

INTERACTIVE TRANSFER OF EXPERTISE: ACQUISITION OF NEW INFERENCE RULES

Randall Davis
Computer Science Department
Stanford University
Stanford, California 94305

Abstract

TEIRESIAS is a program designed to function as an assistant in the task of building large, knowledge-based systems. It embodies a particular model of interactive transfer of knowledge from a human expert to the system, and makes possible knowledge transfer in a high level dialog conducted in a restricted subset of natural language. This paper explores an example of TEIRESIAS in operation, and demonstrates how it guides the acquisition of new inference rules. The concept of *meta-level Knowledge* is described, and an illustration given of its utility and contribution to the creation of intelligent programs.

This work was supported in part by the Bureau of Health Sciences Research and Evaluation of HEW under Grant HS-01544 and by the Advanced Research Projects Agency under ARPA Order 2494. It was carried out on the SUMEX Computer System, supported by the NIH under Grant RR-00785. The views expressed are solely (those of the author).

(1) Introduction

The knowledge base for a high performance, domain-specific program (e.g., DENDRAL [9], MACSYMA[10]) is traditionally assembled by hand, an ongoing task that typically involves numerous man-years of effort. A key element in the construction process is the transfer of expertise from a human expert to the program. Since the domain expert often knows nothing about programming, the interaction between the expert and the performance program usually requires the mediation of a human programmer.

We have chosen this transfer of expertise task as a case study, and have sought to create a program that could supply much the same sort of assistance as that provided by the programmer. That is, we have attempted to create an assistant that will help build intelligent programs. The result is a system called TEIRESIAS [3,5-7] a large INTERLISP program designed to offer assistance in the interactive transfer of knowledge from a human expert to the knowledge base of a high performance program.

Work on TEIRESIAS has two goals. We have attempted first to develop a set of tools and empirical methods for knowledge base construction and maintenance, and sought to abstract from them a methodology applicable to a range of systems. The second, more general goal has been the development of an assistant capable of offering increasingly more sophisticated aid. This involves confronting many of the traditional problems of AI, and has resulted in the exploration of a number of solutions reviewed below.

This paper describes a number of the key ideas in the development of TEIRESIAS and discusses their implementation in the context of a specific task (acquisition of new *inference* rules¹), for a specific system (a rule-based computer consultant system modelled after the MYCIN system [13,4].) While the discussion deals with one particular task, system and knowledge representation, it should become clear that the main ideas are applicable to a number of more general issues.

(2) Meta-level knowledge

A central theme that runs through this and related papers ([3,5-7]) is the concept of *meta-level knowledge*. This takes several different forms as its use is explored, but can be summed up generally by saying that a program can "know what it knows". That is, a program can not only use its knowledge directly, but *may* also be able to examine it, abstract it, reason about it, and direct its application.

To see in general terms how this might be accomplished, recall that one of the principal problems of AI is the question of representation and use of knowledge about the world, for which numerous techniques have been developed. One way to view what we have done is to imagine turning this in on itself, and using some of these same techniques to describe the program itself.

The resulting system contains both object level representations describing the external world, and meta-level representations which describe the internal world of representations. As the discussion of "rule models" in Section [7] will make clear, such a system has a *number* of interesting capabilities.

(3) Perspective on knowledge acquisition

One of the aims of creating TEIRESIAS was to provide a vehicle for developing a particular approach to knowledge acquisition. We describe that approach here; Section {9} contains some comments on its likely range of applicability.

We view the interaction between the domain expert and the performance program in terms of a teacher who continually challenges a student with new problems to solve, and carefully observes the student's performance. The teacher may interrupt to request a justification of some particular step the student has taken in solving the problem, or may challenge the final result. This may uncover a fault in the student's knowledge of the subject, and result in the transfer of information to correct it.

There is an important assumption involved in the attempt to establish this sort of communication: we are assuming that it is possible to distinguish between *basic formalism and degree of expertise*, or equivalently, that control structure and representation in the performance program can be considered separately from the content of its knowledge base. The basic control structure(s) and representations are assumed to be established and debugged, and the fundamental approach to the problem assumed acceptable. The question of *how* knowledge is to be encoded and used is settled by the selection of one or more of the available representations and control structures. The expert's task is to enlarge *what* it is the program knows.

There is a corollary assumption, too, in the belief that the control structures and knowledge representations can be made sufficiently comprehensible to the expert (at the conceptual level) that he can (a) understand the system's behavior in those terms and (b) use them to codify his own knowledge. This insures that the expert understands system performance well enough to know what to correct, *and* can then express the required knowledge, i.e., he can "think" in those terms. Thus part of the task of establishing the link between the expert and the system involves insulating the expert from the details of *implementation*, by establishing a discourse at a level high enough that we do not end up effectively having to teach him how to program.

(4) Design of the performance program

(4.1) Program architecture

Figure 1 shows a picture of the sort of performance program that TEIRESIAS is designed to help construct. (The program described here is modelled after the MYCIN system, which provided the context within which TEIRESIAS was developed. We have abstracted out here just the essential elements of MYCIN's design.) The *knowledge base* is the program's store of task specific knowledge that makes possible high performance. The *inference engine* is an interpreter that uses the knowledge base to solve the problem at hand.

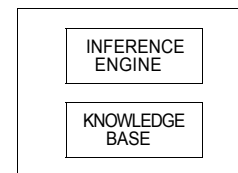


Figure 1 - architecture of the performance program

The main point of interest in this very simple design is the explicit division between these two parts of the program. This design is in keeping with the assumption noted above that the expert's task would be to augment the knowledge base of a program whose control structure (inference engine) was assumed both appropriate and debugged, if all of the control structure information has been kept in the inference engine, then we can engage the domain expert in a discussion of the knowledge base and be assured that the discussion will have to deal only with issues of domain specific expertise (rather than with questions of programming and control structures). The explicit division also offers a degree of domain independence. If all of the task specific knowledge has been kept in the knowledge base, then it should be

possible to remove the current Knowledge base, "plug in" another, and obtain a performance program for a new task.²

In this discussion we assume the knowledge base contains information about selecting an investment in the stock market; the performance program thus functions as an investment consultant. (MYCIN, of course, deals with infectious disease diagnosis and therapy selection, and the rules and dialog shown below dealt with that subject initially. The topic has been changed to keep the discussion phrased in terms familiar to a wide range of readers, and to emphasize that neither the problems attacked nor the solutions suggested are restricted to a particular domain of application. The dialog shown is a real example of TEIRESIAS in action that has been transferred to the new domain by substituting a few words in a medical example: e.g. *Ecoli* became *AT&T*, *infection* became *investment*, etc.)

An example of the program in action is shown in Section {6}. The program Interviews the user, requesting various pieces of information that are relevant to selecting the most appropriate investment, then prints its recommendations. In the remainder of this paper the "user" will be an expert running the program in order to challenge it, offering it a difficult case, and observing and correcting its performance.

{4.2} The knowledge base

* Knowledge in the knowledge base is in the form of a collection of decision rules of the sort shown in Figure 2. (The rule is stored internally in the INTERUSP form shown; the English version is generated from that.) Each rule is a single "chunk" of domain specific information indicating an *action* (in this case a conclusion) which is justified if the conditions specified in the *premise* are fulfilled,

RULE027

If [1] the time scale of the investment is long-term,
 [2] the desired return on the investment is greater than 10%, and
 [3] the area of the investment is not known, and then AT&T is a likely (.4) choice for the investment.

PREMISE (SAND (SAME OBJECT TIMESCALE LONG-TERM)
 GREATER OBJECT RETURNRATE 10)
 (NOTKNOWN OBJECT INVESTMENT-AREA))
 ACTION (CONCLUDE OBJECT STOCKNAME AT&T .4)

Figure 2 - example of a rule

The rules are judgmental, i.e., they make inexact inferences. In the case of the rule in Figure 2, for instance, the evidence cited in the premise is enough to assert the conclusion shown with only a weak degree of confidence (.4 out of 1.0). These numbers embody a model of confirmation described in detail in [14].

Finally, a few points of terminology. The premise is a Boolean combination of one or more *clauses*, each of which is constructed from a *predicate function* with an *associative triple* (*attribute, object, value*) as its argument. For the first clause in Figure 2, for example, the predicate function is *SAME*, and the triple is "*timescale of investment is long-term*." (The identifier *OBJECT* is used (rather than the name of a specific object) for reasons dealing with implementation conventions; see [3] for details.)

{4.3} The inference engine

The rules are invoked in a simple backward-chaining fashion that produces an exhaustive depth-first search of an and/or goal tree. Assume that the program is attempting to determine which stock would make a good investment. It retrieves all the rules which make a conclusion about that topic (i.e., they mention *STOCK-NAME* in their action), and invokes each one in turn, evaluating each premise to see if the conditions specified have been met. For the rule in Figure 2, this means first determining what the timescale of the investment ought to be. This is in turn set up as a subgoal, and the process recurs.

The search is thus depth-first (because each premise condition is thoroughly explored in turn); the tree that is sprouted is an and/or goal tree (because rules may have OR conditions in their premise); and the search is exhaustive (because the rules are inexact, so that even if one succeeds, it was deemed to be a wisely conservative strategy to continue to collect all evidence about the subgoal.)

{5} A word about natural language

Natural language has not been a major focus of this work, and for the most part the simplest techniques have been used that would support the level of performance required. For instance, all questions and responses from TEIRESIAS are either pre-formed, or manufactured by filling in templates with appropriate words (as

evidenced by the occasional appearance of phrases like "a area"). Where answers are free text (rather than multiple choice), the system relies primarily on a simple Keyword-oriented approach, described below. There is no parser in the system, in part to minimize processing time, and in part because users often give ungrammatical responses.

{6} Example of TEIRESIAS in operation

In the interaction shown below, the expert is challenging the investment advisory program with a particular case, observing the program's performance, and then making the necessary correction to the knowledge base. The bug encountered here was created by deleting an existing rule from the knowledge base and removing all references to it. The expert then repairs the bug by tracking down the source of the problem and teaching the system the "new" rule.

The deleted rule was chosen to make possible an acquisition example that displays a large number of TEIRESIAS's capabilities in a single example. Dialogues are thus typically not as long as this one.

In the trace below, computer output is in mixed upper and lower case in this font, user responses follow a double asterisk and are in BOLD CAPITALS, and commentary appears (inside braces). The example is also divided into five phases for later reference.

Investment Advisory Program

```

1) Investor's name?
**FRED SMITH
2) Age?
**34
3) Present income?
**25,000
4) Number of dependents?
**3
5) Occupation?
**ELECTRICIAN
6) Amount of previous investment experience?
{slight = less than one year; moderate =
 2 to 4 years; extensive = more than 4}
**MODERATE
7) Does the investor subscribe to the Wall Street
Journal, or any securities market newsletters?
**Y
:
```

{Most of the consultation with the performance program has been edited out for brevity. The interaction continues with the program printing out its conclusions.}

The following Investments appear to be the most appropriate at this time:

Varian-Corporation
 American-Telephone-and-Telegraph

{At this point TEIRESIAS steps in, and since the interaction is with an expert (rather than an ordinary user), it requests an evaluation of the performance program's results.

From here on the dialog takes place between TEIRESIAS and the expert, with the performance program in the background being manipulated by TEIRESIAS.}

 {PHASE 1: tracking down the bug}

```

Are all these correct?
**** N
Are there any missing?
**** Y
Which one(s)?
**** DIGITAL-EQUIPMENT DATA-GENERAL
Are there any that shouldn't have appeared?
**** Y
Which one(s)?
**** AT&T
Do you want to
1) track down the problem yourself, or
2) allow me to guide the debugging, or
3) ignore it for now and fix it another time?
Type 1, 2 or 3
**** 2
```

{The first option allows the expert to use the explanation facilities described in [3] to track down the problem. In the second, TEIRESIAS uses a simple algorithmic search for the error. The third is for human engineering purposes, and allows the expert to sidestep this whole process if desired.}

We'll start with the AT&T first.

{Heuristics derived from our experience in debugging rule-based systems are built into TEIRESIAS. One such heuristic suggests that the best approach is to (1) track down one problem at a time, and (2) start with a result that shouldn't have appeared (if there is one). Experience indicates that this class of bug is often a manifestation of a deeper bug causing several problems, so curing it will often solve several problems at once.}

The following rules were used in deducing that the Investment should be AT&T:
RULE027

Since [1.1] the area of the investment is not known,
[1.2] the desired return on the investment is greater than 10%, and
[1.3] the time scale of the investment is long-term,
then AT&T is a likely (.4) choice for the investment.

Is this rule correct?

◆◆*** Y

Should its premise have failed for this case?

◆◆*** Y

Which clause of the premise should have been false?

◆◆*** 1

{The expert indicates that this rule should have failed, so the debugging continues by looking further back in the chain. Had he indicated that RULE027 had been invoked properly, then TEIRESIAS would have indicated that the problem was probably a missing rule which tended to rule out AT&T.}

No conclusions were reached concerning investment area
Should it have been possible to conclude an area?

◆◆*** Y

{The error might have been a failure to rule out an investment area.}

Which one?

◆◆*** HIGH TECHNOLOGY

There were no rules successfully invoked in the attempt to deduce that the area of the investment should be high-technology.

The following (potentially applicable) rules were NOT executed because the starred [**] clause failed:

{At this point TEIRESIAS displays rules 50 and 166 (omitted here), which would have made the desired conclusion, but were not invoked because one of their premise clauses failed.}

Are all these rules correct?

◆◆*** Y

Should any of them have been successfully invoked?

◆◆*** N

{This would have been the case if one of the premise clauses that failed was false because it was "stepped on" by yet another rule further back in the chain.}

Is there a rule missing?

◆◆*** Y

Ok, then I need a rule that would allow me to deduce that the area of the Investment should be high-technology.

Would you like to give me such a rule now?

◆◆*** Y

Good...

{To review briefly: the expert is tracking down the reason why the system incorrectly concluded that the stock to invest in was AT&T. Shown below is the sequence of rules involved:

```
= RULE116 ==> I
                    J-> AREA
== RULE050 **> I RETURN    === RULE027 ===> AT&T
                    TIME SCALE
```

AT&T appeared because RULE027 executed successfully; this occurred in part because investment area was unknown. Investment area was unknown because all rules to conclude it (50 and 116) failed.

The expert believes all existing rules are correct, so the knowledge base must be missing a rule that will conclude about investment area, in particular, that the area should be high-technology. Such a rule will cause area to be known, inhibiting rule 27, fixing at least this first problem.

Note how sharply focussed the issue has become, for both the expert and TEIRESIAS, as a result of the large amount of contextual information available. The expert is not being asked to "discuss the stock market"; instead the system can at this point say "I need a rule that would allow me to deduce that the area of the investment should be high-technology," and it must be a rule which is invocable in the current situation.

The contextual information is useful for TEIRESIAS as well. The information provides a number of expectations about the content of the rule which are used both in understanding what the expert types in, and in checking the new rule to make sure that it does indeed solve the problem at hand. Details of this process are explained below.

The trace continues with the expert typing in the rule.}

{PHASE II: interpreting the rule}

The new rule will be called RULE383

```
If: 1 - THE CLIENT'S INCOME TAX BRACKET IS 507.
    and 2 - THE CLIENT IS FOLLOWING UP ON MARKET TRENDS CAREFULLY
    and 3 -
Then: 1 - THERE IS EVIDENCE (.8) THAT THE INVESTMENT AREA SHOULD BE HIGH TECHNOLOGY
      and 2 -
```

This may take a few moments...

{Generally on the order of 10 to 30 seconds of cpu time.}

This is my understanding of your rule:

RULE383

```
If 1) The client's income-tax bracket is 50%,
    2) The market has followed an upward trend recently,
    3) The client manages his assets carefully
```

Then there is evidence (.8) that the area of the Investment should be high-technology

{TEIRESIAS's initial interpretation is not entirely correct (and the example was chosen in part for this reason). The first line of text and the conclusion have been interpreted correctly, but the second line has been turned into two clauses (numbers 2 and 3) in the system's initial attempt at understanding. (The expert occasionally types what should be two separate clauses as one line of text. The system is capable of handling this, but this same capability at times produces spurious clauses like *3).}

This results primarily from the system's simple word-by-word approach to meaning. While the result is in part incorrect, we'll see that the interpretations are both "reasonable" and "consistent" in light of the techniques used.}

{PHASE III: modifying the interpretation}

Okay? (Yes or No)

◆◆*** N

{The expert indicates there is a problem with the interpretation, and uses TEIRESIAS's rule editor to specify the changes needed.}

Change premise?

◆◆*** Y

What clause should be changed in premise?

◆◆*** CHANGE 2

Suppose (2) were changed to read-

2 - The client follows the market carefully

Would that be ok?

◆◆*** Y

{Several simple heuristics are invoked to select the next most likely interpretation, based on the English source of this clause and the set of interpretations already rejected. In this case the second choice is correct.}

Then you'll probably want to delete:

3) The client manages his assets carefully

Correct?

◆◆*** Y

{TEIRESIAS enforces certain consistency constraints on its interpretation of the English text. One of these constraints indicates that if the new interpretation of clause #2 is correct, then clause *3 must have arisen spuriously (since, among other things, it requires a different interpretation for the word "carefully*");}

Anything else?

+◆*** ?

```
1 - the client's income-tax bracket is 50%
and 2 - the client follows the market carefully
and 3 - <deleted>
```

Anything else?

+◆◆*** N

Change action?

+◆*** N

This is my understanding of your rule:

RULE383

If: 1) the client's Income-tax bracket is 50%, and
2) the client follows the market carefully,
Then: there is evidence (.8) that the area of the
Investment should be high-technology.

Okay? (Yes or No)

+♦** Y

{This is the standard test of comprehension: the original English text has been turned into a rule coded in LISP, and is then retranslated back into English for approval. Compare the results with the expert's original input. Now that the expert is satisfied that the system understands the rule as he entered it, the system checks several things about the rule to see if it is satisfied.}

{PHASE IV: "second guessing"}

I hate to criticize, Randy, but did you know that most rules about what the area of an investment might be, that mention-

the income-tax bracket of the client, and how closely the client follows the market
ALSO mention-

A) the amount of investment experience of the client
Shall I try to write a clause to account for [A] ?

+♦** Y

{The mechanism that makes possible this sort of "second guessing" by the system is explained in detail below. Briefly, the system maintains a model of its own knowledge. In checking to see how this new piece of information "fits in" to that model, it discovers that its expectations have been only partially met. It indicates to the expert what expectations were left unfulfilled, and asks if it should attempt to fill in the missing parts.}

How about-

A) The amount of investment experience of the client
is moderate

Ok?

♦♦** Y

{The expert agreed to the inclusion of the new clause and TEIRESIAS is able to construct a plausible guess by relying on the fact that it must be possible to invoke the new rule in the current situation (in order to fix the known bug). Details of this process are explained below.}

Also, most rules about what the area of an investment might be ALSO conclude about-

A) the desired risk classification of the investment
Shall I try to write a clause to account for [A] ?

++++ N

{The system noticed another unfulfilled expectation concerning the action part of the rule, but the expert always has the option of overriding the system's suggestions, and does so.}

{PHASE V: final checkout, bookkeeping, rerunning the consultation}

The rule now reads:

RULE383

If: 1) the income-tax bracket of the client is 50%,
2) the client follows the market carefully,
3) the amount of investment experience of the
client is moderate

Then: there is evidence (.8) that the area of the
Investment should be high-technology.

Please describe briefly why you added (or changed) this rule. Type an empty line when done.

♦♦**THE SYSTEM FAILED TO REALIZE IT COULD CONCLUDE
♦♦*INVESTMENT AREA; THIS ALLOWED RULE027 TO BE INVOKED
♦♦**INAPPROPRIATELY

+ 4**

{Experience suggests that in building any large knowledge base it is very useful to keep fairly detailed records about the source of various contributions. The new rule is tagged with a range of such information.}

RULE383 has now been added to the knowledge base.

The consultation will be re-run to check the effectiveness of your new rule. Get comfortable, since this may take a bit.

Also, new questions may be asked as a result of this new rule. If a question is printed, please answer it with reference to this current client.

{The consultation is then rerun, but all the previous answers are saved, so the expert can in most cases just watch it proceed. This has been edited out here for brevity.}

The following investments appear to be the most appropriate at this time:

Var1an-Corporation
Data-General
Digital-Equipment

Well, it appears that things are straightened out.
{The new conclusions are the ones the expert indicated earlier as the correct answers. This does not go unnoticed.}

Are all these values correct now?

♦ Y

{As always, however, the expert gets the last word, and TEIRESIAS waits for approval before proceeding.}

Good, then the problem's been fixed.

{Three errors have disappeared with the addition of one rule. This is explained schematically below:

RULE116 *

RULE050 *

RULE383 ==> AREA == RULE027 * ==> AT&T

== RULE084 ==> DATA-GENERAL

DIGITAL-EQUIPMENT

As before, 116 and 50 failed (hence the "s"), but now the new rule, 383, has succeeded, supplying a value for investment area. This in turn means that 27 now fails, as hoped.

However, rule 84, which requires a value of area, failed last time because the system was unable to deduce one. Now that 383 supplies one, 84 can fire, and it turns out to yield the missing answers.}

(7.1 How it all works

7.1.1 Overview of the main ideas

Before reviewing the trace in more detail, we describe the ideas which make possible the capabilities displayed. The list below serves primarily to name and briefly sketch each in turn; the details are supplied in reviewing the example.

I. Knowledge acquisition in context

Performance programs of the sort TEIRESIAS helps create will typically find their greatest utility in domains where there are no unifying laws on which to base algorithmic methods. In such domains there is instead a collection of informal knowledge based on accumulated experience. This means an expert specifying a new rule may be codifying a piece of knowledge that has never previously been isolated and expressed as such. Since this is difficult, anything which can be done to ease the task will prove very useful.

In response, we have emphasized knowledge acquisition in the context of a shortcoming in the knowledge base. To illustrate its utility, consider the difference between asking the expert

What should I know about the stock market?

and saying to him

Here is an example in which you claim the performance program made a mistake. Here is all the knowledge the program used, here are all the facts of the case, and here is how it reached its conclusions. Now, what is it that you know and the system doesn't that allows you to avoid making that same mistake!

Note how much more focussed the second question is, and how much easier it is to answer.

II. Building expectations

The focussing provided by the context is also an important aid to TEIRESIAS. In particular, it permits the system to build up a set of expectations concerning the knowledge to be acquired, facilitating knowledge transfer and making possible several useful features illustrated in the trace and described below.

III. Model-based understanding

Model-based understanding suggests that some aspects of understanding can be viewed as a process of matching: the entity to be understood is matched against a collection of prototypes, or models, and the most appropriate model selected. This sets the framework in which further interpretation takes place, as that model can then be used as a guide to further processing.

While this view is not new, TEIRESIAS employs a novel application of it, since the system has a model of the knowledge it is likely to be acquiring from the expert.

IV. Giving programs a model of their own knowledge

We will see that the combination of TEIRESIAS and the performance program amounts to a system which has a picture of

its own Knowledge. That is, it not only knows something about a particular domain, but in a primitive sense it knows what it knows, and employs that model of its knowledge in several ways.

V. Learning by experience

One of the long-recognized potential weaknesses of any model-based system is dependence on a fixed set of models, since the scope of the program's "understanding" of the world is constrained by the number and type of models it has. As will become clear, the models TEIRESIAS employs are not hand-crafted and static, but are instead formed and continually revised as a by-product of its experience in interacting with the expert.

{7.2} Phase I: tracking down the bug

To provide the debugging facility shown, TEIRESIAS maintains a detailed record of the actions of the performance program during the consultation, and then interprets this record on the basis of an exhaustive analysis of the performance program's control structure (see [3] for details). This succeeds because (a) the backward-chaining technique used by the performance program is sufficiently straightforward and intuitive, even to a non-programmer) and (b) the rules are designed to encode knowledge at a reasonably high conceptual level. As a result, even though TEIRESIAS is running through an exhaustive case-by-case analysis of the preceding consultation, the expert is given the task of debugging *reasoning* rather than *code*.

The availability of an algorithmic debugging process is also an important factor in encouraging the expert to be as precise as possible in his responses. Note that at each point in tracking down the error the expert must either approve of the rules invoked and conclusions made, or indicate which one was in error and supply the correction. This is extremely useful in domains where knowledge has not yet been formalized, and the traditional reductionist approach of dissecting reasoning down to observational primitives is not yet well established.

Finally, consider the extensive amount of contextual information that is now available. The expert has been presented with a detailed example of the performance program in action, he has available all of the facts of the case, and has seen how the relevant knowledge has been applied. This makes it much easier for him to specify the particular chunk of knowledge which may be missing. This contextual information will prove *very* useful for TEIRESIAS as well. It is clear, for instance, what the *effect* of invocation of the new rule must be (as TEIRESIAS indicates, it must be a rule that will "deduce that the area of the investment should be high-technology"), and it is also clear what the *circumstances* of its invocation must be (the rule must be invocable for the case under consideration, or it won't repair the bug). Both of these will be seen to be quite useful.

{7.3} Phase II: interpreting the rule

As is traditional, "understanding" the expert's natural language version of the rule is viewed in terms of converting it to an internal representation, and then retranslating that into English for the expert's approval. In this case the internal representation is the INTERLISP form of the rule, so the process is also a simple type of code generation.

{7.3.1} Models and model-based understanding

As a background for reviewing the interpretation process, we digress for a moment to consider the idea of models and model-based understanding, then explore their application in TEIRESIAS.

In the most general terms, a model can be seen as a *compact, high-level description of structure, organization, or content* that may be used both to provide a framework for lower-level processing, and to express expectations about the world. One particularly graphic example of this idea can be found in the work on computer vision by Falk [8] in 1970. The task there was the standard one of understanding blocks-world scenes: the goal was to determine the identity, location, and orientation of each block in a scene containing one or more blocks selected from a known set of possibilities.

The key element of his work of interest here is the use of a set of *prototypes* for the blocks, prototypes that resembled wire frame models. While it oversimplifies slightly, part of the operation of his system can be described in terms of two phases. The system first performed a preliminary pass to detect possible edge points in the scene, and attempted to fit a block model to each collection of edges. The model chosen was then used in the second phase as a guide to further processing. If, for instance, the model accounted for all but one of the lines in a region, this suggested that the extra line might be spurious. If the model fit well except for some line missing from the scene, that was a good hint that a line had been overlooked, and indicated as well where to go looking for it.

While it was not a part of Falk's system, we can imagine one

further refinement in the interpretation process and explain it in these same terms. Imagine that the system had available some *a priori* hints about what blocks might be found in the next scene. One way to express those hints would be to bias the matching process. That is, in the attempt to match a model against the data, the system might (depending on the strength of the hint), try the indicated models first, make a greater attempt to effect a match with one of them, or even restrict the set of possibilities to just those contained in the hint.

Note that in this system, (i) the models supply a compact, high-level description of structure (the structure of each block), (ii) the description is used to guide lower level processing (processing of the array of digitized intensity values), (iii) expectations can be expressed by a biasing or restriction on the set of models used, and (iv) "understanding" is viewed in terms of a matching and selection process (matching models against the data and selecting one that fits).

{7.3.21} Rule models

Now recall our original task of interpreting the expert's natural language version of the rule, and view it in the terms described above. As in the vision example, there is a signal to be processed (the text), it is noisy (words can be ambiguous), and there is context available (from the debugging process) that can supply some hints about the likely content of the signal. To complete the analogy, we need a model, one that could (a) capture the structure, organization, or content of the expert's reasoning, (b) be used to guide the interpretation process, and (c) be used to express expectations about the likely content of the new rule.

Where might we get such a thing? Not surprisingly, there are regularities in the knowledge base: rules about a single topic tend to have characteristics in common. From these regularities we have constructed *rule models*. These are abstract descriptions of subsets of rules, built from empirical generalizations about those rules, and are used to characterize a "typical" member of the subset.

Rule models are composed of four parts. They contain, first, a list of EXAMPLES, the subset of rules from which this model was constructed.

Next, a DESCRIPTION characterizes a typical member of the subset. Since we are dealing in this case with rules composed of premise-action pairs, the DESCRIPTION currently implemented contains individual characterizations of a typical premise and a typical action. Then, since the current representation scheme used in those rules is based on associative triples, we have implemented those characterizations by indicating (a) which attributes typically appear in the premise (action) of a rule in this subset, and (b) correlations of attributes appearing in the premise (action).

Note that the central idea is the concept of *characterizing a typical member of the subset*. Naturally, that characterization would look different for subsets of rules, procedures, theorems, etc. But the main idea of characterization is widely applicable and not restricted to any particular representational formalism.

The two other parts of the rule model are pointers to models describing more general and more specific subsets of rules. The set of models is organized into a number of tree structures. At the root of each tree is the model made from all the rules which conclude about <attribute> (e.g., the INVESTMENT-AREA model), below this are two models dealing with all affirmative and all negative rules (e.g., the INVESTMENT-AREA-IS model), and below this are models dealing with rules which affirm or deny specific values of the attribute.

Rather than being hand-tooled, the models are assembled by TEIRESIAS on the basis of the current contents of the knowledge base, in what amounts to a very simple (i.e., statistical) form of concept formation. The combination of TEIRESIAS and the performance program thus presents a system which has a model of its own knowledge, one it forms itself.

The rule models are the primary example of meta-level knowledge in this paper. This form of knowledge, and its generation by the system itself have several implications illustrated in later sections.

Figure 3 shows a rule model; this is the one used by TEIRESIAS in the interaction shown earlier. (Since not all of the details of implementation are relevant here, this discussion will omit some. See [3] for a full explanation.) As indicated above, there is a list of the rules from which this model was constructed, descriptions characterizing the premise and the action, and pointers to more specific and more general models. Each characterization in the description contains the two kinds of entries noted: one concerning the presence of individual attributes and the other describing correlations. The first item in the premise description, for instance, indicates that "most" rules about what the area of an investment should be mention the attribute *rate of return* in their premise; when they do mention it they "typically" use the predicate functions SAME and NOTSAME and the "strength", or reliability, of

guessing ability. Alternatively, more sophisticated concept formation techniques might be borrowed from existing work.

There is a potential problem in the way the models are used. Their effectiveness in both guiding the parsing of the new rule and in "second guessing" its content is dependent on the assumptions that the present knowledge base is both correct, and a good basis for predicting the content of future rules. Either of these can at times be false, and the system may then tend to continue stubbornly down the wrong path.

In addition, the weakness of the natural language understanding technique presents a substantial barrier to better performance. Once again there are several improvements that could be made to the existing approach (see [3]), but more sophisticated techniques should also be considered.

There is also the difficult problem of determining the impact of any new or changed rule on the rest of the knowledge base. The difficulty lies in establishing a formal definition of inconsistency for inexact logics, since, except for obvious cases (e.g., two identical rules with different strengths), it is not clear what constitutes an inconsistency. Additional work is needed here.

[10] Conclusions

The ideas reviewed above each offer some contribution toward achieving the two goals set out at the beginning of this paper: the development of a methodology of knowledge base construction via transfer of expertise, and the creation of an intelligent assistant.

Knowledge acquisition in the context of a shortcoming in the Knowledge base, for instance, has proved to be a useful technique for achieving transfer of expertise, offering advantages to both the expert and TEIRESIAS. It offered the expert a framework for the explication of a new chunk of domain knowledge. By providing him with a specific example of the performance program's operation, and forcing him to be specific in his criticism, it encourages the formalization of previously implicit knowledge. It also enabled TEIRESIAS to form a number of expectations about the knowledge it was going to acquire, and made possible several checks on the content of that knowledge to insure that it would in fact fix the bug.

In addition, because the system has 8 *model of its own knowledge*, it was able to determine whether a newly added piece of knowledge "fit into" its existing knowledge base.

A second contribution of the ideas reviewed above lies in their ability to support a number of intelligent actions on the part of the assistant. While those actions have been demonstrated for a single task and system, none of the underlying ideas appears to be limited to this task, or to associative triples or rules as a knowledge representation.

The idea of *model-based understanding*, for instance, found a novel application in the fact that TEIRESIAS has a model of both its own knowledge and the information it expects to receive, and uses this to guide acquisition.

The idea of *biasing the set of models* to be considered offered a specific mechanism for the general notion of *program-generated expectations*, and made possible an assistant whose understanding of the dialog was more effective.

TEIRESIAS was able to "second guess" the expert with respect to the content of the new knowledge by using its models to see how well the new piece of knowledge "fit into what it already knew". It was the presence of a partial match between the new knowledge and the system's model of its knowledge that prompted it to make a suggestion to the expert.

The concept of meta-level knowledge made possible *multiple uses of the knowledge in the system*: information in the knowledge base was not only used directly (during the consultation), but was also examined and abstracted to form the rule models.

TEIRESIAS also represents a synthesis of the ideas of model-based understanding and learning by experience. While both of these have been developed independently in previous AI research, their combination produced a novel sort of feedback loop (Figure 5). Rule acquisition relies on the set of rule models to effect the model-based understanding process. This results in the addition of a new rule to the knowledge base, and this in turn prompts the recomputation of the relevant rule model(s).

This loop has a number of interesting implications. First, performance on the acquisition of the next rule may be better, because the system's "picture" of its knowledge base has improved — the rule models are now computed from a larger set of instances, and their generalizations are more likely to be valid.

Second, since the relevant rule models are recomputed each time a change is made to the knowledge base, the picture they supply is kept constantly up to date, and they will at all times be an accurate reflection of the shifting patterns in the knowledge base.

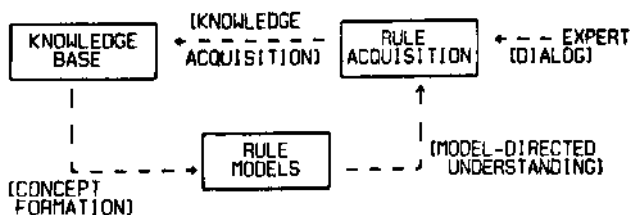


Figure 5

Finally, and perhaps most interesting, the models are not hand-tooled by the system architect, or specified by the expert. They are instead formed by the system itself, and formed as a result of its experience in acquiring rules from the expert. Thus despite its reliance on a set of models as a basis for understanding, TEIRESIAS's abilities are not restricted by the existing set of models. As its store of knowledge grows, old models can become more accurate, new models will be formed, and the system's stock of knowledge about its knowledge will continue to expand. This appears to be a novel capability for a model-based system.

Acknowledgments

The work described here was performed as part of a doctoral thesis supervised by Bruce Buchanan, whose assistance and encouragement were important contributions at every stage. Nancy Martin, Dave Barstow, and Mike Clancy made a number of very helpful comments on earlier drafts of this paper.

Notes

- [1] Acquisition of new conceptual primitives from which rules are built is discussed in [7], while the design and implementation of an explanation capability is discussed in [3] and [11].
- [2] Two experiments of this sort have been performed with the MYCIN system, and suggest that this sort of "plug compatibility" of knowledge bases is a realistic possibility for a range of tasks.
- [3] This description is necessarily abbreviated. See [3] for a more detailed discussion.
- [4] Both (a) and (b) are constructed via simple statistical thresholding operations.
- [5] Note that these tests require the ability to dissect and partially evaluate the rule. The same function template which is used as a pattern for constructing rules is also used as a guide in this dissection and partial evaluation process. See [3] for details.
- [6] The issue is related to work described in [2], on closed vs. open sets.

References

- [1] Baumgart B G. Geometric models for computer vision, Stanford University AI Memo 249, October 1974
- [2] Carbonell J R, Collins A M. Natural semantics in artificial intelligence, *Proc 3rd IJCAI*, pp 344-351, August 1973.
- [3] Davis R. Applications of meta-level knowledge to the construction, maintenance, and use of large knowledge bases, Stanford University HPP Memo 76-7, July 1976.
- [4] Davis R, Buchanan B, Shortliffe E H. Production rules as a representation for a knowledge-based consultation system, *Artificial Intelligence*, 8 (Spring 1977) 15-45.
- [5] Davis R. Generalized procedure calling and content directed invocation, to appear in *Proc ACM Conf on AI and Programming Languages*, August 1977.
- [6] Davis R, Buchanan B G. Meta-level knowledge: overview and applications, to appear in *Proc 4th IJCAI*, Cambridge, Mass.
- [7] Davis R. Knowledge about representations as a basis for system construction and maintenance, to appear in *Pattern-Directed Inference Systems*, Academic Press, (in press)
- [8] Falk G. Computer interpretation of imperfect line data, Stanford University AI Memo 132, August 1970.
- [9] Feigenbaum E A, et al. On generality and problem solving, in *MI 6*, pp 165-190, Edinburgh University Press, 1971.
- [10] Mathlab Group, The MACSYMA Reference Manual, Sept 1974.
- [11] Scott A C, Clancey W, Davis R, Shortliffe E K. Explanation capabilities of production-based consultation systems, *Am Jnl Computational Linguistics*, Microfiche 62, 1977.
- [12] Scott A C, Aikins J, Buchanan B G, Clancey W, Davis R, van Melie W. Issues of dealing with a growing knowledge base, in preparation.
- [13] Shortliffe E H. *MYCIN: Computer-based Consultations in Medical Consultations*, American Elsevier, 1976.
- [14] Shortliffe E H, Buchanan B G. A model of inexact reasoning in medicine, *Mathematical Biosciences*, 23 (1975) pp 351-379.
- [15] Waterman D. Exemplary programming, in *Pattern-directed Inference Systems*, Waterman and Hayes-Roth (eds), Academic Press, (in press).