

THE APPLICATION OF ARTIFICIAL INTELLIGENCE TO OATA BASE MANAGEMENT

Rob Gerritsen
The Wharton School (*)
University of Pennsylvania
Philadelphia, Pennsylvania

Abstract

An application of Artificial Intelligence is discussed: an automatic programming system that generates information retrieval programs and data base structure designs for highly structured or network data bases. It is claimed that these applications are unusual in that they have more practical value than is usually the case, being of immediate utility to commercial data base management. The paper concludes with a short reflection on the problems associated with the representation and acquisition of knowledge for problem solving programs. In particular, representation of knowledge is important for efficiency of the programs, yet an efficient representation may not be "natural". It was difficult to separate representation from content/ and successful application required that the expert providing the knowledge also know how the program uses that knowledge.

CR categories: 3.5, 3.6U, 3.7, U.2, 5.23

Key words and phrases: application of Artificial Intelligence, automatic programming, procedure generation, information retrieval, network data base. Data Base Task Group, data base design, non-procedural languages.

Introduction

Most successful Artificial Intelligence programs exercise their problem solving ability in non-realistic environments or apply it to a game playing task, but this state of affairs does not constitute a basis for criticism of AI. As many of the proponents of AI have maintained, it is necessary for AI to cut its teeth on toy problems, games, and in otherwise limited environments before it can proceed to fulfill its highly touted promises. Perhaps the only fair criticism is that this teeth cutting is taking longer than envisioned by early workers in the field.

AI is entering an age where practical and objective application is feasible. This is not to say that there have not been very valuable spin-offs from AI research that have contributed to

* A major portion of the work described herein was done while the author was at Carnegie-Mellon University, Pittsburgh, Pennsylvania.

advances in fields such as computer science, linguistics/ and psychology. However, there have not been many direct applications of AI technologies (with perhaps the exception of pattern recognition) to problem areas outside the research laboratory.

My background is not primarily in AI, but lies more in the area of commercial data processing, data base management, and information retrieval. My experience in these areas predates by several years my introductory exposure to AI. This fact is relevant because the problems, which will be described shortly, existed for me long before I discovered the technology to solve them. As such, my circumstances provided an amount of practical experience not typically possessed by the AI researcher who is looking for problems to adapt to his technology.

The application described here is a practical one. The problem solving systems that have been developed can supplant or augment some of the tasks currently being performed by programmers and data base administrators.

Data Management Application

The subject of this paper is the successful application of AI, particularly automatic programming technology, to the problems of information retrieval and data base design. This has resulted in a system<1,2,3> that generates information retrieval programs and data structure declarations.

These two problem solving tasks have been around for a long time, but growing commercial usage of larger and more complex network data bases has suddenly increased their significance. CAs used here, the term "network data base" refers to a data base in which a data item may be linked to many other data items. Our subject is not data bases on computer networks.) There has been an attempt within the United States and Canada to develop a standard for network data bases through CODASYL which spawned the Data Base Task Group (DBTG) involving both industry and academia. This group has produced a specification for a general purpose network data base system<U>. Although not yet accepted as a standard. It is believed by many that the DBTG specification will be so adopted, perhaps

by default since many computer manufacturers and software firms are now selling implementations of the DBTG specification.

Both applications described herein function in the data base environment as specified by the DBTG. That is, data structures designed by the system conform to DBTG specifications, and the information retrieval programs generated by it are COBOL Procedure Divisions containing Data Manipulation Language (DML), as defined by the DBTG.

Programming for Network Data Bases

Introduction of the network type of data base brings with it a new level of complexity for the business programmer. Earley^{<5>} points out that instead of the usual logical-physical dichotomy that a programmer traditionally deals with, there is a logical(a)-logical(b)-physical trichotomy facing the network data base programmer. The logical(a) level corresponds to data relationships as perceived by the data user. The logical(b) level concerns the actual links between the data in the data base.

Bachman^{<6>} recognizes this complexity in the 1973 ACM Turing Award Lecture where he likens the programmer's problem (in the network data base environment) to that of a navigator. The programmer is given the user's query requirements; his task is to find the proper access path to the data through the myriad of connective links in the data base.

The programmer's task should be reduced. After all, much of its complexity is caused by an artifact introduced by the data network, the data links. It should be noted here that others in the field, primarily Codd^{<7,8,9>}, have proposed relational models of data bases whereby they hope to avoid this artifact. It is not clear, however, how this can be done without exactly the kind of problem solving program that is the subject of this paper.

The technology used to reduce this task comes from work in the area of automatic programming. A compiler developed by Buchanan and Luckham^(10,11) was used to create the Information Retrieval Program Generator.

Programming Rules

The Information Retrieval Program Generator is compiled from a set of rules. These rules are stated in a formalism that closely corresponds to the formalism developed by Hoare^{<12>} to describe the logic of programs. Each

such rule describes some type of program construct, be it a loop, assignment statement, a DML statement, or other program construct*

One type of rule describes the conditions under which a program statement may be used and the effect of its use. For example, the rule for a "GET" statement says that it is not possible to GET a record unless its location in the data base is known to the data management system, and that the object record will be located in working storage following execution of the GET.

This rule is formally stated (to the compiler) as:

```
(CURRENT R1) | GET R1 | (INCORE R1)
In this rule R1 is a variable (for record name) and (CURRENT R1) and (INCORE R1) are predicates. Henceforth we use the convention that R1/ R2/ X1, I1, etc., are variable names.
```

Other rules do not bear such close resemblance to a description one might find in a programming manual as does the rule for GET. These other rules capture my own programming knowledge, including heuristics, that I had previously acquired and developed as a professional programmer.

One such rule is the rule that describes the top level composition of a program:

```
(OPENED A1)A(CURRENT R1)A(LINKED R1 R2)A
(CLOSED A1)A(STOP X1)
--> (PROGRAM X1)
```

This rule is different than the rule for GET: it does not define a program statement, rather it is an implication stated with predicates. The order of the predicates is obviously important.

A semantic interpretation of this rule is that a program consists of five blocks that: (1) open the required files of the data base, (2) locate the first record of interest, (3) follow the access path from this record to all other records of interest, (4) close the files opened above, and (5) terminate execution.

Of these five blocks, the second and especially the third, are the most complex. They may involve many further levels of refinement, revealing iterations, alternations and compositions. The third block will include the processing of the records on the access path to display values, calculate statistics, etc.

Translation to Micro-Planner Theorems

The rules are translated by the Buchanan-Luckham compiler to Micro-Planner<13,14> theorems that constitute the Information Retrieval Program Generator.

For example, the GET rule is compiled to:

```
(DEFPROP GET
  (THCONSE (R1) (INCORE (THV R1)
    (THGOAL (CURRENT (THV R1)))
    (THSET (CAR (THV ANS>)
      (CONS (CONS (QUOTE GET)
        (LIST (THV R1)))
      (EVAL (CAR (THV ANS))))))
    (THASSERT (INCORE (THV R1))))
  THEOREM)
```

The theorem presented here is quite condensed from the one actually generated, which also contains bookkeeping, tracing, uniqueness, and other special functions. Note that the pre-condition of the rule, (CURRENT R1), is a goal in the theorem. Similarly, the post-condition, (INCORE R1), is defined as the consequence of the theorem and will be asserted when the theorem is used. When the theorem is used 'GET R1' will be catenated onto the ANSWER.

Program construction proceeds in a sub-goaling fashion. For example, if (INCORE PATIENT) is a goal then the GET theorem is tried, and (CURRENT PATIENT) becomes a sub-goal. A sub-goal may invoke other theorems in turn, generating further sub-goals. This process generates a goal tree with leaf nodes which are true or false in a set of assertions known as the "state".

The initial state is the only input to the program generation process and is made up of two major subsets of assertions. One such subset describes the structure of the data base, the other describes the query for which program generation is required.

The description of one more rule will give us a sufficient number of rules for an example. A rule for FINDING records:

```
(HASHKEY R1 K1)A(DESIRE K1 V1)A
(CONTAINS KI V1) IIFIND R1|| (CURRENT R1)
```

This rule is similar to the rule for GET; it defines a program statement. Again, it corresponds closely to a paragraph in a programming manual: To FIND a record, first initialize the key for that record with the value identifying the desired record. Following the execution of FIND, the record is current to the data management system. Current is a DBTG term meaning that a record's exact location in the

data base is presently known to the data management system,

Example

Assume that a data base definition specifies that STUDENT records are (hash)keyed on SNUM. If a query asks for information regarding the student having number (SNUM) 126, then two of the assertions in the initial state will be: (HASHKEY STUDENT SNUM) (DESIRE SNUM 126).

The first assertion belongs to the set defining the data base, the second to the set defining the query.

Micro-Planner is given the goal (PROGRAM EXAMPLE), and program generation commences. The first sub-goal, (OPENED A1), is satisfied with a set of rules not illustrated here resulting in an initial block of procedure to open the files.

The next sub-goal is (CURRENT R1). Since the rule for FIND has this as a consequence, it is tried. The first two sub-goals in this rule are immediately true in the state, and the variables R1, KI, and V1 get bound to STUDENT, SNUM, and 126, respectively. The third sub-goal is satisfied with the assignment statement rule (not illustrated here), and an assignment statement is inserted in the program. The pre-condition now being true, 'FIND STUDENT' is put into the program.

At this stage the program appears as:

```
(Block to open files)
MOVE 126 TO SNUM.
FIND STUDENT.
```

The state (set of assertions) has also been altered and now contains an additional assertion, (CURRENT STUDENT), because of the FIND rule. This may affect subsequent procedure generation. For example, if the query requires that values from the STUDENT record be in working-storage, the system will not generate any procedure for the (CURRENT R1) sub-goal in the GET rule.

Designing Network Data Structures

Another problem attendant to data base management is data base structure design, that is, the design of appropriate links between data elements so that data relationships can be properly captured and reconstructed. Since all such relationships must be implicitly contained in the set of all retrieval requests, it should be possible to derive a data base structure from such a set<3>.

The Automatic Data Structure Designer operates in exactly this fashion. It generates a data structure design that is satisfactory for a set of queries. The resulting structure is not optimal but has the characteristic that data redundancy is minimized (and as a corollary, the number of links are maximized).

The design program has two major tasks: to design a network structure that can accommodate all of the hierarchical relationships referenced in the queries, and to determine the best location in the structure for each data item.

Detection of a hierarchical relationship in a query is not difficult. The system simply looks for quantifiers or summarizing commands such as "all," "any," "average," "total," etc. For example, a query asking for a class list (all values of student-name when class = "history") suggests a hierarchical relationship between class and student-name. Similarly, average grade for student-name = "Anderson", suggests that there may be many values for grade for a given value of student-name.

A designer has a problem if, for example. In addition to the class list discussed above, a second query asks for a student transcript (all values of class when student-name = "Anderson"). The two queries' views of the data base are inverted with respect to each other.

Included in the Automatic Design System are rules that recognize such situations. These rules will insure that a special linking record is included in the data base design. Through this record will pass two linked lists, one for each hierarchy.

The second principal task of the Designer is to determine the location of data items (attributes) in the hierarchical structure. Queries may differ in their views of attribute associations. For example, hospital name may be viewed as a patient attribute in one query (where he's being treated), as a doctor attribute (where he works), or as a hospital attribute in other queries. If hospital-patient and hospital-doctor hierarchies (several patients and several doctors per hospital) have already been constructed by the system, then it will assign the attribute in question to the hospital record. This permits usage of hospital name as a unique attribute of hospital, doctor, and patient.

Using Earley's terminology, we can say that the Automatic Data Structure Designer finds a logical(b) arrangement

that accommodates all logical(a) relationships. As a matter of fact, translation from logical(a) to logical(b) also characterizes the system's activities as a programmer. Further translation to the physical level is not the task of these programs. That activity is performed by the data management system, e.g. an implementation of the DBTG specification.

As was the case with the Information Retrieval Program Generator, the Automatic Data Structure Designer is defined with a set of rules that are translated by a compiler to a program of Micro-Planner theorems. However, in this case the resulting system is not a code generator, it is a declarative generator.

Representation of Knowledge

If a system is to apply knowledge. It must of course first obtain such knowledge. Until generalized learning systems become available that *can* be adapted in a practical way for use within a particular application domain. It is necessary to imbed knowledge of the application either by building a specific learner for that system or by providing the system directly with the knowledge it needs. The trouble with the *former* approach is that building the learner may turn out to be impractically difficult. The latter approach is the one that was taken for the applications discussed here.

Complexity of Knowledge

The rule formalism provided a good framework for capturing network data base programming and design knowledge. Even so, capturing the knowledge involved many trial-and-error iterations because much of it does not have a well understood structure. Furthermore, it is possible to map the knowledge into a structured set of formal rules in several valid ways. These representations may differ with respect to efficiency.

However, some of the knowledge does have a well-understood structure, for example the knowledge regarding GET or FIND as discussed above. In general, the rules for single program statements like GET are similar to usage rules that might be found in a reference manual. It is the rules that control the next higher level of programming, the construction of program blocks, that are complex and not easily derived.

Such rules were further complicated by considerations of efficiency both for the generated programs and for the program generation process itself.

As an example we illustrate the two rules that control code generation for alternation. Alternation occurs in the generated program whenever the query specifies conditional processing. During program generation, the condition that is to be tested in the generated program is contained in a list, FORMS, in disjunctive form.

For example, the list
 <<<GT A B>>(GT B C) ((LT A 26))) means
 (A>B and B>D) or A<26.

The two rules are:

```
(NULL FORLIS)V
(CONJUNCT ACTION (CAR FORLIS) DUM)A
(DISJUNCT ACTION (CDR FORLIS))
  -> (DISJUNCT ACTION FORLIS)
```

```
((NULL CONJ)A(ACT ACTION)A(SETQ DUM 0))V
(SETQ REL (CAAR CONJ))A
(SETQ ITM1 (CADAR CONJ))A
(SETQ ITM2 (CADDAR CONJ))A
(DETV AL ITM2)A(DETV AL ITMDA
(TEST (LIST REL ITM1 ITM2))A
(CONJUNCT ACTION (CDR CONJ) DUM)
  --> (CONJUNCT ACTION CONJ DUM)
```

Those readers who are programmers will agree that these rules bear little correspondence to their own programming knowledge, either viewed introspectively, or as such knowledge might be represented in a programming text.

The first rule CDR's through FORLIS, each time picking off the CAR, which is a conjunctive list of predicates. The CAR of FORLIS will be bound to CONJ in the second rule.

The second rule CDR's through CONJ until it becomes NULL, at which point the processing code can be generated (accomplished by invoking the ACT rule, not shown here).

Each time the second rule is used, code is generated to determine the values of the two items in the test described by the CAR of CONJ. Then the test itself is inserted in the generated program.

For the example of FORLIS given above, the generated procedure has the following form:

```
(block to determine value of B)
(block to determine value of A)
IF A>B THEN
  BEGIN
    (block to determine value of C)
    IF B>C THEN
      BEGIN
        (block of conditional
        processing)
      END
    ELSE PROC2
  END
ELSE PROC3
```

```
PROC2: (and PROC3)
BEGIN
IF A<6 THEN
  BEGIN
    (block of conditional processing)
  END
END
```

The Buchanan-Luckham compiler takes special care with rules that may generate alternations. Unbound variables become important, and it is for this reason that DUM is included in the second rule. For further details see <1,10,11>.

Tests are separated in the generated procedure to increase runtime efficiency. In the example above, the value of C is determined only if A>B. Since such a value determination may require a lot of processing, it is best to avoid it if possible.

Impact of Representation on Efficiency

Since the content of the rules determines the decision tree that controls the search space, the choice of mappings mentioned above can have a dramatic impact on the efficiency of the search for a program. If the use of a particular rule occurs far out in the branches of the tree, and if the rule is a powerful discriminator, then it would be appropriate to consider an alternative rule structure that would permit earlier use of the rule in question.

It is also important to direct the Information Retrieval Program Generator as much as possible. Consider the FIND rule previously discussed. This rule states that FIND has the effect of making a record current. There are in fact four other rules describing statements that have the same effect. These rules all differ in the way in which they make records current: some use pointers, another utilizes sequential access, etc.

Since there are five ways of making a record current/ Micro-Planner may try several inappropriate rules to satisfy a CURRENT goal before the proper one is attempted* This can be very expensive since an inappropriate rule may not be discarded until a very large subtree has been evaluated.

Frequently, such searches can be eliminated because the programming context will determine which type of FIND is needed. In the iteration step of a loop, for example, we would not expect a direct FIND, rather we would expect a FIND for the next record on a list or the next record in sequence.

By changing the rules so that the system is directed to a specific rule whenever possible, we were obviously able to avoid a lot of unnecessary searching.

It is important to note that the rules were correct without these specific directions. The way in which the rules were used required additional knowledge about the application so that the rules could be used efficiently.

For a particular set of programs that the system generated, it tried (on average) about 94 rules per program. Of these (on average) only 13 rules were inappropriately tried, resulting in backtracking to use alternative rules.

Although we have no exact figures to compare with, the system was trying 10 to 20 times as many unnecessary rules prior to the addition of more specific direction. The difference is multiplicative rather than additive because the sub-tree associated with a single inappropriate rule may be quite large.

Conclusion

As a user of AI technology, I am somewhat disappointed that it was not sufficient to directly transfer my knowledge to the machine. In itself a difficult task. It was also necessary to observe how that knowledge was being used by the machine so that I could change and augment the representation for more efficient use of the knowledge.

These problems were also encountered while building the Automatic Data Structure Designer. Another difficulty arose as well. As the Designer took shape it became apparent that my design knowledge contained many ad hoc techniques that were not easily captured in a set of general rules. This realization led me to develop the rules for the more direct, general purpose algorithm that is embedded in the

Designer.

Although the difficulties of transferring knowledge to the machine were somewhat frustrating, the process (and the machine as a mirror of myself) also led to new insights into the problem domain.

Acknowledgements

The assistance and encouragement of my thesis advisor, Jack R. Buchanan, is gratefully acknowledged, as are the suggestions of one of the referees.

This research was supported in part by the Advanced Research Projects Agency of the Office of the Secretary of Defense under contract Fi44620-73-C-0074, and in part by the Office of Naval Research under contract M0001U-75-C-0W62. I am also indebted to the William Larimer Mellon Fellowship, which provided personal funding during the course of this research.

References

1. Gerritsen, Rob, "Understanding Data Structures," PhD Thesis, Carnegie-Mellon University, Pittsburgh, Pennsylvania, 1975.
2. Gerritsen, Rob, "The Relational and Network Models of Data Bases: Bridging the Gap," Proceedings 2nd USA-Japan Computer Conference, 1975.
3. Gerritsen, Rob, "TS Preliminary System for the Design of DBTG Data Structures," will appear in CACM, October 1975.
4. CODASYL, CODASYL Data Base Task Group April 71 Report, available from ACM, New York City.
5. Earley, J., "On the Semantics of Data Structures," in Rustin, R. (Ed), Database Systems. Prentice Hall, 1972.
6. Bachman, C. W., "The Programmer as Navigator," CACM 16. 11, November 1973, pp653-658.
7. Codd, E. F., "A Relational Model of Data for Large Shared Data Banks," CACM 13, 6, June 1970, pp377-387.
8. Codd, E. F. and C. J. Date, "Interactive Support for Non-programmers: the Relational and Network Approaches," IBM Research RJ-1U00, June 6, 1971*.
9. Date, C. J. and E. F. Codd, "The Relational and Network Approaches: Comparison of the Application Programming Interfaces," Proceedings 1971 ACM-SIGFIDET Workshop.
10. Buchanan, J. R. and D. C. Luckham, "On Automating the Construction of Programs," Stanford AI Project Memo, Stanford University, California, 1970.
11. Buchanan, J. R., "A Study In

- Automatic Programming," PhD Thesis, Stanford California, 1974.
12. Hoare, C. A. R., "An Axiomatic Basis for Computer Programming," CACM 12, 10, October 1969, pp576-580.
 13. Hewl tt, C, "Description and Theoretical Analysis of Planner," PhD Thesis, Massachusetts Institute of Technology, Cambridge, Massachusetts, 1971.
 14. Sussman, Gerald J. and Terry Winograd, "Micro-Planner Reference Manual," Project MAC Report, Massachusetts Institute of Technology, Cambridge, Massachusetts, 1972.