set:     (yes no through-pointer)
multiple-operands: (true false)
position-A:     (normal left-lazy right-lazy)
operator-A-lazy: (AND OR)
lazy-A-left-value: (non-zero zero)
on-operator-A-side: (left right)
on-operator-B-side: (left right)
operator-A:      (assign modify-assign mathematical logical comparison
                  auto-incr)
operator-B:      (assign modify-assign mathematical logical comparison
                  auto-incr)

```

Correct Domain Theory

```

compile-error < constant-not-init
compile-error < constant-assigned
constant-not-init (pointer constant) (pointer-init false)
constant-not-init (integer constant) (integer-init false)
constant-assigned (integer constant) integer-init (integer-set yes)
constant-assigned (integer constant) integer-init (integer-set through-poir
constant-assigned (pointer constant) pointer-init pointer-set
ambiguous < multiple-operands operands-linked
operands-linked < operand-A-uses operator-B-sets
operands-linked < operand-A-sets operator-B-uses
operand-A-uses < operand-A-evaluated operator-A-uses
operand-A-sets < operand-A-evaluated operator-A-sets
operand-A-evaluated (position-A normal)
operand-A-evaluated (position-A left-lazy)
operand-A-evaluated (position-A right-lazy) lazy-A-full-eval
lazy-A-full-eval < (operator-A-lazy AND) (lazy-A-left-value non-zero)
lazy-A-full-eval < (operator-A-lazy OR) (lazy-A-left-value zero)
operator-A-uses < (on-operator-A-side right)
operator-A-uses < (on-operator-A-side left) (not (operator-A assign))
operator-A-sets < (operator-A auto-incr)
operator-A-sets < (on-operator-A-side left) (operator-A modify-assign)
operator-A-sets < (on-operator-A-side left) (operator-A assign)
operator-B-uses < (on-operator-B-side right)

```

operator-B-uses < (on-operator-B-side left) (not (operator-B assign))
operator-B-sets < (operator-B auto-incr)
operator-B-sets < (on-operator-B-side left) (operator-B modify-assign)
operator-B-sets < (on-operator-B-side left) (operator-B assign)

References

Anderson, J. R. (1983). Architecture of Cognitive and University Press, Cambridge, MA.

Anderson, J. R., Boyle, C. F., and Reiser, B. J. (1985). Proceedings of the Ninth International Joint conference on Artificial Intelligence, Los Angeles, CA.

Baffes, P. (1994). Automatic Student Modeling and Bug Library Construction using Refinement Ph.D. thesis, Austin, TX: University of Texas.

Baffes, P. and Mooney, R. (1993). Symbolic revision of theories with M-of-N proceedings of the Thirteenth International Joint Conference on Artificial Intelligence, Chambery, France.

Baffes, P. and Mooney, R. J. (1992). Using theory revision to model students' reotypical errors. Proceedings of the Thirteenth Annual Conference of the Cognitive Science Society, pages 617-622. Bloomington, IN.

Brown, J. S. and Burton, R. R. (1978). Diagnostic models for procedural bug ematical skills. Cognitive Science, 2, 155-192.

Brown, J. S. and VanLehn, K. (1980). Repair theory: A generative theory of dural skills. Cognitive Science, 4, 379-426.

Burton, R. R. (1982). Diagnosing bugs in a simple procedural skill. In Sleeman, J. S., ed. Intelligent Tutoring Systems 8. London: Academic Press.

Burton, R. R. and Brown, J. S. (1976). A tutoring and student modeling paradigm. Computer Science and Education. ACM SIGCSE Bulletin, 1, 246.

Carbonell, J. R. (1970). AI in CAI: an artificial intelligence approach to instruction. IEEE Transactions on Man-Machine Systems, 1, 190-202.

Carr, B. and Goldstein, I. (1977). Overlays: a theory of modeling for computation. Technical Report A. I. Memo 406, Cambridge, MA: MIT.

Costa, E., Ducheno, S., and Kodratoff, Y. (1988). A resolution-based method for resolving student misconceptions. In Artificial Intelligence and Human Learning New York, NY: Chapman and Hall.

Craw, S. and Sleeman, D. (1990). The flexibility of spaced practice. In the Eighth International Workshop on Machine Learning, Evanston, IL.

- Dick, W. and Carey, L. (1990). Systematic design of instruction. IL: Scott, Foresman/Little, Brown Higher Education. Third edition.
- Finin, T. (1989). GUMS: A general user modeling system. In Kobsa, A. and V. editors, User models in Dialog, *System* 15, pages 411-430. New York, NY: Springer-Verlag.
- Gilmore, D. and Self, J. (1988). The application of machine learning to intelligent systems. In Self, J., Artificial Intelligence and Human-Computer Interaction, NY: Chapman and Hall.
- Ginsberg, A. (1990). Theory reduction, theory revision, and translation. In the Eighth National Conference on Artificial Intelligence, Detroit, MI, pages 777-782.
- Hoppe, H. U. (1994). Deductive error diagnosis and inductive error generalization in intelligent tutoring systems. *Journal of AI and Educational Systems*, March issue.
- Ikeda, M. and Misoguchi, R. (1993). FITS: A framework for ITS. Technical Report 93-5, Osaka, Japan: ISIR Osaka University.
- Langley, P., and Ohlsson, S. (1984). Automated cognitive modeling. In the National Conference on Artificial Intelligence, Austin, TX, pages 193-197.
- Langley, P., Wogulis, J. and Ohlsson, S. (1990). Rules and principles in cognitive modeling. In Frederiksen, N., Glaser, R., Lesgold, A. and Shoben, M. (editors), *Acquisition of Skill and Knowledge*, chapter 10, pages 217-250. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Liang, H. and Taotao, H. (1991). A uniform student model for intelligent tutoring systems: Declarative and procedural aspects. *Technical Report 11*:44-48.
- Miller, M. and Goldstein, I. (1977). An automated consultation procedure for MACSYMA. In the Proceedings of the Fifth International Joint Conference on Artificial Intelligence, Cambridge, MA, page 789.
- Murray, W. (1991). An endorsement-based approach to student modeling for intelligent tutoring systems. In the Proceedings of the Twelfth International Joint Conference on Artificial Intelligence, Sydney, Australia, pages 1100-1106.
- Newell, A. and Simon, H. (1972). *Human problem solving*. Chapter 6. Englewood Cliffs, NJ: Prentice-Hall.
- Nicolson, R. I. (1992). Diagnosis can help in intelligent tutoring systems. In the thirteenth Annual Conference of the Cognitive Science Society, Bloomington, IN, pages 635-640.

- Ourston, D. (1991). Using Explanation-Based and Empirical Methods in Theory Revision. Ph.D. thesis, Austin, TX: University of Texas. Also appears as Artificial Intelligence Laboratory Technical Report AI 91-164.
- Ourston, D. and Mooney, R. (1994). Theory refinement combining analytical and methods. *Artificial Intelligence* 66:371-394.
- Ourston, D. and Mooney, R. (1990). Changing the rules: A comprehensive approach to theory refinement. *Proceedings of the Eighth National Conference on Artificial Intelligence* pages 815-820. Detroit, MI.
- Plotkin, G. D. (1970). A note on inductive generalization. In Meltzer, B. and others. *Machine Intelligence*. (New York: Elsevier North-Holland.
- Quilici, A. (1989). Detecting and responding to plan-oriented misconceptions and Wahlster, W., editors. *Models in Dialog Systems*, chapter 5, pages 108-132. New York, NY: Springer-Verlag.
- Quinlan, J. (1990). Learning logical definitions from examples. *Artificial Intelligence* 33:239-266.
- Reiser, B. J., Anderson, J. R., and Farrell, R. G. (1985). Dynamic student intelligent tutor for LISP programming. *Proceedings of the Ninth International Joint conference on Artificial Intelligence* 1:8-14. Los Angeles, CA.
- Rich, E. (1989). Stereotypes and user modeling. In Kobsa, A. and Wahlster, V. editors. *Models in Dialog Systems*, chapter 2, pages 35-51. New York, NY: Springer-Verlag.
- Sandberg, J. and Barnard, Y. (1993). Education and technology: What do we have in store? *Artificial Intelligence Communication* 6(1):47-58.
- Self, J. A. (1990). Bypassing the intractable problem of student modeling. Gauthier, G., editor. *Intelligent tutoring systems: at the crossroads of artificial intelligence and education*, chapter 5. Ablex Publishing Corp.
- Sleeman, D. (1983). Inferring student models for intelligent computer-aided instruction. Michalski, R. S., Carbonell, J. G., and Mitchell, T. M., editors. *Machine Learning*, Artificial Intelligence, pages 483-508. Tioga.
- Sleeman, D. (1987). Some challenges for intelligent tutoring systems. In *Proceedings of the Tenth International Joint conference on Artificial Intelligence*, Milan. pages 1166-1171.
- Sleeman, D. H. and Smith, M. J. (1981). Modelling student problem solving. *Artificial Intelligence* 16:171-187.
- Sleeman, D., Hirsh, H., Ellery, I., and Kim, I. (1990). Extending domain tutoring studies in student modeling. *Machine Learning* 5:11-37.
- Soloway, E., Rubin, E., Wolf, B., Bonar, J., and Johnson, W. (1983). MENO-I: a programming tutor. *Journal of Computer-Based Instruction* 10(1):20-34.

- Soloway, E. and Johnson, W. (1984). Remembrance of blunders past: a retrospective development of PROUST. Proceedings of the Sixth Annual Conference of the Cognitive Science Society 57. Boulder, CO.
- Tennyson, R. (1971). Instructional variables which predict specific learner error and errors. Paper presented at the annual Meeting of the American Psychological Association. Washington, D.C.
- Tennyson, R. D. and Park, O. (1980). The teaching of concepts: A review of design research literature. Educational Research 5:55-70.
- Towell, G. and Shavlik, J. (1991). Refining symbolic knowledge using neural networks. Proceedings of the International Workshop on Multiple Learning Environments. Harper's Ferry, W.Va. pages 257-272.
- Tukey, J. W. (1953). The Problem of Multiple Comparisons. Mimeograph, Princeton University. Spence, J., Cotton, J., Underwood, B. and Dunham, R. Elementary Statistics fourth edition, pg 215, note 4. Prentice-Hall.
- Wenger, E. (1987). Artificial Intelligence and Tutoring Systems. Morgan Kaufmann.
- Winston, P. H. and Horn, B. KLIP Reading, MA: Addison-Wesley.
- Young, R. M. and O'Shea, T. (1981). Errors in children's subtraction. Cognitive Science 5:153-177.