

Combining Explanation Types for Learning by Understanding Instructional Examples

Michael Redmond

School of Information and Computer Science
Georgia Institute of Technology

Abstract

Learning from instruction is a powerful technique for improving problem solving. It is most effective when there is cooperation between the instructor and the student. In one cooperative scenario, the instructor presents examples and partial explanations of them, based on the perceived needs of the student. An active student will predict the instructor's actions and then try to explain the differences from the predictions. This focuses the learning, making it more efficient. We expand the concept of explanation beyond the provably correct explanations of explanation-based learning to include other methods of explanation used by human students. The explanations can use deductions from causal domain knowledge, plausible inferences from the instructor's actions, previous cases of problem solving, and induction. They involve the goal being pursued and the action taken in support of the goal. The explanations result in improved diagnosis and improved future explanation. This combination of explanation techniques leads to more opportunities to learn. We present examples of these ideas from the system we have implemented in the domain of automobile diagnosis.

INTRODUCTION

People learn much of what they know from instruction. Presentation of examples can be an important part of instruction. LeFevre and Dixon [1986] found that students prefer examples to written text in learning a procedural task. Reder, Charney and Morgan [1986] found that instruction that included examples was more effective. What is it that makes examples effective teaching instruments?

One characteristic that makes them effective is that active students that try to explain the examples learn through the process of explanation. Lancaster and Kolodner [1988] and Chi, Bassok, Lewis, Reimann, and Glaser [in press] have both observed this in protocol studies. This has been our focus learning from understanding how a teacher solves an example problem.

Figure 1 summarizes the general process. Essentially, the instructor presents the problem, and appropriate actions or solutions. The student uses various types of knowledge to predict the instructor's actions, and then to understand or explain why the instructor's action or solution is appropriate.

The student is testing her ability to diagnose when she predicts what the instructor will do. The same techniques she would use if she were actually diagnosing are used to set up the prediction. In this way, when an opportunity to learn occurs, what is learned will be useful when the student actually goes about diagnosing. The example helps focus the learning.

We have constructed a system that creates explanations using deductions from causal domain knowledge, plausible inferences from the instructor's actions, previous cases of problem solving, and induction. The explanations involve the goal being pursued and the action taken in support of the goal. The explanations result in improved diagnosis and improved future explanation. This combination of explanation techniques leads to more opportunities to learn. This paper discusses the different types of explanations, and how they improve future problem solving and explanation.

-
1. The *instructor* states the problem description.
 2. The *student* attempts to generate an appropriate action for the problem and current context.
 3. The *instructor* generates a correct action or solution for the problem and current context.
 4. The *student* then attempts to explain this action, learning if possible.
 5. Continue with step 2 if the problem is not solved.
-

Figure 1: General Algorithm.

REDMOND

EXPLANATION

In our approach, explanation follows prediction and observation. The first step, therefore, is to compare the prediction with the expert's problem solving. This includes whether the instructor appears to be pursuing the predicted goal, and whether pursuit of the goal leads to the predicted action.

A correct prediction is essentially a successful explanation. Further explanation is required where the prediction isn't met. There can be many different ways of explaining differences. In this paper we discuss explanations involving:

- Inferring the instructor's current goal, and when necessary learning a new goal.
- Inferring the place of the current goal and actions in the diagnosis episode.
- Adjusting the saliency of features for future case retrieval.
- Trying to causally explain actions.

We have also begun to deal with a few other types of explanation that we will not discuss here. For example, explaining differences in implementation detail may rely on differences in car models, available tools, or in the current state of the car.

The types of explanations we make use of overlap with the types of explanations observed by Chi et al [in press]. They observed explanations that:

1. Refine or expand the conditions of an action
2. Explicate or infer different consequences of an action
3. Determine a goal or purpose for an action
4. Give meaning to a set of quantitative expressions.

Their first type of explanation is not a type that we have explored as yet. Our causal chaining explanation type corresponds to their second type, and our inferring the instructor's goal explanation type corresponds to their third type. Their fourth type is not applicable to our domain, though really it is a more specific version of inferring a goal. At a different level, Chi et al [in press] note explanations relating example actions to domain principles and to other example actions. Causal chaining can be seen as relating the observed actions to the domain principles. Inferring the place of the current goal and actions in the current diagnosis episode is one part of relating actions to each other.

In the following sections we will discuss in more detail how explanation of instruction is done, and how it improves the system through what is learned.

INFERRING INSTRUCTOR'S GOAL

Since the instructor's goal is usually not explicitly stated, it must be inferred from her actions. Different goals result in different types of actions being done. The instructor's goal must be inferred so that it can be compared to the predicted goal. The process is focused by the student's prediction of the instructor's goal. The predicted goal is the first goal considered as a possibility. If the instructor's actions are consistent with that goal then it is inferred that that is the goal being used. Otherwise, the goal must be inferred bottom up, with all possible goals being possible. This means that if the student gets lost in the example, she can find actions that make sense and get back to following along from there, and salvage something from the instructional episode.

```
(test (low ^fast-idle-speed))
(do (remove ^air-cleaner))
(do (disconnect ^radiator-fan))
(do (connect ^tachometer ^engine))
(do (plug ^vacuum-advance-hose))
(use (c-4812-2c))
(do (connect c-4812-2c ^choke-cam-follower-pin))
(do (release ^throttle-lever))
(ask ((rpm ^engine-system) nil) ^tachometer (reply 1600))
```

Figure 2: Instructor's Actions. The instructor's actions, entered into the system either interactively or by batch in a variable, are predicate forms specifying the type of action, and the action.

REDMOND

Some possible goals in a diagnostic domain include generating a hypothesis, testing a hypothesis, interpreting a test, fixing a fault, verifying a complaint, and clarifying a complaint. Figure 2 shows a portion of the instructor's actions in a given example. The complaint had been that the engine stalls, and the instructor has just hypothesized that the fast idle speed is set too low. This hypothesis must be tested. The instructor says that she is going to test whether the fast idle speed is low. Then she removes the air cleaner. She disconnects the radiator fan and connects a tachometer, and otherwise prepares for the test. Then using a specific tool specified in a reference book, she carries out the test, reading the value from the tachometer and comparing it to the specifications.

The process of inferring the instructor's goal uses knowledge about the goals stored in their representation. Some goals require particular types of actions. Some action types are inappropriate for some goals. Some action types can occur multiple times in the pursuit of a particular goal, some can only occur once. To give one example of the type of inference involved, testing a hypothesis *must* include an *ask* type action in order for results to be obtained. When it is determined that the predicted goal was not pursued, the other known goals are considered. Once the system knows what goal is being pursued, then the same explaining is done as if the goal had been correctly predicted. The student can recover and resume following the instructor.

If none of the diagnosis-specific goals are appropriate a more general goal can be considered, which could result in a diagnosis-specific specialization of the goal being learned. Figure 3 shows an annotated run of our system CELIA (Cases and Explanations in Learning: an Integrated Approach), reasoning as a student would, realizing that it needs to learn a new goal. For this run of the

```
...
Next Task
G-PREDICT-EXPERTS-ACTION

Next predicted goal
G-REPLACE-FIX
Mentally Simulating strategy S-RETRIEVE-MEMORY-PIECE for goal G-REPLACE-FIX
retrieve a piece from memory now
Matches fragments (pieces) -
  (GEN-REPLACE-FIX-LOW-IDLE          7.6000004)
  (GEN-REPLACE-FIX-THERM-COIL-CHOKE  5.6)
  (GEN-REPLACE-FIX-LEAN-CHOKE       5.6)
  (GEN-REPLACE-FIX-TOO-RICH         1.3)
Simulating based on retrieved piece GEN-REPLACE-FIX-LOW-IDLE
The fault has been determined to be: (LOW IDLE-SPEED)
The fix usually done in previous similar experiences was: (INCREASE (POSITION IDLE-SPEED-SCREW))
The method of doing the fix in previous similar experiences was: ...

Next Task
G-OBSERVE-EXPERTS-ACTION

Expert's next action ***** NOTE - test if engine is cold when it stalls *****
(TEST (TEMPERATURE ENGINE-SYSTEM (WHEN (STALLS ENGINE-SYSTEM)) COLD))
Expert's next action :
(DO (DRIVE CAR) UNTIL (STALLS ENGINE-SYSTEM))
Expert's next action ***** NOTE - read engine temperature gauge when car stalls *****
***** engine is cold when it stalls *****
(ASK ((TEMPERATURE ENGINE-SYSTEM) NIL) ENGINE-TEMP-GAUGE (REPLY (COLD)))

Next Task
G-EXPLAIN-DIFFERENCE

Comparing instructors actions to predicted actions
***** He's using a different goal than expected *****
...
***** don't know the goal being used or know it incorrectly *****
...
He's probably pursuing a specialization of the goal: G-TEST-DECISION
***** Create that specialization *****
NEW GOAL: G-DIAG-TEST-DECISION
**** Add new goal to tables ****
  modify goal-action table
  modify feature-saliency table
  modify goal hierarchy
  modify goal-slot table
  modify slot-action table
  modify slot-context table
reacting to observing learned goal G-DIAG-TEST-DECISION
making new case piece ... CASE-DIAG-TEST-DECISION-1
...
```

Figure 3: Realizing the need to Learn a Goal.

REDMOND

program we removed knowledge of the goal G-TEST-HYPOTHESIS from the student. This is equivalent to the novice student observed by Lancaster and Kolodner [1987], who came up with a reasonable hypothesis, then proceeded directly to trying to fix it without testing to see if it was a correct hypothesis. The example picks up after the instructor has made the hypothesis that the idle speed is low. The student retrieves a case piece suggesting the repair to do as a prediction of the instructor's actions. The instructor, however, correctly tests the hypothesis. These actions do not match expected action types for carrying out a repair, and in fact are not consistent with action types expected for any of the student's known diagnostic goals. It does, however, on further inspection, fit with expectations for a more general, cross-domain goal, of testing a decision. This enables learning a new diagnostic goal which will be a specialization of the more general goal.

There seems to be a difference between the goals that Chi et al [in press] talk about being inferred and the goals that our system infers. Specifically, if one looks at a goal as a goal type plus a parameter, our main effort is in inferring the goal type. The goal type would be our goal, for example, G-REPLACE-FIX, and the parameter would be the specific instantiation, for example (INCREASE (POSITION IDLE-SPEED-SCREW)). The parameter comes pretty easily for our system due to the input representation. Chi et al [in press] observed students trying to infer fully instantiated goals where the parameter could be less than obvious. However, the key point is that the student must understand what goal is being pursued in each part of the example as part of explaining the example. Future work can be directed towards inferring the parameter from less well-tailored input.

INFERRING PLACE IN CURRENT DIAGNOSIS

Inferring the place of the current goal and actions in the diagnosis episode is another step toward understanding observed problem solving. It is not only important in understanding what the instructor is doing, it is also necessary for saving the episode in a useful form as a case for case-based reasoning (CBR) [Kolodner and Simpson 1984]. A case will be more useful in the future if it reflects the problem solving done in the episode.

The instructor in most cases diagnoses hierarchically. People doing diagnosis don't hop around between unrelated hypotheses. The experienced mechanic considers a system as a potential source of the problem, then narrows the hypothesis down until a replaceable or fixable unit is determined to be malfunctioning. To a naive observer the hierarchy is not seen, the instructor's actions are sequential, a straight line instead of a tree. The rank novice observed by Lancaster and Kolodner

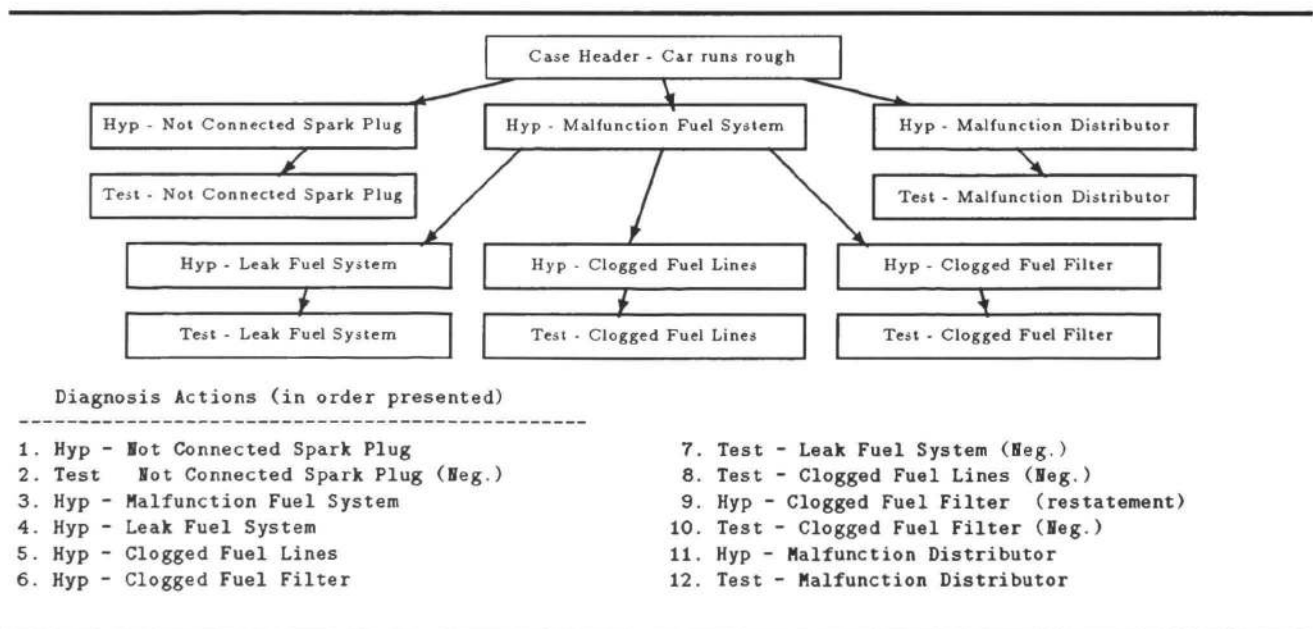


Figure 4: Inferred Diagnosis Structure.

[1987] did not diagnose hierarchically, but the other students, even the one with just six months more experience, did. The ability to diagnose hierarchically requires knowledge of the hierarchy involved. A system cannot rely on a given pattern of actions from the instructor, but must actually explain or understand what is going on. Figure 4 demonstrates this with an example diagnosis sequence. The top part of Figure 4 shows the structure of the instructor's actions which are shown in the bottom part of Figure 4.

Note that a test does not necessarily follow the hypothesis it relates to. Another complication is that there are at least two different reasons that a hypothesis can directly follow another hypothesis — it is a refinement as with the *'fuel system leak'* hypothesis following *'malfunction fuel system'*, or it is another possibility at the same level, such as with the *'clogged fuel lines'* hypothesis directly following the *'leak fuel system'* hypothesis. Also note that there is no 'syntactic' cue that the *'clogged fuel filter'* hypothesis is a refinement of the *'malfunction fuel system'* hypothesis and that the *'malfunction distributor'* hypothesis is not.

Knowledge is necessary to understand the hierarchy being used. Causal knowledge and structural relationships from the model are both useful for this process. A hypothesis can go under a previous hypothesis in the hierarchy if it causes the previous hypothesis, if the component involved is part of the previous component, or if the predicate is more refined.

Chi et al [in press] noted that one type of explanation is relating an action to another action. This process is one way of doing that. It is basically a linking of an action to the action that it follows from, which may *not* be the most recent previous action. The heuristics we use are geared for diagnosis. They were drawn from task analysis of Lancaster and Kolodner's [1987] protocols. They are the set that were necessary to establish the relationships between actions that we saw in the instructor's examples. We don't have any indication whether human students use heuristics such as these to recognize the relationships. Further analysis is required in order to come up with heuristics that would prove useful across domain types, such as for design or planning.

A partial list of heuristics used by our system to explain the instructor's actions in terms of hierarchical diagnosis is shown in Figure 5. The default expectation is that a hypothesis or test will be related to what immediately preceded it. However, as has been noted, this isn't always the case, and the third, fourth, and fifth heuristics are controls on that. The new action must actually

-
1. Try to put new hypothesis under most recent previous hypothesis or test.
 2. Try to put new test under most recent previous hypothesis.
 3. New hypothesis can go under a previous hypothesis if
 - its component is below the previous hypothesis's component in partonomy,
 - if the component is the same and the new predicate is more specific,
 - if the new hypothesis could cause the previous hypothesized fault
 4. New hypothesis can go under a previous test if
 - the test showed results indicating abnormal function and
 - the hypothesis is more refined than the test result (component is below the test's component in partonomy or if the component is the same and the predicate is more specific, or if the hypothesis could cause the test result)
 5. New test can go under a previous hypothesis if
 - the tested component is the same or below the hypothesis's component in partonomy and the test predicate is the same or more refined than the predicate in the hypothesis,
 - no component in the test is higher than any component in the hypothesis in partonomy
 - or if the tested clause could be a result of the hypothesis")
 6. Don't add anything directly under a hypothesis that has already been tested
 7. Don't add anything under a test whose results indicated normal function, this should be followed by backtracking
 8. Don't add a new test directly under a hypothesis that already has subhypotheses
-

Figure 5: Heuristics for inferring the structure of a diagnosis.

REDMOND

be related to the previous one, by being more specific or causally related. For example, in Figure 4, the hypothesis '*leak fuel system*' is more specific than the hypothesis '*malfunction fuel system*' because the predicate is more specific and the involved component is the same. The hypothesis '*clogged fuel lines*' is more specific than the hypothesis '*malfunction fuel system*' because fuel lines is below fuel system in the partonomy in memory. However the hypothesis '*malfunction distributor*' did not qualify on either count so it had to go in a different place. The third way to satisfy heuristic 3 is for the later hypothesis to be causally related to the previous hypothesis. The necessity of this is shown by an example. If the hypothesis '*clogged spark plug gap*' follows the hypothesis '*no spark from spark plug*' it would not be placed beneath it because 'clogged' is a different predicate than 'no spark', and isn't more refined. This could easily be a different problem. However, causal knowledge allows linking the one to the other so that the system knows, as a person would, that '*clogged spark plug gap*' is a refinement of the hypothesis '*no spark from spark plug*'.

If the action cannot go after the most recent action then the system must search for its proper place. Many of the other heuristics are limitations on this process, either avoiding potential incorrect placements, or cutting off search that will prove to be unfruitful.

For example, Heuristic 7 allows cutting off search when the instructor would be backtracking. If in Figure 4 the *malfunction fuel system* hypothesis had been followed by a test that showed normal function for the fuel system, then future hypotheses from the instructor should involve other hypotheses that aren't refinements of a fuel system malfunction, and the system can avoid wasted effort by not trying to see if they fit under that hypothesis.

Once the structure of the observed diagnosis has been determined, the case can be stored in memory for use in future problem solving and explanation. The case is stored in pieces so that the particular pieces can be accessed as necessary, and so the representation is flexible enough to handle diagnosis that doesn't have a set pattern of hypotheses and tests. There are pieces for each instance of each goal pursued in the episode. That is, for each hypothesis made, for each test of a hypothesis, for each interpretation of a test, for each fix attempted, there will be a piece. These pieces are linked together to preserve the structure of the case, as inferred in this step. This allows a future diagnosis using the current case to follow the links as long as the findings are the same. The diagnostician following such a hierarchically organized case will diagnose hierarchically rather than haphazardly like a novice. The case pieces, once correctly linked, are stored beneath general knowledge in the model for the car, under related components.

ADJUSTING THE SALIENCE OF FEATURES

Another important explanation type is adjusting the saliency of features for future case retrieval. It may not seem like adjusting the saliency of features is really explanation. However, when two or more hypotheses are both correct hypotheses, in that they can both cause the observed symptom, causal EBL-like explanations do not provide a way of distinguishing between them. The instructor chooses one of the hypotheses to pursue first. The student predicts a particular hypothesis will be pursued first. If the student's prediction is made based on case based reasoning, then the hypothesis predicted first depends on the matching function. Retrieval of previous cases involves searching for a case or generalization piece which served the goal currently being pursued. The retrieved case piece is selected from the candidate pieces based on a comparison of the feature values of the current problem solving context with the feature values of the problem solving context at the time of the previous case pieces. So adjusting the matching function by adjusting the importance of features in the problem solving context will lead to the prediction being correct in the future. This is an implicit way of explaining the choice between the hypotheses without having reason to say that one is more likely than the other. The intuition is that such weighting of competitive hypotheses in diagnosis is generally inductive, the mechanic doesn't know for a fact that x fails more often than y, statistics aren't readily available or used, nor can such preference be explained deductively. The weighting is inductive from experience, and from instruction. There is no evidence of this type of explanation in Lancaster and Kolodner's and Chi et al's observations. However, it isn't the sort of thing that would be amenable to study through protocols.

The method of adjusting the saliency of features is fairly simple. It is based on the idea of making

REDMOND

Goal - G-GENERATE-HYPOTHESIS
 Piece retrieved - (CASE-HYP-CHOKE-THERM 14.280001) Hypothesis (MALFUNCTION CHOKE-THERMOSTAT)
 Piece Expert's with hypothesis = (LOW IDLE-SPEED) - (GEN-HYP-ENGINE-STALLS 11.6)

Feature	Student's piece	Piece matching Instructor	Feature Importance
-----	-----	-----	-----
CAR-TYPE	Partial match	No Match	less important
CAR-OWNER	Match	No Match	less important
COMPLAINT	Match	Match	no change
FREQUENCY	Partial match	No Match	less important
HOW-LONG	Match	Partial match	less important
OTHER-SYMPT	Match	Match	no change
RULED-IN	Partial match	Match	more important
RULED-OUT	Partial match	Match	more important
TESTS-DONE-N-RESULTS	Partial match	Match	more important
FIXES-DONE	Partial match	Match	more important
CURRENT-HYPOTH	Match (none)	Match (none)	no change
PARTICIPANTS	Partial match	No Match	less important
LOCATION	Match	Match	no change
WHEN	Partial match	No Match	less important

Figure 6: Example Blame Assignment.

features that match when the problem solver is successful more important, and features that match when the problem solver is unsuccessful less important. Since the salience of various features varies depending on the goal being pursued by the problem solver, separate measures of feature importance are maintained for different goals. When the student predicts the same action the instructor makes, the student has been successful. The features of the current problem solving context that matched the features in the previous case are made slightly more important. When the student predicts a different action than the instructor, presumably the student has been unsuccessful. The blame assignment is best made by retrieving another case piece in which the instructor's action was the one done. Figure 6 shows how the blame assignment is done on an example incorrect prediction of a hypothesis. Those features of the current context that more closely match the context of the newly retrieved case piece than the context of the originally retrieved case piece will be made more important. Those features of the current context that more closely match the context of the originally retrieved case piece than the context of the 'correct' piece are made less important. Thus instruction with examples helps deal with the feature saliency problem, by giving feedback on the correctness of case retrieval, allowing comparison of the matching features.

This will lead to the correct piece being retrieved in the same situation in the future. A combination of instruction, case retrieval, and induction has been used to improve the performance of the CBR part of diagnosis.

CAUSAL EXPLANATION OF ACTIONS

Causal explanations of actions enable filling gaps in the causal domain knowledge through the basic LBUE methods described in Redmond and Martin [1988]. These were an extension of explanation-based learning (EBL) [DeJong 1983; DeJong and Mooney 1986; Mitchell, Kellar, and Kedar-Cabelli 1986], to allow learning without a complete and consistent domain model. An example will illustrate the ideas. An instructor may present the student with a malfunctioning car in which the engine cranks but does not start. She may suggest a hypothesis that the distributor cap is cracked. A complete causal explanation would be:

```

(cracked distributor-cap) causes
  (contains distributor-cap moisture) causes
    (low (input spark-plug electricity)) causes
      (not (ignite spark-plug)) causes
        (not (combustion cylinder)) causes
          (not (start engine))
  
```

REDMOND

If the student can complete the explanation, she can learn that a cracked distributor cap causes the symptom of the engine cranking but not starting. If the student was missing some knowledge, it is possible that the knowledge could be inferred as plausible. For example, if the student was missing the fact that moisture in the distributor cap can cause less electricity to reach the spark plug, she may still be able to infer that fact based on the partial explanation having been given by the trusted expert, in conjunction with the partial explanation formed by the student and general knowledge possessed by the student about water's effect on electricity.

In addition to enabling filling gaps in the causal domain knowledge, trying to causally explain actions can make causal explanations available as indices to the new case containing the action. Hammond and Hurwitz [1988], and Barletta and Mark [1988] both use this approach, which hasn't yet been implemented in the current system.

CONCLUSION

Explanation of solved example problems is an effective way of learning. A system has been constructed that uses EBL-like deduction, induction, and retrieval of previous cases in creating explanations, improving future diagnoses and future explanations of observed problem solving. The use of multiple types of explanation of examples follows the lead of the studies by Lancaster and Kolodner [1987, 1988] and Chi et al [in press]. Their observations suggest further types of explanation that could be exploited in making our system a better student. The exploitation of instruction turns out to be a powerful way of learning, and integrates several learning techniques.

ACKNOWLEDGEMENTS

This research was supported by the Army Research Institute for the Behavioral and Social Sciences under Contract No. MDA-903-86-C-173. The author wishes to thank Janet Kolodner for her advice and guidance, and Joel Martin, Louise Penberthy, and Chris Hale for helpful comments on earlier versions of the paper.

REFERENCES

- Barletta, R. & Mark, W. (1988). Explanation-based indexing of cases. In *Proceedings of a Workshop on Case-Based Reasoning*.
- Chi, M., Bassok, M., Lewis, M., Reimann, P., & Glaser, R. (in press). Self-explanations: how students study and use examples to solve problems. *Cognitive Science*, in press.
- DeJong, G. & Mooney, R. (1986). Explanation based learning: an alternative view. *Machine Learning*, 1, 145-176.
- DeJong, G. (1983). Acquiring schemata through understanding and generalized plans. In *Proceedings of the Eighth International Joint Conference on Artificial Intelligence*.
- Hammond, K. J. & Hurwitz, N. (1988). Extracting diagnostic features from explanations. In *Proceedings of a Workshop on Case-Based Reasoning*.
- Kolodner, J. & Simpson, R. Jr. (1984). Experience and problem solving: a framework. In *Proceedings of the Sixth Annual Conference of the Cognitive Science Society*.
- Lancaster, J. & Kolodner, J. (1987). Problem solving in a natural task as a function of experience. In *Proceedings of the Ninth Annual Conference of the Cognitive Science Society*.
- Lancaster, J. & Kolodner, J. (1988). Varieties of learning from problem solving experience. In *Proceedings of the Tenth Annual Conference of the Cognitive Science Society*.
- LeFevre, J. & Dixon, P. (1986). Do written instructions need examples?. *Cognition and Instruction*, 3, 1-30.
- Martin, J. & Redmond, M. (1988). The use of explanations for completing and correcting causal models. In *Proceedings of the Tenth Annual Conference of the Cognitive Science Society*.
- Mitchell, T. M., Kellar, R. M., & Kedar-Cabelli, S. T. (1986). Explanation based learning: an unifying view. *Machine Learning*, 1, 47-80.
- Reder, L., Charney, D., & Morgan, K. (1986). The role of elaborations in learning a skill from an instructional text. *Memory and Cognition*, 14, 64-78.
- Redmond, M. & Martin, J. (1988). Learning by understanding explanations. In *Proceedings of the 26th Annual Conference of the Southeast Region ACM*.