

**EXPERIMENTS WITH AN INTEGER PROGRAMMING
FORMULATION OF AN EXPERT SYSTEM**

by

Vasant Dhar

Leonard N. Stern School of Business
Information Systems Department
New York University
90 Trinity Place
New York, NY 10006

and

Nicky Ranganathan

Leonard N. Stern School of Business
Information Systems Department
New York University
90 Trinity Place
New York, NY 10006

March 1989

Center for Research on Information Systems
Information Systems Department
Leonard N. Stern School of Business
New York University

Working Paper Series

CRIS #205
STERN #89-39

This paper also appears as
MCC AI Laboratory Technical Report ACA-AI-022-89

Abstract*

We present expert system (ES) and Integer Programming (IP) formulations of an NP-complete constraint satisfaction problem (CSP). The problem involves generating a plan for assigning faculty to courses given a variety of constraints and preferences and other tentative data. The expert system consists of a heuristic rule-based problem solver and a truth maintenance system. The IP model consists of about 700 zero-one decision variables and 300 constraints. We describe and contrast the expert system and IP models in terms of behavior, quality of results, and computational performance. We find that the expressiveness of the IP model is hampered by its single objective function, inability to encode various types of complex preferences, the lack of useful output when it fails to find a feasible solution, and a general lack of control over inference. It is also difficult to make incremental revisions to the plan produced by the IP model. In contrast, the truth maintenance system maintains justifications for assignments, which makes it possible to reason about incremental modifications to a plan. In terms of performance, we found that whenever the IP approach finds a solution, it does so quickly using the Pivot and Complement heuristic of Balas and Martin (1980). The branch and bound *a/ways* failed to find a feasible integer solution when the heuristic failed to find one.

*This work was supported by the following shareholders: Bellcore, CDC, DEC, Harris, Kodak, and NCR. This technical report also appears as: Working Paper 205, Dept. of Information Systems, Leonard N. Stern School of Business, New York University, 90 Trinity Place, New York, NY 10006.

Table of Contents

1. Introduction	2
2. The Integer Programming Formulation	3
2.1. Subscripts	4
2.2. Variables	4
2.3. Coefficients	4
2.4. Constraints	6
2.5. Objective Function	7
2.6. Problem Structure	7
3. The Expert System Architecture	8
3.1. The Truth Maintenance System	9
3.2. Rules	11
3.3. An Example	13
4. The IP Architecture	14
5. Results	18
5.1. Performance	18
5.2. Control of Reasoning	18
5.2.1. Single Objective Limitations	18
5.2.2. Compiled Knowledge Limitations	19
5.2.3. Global Optimization Limitations	21
5.3. The Need for Partial Solutions	22
5.4. The Division of Responsibility	23
6. Plan Revision	23
7. Directions for Future Work	25
8. Acknowledgments.	26

1. Introduction

Virtually all decision making situations involve constraints. What distinguishes various types of problems is the form of these constraints. Depending on how the problem is visualized, they can arise as rules, data dependencies, algebraic expressions, or other forms.

Constraint Satisfaction Problems (CSPs) have been studied extensively in the Operations Research (OR) and Artificial Intelligence (AI) literature. In OR formulations constraints are quantitative, and the solver (such as the Simplex algorithm) optimizes (maximizes or minimizes) the value of a specified objective function subject to the constraints. In contrast, the interest in AI has been more with inference based on mostly symbolic constraints. The inference mechanisms employed include theorem provers, production rule interpreters, and various labeling procedures such as those used in truth maintenance systems.

It has been apparent for some time that there is a close relationship between logical and quantitative inference. Dantzig (1963) showed how logical propositions could be expressed quantitatively using Boolean variables together with algebraic operators. In this way, symbolic constraints can be modeled in terms of an integer programming formulation. For example, the constraint $\neg x_1, x_2 \rightarrow x_3$ is equivalent to the clause " x_1 or not- x_2 or x_3 " where each x_i is a propositional variable. A clause can be expressed as an inequality; the above clause is equivalent to

$$x_1 + (1 - x_2) + x_3 \geq 1$$

where the truth values true and false are denoted by 1 and 0 respectively. In general, as has been noted independently by Hinton (1979) and Hooker (1988), a clause can be expressed as the following inequality:

$$c_1 x_1 + \dots + c_n x_n \geq 1 - n(c)$$

where c is a row vector and x is a column vector, and $n(c)$ is the number of negative elements in the vector c . Each c_i is 1, 0, or -1, indicating whether x_i appears, does not appear, or $\neg x_i$ appears in the clause respectively. The above notation is due to Hooker (1988).

In concrete terms, the symbolic/quantitative equivalence means that many constraint satisfaction problems can in principle be modeled using symbolic inference techniques employed in Expert Systems or the quantitative techniques of OR. In practice, there are pros and cons to both approaches. One of our objectives has been to study the tradeoffs involved in detail.

Over the last year, several researchers of the Reasoning Architectures group at the MCC AI Laboratory have collaborated with an experienced department head at the University of Texas in analyzing the problem of planning the assignment of faculty to

courses, and revising plans as assumptions (about faculty availability and enrollments) change. The problem is of interest for several reasons. It characterizes many types of planning problems where solutions are assumption based or defeasible and must be revised in light of a changing reality. The problem also involves a diverse variety of constraints, making it a good canonical problem for analyzing different modeling approaches.

Based on an extensive knowledge engineering effort over the last year, an expert system has been built and implemented in PROTEUS (Petrie et. al, 1987). The expert system consists of a heuristic-rule-based problem solver which contains the expert's knowledge, and a truth maintenance system (TMS) which maintains a contradiction-free database of assignments.

In this paper, we present an OR (integer programming) formulation corresponding to the expert system and report on our experiences with running this model. We are interested in studying whether the integer programming model can provide the results that derive from the reasoning process of the expert, whether it can be modified easily to accommodate changes, and how general purpose integer programming techniques do in terms of computational performance.

The remainder of this paper is organized as follows. In the next section we provide a mathematical formulation of the problem and describe its general structure. This is followed by a description of the problem-solver/TMS architecture that constitutes the expert system. In section 4 we describe the architectural components of the IP implementation. We then analyze the results from running this model, comparing them with an expert's decisions when presented with the same problem data. In doing so, we classify some of the limitations of the mathematical model and describe why the problem-solver/TMS combination does not suffer from these problems (although it has other drawbacks). We conclude with a list of short and long term issues that must be addressed in order to design more expressive and flexible reasoning systems.

2. The Integer Programming Formulation

Generating a course schedule requires assigning faculty to courses while taking into consideration a variety of constraints and preferences. Constraints pertain to bounds on the number of sections of courses, minimal faculty teaching requirements, dependencies (i.e. whoever teaches *Introduction to Logic* must also teach *Advanced Logic*), etc. Preferences are expressed by teachers for desired courses and alternatives. In addition, the scheduler (chairman) expresses preferences in deciding how to deal with conflicts that arise in synthesizing the plan. After courses have been assigned, class rooms and times are scheduled. In this paper we do not deal with this latter problem.

The structure of the problem data is presented in figure 1. There are three *classes* of data involved: teacher, course, and term. There are three types of teachers (three sub-classes), each with specific instances. Likewise, there are two types of courses and instances of these two types, and two terms -- fall and spring. Constraints are defined at all levels of abstraction, in terms of the various classes and instances.

The formulation of the problem we present is based on a detailed problem description that appears in Petrie et.al (1989). Many of the constraints in the IP formulation have been derived from the rules in the expert system. The IP formulation is best presented by first defining all subscripts, variables, and coefficients, followed by the constraints and objective function.

2.1. Subscripts

$i = 1, 2, \dots, n$ instructors

$j = 1, 2, \dots, m$ courses

$t = 1, 2$ terms (fall or spring)

$k = 1, 2, \dots, p$ categories of courses (for balancing curriculum)

2.2. Variables

$$x_{ijt} = \begin{cases} 1 & \text{if faculty } i \text{ is assigned to course } j \text{ in term } t \\ 0 & \text{otherwise} \end{cases}$$

2.3. Coefficients

ρ_j = load factor of course j (depending on its size and type)

T_i = teaching load requirement (TL) for teacher i

lb_{jt} = lower bound on sections of course j to be taught in term t

ub_{jt} = upper bound on sections of course j to be taught in term t

lb_j = lower bound on total sections of course j to be taught in year

ub_j = upper bound on total sections of course j to be taught in year

y_j = number of sections of course j to be taught in the entire year

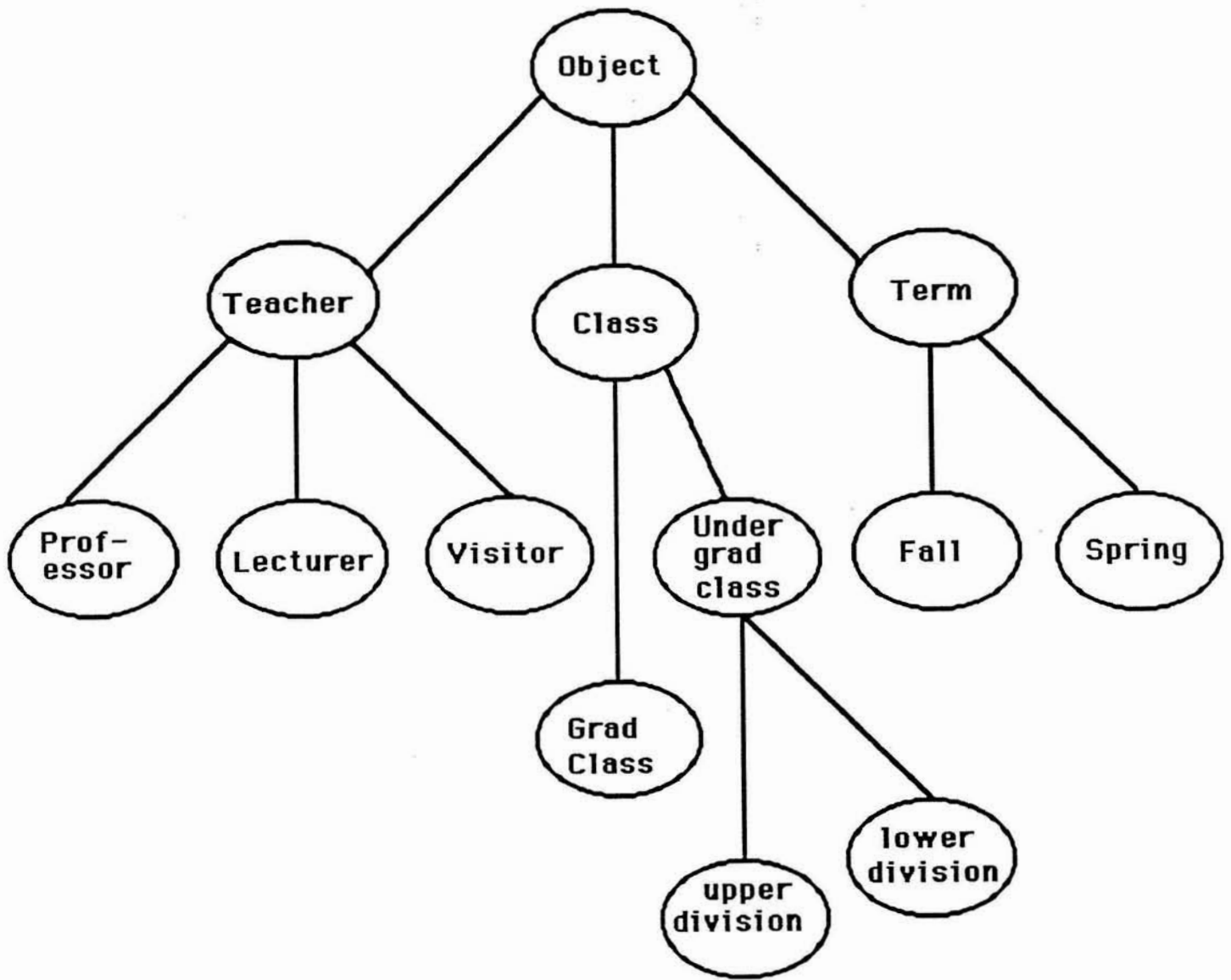
G_t = the number of graduate level courses normally offered during term t .

c_{ijt} = "cost" of assigning teacher i to course j in term t .

This can take on three values:

$$c_{ijt} = \begin{cases} 0 & \text{if course } j \text{ is desired by faculty } i \text{ in term } t \\ M, & \text{a small number if } j \text{ is an explicit alternative} \\ Z, & \text{a large number if } j \text{ is an implicit alternative or} \\ & \text{if faculty } i \text{ can teach course } j \text{ in term } t_1 \neq t \end{cases}$$

Actually, another number is added to the value of c to incorporate the "flexibility" of a teacher in terms of his capability and propensity to teach alternative courses.



Data Hierarchy

Figure 1

Specifically, teachers are ranked on a 1-5 scale (both inclusive) where a 1 expresses least flexibility. The combined number is the cost coefficient and expresses a "penalty" associated with assigning faculty i to course j in term t . We shall discuss this further in section 5.

$$g_j = \begin{cases} 1 & \text{if } j \text{ is a graduate level course} \\ 0 & \text{otherwise} \end{cases}$$

$$u_j = \begin{cases} 1 & \text{if } j \text{ is an upper division course} \\ 0 & \text{otherwise} \end{cases}$$

$$l_j = \begin{cases} 1 & \text{if } j \text{ is a lower division course} \\ 0 & \text{otherwise} \end{cases}$$

$$w_j = \begin{cases} 1 & \text{if } j \text{ is a writing course} \\ 0 & \text{otherwise} \end{cases}$$

$$f_i = \begin{cases} 1 & \text{if } i \text{ is a faculty member} \\ 0 & \text{otherwise} \end{cases}$$

$$C_{jk} = \begin{cases} 1 & \text{if course } j \text{ is in category } k \\ 0 & \text{otherwise} \end{cases}$$

$$t_{ijt} = \begin{cases} 1 & \text{if course } j \text{ is proposed as a tutorial by } i \text{ in term } t \\ 0 & \text{otherwise} \end{cases}$$

$$AB_j = \begin{cases} 1 & \text{if course } j \text{ has the A-B sequence} \\ 0 & \text{otherwise} \end{cases}$$

2.4. Constraints

1. number of teachers assigned to a course in each term should be between the lower and upper bounds on the number of sections of that course:

$$lb_{jt} \leq \sum_{i=1}^n x_{ijt} \leq ub_{jt} \text{ for } 1 \leq j \leq m, \text{ and } t=1,2$$

2. the total number of sections of course j taught in the year should be:

$$lb_j \leq \sum_{t=1}^2 \sum_{i=1}^n x_{ijt} \leq ub_j \text{ for } 1 \leq j \leq m, \text{ and } t=1,2$$

3. each teacher must satisfy some minimal teaching load:

$$\sum_{j=1}^m \rho_j x_{ijt} \geq T_i \text{ for } i=1,2,\dots,n, \text{ and } t=1,2$$

4. Only professors can teach graduate courses:

$$x_{ijt} \leq 1 - (1 - f_i)g_j \text{ for } i=1,2,\dots,n, \text{ for } j=1,2,\dots,m \text{ and } t=1,2$$

5. No professor can teach more than G graduate courses per year:

$$\sum_{j=1}^m \sum_{t=1}^2 x_{ijt} g_j \leq G \text{ for } i=1,2,\dots,n$$

6. Course sequence continuity (A-B) constraint:

$$AB_j(x_{ij2} - x_{ij1}) \leq 0 \text{ for } i=1,2,\dots,n, \text{ and } j=1,2,\dots,m$$

The above constraint says that someone can only teach a "B" course in spring if they taught the "A" course in the fall. If someone MUST teach the B course if they taught the A course in fall, the inequality must be changed to an equality.

7. There must be at least U upper division writing course offered each term:

$$\sum_{j=1}^m x_{ijt} w_j u_j \geq U \text{ for } i=1,2,\dots,n \text{ and } t=1,2$$

8. Similar constraints are formulated for lower division writing courses and the number of graduate courses that must be offered each term.

2.5. Objective Function

A variety of objective functions can be formulated for this problem depending on the primary goal. One reasonable goal is to maximize fit, or minimize deviations from teachers' desired courses. To express this, the objective function is:

$$\text{Minimize } z = \sum_{i=1}^n \sum_{j=1}^m \sum_{t=1}^2 c_{ijt} x_{ijt}$$

The above objective function states that the good of the many outweighs the good of the few.

2.6. Problem Structure

If one considers courses as "producing" sections for which there are maximum and minimum supplies (capacities) and teachers as "demand" points, the problem can be viewed as a transportation problem. In fact, with constraints 1 and 2 and the objective function, the problem is a classical transportation problem. This problem has a nice property in that it produces an integer solution for integer inputs. Thus the problem of fractional assignments does not arise.

Unfortunately, the third constraint, relating course load factors to teaching load requirements, is a knapsack problem constraint. This destroys the integer solution property. More importantly, the problem becomes NP-complete and its solution time varies extremely unpredictably.

Actually, the real problem is more complex in one other way. In our formulation, all sections of a course have been assumed to be of equal size (they have the same load factor). In fact, a few courses have sections of varying sizes. If size is also considered variable (another subscript added to the decision variable), the number of decision variables can increase significantly thereby increasing the overall complexity of an already difficult to solve problem. Specifically, the problem becomes a bin packing problem which is known to be NP-hard (Garey and Johnson, 1979). The expert deals with this complexity by assuming default section sizes in making assignments (ignoring TLs), and then increasing section sizes in order to satisfy the various teaching load requirements. In the IP formulation, we have dealt with this problem in a similar way by assuming default section sizes but lowering the TLs a little (about 10 percent). The expert then has the latitude in increasing section size and fine tuning the loads.

3. The Expert System Architecture

In contrast to the IP model, the expert system has no global objective function that guides it towards an "optimal" solution. Rather, the process of finding a solution involves a series of local decisions that are most preferred at each point in the problem solving process.

The expert system consists of two components, the problem-solver and the TMS. The problem solver is essentially a production system consisting of rules whose patterns are matched by data (assertions) in a global database. Each assertion is a problem-solver datum such as an assignment of a value to a variable (i.e. a decision) or an operation that can result in an assignment. The problem data is represented in terms of an ISA object hierarchy corresponding to the structure in figure 1. Problem constraints are stated in terms of the objects in the hierarchy.

Each problem-solver datum has a justification associated with it. The justification encodes the reasons for belief in the datum. The task of the TMS is to ensure that the global database is contradiction-free, that is, the data are logically consistent (their truth values do not result in logical contradictions).

Before describing the types of problem-solver rules it is expedient to explain how the TMS works. For an overview of truth maintenance systems, the reader should refer to survey articles by Reinfrank (1988), McAllester (1982), or Dhar (1989). However, the fundamental concepts and mechanics underlying the TMS should be clear in the following description.

3.1. The Truth Maintenance System

There are two types of TMSs, *justification-based* (Doyle, 1979) and *assumption-based* (de Kleer, 1986). We shall concern ourselves with the former type. Further, the discussion is in terms of a Doyle-style TMS since this is the type of TMS used in the expert system.

In a TMS, every datum has a *support status* associated with it; a status of IN indicates that the datum is currently believed whereas an OUT indicates disbelief. The values IN and OUT are computed via *justifications* associated with the datum. Each justification has two parts, an *inlist* and an *outlist*. A justification is considered *valid* if it evaluates to *true*, that is, if each datum in its inlist is IN and each datum in its outlist is OUT. Ultimately, all data depend on "ground" level justifications of two types, *premises* and *assumptions*. In a premise justification, the inlist and outlist are both empty. This type of justification is always valid. An assumption justification has a non-empty outlist, that is, its belief is justified by a lack of belief in some other datum. Such a justification is considered *non-monotonic* since its validity depends on a lack of belief in some other datum. Finally, deductive justifications are those in which the outlist is empty. A datum can have more than one justification associated with it. A datum is IN if it has at least one valid justification, otherwise it is OUT.

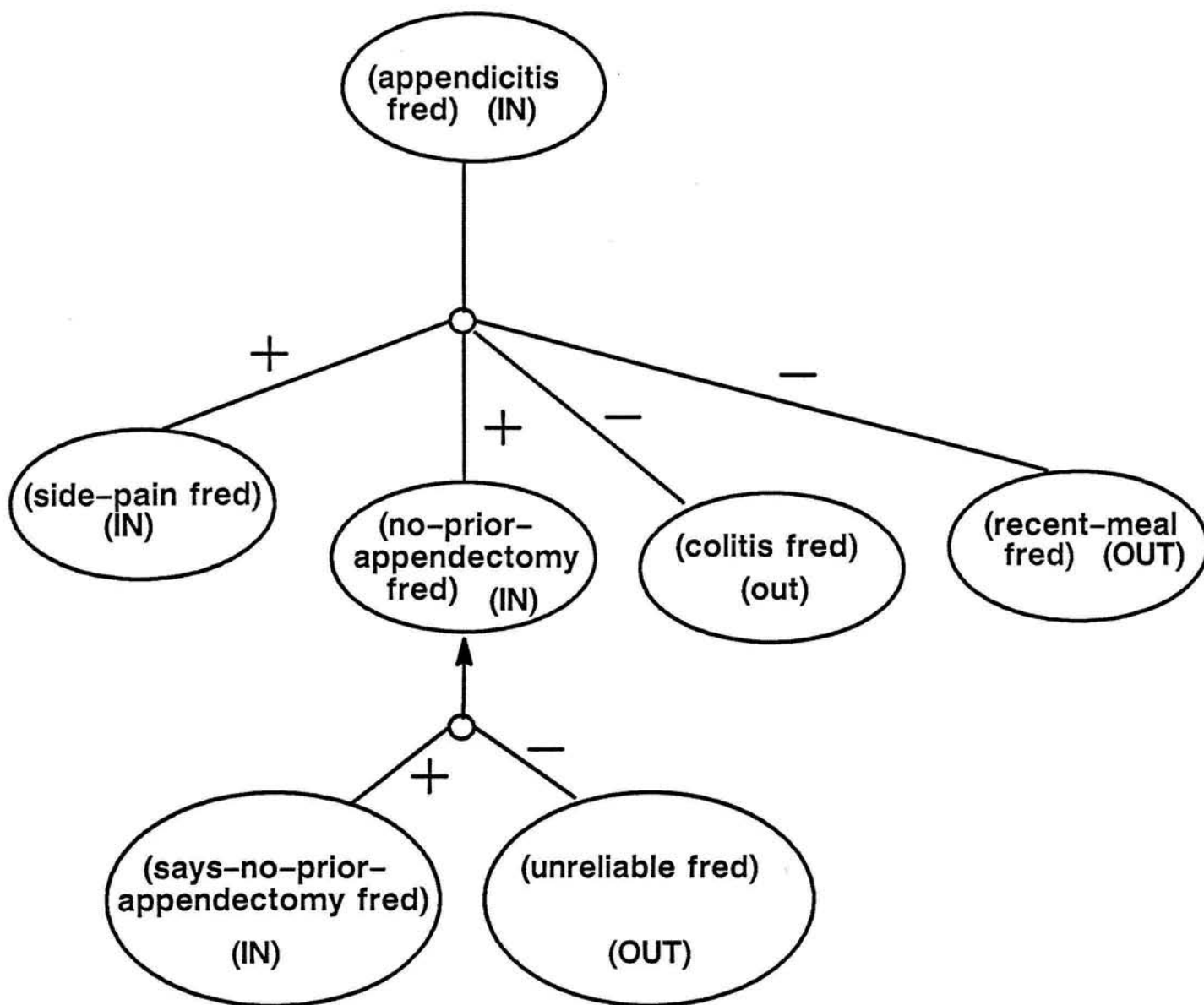
To illustrate, consider the following simple example which is taken from Petrie et.al (1987). Upper case symbols denote functions and the lists of lower case symbols refer to data:

```
Datum: (appendicitis fred)
Justification: (AND (INLIST (side-pain fred) (no-prior-appendectomy fred))
                (OUTLIST (colitis fred) (recent-meal fred)))
```

```
Datum: (no-prior-appendectomy fred)
Justification: (AND (INLIST (says-no-prior-appendectomy fred))
                (OUTLIST (unreliable fred)))
```

The example can be visualized in terms of the graph in figure 2 which shows the justifications of data referenced in the above two justifications. Each circle corresponds to a justification, with an arrow pointing to the justified datum, positive arcs connected to the elements of the inlist and negative arcs connected to the elements of the outlist. The datum "(appendicitis fred)" has a two element inlist and two element outlist. The datum "(side-pain fred)" has an empty inlist and empty outlist and is therefore a premise. The datum "(unreliable fred)" has an empty list of justifications and is therefore OUT. If this datum were to acquire a valid justification, its support status and that of those that depend on it must be reevaluated. Specifically, the "(no-prior-appendectomy)" would go out, also causing "(appendicitis fred)" to go out.

Actually, in reevaluating the belief status of the data, referred to as *reason*



A Consistent Well-Founded State

Figure 2

maintenance, what the TMS really does is to ensure that the data as a whole satisfies two properties, *consistency* and *well-foundedness*. In a stable state each datum with at least one valid justification is IN and each one without a valid justification is OUT. A state is well-founded if no set of beliefs is mutually monotonically dependent; in terms of a labeled network such as that in figure 2, this means that there is no set of arcs from a node to itself all of which are labeled positively.

In the expert system, the TMS works in conjunction with the problem solver as follows. Each problem-solver action is communicated to the TMS at which point the TMS executes a constraint satisfaction procedure to ensure consistency and well-foundedness. Essentially, this involves updating the justifications associated with each datum such that the two conditions are satisfied. A problem-solver action can lead to a constraint violation which is recorded in the network as a contradiction, that is a special node called "contradiction" becomes IN. When this is detected by the TMS, it tries to compute a new labeling that makes the contradiction OUT. Control then passes back to the problem solver, and the cycle repeats. We shall describe in more detail the workings of the TMS in the context of an example later in this section.

3.2. Rules

In the following paragraphs we describe the expert system architecture using a simplified set of rules that play a role in one part of the problem solving process. The rules employ a LISP-like prefix notation that is used in PROTEUS (Petrie et.al, 1987). In the following rules a question mark followed by a symbol denotes a variable, i.e. ?xxx is the variable xxx. Each list in a rule consists of a pattern (a form) that is matched against a datum in the global database. There are two types of rules, forward and backward. When the problem solver has a proposition that it is trying to *prove* to be true, it attempts to do so via a backward rule. A "<--" symbol denotes a backward rule; the form appearing before it is the consequent and the ones after it are antecedents. In trying to prove a goal, the assertions that match the antecedents become subgoals. Proving all of these subgoals completes the proof of the original goal. This process is known as backward chaining.

The symbol "-->" denotes a forward rule. The forms appearing before it are antecedents, the others are consequents. The forward rule "fires" when assertions match its antecedent. When this happens, the datum (instance) corresponding to its consequent is added automatically to the global database. This process is known as forward chaining.

In the following rules, the variables used in the patterns are as follows:

?*prof* refers to the current professor under consideration.

?*sem* refers to a semester (either fall or spring).

?*course* refers to the course that is under consideration for assignment.

?p-cycle refers to the list of professors who have already been considered for *?course*.
?c-cycle refers to the list of courses for *?prof*.

The following backward rule states that the goal of making a course *?course* a chosen course for professor *?prof* in semester *?sem* can be achieved if that course is actually desired by the professor unless that course has already been considered (and presumably failed) for that professor:

```
(Mark-desired-course-as-chosen ;;name of this rule, rule 1
  (chosen-course ?prof ?course ?sem ?c-cycle)
  <--
  (desired-course ?prof ?course ?sem ?c-cycle)
  (unless (element ?course ?c-cycle)))
```

The following forward rule would fire whenever there is an assertion stating that courses be assigned to a professor:

```
(Select-Course ;; rule 2
  (assign-courses-to-prof ?prof ?sem ?p-cycle ?c-cycle)
  (chosen-course ?prof ?course ?sem ?c-cycle)
  -->
  (attempt-to-satisfy ?prof ?course ?sem ?p-cycle ?c-cycle))
```

In the rule above, the datum corresponding to the consequent would have the two data items matching the first two antecedent forms on its inlist.

The remaining rules can be interpreted similarly. Each form is accompanied by a comment whenever necessary.

```
(Semester-Succeeds ;; rule 3
  (attempt-to-satisfy ?prof ?course ?sem ?p-cycle ?c-cycle)
  (unless (unacceptable-for-semester ?sem ?prof ?course ?p-cycle ?c-cycle))
  -->
  (successful ?prof ?course ?sem ?p-cycle ?c-cycle)) ;;mark as successful
```

In the above rule, a datum corresponding to the consequent would have the datum matching the first form on its inlist and the datum matching the form in the "unless" part on its outlist. In effect, it would have a non-monotonic justification.

```
(Semester-Fails ;; rule 4
  (attempt-to-satisfy ?sem ?prof ?course ?p-cycle ?c-cycle)
  ;; If this semester is not acceptable
  (unacceptable-for-semester ?sem ?prof ?course ?p-cycle ?c-cycle)
  -->
  (failed-semester ?sem ?prof ?course ?p-cycle ?c-cycle))
```

```
(contradiction-detection-type-1 ;; rule 4
  (failed-semester ?sem ?prof ?course ?p-cycle ?c-cycle)
  -->
  (CONTRADICTION))
```

There are many rules used to determine why a course should not be assigned to a prof in a particular semester (unacceptability of a semester). The following rule, for

example, states that a professor may not teach a course if the maximum number of sections for it have already been assigned:

```
(course-section-check    ;; rule 6
 (unacceptable-for-semester ?sem ?prof ?course ?p-cycle ?c-cycle)
 <--
 (semester-max-constraint-satisfied ?course ?sem))
```

In effect, the rule above encodes constraint knowledge. In addition to the assignment and constraint checking types of rules, the system also contains rules for contradiction resolution. When a contradiction arises, the TMS first finds a *culprit* to be made OUT to resolve the contradiction. The process of finding a culprit is a recursive one that involves determining the ground level data (the assumptions) that support the contradiction. The contradiction is resolved by invalidating all valid justifications of one of its supporters, called the *culprit*. This involves justifying some belief on the outlist of the justification being invalidated. For example, if an attempted assignment fails, one way to resolve the problem is to exchange that (chosen) course with an alternative course that that professor can teach. The following rule shows the use of a predicate, FIX, that is designed to do contradiction resolution as described above. In Petrie's (1987) terminology, a datum that unifies with its first argument is referred to as the target, one that unifies with its second as the fix-culprit, and the third as the fix-elective:

```
(fix-via-alternative-course    ;; rule 7
 (fix (failed-semester ?sem ?prof ?course ?p-cycle ?c-cycle)
      (chosen-course ?prof ?course ?p-cycle ?c-cycle)
      (exchange ?alt-course ?course ?prof ?sem ?p-cycle ?c-cycle))
 <--
 (alternative-course ?prof ?course ?alternatives)
 (element ?alt-course ?alternatives))
```

If an exchange is successful in resolving a contradiction, the following rule makes the alternative course a desired course, which in turn causes the first rule to make it chosen, thereby repeating the cycle.

```
(Try-An-Alternative    ;; rule 8
 (exchange ?alt-course ?course ?prof ?sem ?p-cycle ?c-cycle)
 -->
      ;; Mark the alternative course as desired
      ;; (this will cause the Select-Course-Rule to fire)
 (desired-course ?prof ?alt-course)
      ;; Reject the original course.
 (rejected-course ?prof ?course ?p-cycle ?c-cycle))
```

3.3. An Example

In order to illustrate how the rules above work, let us consider a scenario where among other data, the following are in the database:

```
Sartre wants to teach phi304 in fall
(desired-course sartre phi304)
```


Sartre has as alternatives to phl304 in the fall, phl310 and phl318:
(alternative-course sartre phl304 fall (phl310 phl318))

Suppose that the maximum allowable sections of phl304 have already been allocated for fall:
(semester-max-constraint-satisfied phl304 fall)

Suppose Sartre is now up for consideration (Anselm has already been considered for phl304):
(assign-courses-to-prof Sartre fall (Anselm) ())

The rules try assigning phl304 to Sartre in the fall and fail because the maximum constraint has been satisfied for the term. This creates a contradiction as shown in figure 3. For simplicity we have truncated the data, assuming that they refer to Sartre in the fall term.

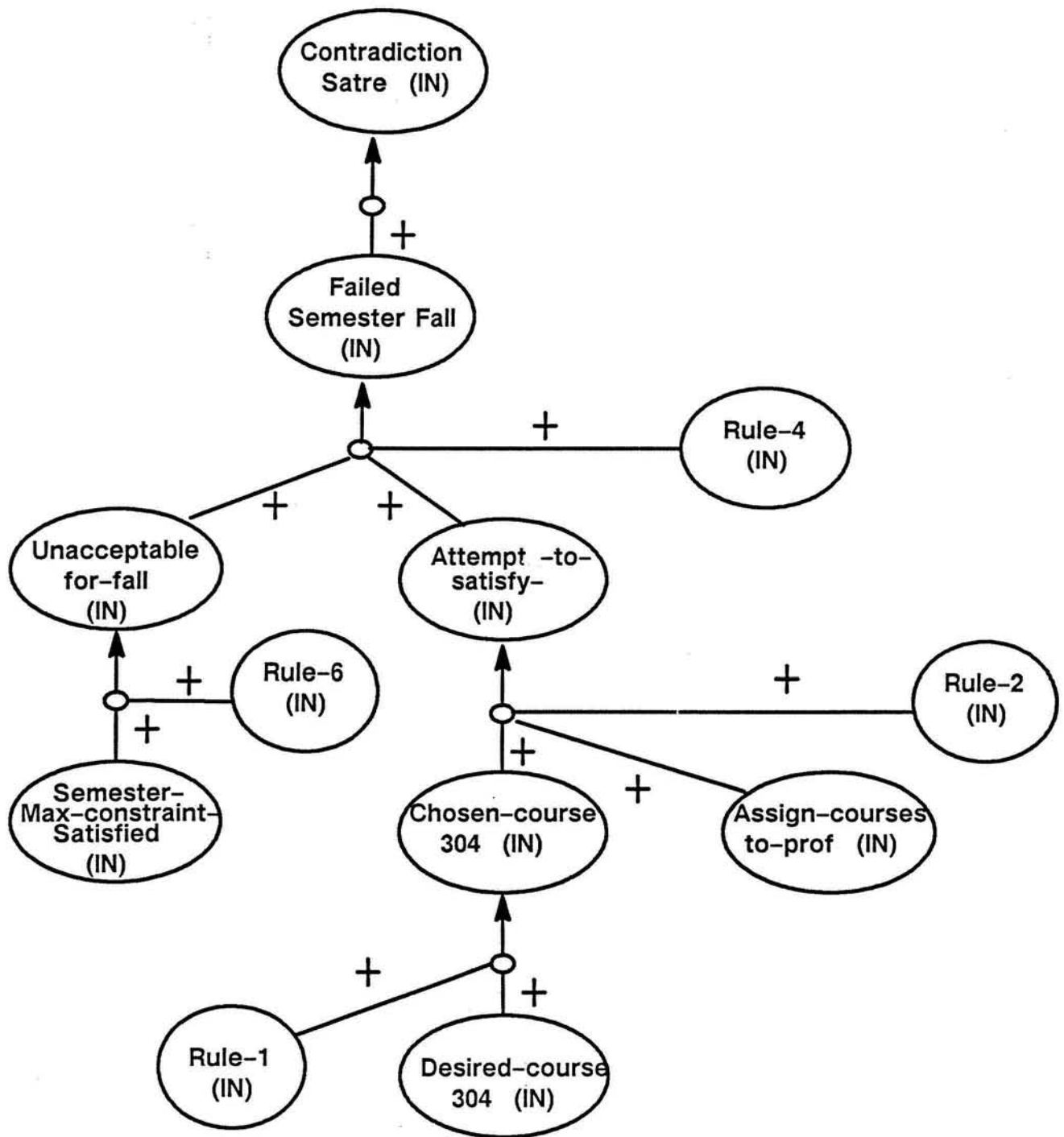
In the dependency network in figure 3, the contradiction is IN. The TMS attempts to make the contradiction OUT by finding a culprit and invalidating one of its in-supporters. In this case, the culprit is "(chosen-course 304)" (it unifies with the fix-culprit in rule 7) and its in-supporter "(desired-course 304)" is invalidated by putting the "(exchange 310 304)" datum on its outlist, and "(desired-course 310)" gets a valid justification by the exchange datum (via rule 8) on its inlist. Finally, the contradiction is put on the outlist of exchange which makes it (the contradiction) OUT. Part of this new stable and well-founded state is shown in figure 4.

In summary, the dependency network maintains the reasons for assignments. This includes desired as well as unexpected assignments such as the one resulting from the exchange effected by rule 8. Whenever a justification for some datum becomes invalid, the TMS computes what beliefs must be revised in order to restore consistency and well-foundedness. It is important to recognize that the problem of contradiction resolution is often an underconstrained one, that is, there are many possible labeling that satisfy consistency and well-foundedness. The more the application specific knowledge provided, the more *reasoned* is the behavior of the system.

4. The IP Architecture

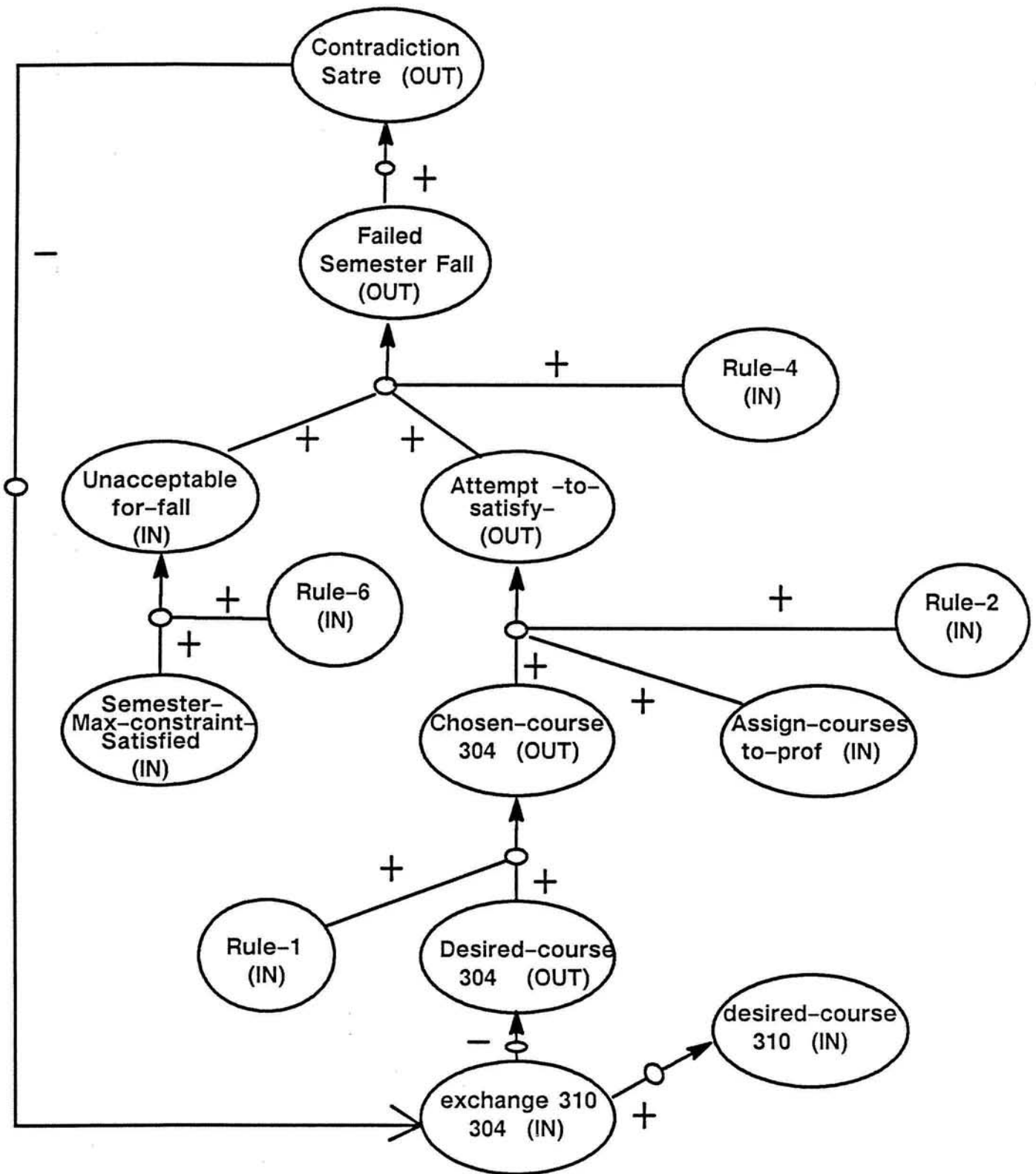
The integer programming formulation consists of about 700 binary variables and 300 constraints. The model has been implemented using the ZOOM (Zero One Optimization Model) library of the XMP package (Marsten, 1987). XMP includes a modeling language, XML, for expressing the problem.

ZOOM solves an integer program as follows. The LP relaxation (ignoring integrality) is



Contradiction State

Figure 3



Contradiction-Free State

Figure 4

first solved. After an LP solution is obtained, a heuristic procedure, called Pivot and Complement (Balas and Martin, 1980) attempts to find an integer solution. Basically, this involves a sequence of pivot operations which put all slacks into the basis at minimum cost. If a feasible integer solution is found, it is then improved by flipping variables to their opposite bounds.

As shown in figure 5, we have implemented a preprocessor that translates the input data used by our expert system into a form acceptable by XML which produces an MPS file that is used by ZOOM. The preprocessor takes as input problem data expressed in terms of the object hierarchy of figure 1. Expressing the constraints in XML requires the preprocessor to translate them into algebraic expressions stated in terms of decision variables. For example, a constraint such as "A faculty member can teach at most 1 graduate course per year" involves searching the ISA hierarchy to locate all instances of faculty, graduate courses and terms, defining the decision variables, and writing out the constraint. In this way, only decision variables essential to the formulation are defined. The number of decision variables can be further reduced by analyzing each teacher's course repertoire and excluding variables corresponding to impossible assignments. For example, if Frege's repertoire slot does not include phi381, he can never be assigned this course, and there is no point in defining a decision variable for this assignment.



Figure 5

The results produced by ZOOM are translated into a schedule for the user. Although not implemented by us, it is also possible to produce other summaries and answer questions using the generated schedule and the object hierarchy. For example answers to questions like "how many faculty are teaching undergraduate courses that meet twice a week", can be very useful to the decision maker.

The preprocessor and translator are both implemented in Common Lisp. The experiments were carried out on a SUN-3 workstation, which is also one of the platforms on which the expert system has been developed.

5. Results

In this section we report the computational performance of the IP model and an analysis of the results it produces relative to the expert model.

5.1. Performance

Like many integer programs, we found the solution time to be highly unpredictable, varying from a few minutes to a few days. It was extremely sensitive to the TL constraint. Changing the TL requirement by as little as 2% resulted in orders of magnitude variations in solution time. We varied the TL in roughly an 8 percent range and could not find feasible solutions in the top quarter of this range.

In all cases, when a solution was found, it was by the pivot and complement heuristic (Balas and Martin, 1980) which is incorporated in the ZOOM code. It always took under 5 minutes on a dedicated SUN-3. The heuristic never found a solution once it started executing the rounding procedure. In effect, if a feasible solution was not found quickly, it was not found at all. Further, the branch and bound algorithm, which ZOOM resorts to if the heuristic fails, *never* found a solution even after many hours of running time. This is because the number of fractional valued variables remaining after the LP solution, typically between 300 and 500, gives rise to an enormous search tree.

In contrast, the expert system's solution time was less volatile. It usually takes between 1 and 2 hours for it to generate a solution, or in cases it cannot find one, to generate a partial solution.

5.2. Control of Reasoning

In general, the integer program generates a plan that has about three-quarters of the same assignments as those made by the expert. In some cases the expert was pleasantly surprised by its decisions, but in a larger proportion of the cases the differing assignments were judged to be undesirable. The differing assignments arise for three reasons which we call single objective limitations, compiled knowledge limitations, and global optimization limitations.

5.2.1. Single Objective Limitations

The objective function used in the formulation expresses one goal of the decision maker: to give all teachers their desired courses to the extent possible. However, in reality there are other goals that the expert tries to satisfy simultaneously. One of these is to ensure that as far as possible, each teacher's load is as close as possible to the minimum required. What is required, therefore, is a multi-objective formulation. Unfortunately, this makes solution extremely difficult. Even if the problem is solvable, a "frontier" of optimal solutions is generated which must be evaluated by the decision

maker. For large problems involving discrete choices, analyzing the goodness of alternative sets of assignments can be difficult.

While the single objective formulation works well in most cases, it sometimes produces undesirable results. For example, consider two teachers, Helga and Hegel, who have been assigned courses that result in TLs of 9 and 12 respectively and that 9 is the minimal requirement. Assume that Hegel is more flexible than Helga. Now, if a section of Ethics needs to be taught and Helga and Hegel are the only qualified teachers, Hegel will be assigned since his cost coefficient is lower than Helga's. Clearly, if balancing loads is a concern, the course should be assigned to Helga.

It is possible to alleviate this problem to some extent by expressing the objective function as a constraint, but this too has severe limitations. One must decide on a reasonable TL upper bounds for each teacher. In the above example, if Hegel's upper bound is 14, the course will be assigned to Helga, as desired. However, setting such bounds is difficult. If an inappropriate value is used, feasible solutions can be excluded.

In contrast, it is relatively easy to encode knowledge about multiple objectives in the expert system rules. For example, in the Hegel/Helga case, it is possible to have in the antecedent of a rule that attempts to assign courses (such as the select-course rule) a form such as "(least-loaded ?prof)" which would result in the rule condition becoming true (unification succeeding) only with the least loaded professor. Such rules ensure that "locally good" decisions will be made but they do not guarantee a global optimal solution. We shall elaborate on this point in the following subsection.

5.2.2. Compiled Knowledge Limitations

The behavior of the system and the solution are very sensitive to the cost coefficients. Each teacher has three cost values, corresponding to the penalty associated with assigning a desired course, an explicit alternative, and an implicit alternative (a course that the teacher is capable of teaching but did not ask for). The values are organized as follows:

	Desired	Explicit alt.	Implicit alt.
CFI=1	1	10	19
CFI=2	2	9	18
CFI=3	3	8	17
CFI=4	4	7	16
CFI=5	5	6	15

Observing the first column, we see that a teacher with less flexibility (CFI=1) will get preference for a desired course over a teacher that has more flexibility. Moreover, the penalty associated with giving the former an explicit alternative instead of a desired course (9) is higher than doing it for the latter (1). For conflicting desired courses, if the

less flexible teacher gets the desired and the less flexible one gets an explicit alternative, the cost is 1+6, otherwise it is 5+10. Similar reasoning applies to the situation where there are conflicts in the explicit alternatives. Finally, note that the coefficients are designed so that a desired plus explicit alternative combination is always preferred to a desired plus implicit alternative combination.

In general, the above cost coefficients work well. However, they incorporate a lot of compiled knowledge about preferences, flexibility and trading criteria which makes the behavior of the system somewhat unpredictable. We have found that changing the cost values and differentials can have significant unforeseen (desirable or undesirable) consequences on the assignments in the following ways. When the differentials between explicit and implicit alternatives are made larger, the penalty associated with assigning implicit alternatives is high, hence fewer of these are assigned. However, this also has the effect of reducing the desired courses assigned since in conflicting situations the penalty associated with assigning explicit alternatives relative to the desired courses is low. When both types of differentials, desired-explicit and explicit-implicit, were increased, there were a few changes in assignments. However, we were not able to determine a general pattern underlying these changes nor the reasons for the changes.

Another related difficulty with compiled knowledge is that of encoding complex preferences in it. For example, it is usually preferable to teach two instead of three courses but if someone needs to teach three (to make the required TLs) it is preferable to have two sections of one course and one section of another instead of one section each of three different courses. This is because there is a diminishing marginal effort associated with teaching an additional section of a course. Such knowledge is not expressed in the coefficients of the objective function, but it is important in matching the expert's behavior.

Another type of knowledge that is difficult to express either in the constraints or the objective function is one involving unusual situations, that is, an action that is rarely undertaken but is a good one under some circumstances. For example, a seminar required in the spring term would not normally be swapped into the fall term except under certain unusual circumstances. Clearly, it is not appropriate to exclude such actions via constraints. Nor is it appropriate to attach heavy penalties to such actions, since this would prevent them from being taken even under appropriate conditions.

The examples above bring out some of the differences in how the IP and ES models incorporate problem solving knowledge. The ES proceeds locally in a GPS-like manner, reducing the differences between an evolving partial solution and the requirements based on preference information encoded as heuristic rules. In effect, the rules are

sensitive to the state of the evolving solution, giving the system builder considerable control over inference.

It is easy to define rules that encode knowledge about the types of preferences discussed above, that is, preferences for assigning additional sections of a course before considering a new course, preference in unusual situations, or for assigning courses to less loaded teachers as discussed in the previous section. In fact, it is such rules that are "compiled" into the coefficients in the IP formulation. Given the inevitable loss of information resulting from this transformation it is not surprising that there is a general loss in control over reasoning with the IP model. In other words, in the IP formulation, all such preference knowledge is usually compiled into one global objective function, which controls how the search space is explored. In the next section, we comment on some of the consequences of global optimization.

5.2.3. Global Optimization Limitations

While the minimization of cost is designed to maximize the extent to which teachers are assigned desired courses as a whole, the system has a tendency to schedule sections that are as close as possible to the lower bound since this also minimizes cost. This can have the effect of not assigning such courses to teachers that desired them. To illustrate, if at most one section on Ethics should be scheduled, the system has a tendency to schedule none. This turns out to be undesirable in situations where a professor requested Ethics but was assigned an alternative instead.

Essentially, the above problem would be avoided if in the case of the senior professor asking for that course, the lower bound is set to one instead of zero. In effect, the constraint is conditional on the data. However, deciding the appropriate bound based on the input data is tricky since tightening it could rule out feasible solutions or interfere negatively in unforeseen ways with others' assignments.

The only way to express constraints as conditional on the data is to use non-monotonic justifications. For the example above, a rule with an "(unless ...)" form could be used to set the bound. Then, an assignment would have the datum matching the condition in the *unless* in its outlist, making it clear that one section was scheduled since a senior professor desired it. If that senior professor were to go on leave, invalidating the desire for that course, the TMS could automatically validate the datum specifying the bound of zero. In general, a TMS is a natural mechanism for modeling default reasoning of this type.

Another consequence of global optimization is the lack of explanation for its decisions. In analyzing the results of the IP we were generally able to infer after analyzing the data in detail, why a teacher had not been assigned a desired course. For the most part, this

happened when there was competition for a limited number of sections. Often, however, the rationale for assignments could not be determined even after considerable analysis of the constraints, preferences and flexibility indices. In such cases, global optimization essentially obscures the reasons for assignments.

In contrast, the ES attempts to make good local assignments the justifications for which are recorded by the TMS. The justifications are extremely useful for purposes of explanation and for incremental revision of existing decisions in an evolving solution. Although the ES has no notion of a global optimum, the more the knowledge provided to the system for resolving conflicts, making choices, etc, the better the quality of its solution. The factors that determine whether the expert system will work better than the IP model are the extent to which it is important to specify complex preferences of the type described in the preceding discussion, and the ease with which such knowledge can be specified by the expert. These two factors interact in a complex manner. If complex preferences need to be specified, it is usually indicative of the complexity of defining optimal solutions, and hence the limitations of the objective function as the mechanism driving the search. However, in such situations the expert also finds it increasingly difficult to specify the preference knowledge in terms of abstractions. Over time, this can result in a situation where interactions among the various pieces of knowledge become very complex, thereby eroding the modularity of the knowledge base.

5.3. The Need for Partial Solutions

A major problem during the early runs of the IP model was its inability to find a feasible integer solution even after many hours of running time. Two factors contributed to this situation. Firstly, the data were often such that the constraints were not satisfiable. For example, if three sections of a course were required and only two people were qualified to teach it, clearly, no solution would be found. Secondly, we found that the knapsack constraints, on TLs, were often too tight for it to find a feasible integer solution.

We solved the first problem by analyzing the behavior of the expert system. While that system is also often unable to find a feasible solution, it generates a useful partial solution indicating those "holes" in the schedule that still need to be filled.

Partial solutions are extremely useful to the decision maker. There is no guarantee that the skills of the teachers will cover the requirements (particularly if a significant number of faculty are on leave). In such cases, and the decision maker needs to know what the holes are. Typically, there are several holes in the schedule in the initial draft. These are patched by hiring visiting professors, lecturers, or graduate students. In this respect, the integer program is deficient because if it fails to find a feasible integer

solution, it runs virtually forever and finally provides no useful information to the decision maker. What is required in such cases is an integer solution which even though not feasible, is "close" to feasible.

Even though the results of the expert system pointed us to the section requirements constraints to loosen, we still did not manage to obtain a solution until we simplified the problem by excluding the TL constraints, solving it as a transportation problem. We then introduced lowered TL requirements. As we mentioned earlier, we found that the ability of the system to find an integer solution was extremely sensitive to the TL value. On analyzing the expert's behavior, we found that he dealt with the TL requirements by first ignoring them, making assignments that put teachers' loads in the "ballpark", and then massaging the schedule to satisfy the requirements. For the most part, this massaging consists of giving teachers larger sections (which have higher credits) or assigning them light administrative responsibilities for which they earn small amounts of credit.

5.4. The Division of Responsibility

We have come to believe that it is neither practical nor reasonable for a system to always generate a complete plan. There are simply too many factors to be accounted for in generating a plan that are not reasonably enumerable a priori. For example, the expert might feel that a teacher has been assigned too much of a certain topic, which is a judgment call, and wish to alter the plan accordingly.

The user might also use creative tricks in massaging a plan. In one schedule, for example, a course was cross-listed with another department and therefore excluded from participating in the constraints defined over courses in its class. In that case, the expert added another course from that class to the schedule in order to satisfy the desires of a professor.

In summary, from a practical standpoint, a support system must function as a decision support system instead of a fully automated planner. In the next section we discuss the types of levers that are available to the user for tuning and revising plans in the IP model and the directions in which it needs to be enhanced.

6. Plan Revision

Since planning problems are based on assumptions that may change, it is important that a decision maker be able to modify a model incrementally. Planning course assignments is based on assumptions about enrollments (which determine the numbers of sections planned), and faculty availability. If enrollments turn out to be higher than expected, additional sections must be scheduled. If a faculty member gets a grant or goes on leave, substitute teachers must be found. In all such cases, it is important that the overall plan be perturbed as little as possible.

Several types of actions can be taken by the decision maker when assumptions change. Depending on the change, the action can include hiring assistant instructors, hiring visiting professors, and swapping assigned courses among faculty. Hiring a visitor is usually feasible only if there is sufficient time to negotiate a contract.

Some changes are easy to manage. For example, if a faculty that goes on leave was teaching low level undergraduate courses, instructors can be hired to fill in. However, if graduate courses are involved revision becomes more complex since such courses can be highly specialized, making it difficult to find faculty qualified to teach them. Also, faculty are limited to a maximum of one graduate course per year. In summary, determining who to assign to an unassigned graduate course can be difficult, depending on how many faculty are qualified to teach the course, how flexible they are, whether are already scheduled to teach a graduate course, and whether substitutes can be found for such courses.

The IP model supports revision as follows. The decision maker must specify what subset of the existing set of assignments can be changed. The complement set is considered fixed; the assignments in this set are therefore added as constraints and the problem is re-solved. In effect, the decision maker must specify what part of the plan is fixed and what is variable. Further, the decision maker must make a judicious choice in specifying what combination of lecturers, Assistant Instructors and visiting professors should be considered to fill the newly created holes in the plan.

Unfortunately, it is usually difficult for the decision maker to specify in advance with any degree of confidence which parts of the plan should be considered variable. Rather, the process of figuring out what to change requires negotiation, the results of which serve as further input in determining what parts of the model can be considered changeable. In effect, figuring out what can be changed is where support is most needed. In this sense, the IP model requires the decision maker to do too much. Specifically, if changes are to be minimal, he must keep the variable set small, otherwise the new solution can contain too many changes. However, this small set of changeable assignments that should also result in a solution being found, otherwise the exercise is of little use. Determining the appropriate set can require considerable trial and error; finding it is therefore the really difficult part of the problem where support is most needed.

Another aspect of revision is that certain changes actually require that changes be made to the constraints. For example, if an instructor who goes on leave was teaching a course that is not strictly required, it is usually not re-assigned to anyone when the IP model is run with the new data since the objective is to minimize cost. Clearly, this is problematic if that course has already been listed as an offering. In such cases, what is

really required that the constraints be modified to state that the listed course *must* be offered.

In order to support revision, it is necessary that a system actually *suggest* alternative courses of action to the decision maker. For this to be possible, it is necessary for it to record the rationales for existing assignments. For example, if an instructor I_1 who goes on leave was teaching X_1 which was desired by but not assigned to I_2 (perhaps due to an exchange like the one illustrated in the example), it might make sense to consider whether I_2 should now teach X_1 . This type of reasoning requires a truth maintenance system. In a TMS, for example, one of the justifications for the assignment $I_2 \leftarrow X_2$ would be the fact that I_1 was assigned X_1 . When this latter proposition is no longer true, the justification for $I_2 \leftarrow X_2$ becomes invalid, making it possible to assign X_1 to I_2 . This process can be repeated recursively until a consistent set of assignments is found.

In practice, however, it would probably be undesirable for the above process to happen automatically for two reasons. Firstly, it is conceivable that taking away X_2 from I_2 could create an even bigger problem if it is difficult to find someone to teach X_2 . Secondly, it would not make sense to take such a course away if a significant amount of preparation has already gone into preparing for that course or if the teacher assigned to it is inflexible.

Regardless of the above limitations, a TMS should prove to be useful. Even if it fails in repeated attempts at finding a solution, the reasons for failure can be recorded and presented to the decision maker. In this way, even if the TMS does not find a solution, it provides useful information that can be used to find a solution. Since it is desirable to try several avenues simultaneously, an ATMS (de Kleer, 1986) might be suitable even though in principle any TMS could be used.

The problem of determining how breakable an assignment is is a much more difficult problem. It depends on the flexibility of the teacher, how suited substitutes are for a course, and how much time has already gone into preparation. We are currently in the process of trying to formalize these concepts so that they might be representable and used by a TMS.

7. Directions for Future Work

It is clear that the optimization problems solvers such as ZOOM must be made much more flexible if they are to prove useful as decision support tools for practical problems. Based on our experience with the experiments and analysis of expert behavior, we feel that there are at least two directions that are worthwhile pursuing. First, a more expressive interface is required which allows the problem the problem data to be specified as "naturally" as possible. Second, it should be possible to augment ZOOM so that it provides more useful output to the user.

On the interface front, we have found that specifying problem data in terms of the object hierarchy and defining constraints in terms of such objects is a very useful functionality for the user. In effect, we have built a layer on top of XML. Currently languages like XML and GAMS are used directly as specification languages. While they (especially GAMS) allow the user to specify the problem in terms of a compact notation, they are still relatively low-level utilities. We feel that there are significant productivity gains possible by providing a higher level modeling environment where a user can specify an arbitrary class hierarchy and constraints over it, and have the mathematical formulation generated automatically. We are currently working on the primitives that such an environment must have for it to be able to work for all mathematical programming problems.

We also feel that it should be able to build in a truth maintenance functionality into the optimization package, making it more flexible and useful to the user. In ZOOM, for example, considering that the branch and bound is usually unsuccessful (in our case it was *always* unsuccessful), it makes sense to try and generate an "almost feasible" solution in cases when the Pivot and Complement heuristic fails.

When the heuristic fails, the system knows which constraints are being violated (in our problem these were the knapsack constraints). At this point, if the number of non-integer variables is not too large (i.e. less than about 30), the system could perform a local search around these variables only (keeping values of all integer variables fixed), that is, attempt different combinations of integer values for the non-integer variables, and keeping a record of the contradictions. Such a history could be used by the system (or the user) to decide which constraints to loosen in order to generate a feasible solution. In effect, this boils down to augmenting the Pivot and Complement heuristic with a primitive truth maintenance system. Actually, this is quite similar to what the expert does in overconstrained situations (i.e. ignore TLs, obtain a solution, and gradually introduce the TLs). By enabling the system to loosen constraints in situations where no solution is in sight, it can begin to approximate the problem solving behavior of experts. As a next step in this research we are considering ways of incorporating a truth maintenance functionality into optimizers to handle such situations.

8. Acknowledgments.

We are extremely grateful to Ted Stohr, Charles Petrie and Elaine Rich for comments that have greatly improved this paper.

The expert system rules were formulated by Charles Petrie, and the system was implemented by Donald Steiner and Petrie. We thank Donald Steiner for providing us with the rules used for the belief revision example in this paper. Robert Causey provided extensive critiques of the solutions produced by the integer programming model. Finally, thanks to Roy Marsten for ZOOM.

References

- Balas, E., and Martin, C.H., Pivot and Complement -- A Heuristic for 0-1 Programming, *Management Science*, vol. 26, no. 1, January 1980, pp. 86-96.
- Dantzig, G., *Linear Programming and Extensions*, Princeton University Press, 1963.
- de Kleer, J., 1986. An Assumption-based TMS, *Artificial Intelligence*, vol 28, no.2, March 1986.
- Dhar, V., A Truth Maintenance System for Supporting Constraint-Based Reasoning, *Decision Support Systems*, Fall 1989.
- Doyle, J., A Truth Maintenance System, *Artificial Intelligence*, vol 12, no. 3, 1979.
- Garey, M., and Johnson, D., *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W.H. Freeman and Company, New York, 1979.
- Hinton, G.E, Relaxation and its Role in Vision, Ph.D Thesis, University of Edinburgh, 1977.
- Hooker, J.N., A Quantitative Approach to Logical Inference, *Decision Support Systems*, vol 4, no. 1, March 1988.
- Marsten, R., ZOOM/XMP User's Manual, Release 4.0, XMP Optimization Software Company, Tucson Arizona, July 1987.
- McAllester, D., Reasoning Utility Package, MIT-AI Lab Memo 667, April 1982.
- Petrie, C., Russinoff, D., and Steiner, D., Proteus 2: System Description, MCC Technical Report AI-136-87, May 1987.
- Petrie, C., Revised Dependency-Directed Reasoning for Default Reasoning, Proceedings of the Sixth National Conference on Artificial Intelligence (AAAI-87), July 1987, pp. 167-172.
- Petrie, C., Causey, R., Steiner, D., Dhar, V., A Planning Problem: Revisable Academic Course Scheduling, MCC Technical Report.
- Reinfrank, M., Lecture Notes on Reason Maintenance Systems, Technical Report INF2 ARM-5-88, Siemens AG, Munich, West Germany, 1988.