

Towell, G. G., Shavlik, and Worden O. (1990). Refinement of approximate domain theories by knowledge-based artificial neural networks. In the Eighth National Conference on Artificial Intelligence, pages 861-866. Boston, MA.

Tukey J..W(1953). The Problem of Multiple Comparisons. Mimeograph, Princeton Spence, J., Cotton, J., Underwood, B. and Dunham, J. Elementary Statistics fourth edition, pg 215, note 4. Prentice-Hall.

Young, R. M. and O'Shea (1981). Errors in children's cognitive science 5:153-177.

- Ourston, D. and Money (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. *Ed. by Michèle Meltzer*
Machine Intelligence 5 (New York: Elsevier North-Holland).
- Quilici, A. (1989). Detecting and responding to plan-oriented misconceptions
 and *Wahlster*, editor. *User models in Dialog Systems*, chapter 5, pages 108-132.
 New York, NY Springer-Verlag.
- Quinlan, J. (1990). Learning logical definitions. *Machine Learning* (3): 239-
 266.
- Reiser B. J., Anderson, J. R., and Farrell, R. G. (1985). Dynamic student model
 intelligent tutor for LISP programming. *Proceedings of the Ninth International Joint
 conference on Artificial Intelligence*, pages 814. Los Angeles, CA.
- Rich, E. (1989). Stereotypes and user modeling. *Ed. by Stuart A. Anderson*
models in Dialog Systems, chapter 2, pages 35-51, New York, NY Springer-Verlag.
- Richards, B. and Money (1995). Refinement of first-order horn-clause domain theories.
Machine Learning (2): 95-131.
- Sandberg, J. and Barnard (1993). Education and technology: What do we know? And
 where is Artificial Intelligence. *Communication* (4): 45-58.
- Sleeman, D. (1987). Some challenges for intelligent tutoring systems. In
Tenth International Joint conference on Artificial Intelligence, pages 166-168. Milan.
- Sleeman, D. H. and Smith, M. J. (1981). Modelling student problem solving.
Intelligence (6): 171-187.
- Sleeman, D., Hirsh, H. J. and Kim, I. (1990). Extending domain theories: two
 studies in student modeling. *Machine Learning* (5): 1-37.
- Soloway E., Rubin, E. F., Bonar, and Johnson (1983). MENO-II: an AI-based
 programming tutor. *Journal of Computer-Based Instruction* (1): 20-34.
- Soloway E. and Johnson (1984). Remembrance of blunders past: a retrospective on
 development of PROUST. *Proceedings of the Sixth Annual Conference of the Cogni-
 tive Science Society* 57. Boulder
- Tennyson, R. (1971). Instructional variables which predict specific learner con-
 tention 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. *Review of Educational Research* (50): 55-70.

- Costa, E., Duchesoy and Kodratoff (1988). A resolution-based method for discovering student misconceptions. In *Artificial Intelligence and Human Learning* New York, NY Chapman and Hall.
- Craw, S. and Sleeman, D. (1990). The flexibility of specialized refinement. In *Proceedings of the Eighth International Workshop on Machine Learning* pages 28-32. Evanston, IL.
- Dick, W and Carey. (1990). *The systematic design of instruction*: Scott, Foresman/Little, Brown Higher Education. Third edition.
- Finin, T. (1989). GUMS: A general user modeling system. In *Knowledge, A. and W. editors, User models in Dialog Systems*, volume 15, pages 430-441. New York, NY Springer-Verlag.
- Gilmore, D. and Self, J. (1988). The application of machine learning to intelligent systems. In *Artificial Intelligence and Human Learning*, NY Chapman and Hall.
- Ginsberg, A. (1990). Theory reduction, theory revision and translation. In *Proceedings of the Eighth National Conference on Artificial Intelligence* pages 777-782. Detroit, MI.
- Hoppe, H. U. (1994). Deductive error diagnosis and inductive error generalization in intelligent tutoring systems. *Journal of Artificial Intelligence in Education* 5:27-49.
- Ikeda, M. and Misoguchi, R. (1994). FITS: A framework for ITS - a computational tutoring system. *Journal of Artificial Intelligence in Education* 5:319-348.
- Langley P, and Ohlsson, S. (1984). Automated cognitive modeling. In *Proceedings of the National Conference on Artificial Intelligence* pages 193-197, Austin, TX.
- Langley P, Gulis, J. and Ohlsson, S. (1990). Rules and principles in cognitive modeling. In *Frederiksen, N., Glasersfeld, A. and Shafto, D. editors, Monitoring of Skill and Knowledge Acquisition*, volume 10, pages 217-250. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Miller M. and Goldstein, I. (1977). An automated consultation procedure. In *Proceedings of the Fifth International Joint Conference on Artificial Intelligence* pages 789. Cambridge, MA.
- Murray W (1991). An endorsement-based approach to student modeling for planned tutoring. *Proceedings of the Eighth International Joint Conference on Artificial Intelligence* pages 100-106. Sydney, Australia.
- Nicolson, R. I. (1992). Diagnosis can help in strengthening. In *Proceedings of the Tenth Annual Conference of the Cognitive Science Society* pages 635-640. Bloomington, IN.
- Ourston, D. and Moore (1994). Theory refinement combining analytical and empirical methods. *Artificial Intelligence* 66:1-13.

operator-uses < (on-operator-side right)
operator-uses < (on-operator-side left) (not-~~Aopassign~~)
operator-sets < (operator-auto-incr)
operator-sets < (on-operator-side left) (operator-modify-assign)
operator-sets < (on-operator-side left) (operator-assign)
operator-uses < (on-operator-side right)
operator-uses < (on-operator-side left) (not-~~Bopassign~~)
operator-sets < (operator-auto-incr)
operator-sets < (on-operator-side left) (operator-modify-assign)
operator-sets < (on-operator-side left) (operator-assign)

References

Anderson, J. R. (1983). *The Architecture of Cognition*. Harvard University Press, Cambridge, MA.

Anderson, J. R., Boyle and Reiser, J. (1985). The geometry of the Ninth International Conference on Artificial Intelligence. Los Angeles, CA.

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

Baffes, P. and Mooney, R. (1993). Symbolic revision of theories with M-of-N rules. Proceedings of the Thirteenth International Conference on Artificial Intelligence pages 135-140. Chambery, France.

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

Brown, J. S. and Burton, R. R. (1978). Diagnostic models for procedural bugs in mathematical skills. *Cognitive Science* 2:155-192.

Brown, J. S. and Lehman, K. (1980). Repair theory: A generative theory of bugs in procedural skills. *Cognitive Science* 4:379-426.

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

Burton, R. R. and Brown, J. S. (1976). A tutoring and student modeling paradigm in an environment. *Computer Science and Education*. ACM SIGCSE (Bulletin) 12:6:1246.

Carbonell, J. R. (1970). AI in CAI: an artificial intelligence approach to computer instruction. *IEEE Transactions on Man-Machine, Systems* 10:190-202.

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

which operate in concept learning domains. It is able to construct student models and automatically detecting both expected and novel student behaviors. The first modeling system which can construct bug libraries automatically using the interactions of students, without requiring input from the teacher, and incorporate the results into future modeling efforts. Finally, empirical evidence presented supports the two principal claims of this 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 constructed automatically using multiple student models as input.

Acknowledgments

This research was supported by the NASA Graduate Student Researchers Program, grant number NGT50732.

Appendix A C++ 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}
operator-B:    {assign, modify-assign, mathematical, logical, comparison, auto}

```

Correct Domain Theory

```

compile-error < constant-not-init
compile-error < constant-assigned
constant-not-init (pointer constant) - (pointer false)
constant-not-init (integer constant) - (integer false)
constant-assigned (integer constant) - (integer set yes)
constant-assigned (integer constant) - (integer set through-pointer)
constant-assigned (pointer constant) - (pointer set)
ambiguous < multiple-operands operands-linked
operands-linked < operand-A-uses operator-A-sets
operands-linked < operand-A-sets operator-A-uses
operand-A-uses < operand-A-evaluated operator-A-sets
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)

```

designed for use in classification domains. As an example, several previous student modeling systems have focused on the domain of writing computer programs and Goldstein, 1977; Soloway et al., 1983; Soloway and Johnson, 1984), while research was tested using a classification task where students were asked to correct program segments. This tie to classification domains is fast that the most mature theory-refinement algorithms developed thus far are designed for classification and is not a limitation of the general framework. For instance, as first order logic refinement methods are enhanced (Ridley, 1985), and can be updated accordingly to address a wider range of applications. However, it is not immediately clear how easy it would be to port to a procedural domain.

Shifting the focus of modeling to a concept-learning emphasis is not unusual. Other researchers, most notably Gilmore and Self (Gilmore & Self, 1988), have shown the potential of using machine learning for tutoring conceptual knowledge. Concept learning also has a fairly well explored pedagogy (Dobson, 1990) and effective techniques for responding to incorrect student classifications are already present in the literature (e.g., & Park, 1980). Thus a general technique for modeling in concept domains has been developed and can be coupled with instructional techniques to show the potential benefit of conceptual modeling (Laird, 1971).

The second issue of importance is the comparison of the results to previous studies performed on the utility of student modeling. Much of the controversy in student modeling stems from a popular belief that detailed student models are difficult to build and yet result in little or no practical benefits. In truth, there are few studies and have reached disparate conclusions: one shows a positive effect (Sleeman, 1987), another shows that modeling can indeed have a positive effect (Nicolson, 1992). Both of these studies used a bug library approach. In the extensive library was built by hand rather than automatically. Furthermore, in both previous studies there was a form of reteaching that was used. In the C⁺⁺Tutor, thus the C⁺⁺Tutor cannot be seen as a direct comparison to previous experiments and should, however, be seen as evidence that student models can be constructed automatically which will positively impact student performance.

Which leads to the third important issue, the kind of modeling done and the way 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 whether 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 modeling system 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. Further, there need not restrict the modeling system to be used only on remediation tasks. Another way to interpret a "misconception" is as a far extension of a valid given problem. The most common bugs in the bug library may actually indicate that the current "correct" rule base does not, in fact, cover all correct solutions. Thus, to use ASSERT as a knowledge acquisition tool which can learn from "creative" solutions. In any event, the significant feature is that it is a general-purpose method, that works automatically and that it can significantly enhance student performance.

7 Conclusions

In conclusion, ASSERT is a general-purpose method for constructing student modeling

Part (a) of [Table 1](#) compares the three libraries based on size and on how many of standard bugs were found. The 20-180 library performed the best, finding all bugs. That result is readily justified—by drastically reducing the amount of data for each student, individual bugs are much less likely to be found and thus much more likely to end up in the bug library. The 100-12 library outperformed the 20-180 library though having smaller amounts of data per student reduces the chances that a bug will be found for any given individual, increasing the number of students improves the likelihood that the bug will be present in some student.

However, note that the total overall size of the 20-180 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. This dilution is a potential problem when the bug library is used for modeling efforts. Recall from [Section 4.2](#) that the bug library is treated as a search space for refinements which is traversed iteratively to improve the accuracy of the rule base before it is passed to the solver. Increasing the size of the bug library widens the search, potentially less likely that the common bugs will be selected, which could, in turn, reduce the modeling accuracy.

This concern is addressed in [Table 2](#) which shows results from more ablation tests aimed at determining if the ratio of common bugs to library size is detrimental to modeling accuracy. For these tests, a new crop of simulated students was created for each ablation test. For each test, only 10 of the 180 examples were used for training, leaving the remainder for testing. As the numbers in [Table 2](#) show, the 100-12 library results in almost no improvement for ASSERT-Full and ASSERT-BugOnly (which used just the library) over the ASSERT-NoBugs (which used no bug library). By contrast, when the “better” libraries are used, the ASSERT-Full and ASSERT-BugOnly is significantly better.

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 modest improvements in modeling. And as the data in [Table 2](#) shows, more accurate models lead to better remediation and improved student performance.

5.2 Summary of Results

To recap, the main result presented in [Table 1](#) is that the use of a bug library significantly improves student performance in a test involving 100 college level student-developed C++ Tutor questions. Furthermore, it was shown that those students for whom ASSERT 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 the likelihood 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 this research: (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 knowledge modeled. Unlike the previous modeling research which focus on procedural tasks,

Library	Total Students	Examples per Student	Common Bugs Found	Total Bugs in Library
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 Library		
	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 the total number of bugs found, part (b) compares accuracy

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. Those 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, they were combined to build a bug library.

Table 3 is a comparison of three libraries constructed using this technique, varying numbers of simulated students and varying amounts of example answers per student. The first library, denoted 20-180, was formed from 20 student models built with all 180 example answers per student. The second library, denoted 20-12, was formed from 20 student models built with only 12 example answers per student chosen randomly for each student from among the 180 possible. The 100-12 library was built from 100 students answering 12 randomly selected questions. Conceptually, these three libraries were designed to compare three different conditions: a few students answering lots of questions, a few students answering a few questions, and a large number of students answering a few questions. This comparison was done at answering the following two-part question: (1) how many bugs can be found in a library when students answer a large number of questions, or (2) can a more accurate bug library be constructed from a smaller number of students answering a more reasonable number of questions. If useful bug libraries cannot be constructed from small student models, 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 a library. While more student data will always result in more accurate models, it is important to show that a good bug library can still be built over time using less accurate models as in

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

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

so as to be equivalently representative across the correct domain rules. Such quality is important to maintain test set from modeling with the training set and manifested in the test set. Therefore, the 20 examples from the pre-test are grouped into 10 pairs, where each pair consisted of the two examples (one from and one from the post-test) which covered the same domain rule. Then, training splits were generated by randomly dividing each pair

The result was possible training-test set splits. For each of the 25 No Feedbacks, 25 training-test splits were generated, yielding 1250 examples for comparison Full and ASSERT-NoBugs. Each system was trained with the training set and accuracy measured on the test set by comparing what the system predicted with what the the No Feedback group actually answered. The results are shown in Table 2. For comparison purposes, we also measured the accuracy of both, using the same training and test set splits, and the correct domain rules. The inductive learner NEITHER with no initial theory which NEITHER builds rules by induction over the training examples using a propositional version of the FOIL algorithm (Quinlan) the correct theory no learning was performed, i.e., the correct domain rules out modification to predict the answers. Statistical significance was measured using a two-tailed Student t-test for paired means at the 0.05 level of confidence.

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 accurately remediation reinforces the finding of other studies (Our 1994) and Mooney's inductive methods are simply not as effective as theory refinement in terms of accuracy

5.1.3 Bug Library Utility Test

However, note that the differences in Table 2 between ASSERT-Full and ASSERT-NoBugs are not significant. This means the use of the bug library did not significantly affect the performance of the student as expected, casting doubt on the utility of 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 useful to know more about the conditions under which, as a bug library, used automatically by ASSERT, might be expected to impact the modeling process.

A series of tests designed to address this question, (see 1994) in detail, can be summarized by the results shown in Table 3. This data was generated using simulated

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 5: 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.

cluded from the model is the ASSERT-style feedback based on a model of the student significantly increase performance. There are no assumptions, however, one will get, whether the increase will always arise for every domain, how much depends upon the size of the pre-test score, nor what the performance will be for other forms of modeling or reteaching. What has been illustrated is that automatic modeling and feedback performed by the system 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 the system can be used to build a tutorial that significantly impacts student performance. Bug models and bug libraries are automatically constructed using only correct knowledge in the domain. Furthermore, it is another example in favor of the use of student models since it shows (1) that they can have significant impact over not modeling at all and (2) that they can be constructed automatically without resorting to the time-consuming task of manually constructing a library of bugs.

5.1.2 Modeling Performance using the C domain

The second important question to answer is whether or not there is a correlation between the ability of ASSERT 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 pieces of the system are disabled and the resulting systems compared based on the accuracy of the models they produce. There are two configurations in which ASSERT was used for modeling. The first, which is labeled "ASSERT-Full," uses everything available to construct the model. This means referencing a bug library to create a modified theory which is then fed back into the system for further refinement. This method should produce the most accurate models. The second, labeled "ASSERT-NoBugs," skips the bug library and uses only the information available in the domain. We expect ASSERT-Full to outperform ASSERT-NoBugs because of the additional information in the library.

In the C domain, only the data from the No Feedback group is useful for analysis. 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

is meant by "reteaching" is extremely important, as it can have a profound results of the experiment. Furthermore, reteaching methods important to clarify the exact approach used. For this experiment, the essential point was to feedback based on modeling made any other feedback based on no modeling at all. To that end, we chose an automated form of reteaching which used no information 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, selected four rules at random from the rule base, and an explanation and example was generated for each rule.

The other two groups received feedback based on the models constructed from their or her answers to the pre-test questions. For the "ASSERT-Full" group, the ASSERT algorithm was used to build the model and for the "ASSERT-NoBugs" group, NOBUGS was used, i.e., no bug library information was given to the system. For both these groups, bugs were selected for remediation based on those found by the system.⁷ For the ASSERT-Full group, bugs from the bug library were given preference to those found by the system in order of their stereotypicality value. In both 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. The ASSERT-Full group required a bug library. The first 45 students to take the tutorial were randomly assigned to the ASSERT-NoBugs, Reteaching and No Feedback groups. The models from these first 45 students were then used to construct the bug library. The next 55 students were randomly assigned to all four groups but at three times the rate of the other groups. 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 on the pre-test and post-test, they were compared using an analysis of variance to determine if there was a significant difference in accuracy between pre-test and post-test. Also because each group consisted of students with no pairing between groups, significance was measured using a one-way ANOVA. The only variable between groups was the feedback received, the significance test used was a 1-way ANOVA test at the 0.05 level of confidence using a multiple comparison method (Tukey 1953). The average improvement in performance for the four groups is shown in Table 4.

The results of the experiment confirmed most of our expectations. As predicted, average performance decreased as the feedback went from bug library to reteaching to nothing. Moreover, the ASSERT-Full and ASSERT-NoBugs students improved significantly more than students in the Reteaching and NoBugs groups. For the ASSERT-Full group, the average improvement is more than 12 percent and for the ASSERT-NoBugs group, the average improvement is even greater.

It is important to be very clear about the results that there is a great deal of variance among the mean pre-test scores in the four groups. Though none of the differences among mean pre-test scores is significant, their variance is a concern because the significance of the difference in average increase from pre-test to post-test. This is precisely why the ANOVA was run to compare the significance. What can be

6. Note that we specifically avoided comparing reteaching against reteaching by a human tutor. This comparison represents an entirely different experiment. Here we were concerned with determining whether automatic reteaching was a useful feature of the algorithm.

7. NEITHER orders its refinements by preferring those which increase accuracy the most with the system.

formance over a control group which received no feedback was expected that students who were modeled with the benefit of a bug library would see performance increases over students who were modeled without as in previous student modeling studies (Sleeman, 1987; Nicolson, 1992) we wanted to test how receiving feedback based on student models would compare against students simple form of reteaching feedback. In this case, the expectation was that re on modeling would result in greater post-test performance than simple reteach

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

5.1.1 Remediation with the Tutor

For the remediation test, students who were divided into four groups. One group received the full benefit of 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 de performance on a post-test as the remediation was targeted from bug 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 background characteristics. This must be collected 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 students' capability and must be collected in similar ways. The only way to detect information from the tutoring program to the student is to have both tests on 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 database 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 problems, three were questions with no errors. This process was used to generate two sets of questions, both of which covered the same subset of the correct rule base. This two sets of questions covered the same concepts at different levels of difficulty. 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 for students, 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 presented was 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 Tutor. 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 3.2. To ensure that the difference between feedback groups was the type of feedback received, each group was given the same amount of feedback, four examples and four explanations for each student.

One feedback group received a form of automated reteaching. Specifying pre

```

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 8: Example C++ Problem

tutoring college-level freshmen taking an introductory University of T at Austin. In addition to this evidence, experiments are presented from an ar which student responses were simulated. The advantage of this simulation dom 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 University of T at Austin. The tutorial covered two concepts most basic to C++ students: ambiguity involving statements with lazy operators and the prop and use of constants. These two concepts plus examples of correct programs categories into which example programs could be classified. A set of 27 doma developed to classify problems, using a set of 14 domain features, as being ous, a compile error (for incorrectly declared or used constants) or category was the default category assumed for any example which could not be ambiguous or a compile error. Figure 8 shows an example question from the for the complete listing of the rule base see Appendix A).

Students who used the tutorial did so on a voluntary basis and received e their participation. As an added incentive, the material in the tutorial cove would be present on the course final exam. This established a high level of mo the students who participated in the test. The number of students involved, the 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 in

Three major questions were the focus of the C First, it was important to establish whether a model could be a useful modeler for students in a realistic s ting. This was measured by testing the model produced for a student on a se taken from the student which had not been seen. The 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 clear how the conclusion was reached. The only way to determine the utility of a model is to provide the student with f that model and measure any change in the student. Our hypothesis was that remediation generated using models would result in increased student per-

```

> (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 7: Trace of bug selection from the bug library

5 Experimental Results

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

In this section, evidence is presented in support of the claim that the be used to construct tutoring systems that significantly impact student performance. The bulk of this evidence comes from a test using 100 students who interacted with a C

```

function ModifyRules (CR:correct rule base,
                    E:labeled student examples,
                    L:bug library): modified rule base;
begin
  R = CR;
  repeat as long as R is updated do begin
    A = ~;
    for b . L do begin
      if Accuracy(R+b, E) > Accuracy(R, E) then add b to A;
    end;
    if A „ ~ then begin
      best = x . A with best accuracy value;
      A¢ = best;
      for x . A do begin
        if Paired-t-Test(best, x) not significant then add x to
A¢
      end;
    end;
    if A „ ~ then update R with x . A¢ with highest
stereotypicality;
  end;
end;

```

Figure 6: Pseudocode for bug library use.

breaker). When no bugs can be found that increase the accuracy of the rule base, the routine quits returning the most current version of the rules.

As an example of bug library selection, refer to a trace of the execution of the "student" routine shown in Figure 7. This function is the implementation of the pseudocode used by the modeler mentioned in Section 5. The trace shown is taken from a student who interacted with the system as a part of a design described in detail below in Section 6. The bug library used in this trace consisted of 34 bugs taken from the students who used the tutor. However, only a portion of the bug library is actually shown in the trace, corresponding to those bugs which were applicable to the mistakes made by the student. For a complete listing of the bug library, refer to the "student" routine trace in the appendix, Figure 9.

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 based methods, the ability to use a cache of expected errors gives the modeler an advantage in domains where a large amount of data would otherwise be required for an accurate diagnosis. But because the bug library here is used as a precursor to theory refinement, it is not restricted to using only those bugs present in the student's problem domain. Bugs not in the bug library can still be captured by the theory refinement component of the algorithm. Thus the final theory of a student is partially accounted for by the bugs in the bug library, and the rest gets done by the theory refinement component.

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

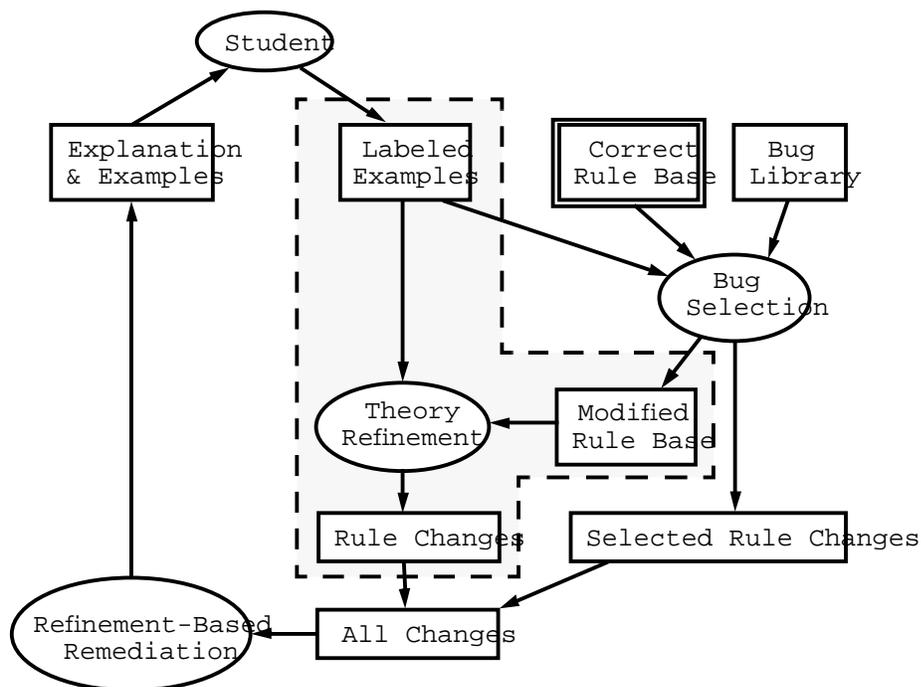


Figure 15: Overview of extended algorithm. The area within the dashed line is the theory refinement component.

the library to be added to the rule base using a hill-climbing search. Bugs upon the predictive accuracy of the rule base are then returned as a modified rule base which resembles student behavior more closely than the correct rules which may still be incomplete. The bugs which were selected are returned as a modified rule base since they must be included with the final model of the student. The modified rule base is constructed, it is passed to theory refinement along with examples to determine any additional refinements necessary for reproducing the student. All rule changes, whether selected from the library or constructed by theory refinement, are returned as the final student model.

The pseudocode for constructing the modified rule base is shown in Figure 16. ModifiedRuleBase starts with the correct rule base, and loops as long as a bug can be found to increase the accuracy of the rule base on the set of labeled examples. The first bug to find the accuracy of the rule base when the bug in question is added to the set of bugs which result in improved accuracy are saved. Next, the bug which increases accuracy the most is found, and an inner loop is entered to pare down the list to only those bugs in which the improvement in accuracy is "statistically equivalent" to the best bug, using a paired Student t-test. Finally, if there are still multiple bugs left, then the one with the greatest improvement in accuracy is picked to be added to the current rule base (random selection is used

4. more precisely, the improvement in accuracy is less, but not statistically significantly less.
 5. 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 the probability that the difference in accuracy between any two bugs is statistically significant (using the standard 0.05 level of confidence to indicate significance).

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

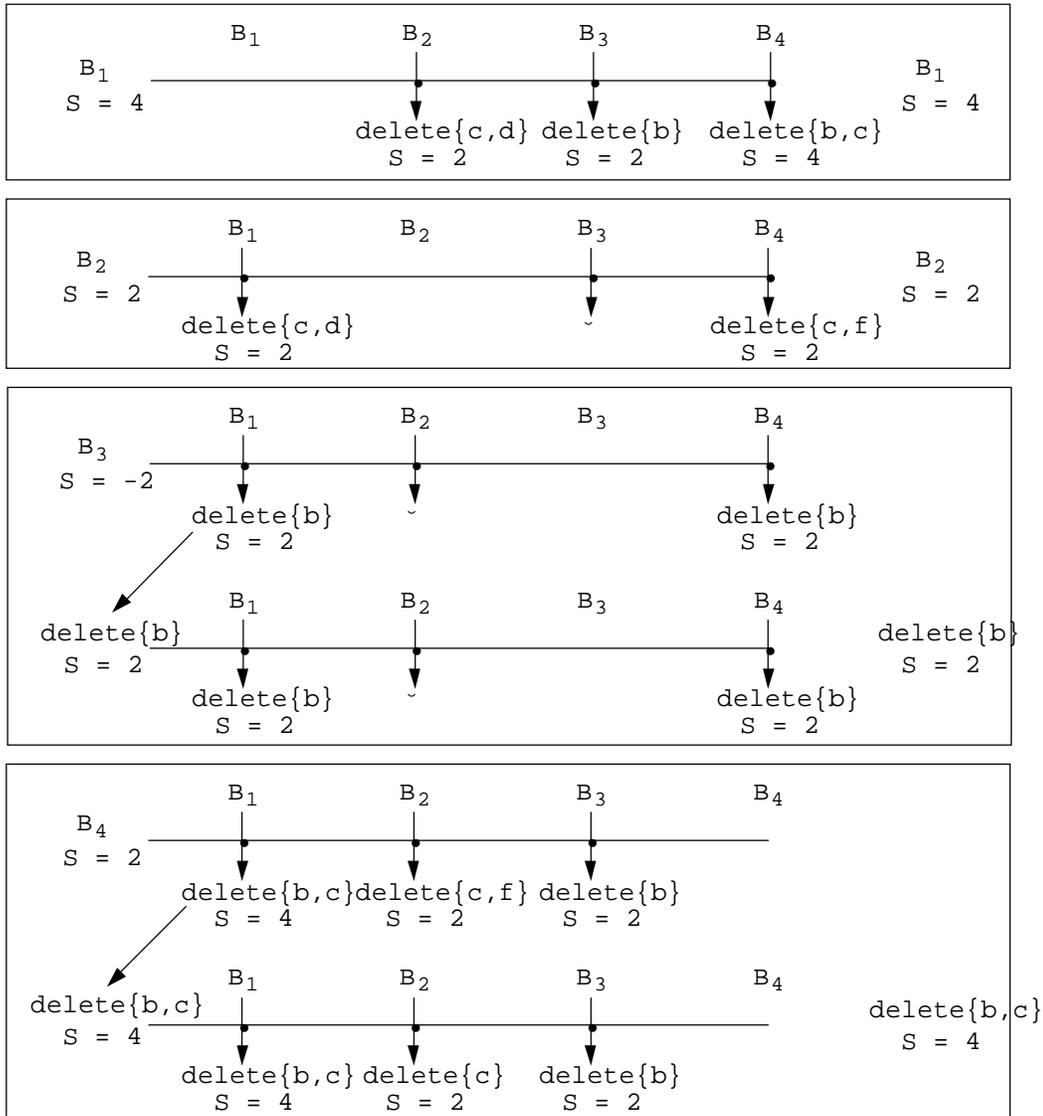


Figure 14: Bug library construction example. "S" stands for:

mistakes made by a student. The disadvantage of such an approach is that it is not modular as is the design of the refinement theory. Theory refinement would no longer be an interchangeable component which could be swapped for different refinement algorithms. A simpler approach, and the one used here, is to modify one of the input rules given to the refinement algorithm intact. Specifically, the rule base is modified before the refinement process begins by incorporating elements of the bug library which are the current student's mistakes.

Figure 15 shows a schematic for how this is accomplished. The bug library rule base, and the student's rule base, are input to a process which selects bugs

```

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;

```

Figure 13: Pseudocode for bug library construction

Figure 14 shows a complete example for constructing a bug library from Figure 1 plus one extra bug to highlight the hill-climbing nature of the algorithm. Four bugs are a series of boxes, each representing one iteration of the outer loop. The first box is the iteration which computes the bug to be added to the library as a seed, the second is the seed, et cetera. After saving the most stereotypical bug, the inner loop is entered and an LGG is found. In the first iteration three bugs are considered, and the best LGG is found, B_1 , with a stereotypicality of 4. This is compared with the current best stereotypicality, which is 0. Improvement in the inner loop has been found, so the bug library is updated. In the second box, for bug B_2 , no improvement from generalization is found, so the bug library is unchanged to the bug library.

The next two boxes representing the iterations B_3 and B_4 are more interesting. For the best LGG results from B_3 , combining the resulting generalization B_3 with the bug library whose stereotypicality value of 2 is an improvement of 4 points over the current value B_2 alone. Consequently the inner loop repeats. A second iteration produces an LGG with a stereotypicality of 6, which is an improvement over the current value of 4. No further improvement, resulting in the addition of B_3 to the bug library. For B_4 the process is similar. A round of LGG produces an improvement which cannot be surpassed by a second iteration of the inner loop. Note that the bug B_1 , which was computed and rejected in the first iteration, turns out to be a useful improvement over B_4 alone. The final bug library consists of the following four bugs, sorted in the following order: $\{b, B_2\}$ and $\{c\}$.

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. Perhaps the most obvious way to incorporate the bugs in the theory-refinement algorithm is to use the bugs as a means for selecting

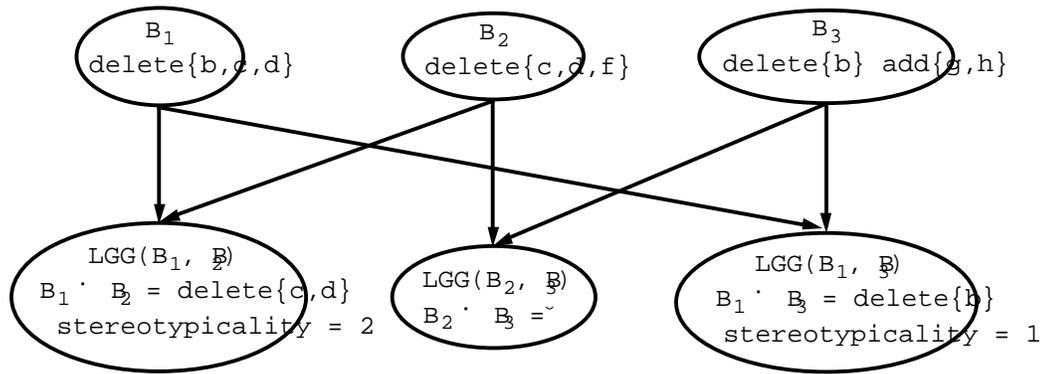


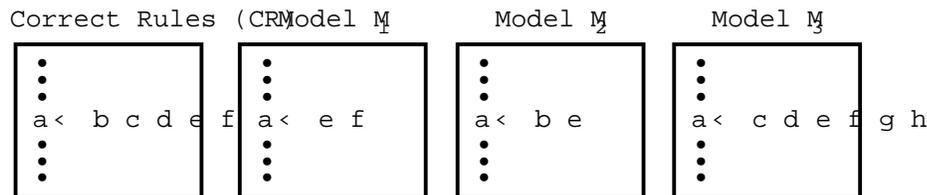
Figure 2: Bug generalization using the LGG operator

bug decreases the 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 model each of the three bugs is added. Calculating the distance between two rule sets amounts to counting the number of literal changes required to convert one to the other. 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 for each of the 3 bugs.

Figure 2 illustrates how forms generalizations among the bugs. Since any refinement to a propositional theory can be expressed as a logic clause, we can compute generalizations using the general(LGG) operator (Plotkin, 1970). For propositional logic literals, the LGG of two refinements is simply their intersection. The LGG will often form a generalization that has better stereotypicality than from which it was taken. For instance, LGG(B₁, B₂) has a value of 2. Likewise, LGG(B₁, B₃) is better than B₃ alone. 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. Note that the LGG is not beneficial in all cases, as mentioned above. LGG(B₂, B₃) is worse than B₂ alone, even though it is an improvement. 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 3 shows the pseudocode for constructing the fundamental idea is to perform a hill-climbing search using successive LGG operations. Starting with each refinement as a seed, multiple calls are made to the LGG operator to combine the seed with all other refinements from the models. As long as this continues to result in a better generalization, 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 continued refinement between refinements will eventually produce no change or the null set. The best generalization found starting with each refinement as the initial seed is kept and inserted into the library. Duplications in the library are eliminated and the results sorted by stereotypicality.

3. Note that a bug may have a negative stereotypicality if it is present in more than half of the student models, it will have a negative stereotypicality since the distance between the bug and the majority of the models. What is important is the relationship between stereotypicality values.



$B_1 = \text{bug from } \mathbb{M} \text{ delete}\{b,c,d\}$
 $B_2 = \text{bug from } \mathbb{M} \text{ delete}\{c,d,f\}$
 $B_3 = \text{bug from } \mathbb{M} \text{ delete}\{b\} \text{ add}\{g,h\}$

Distance from Correct Rules to Models:

$\text{distance}(CR_1, \mathbb{M} \text{ delete}\{b,c,d\}) = 3$
 $\text{distance}(CR_2, \mathbb{M} \text{ delete}\{c,d,f\}) = 3 \text{ total} = 9$
 $\text{distance}(CR_3, \mathbb{M} \text{ delete}\{b\} \text{ add}\{g,h\}) = 3$

Distance from Correct Rules with ~~bug~~ Models \mathbb{M}

$\text{distance}(CR_1 + \mathbb{M}_1) = 0$
 $\text{distance}(CR_1 + \mathbb{M}_2) = \text{delete}\{f\} \text{ add}\{b\} \text{ total} = 6$
 $\text{distance}(CR_1 + \mathbb{M}_3) = \text{add}\{c,d,g,h\} = 4$

Distance from Correct Rules with ~~bug~~ Models \mathbb{M}

$\text{distance}(CR_2 + \mathbb{M}_1) = \text{delete}\{b\} \text{ add}\{f\} = 2$
 $\text{distance}(CR_2 + \mathbb{M}_2) = 0 \text{ total} = 8$
 $\text{distance}(CR_2 + \mathbb{M}_3) = \text{delete}\{b\} \text{ add}\{c,d,f,g,h\} = 6$

Distance from Correct Rules with ~~bug~~ Models \mathbb{M}

$\text{distance}(CR_3 + \mathbb{M}_1) = \text{delete}\{c,d,g,h\} = 4$
 $\text{distance}(CR_3 + \mathbb{M}_2) = \text{delete}\{c,d,f,g,h\} \text{ add}\{b\} = 6$
 $\text{distance}(CR_3 + \mathbb{M}_3) = 0$

$\text{Stereotypicality}(B_1) - 6 = 3$
 $\text{Stereotypicality}(B_2) - 8 = 1$
 $\text{Stereotypicality}(B_3) - 10 = -1$

Figure 1: Stereotypicality computation.

els, called stereotypicality of the rule change. Third, in the process of ranking each change, ASSERT tests generalizations of the change to see if they result in better stereotypicality. If a generalization is found which has superior stereotypicality, it is replaced by the generalization. The final bug library contains the best generalization for each rule change with any duplicates removed.

Figures 4 through 14 illustrate how a bug library is constructed for a contract. Figure 4 shows how stereotypicality is computed (for details on how stereotypicality can be computed for a line of text, see section 4.1). Three models are 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. For example, the first model changes the rule $b < c < d$ by deleting the set of $\{b, c, d\}$. The refinement for model 1, labeled B_1 , is to delete $\{b, c, d\}$. Below the models are the calculations for determining the stereotypicality of each of the three bugs by computing

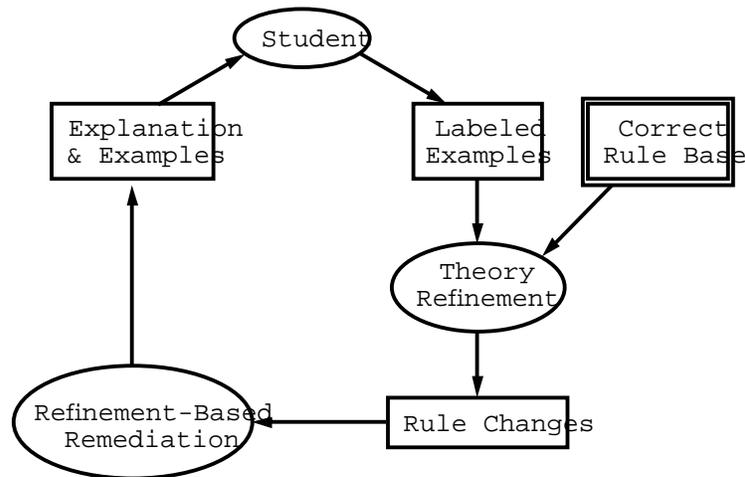


Figure 10: Basic design of the algorithm.

conclusion.

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

4 Extending ASSERT's Modeler

The previous three sections have described the algorithm, showing how the flow of information between student and system can be implemented as a refinements that highlight the differences between how the system and the student evaluate the same set of problems. As such, this explains how student models are traced by tracing the student behavior against a known standard rule base. Nothing, though, is said about how multiple student models are mined to construct a bug library. A bug library is incorporated back into the modeling process.

4.1 Building a Bug Library

ASSERT uses the rule changes resulting from theory refinement for each student for constructing its bug library. This library has two advantages. First, the rule changes are closely related to the type of input generated by the author of the tutor. The rule base must be supplied as input, expressing the bug library in terms of rule base in a way to communicate buggy information. Second, the author of the rule-change format is precisely what to simulate the behavior of the student. A bug library built of rule changes is thus already in a form which can be 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 into a single library, removing any duplicates. Second, it ranks each rule change by a measure of how frequently it occurs among the various

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.

ASSERT can thus generate an example which is correct for every added or missing conditions in the refinement. The result is then presented as a counterexample to the student, and the various added or missing conditions highlighted. Note that this is very closely to tutorial methods outlined for conceptual learning (Touretzky, 1997).

Figure 9 shows an explanation and example regarding to one of the refinements depicted previously in Figure 8. The last rule of Figure 8 is remediated by removing one of its antecedents. The feedback generated for this missing antecedent illustrates 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 rule fits in the overall rule base. As part of the explanation, the three conditions of the rule and its three correct antecedents, are itemized at the end of the explanation. The feedback generated which highlights how the "(pointer constant)" condition bears on the answer to the example, showing how the truth or falsity of the condition leads to a counterexample.



Figure 8: System response diagram.

How the correct rule base is constructed is crucial since it becomes the 1 which ASSERT interprets the student's actions. If the correct rules are expressed at too or too low a level of detail, the ability of the system to form accurate models is diminished. Of course, this type of knowledge representation problem exists for all systems. However, ASSERT gains an advantage by purposely isolating the correct domain knowledge as a separate component: the author can easily change the focus of the system by altering the correct rule base. Moreover, since students possessing different levels of understanding will use the multiple rule bases can be written to give the system more flexibility.

3.2.3 Refinement-Based Remediation

The last component of the system response, is outlined in Figure 8. Using the refinements produced by ASSERT, ASSERT generates explanations and examples to reinforce the correct form of the rule or rules modified. The underlying approach, called remediation, is based on fundamental units of explanation called *explananda*. Rather than implementing any particular pedagogical strategy, ASSERT applies the most elementary information required for explanation with one or more examples. For each refinement detected by ASSERT, ASSERT 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, and whether the system or the student controls which explanations and examples are given. By providing such explanation-example units, ASSERT supplies the raw materials for a variety of remediation techniques.

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

Generating examples is a bit more complicated since they are dynamically constructed rather than being drawn from storage. Recall that each refinement is made by adding or deleting literals from whole rules rather than the original or deleted as well, but this is the same as adding or removing blocks of literals. Using deductive methods, the added and removed literals can be traced down to the original rule. The result is a set of conditions in the feature vector which the student is missing, and extra conditions not present in the feature vector which the student thinks are present.



Figure 6: Student behavior diagram.

that this knowledge base can be modified to replicate the solutions furnished

3.2.1 The Student as a Classifier

Figure 6 depicts how the student behaves. It is assumed that all actions taken by a student can be broken down into classification. That is, given a set of inputs, called problems, the student will produce a set of 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 each feature vector selected from among some predetermined set of legal labels given to the student. The resulting set of labeled examples associates 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 taken from a list of possible categories. Thus, for example, the classic diagnosis problem asks a patient's symptoms (the feature vector) and ask the student to select a diagnosis (the label). This model is used in concept learning domains, which are common applications for automated training systems. It also means that the student's behavior will translate directly into a form usable by theory refinement, which requires labeled examples as one of its inputs.

3.2.2 Modeling by Theory Refinement

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

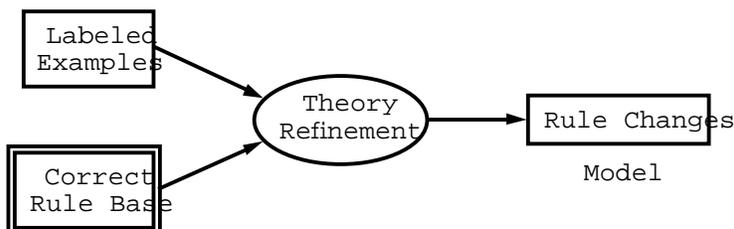
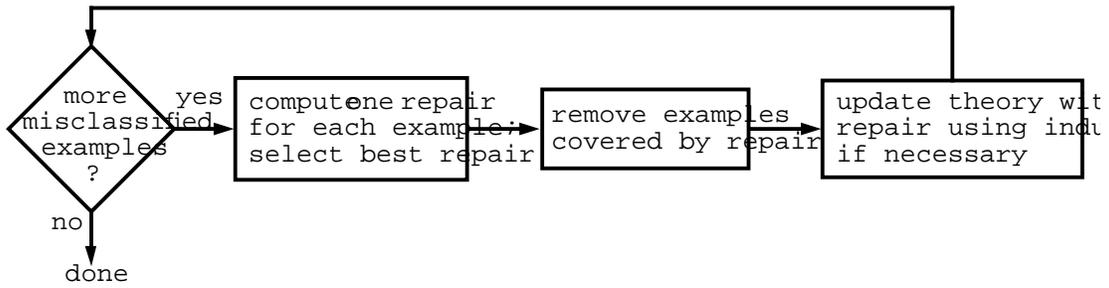


Figure 7: The student simulation model.



Criterion for computing a repair for ONE example
 Find the deepest, shortest, repair which causes the fewest new
 Criterion for selecting a repair AMONG examples
 Select the shortest repair fixing the most examples with the few

Figure 4: NEITHER main loop.

to modify the repair to avoid new misclassifications. The whole process is a loop which continues until all misclassified examples have been accounted for. The NEITHER runs very quickly (see & Madoney 1993), giving response times that are on the order of a few seconds. This is critical to an interactive tutoring system where feedback must be generated for the student in a timely fashion.

3.2 Overview of ASSERT

Having reviewed the basics of theory refinement, we can now return to the details. ASSERT views tutoring as a process of communicating knowledge to a student, and the contribution of the modeling subsystem is to pinpoint elements of the interaction base to be communicated. At its most abstract level, such a tutorial can be viewed as a process between the student and the system as shown in Figure 5. The focus is on how to construct a useful interpretation of the student's focus, which is depicted as a new component inserted into the diagram as shown in the right-hand side of Figure 5. The student simulation model is that the system contains a knowledge base that can be used to solve problems in the same context as the student.

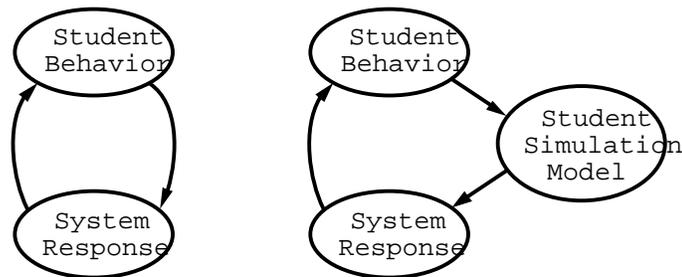


Figure 5: Abstract view of student-tutor interaction.

```

R1: compile-error constant-not-init
R2: compile-error constant-assigned
R3: constant-not-(pointer constant) pointer-init false)
R4: constant-not-integer constant integer-init false)
R5: constant-assigned integer constant integer-init integer-set yes)
R6: constant-assigned integer constant integer-init integer-set through-pc
R7: constant-assigned pointer constant pointer-init pointer-set
-----
R1: compile-error constant-not-init
R2: compile-error constant-assigned
R3: constant-not-(pointer constant) pointer-init false) integer-set no)
R4: constant-not-integer constant integer-init false)
R5: constant-assigned integer constant integer-init integer-set yes)
R6: constant-assigned integer constant integer-init integer-set through-pc
R7: constant-assigned pointer constant pointer-init pointer-set

```

Figure 2: Example of THEHER refinement. Above the dashed line are the rules in the theory base before refinement; below are the rules after refinement. Rules that are shown in boldface are the rules that were added or modified during refinement.

rule. Passing examples 2 and 4 to an induction algorithm would return "(integer constant through pointer)" as the condition which can discriminate between the examples. The final revision which correctly classifies all four examples is shown in Figure 3.

Notice that the repairs chosen for examples 3 and 4 are not the only possible repairs for these examples. For instance, example 3 could have been classified as a compile-error by removing the conditions "(integer constant through pointer)" from rule R6, or by removing the conditions "(integer constant)" and "(integer constant)" from rule R5. For that matter, removing all of the antecedents from any one of the rules through R4 would also have repaired the theory for example 3, by making either "compile-error" or "constant-not-init" concepts provable by default. In fact, computing all possible repairs for an example in the general case is exponential in the size of the theory. Consequently, the way in which repairs are calculated, as well as when a repair is applied to the theory in relation to computing repairs for other examples, can have a profound impact on the accuracy and performance of the theory refinement algorithm.

A summary of the THEHER algorithm is provided in Figure 4. The algorithm focuses on quickly finding one good repair for each example. The algorithm is to find the smallest repair in the deepest possible part of the theory. After converting the theory into a graph, the algorithm uses a recursive routine which starts at the leaf rules and works back up the graph. Failing conditions are collected at each rule and passed up to parent rules. The choice is possible. THEHER always chooses the smallest repair randomly to break ties. Each rule is visited only once, making the repair computation linear in the size of the theory.

Once a repair has been calculated for each example, one repair from the set is chosen to apply to the theory. The choice 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 smallest repair. The repair which fixes the most examples with the fewest new misclassifications is chosen. This repair is then tested against the rest of the examples, and induction is performed.

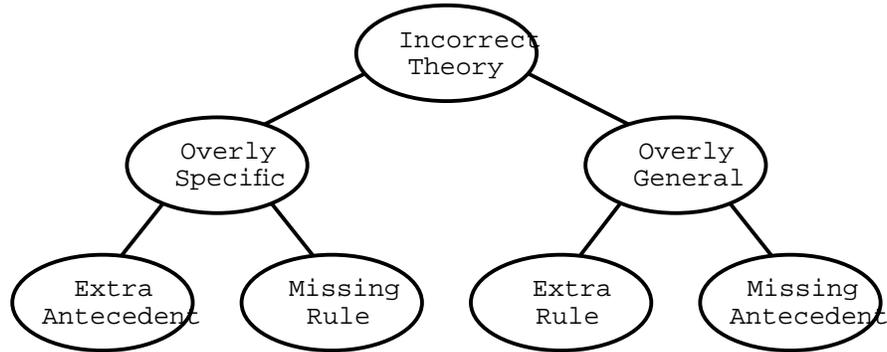


Figure 2: Theory error taxonomy for propositional Horn-cl

rule base is "incorrect" since it does not produce the desired classification examples. Propositional Horn-clause theories can have four types of semantics depicted in Figure 2. An overly-general theory is one that causes an example to be in categories other than its own, i.e. NEITHER adds new antecedents and deletes rules to fix such problems. An overly-specific theory causes an example classified in its own category false negative retracts existing antecedents and learns new rules to fix these problems. By making these four kinds of semantic changes, NEITHER can correct the semantics of the theory by altering the conditions which rules are satisfied.

NEITHER use a combination of three computations to determine how to modify the theory. The first step is repair for a single failing example by analyzing the rule base to determine what rules need to be changed to fix the theory for the example. For a positive example, a set of rule antecedents is found which, if deleted, will fix the example. Alternatively, for a failing negative example, a set of rules is computed which, if deleted, will repair the theory. The second step is to test the repair against all the other input examples to see if the repair will cause new misclassifications. If not, then the repair can be applied to the theory. The third step is taken using induction to learn a set of additional conditions which will separate the examples by the repair from the examples for which the repair causes new misclassifications. The additional conditions are then used to modify the repair.

As an example, in Figure 3 note that both example 3 and example 4 are false negative examples since neither is classified as "pointer constant". This indicates that the theory is overly specific and must be generalized. One way to repair the theory for example 3 is to delete the "(pointer constant)" condition from rule R7. This allows rule R7 to be satisfied by the example, without hindering the classification of example 1, and without causing the theory to become so general that it would be satisfied for example 2. A repair against example 3 against examples 1, 2 and 4 yields no new misclassifications, and it can be applied to the theory.

Finding a repair for example 4 yields a different result. The simplest repair is to delete the "(pointer constant)" condition from rule R3. However, if this repair is tested against examples 1, 2 and 3, example 2 is erroneously classified as "pointer constant". The only way to fix the repair is to remove the "(pointer constant)" condition from rule R3 and add a new condition to the rule which keeps example 2 from being classified as "pointer constant".

- R1: compile-errorconstant-not-init
- R2: compile-errorconstant-assigned
- R3: constant-not-(pointer constant(pointer-init false))
- R4: constant-not-(integer constant(integer-init false))
- R5: constant-assigned(integer constant(integer-init integer-set yes))
- R6: constant-assigned(integer constant(integer-init integer-set through-pc))
- R7: constant-assigned(pointer constant(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 2: A Theory and Examples. The desired classification is shown in italics (thus, Examples 3 and 4 are misclassified).

is repaired using a set of examples. The examples are assumed to be lists of feature value pairs chosen from a base set of main features. Each example has an associated label category which should be provable using the theory with the feature values of that example. NEITHER can generalize or specialize without user intervention, and is guaranteed to produce a set of refinements which are consistent with the input examples.

Figure 2 shows an example theory and four input examples. The top of the figure is part of a rule-based task solver that has been built for teaching ² a subset of C (for a complete listing of the rules, see Appendix A, numbered R1-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. Predicate values represent either intermediate concepts or are shorthand for binary feature values. "true" as a value. This simplified theory has only one category which it can classify examples. The input examples, shown in the table below the rule, 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 the closed-world assumption is used to classify the example 3 as, for compile-error instance, example 1 is correctly classified as a compile-error because it can satisfy either R6 or R7. Likewise, example 2 fails to satisfy any of the rules and is not correctly classified as non compile-error.

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

2. The source code for the task solver and the C tutor is available from the authors by anonymous FTP.

behavior by incorporating knowledge from a library of expected misconceptions. To be truly adaptive and to avoid the costs of bug library construction, one must use some sort of dynamic modeling or learning algorithm. And third, tracing student behavior in comparison to expected correct behavior can be a good method for detecting faulty behavior without requiring a great deal of searching. Also, we can combine these ideas by using a machine-learning technique called theory refinement (Ginsberg 1990; Ourston and Mooney 1994; Crow and Sleeman, 1991; Tet et al., 1990). Theory refinement is a systematic method for revising a knowledge base to be consistent with typical examples. If a knowledge base is considered incorrect or incomplete, and the examples represent behavior which the knowledge base should be able to handle, then the procedure itself is blind to whether or not the input knowledge base is "correct" in any sense; the theory-refinement process merely modifies the knowledge base until it is consistent with the examples. Thus, one can also use theory refinement by inputting a knowledge base and examples of behavior, and theory refinement will introduce whatever modifications are necessary to cause the knowledge base to simulate the 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, refinement will revise the knowledge base to make it consistent with the student, i.e., introduce "faulty" knowledge to account for the student's mistakes. The refinements made to the knowledge base then represent a model of the student, and can be used directly to guide tutorial feedback by comparing the student's behavior 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 MINSKY) and empirical (as in ACM) learning techniques in an integrated, domain-independent way. ASSERT can model any misconception consistent within the primitives that define the domain. And, ASSERT provides an extension to theory refinement that combines the results of different student models. This mechanism allows ASSERT to construct a bug library, without the necessity of intervention on the part of the instructor. Section 3.1 describes this algorithm in detail. We first turn our attention to the mechanism of theory refinement and ASSERT's role in the design.

3.1 Outline of Theory Refinement

Having outlined the philosophy, we can now turn our attention to the theory-refinement algorithm around which ASSERT is constructed. It is important to point out at the start that ASSERT is not tied to a particular theory-refinement algorithm. Other theory-refinement systems, which could be used to provide ASSERT with different or enhanced capabilities.

ASSERT uses the NEITHER algorithm (Ea, 1994; Ea & Mooney 1993) which is based on the NEITHER theory-refinement system (Ourston 1990). NEITHER was chosen as a starting point because it was the most complete symbolic theory-refinement system available. NEITHER is designed to work with a propositional Horn-clause knowledge representation. It takes two inputs, a propositional Horn-clause knowledge base called

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.

with a number of proposed mal-rules and must decide which ones are the "keepers"

2.3.2 Modeling by Induction

In an effort to avoid the cost associated with hand-constructed bug libraries, we turned to machine learning. Their ACM system (Langley & Ohlsson, 1984; Langley & Ohlsson, 1990) was the first to harness machine learning techniques for diagnosis of misconceptions through the use of ACM uses a domain-independent induction algorithm to induce control knowledge representing how students apply operators in a given problem. The output of induction is a set of conditions which are used to generate a model which lies on the path connecting the input for a given problem to the student's solution. The conditions found by induction are then used to specialize the model. The result is a procedure that models a student's problem solving behavior.

By using induction, ACM can operate automatically on models that capture both correct and buggy knowledge, because the operators must initially be general enough to model many kinds of behavior correct and incorrect, the potential search space is huge. Langley et al. note this, and suggest various "psychologically-motivated" conditions which can be applied to the operators to limit the search. However, the system is fundamentally limited by the complexity of having to construct a model from scratch. This can only be remedied by collecting data on each student or by imposing further constraints on the search space, which would require finding the constraints by using the very human-intensive methods this technique is designed to avoid.

2.4 Tracing Techniques

One final style of student modeler bears mentioning because it represents a blend of the previous techniques described to this point. In what might be loosely termed "tracing" 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. As such, tracing techniques must have knowledge of both correct and incorrect actions like bugs and must also have a mechanism for reproducing the steps followed by the student. ACM's solution paths. Using the correct path as a bias, tracing systems can operate efficiently.

The pioneering figures in this area are Anderson and his colleagues (Anderson, 1983; Anderson et al., 1985; Reiser et al., 1985), which follow student behavior through interaction with the tutor to occur through menu selection. Other tracing systems use a logic-based representation (Costa et al., 1988; Ikeda & Misoguchi, 1993; Misoguchi & Ikeda, 1993) where the idea is to use an analytical approach, such as deduction or resolution, through a rule-based system to determine where a misconception lies. Whenever the rule-based system produces a "proof" which mimics the student's solution, the points where the proof fails become candidates for querying the user about his or her beliefs.

Unlike the previous methods, tracing techniques do not dynamically construct models. Instead, they rely upon either the assumption that the student can be followed 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 future research presented here. First, modeling systems can increase their coverage by

directly on to the knowledge used to engineer the system. The disadvantage, a notational restriction placed on the model-only missing elements of the correct model can be modeled. Alternative notions which a student might have cannot be captured. This means that misconceptions cannot be modeled. Thus, overlay models can capture the notion of a student's knowledge, but they cannot be used to model the student who knows of a topic but misunderstands it.

2.2 Bug Libraries

To address the limitation of overlay models, other researchers focused on constructing bases of student misconceptions typically called classic bug-library work was done by Brown, Burton and Lehn (Brown & Burton, 1978; Burton, 1982; Brown & VanLehn, 1980), Sleeman and Smith (Sleeman & Smith, 1981), O'Shea and Young (Young & O'Shea, 1981), but a number of other systems can be said to incorporate a form of stored misconceptions (Rich, 1989; Miller & Goldstein, 1977; Quilici, 1984; Way & Johnson, 1984). Bug libraries are built by matching student behavior against a catalog of expected bugs which they hand through an analysis of student errors.

The idea is a very powerful one, especially if specific responses can be tied to buggy structures are encoded. However, important problems remain with this simple bug-library approach. First, the construction of such a catalog is a costly and 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 student misconceptions. That is, a student may exhibit a misconception which was not anticipated by the author of the library. As with overlay models, the static nature of bug libraries makes them 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 extend 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 consulting buggy information dynamically by data from a student to bound the search.

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. In INFER*, libraries, misconceptions are encoded as fault-rules. In MALGEN* and MALGEN attempt to generate new mal-rules when the student exhibits a problem not be modeled using the mal-rules already in the library. The difference between the two extensions is that INFER* builds new mal-rules by finding representations of a student's solution, whereas MALGEN uses a generate-and-test method to create candidate mal-rules.

The disadvantage of both systems is their reliance upon a user to decide which rules are appropriate extensions for the bug library. The authors do raise this issue and discuss potential filters that might be used to cut down on the number of rules presented to the student. Unfortunately, by this point no general-purpose filtering mechanism which might be usable across domains has been found. In the end, the use

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. In a 1993 interview of ten well-known ITS researchers which appears in the 1993 issue of *Artificial Intelligence Magazine* came to the conclusion that "most of the researchers no longer believe in on-line student modeling (Barnard and 1993). The article went on to conclude that "instead of becoming more integrated, the field has become more divergent in the last few years. It appears that scientists in the field of educational technology 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 Stereotypical Student Errors by Refining Theories*. ASSERT was designed to show that student modeling is a viable tool for an effective tutoring system. By taking advantage of some of the latest techniques in machine learning, ASSERT is able to construct student models and automatically catching both expected and novel student misconceptions. Also, it is the first system which can construct bug libraries automatically using the interactions with students, without requiring input from the instructor. The results so far indicate that future modeling efforts. In this paper, we describe a self-improving student modeler. ASSERT can be used to significantly improve student performance, as will be seen in sections which follow.

The remaining sections are organized as follows. Section 2 reviews previous work in student modeling as a motivation underlying the design. 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 students are automatically collected in 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 earliest AI-based student models, embodied in systems such as SCHOLAR (Carr, 1970), WEST (Burton & Brown, 1976) and WUSOR (Carr & Goldstein, 1977), used a form of modeling which is 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 knowledge. As the student performs actions which illustrate that he or she understands certain elements of the domain knowledge, these are marked in the overlay model. More sophisticated overlay models can express a range of values indicating the extent to which a student believes a student understands a given topic using some form of truth-maintenance (Finnin, 1989; Murray). However the marking is achieved, typically the unmarked elements of the model are used to focus tutoring on new problem areas for the student to ensure full coverage of the domain.

The advantage of the overlay is its simplicity; the elements of the model

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 new modeling technique which can automatically capture novel student errors using only domain knowledge, and can automatically compile trends across multiple student errors. This approach has been implemented as a computer program on 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 base that are necessary to account for the student's behavior. The efficacy of the approach is demonstrated by evaluating using 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 the system performed significantly better on a post test than students who received simple reteaching.

1 Introduction

Student modeling has a long and interesting history, well into the earliest efforts to produce intelligent tutoring systems. The best method for constructing and using a student model is still the subject of much debate. Most student modeling techniques ever 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 system, capturing misconceptions specific to each student and 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.