# A NEW SYSTEM FOR INDUCTIVE LEARNING IN ATTRIBUTE–BASED SPACES

Cezary Z. Janikow

University of North Carolina at Chapel Hill,
Department of Computer Science, CB# 3175
Chapel Hill, NC 27599, USA
and
University of Missouri at St. Louis
Department of Mathematics and Computer Science
St. Louis, MO 63121, USA

### Abstract

Despite the fact that the theory and methodology of inductive learning has previously been described, actual symbolic systems for learning descriptions in attribute–based spaces use other algorithms. The AQ family of programs uses properly guided logic–based negation, union, and intersection operators. The ID family of programs uses the idea of iteratively partitioning the search space. The inductive methodology, on the other hand, describes a number of generalization and specialization operators, some of which can be restricted to the attribute–based spaces. Why, then, they are not used in practice? The reasons are the immense search spaces of inductive learning and its weak heuristics.

This paper describes an approach that uses such task–specific inference rules directly on formulas of a rule–based framework. To control such a search, in the absence of strong heuristics, we use the very liberal and robust search mechanism of genetic algorithms. A simple implementation is tested on two standards data sets. The results suggest the feasibility of this approach, not only in terms of numerical qualities, but also in terms of understanding the output knowledge and the processing mechanisms.

## 1  Introduction and Preliminaries

Concept learning from examples in an attribute–based description space is one of the most heavily explored areas of machine learning. Such popularity can be attributed to both the existence of many practical problems sufficiently described in a feature–based world and a wide availability of such data. Over the last few years there have been quite few different approaches to this kind of learning including rule–based, decision tree, connectionist, statistical, and even genetic algorithms.

1

The goal of any system implementing learning from examples is, given initial set of events, to produce classification "rules" for the concept presented in the input set. Such rules might be quite implicit (*e.g.* a network's topology or some statistics) as seen in most non–symbolic systems, or they might be statements of a high level descriptive language in the case of symbolic systems. The purpose of such acquired knowledge is two–fold: to predict classifications of new events, and to make such knowledge available to humans (or other entities of an intelligent system). Of those two different classes of methods, the symbolic one is more appealing since its high level output language allows for an easier understanding of the generated knowledge. This issue becomes more and more important with the increasing popularity and applicability of artificial systems. For example, Michalski wrote in [8]:

> "An important implication [...] is that any new knowledge generated by machines should be subjected to close *human scrutiny* before it is used. This suggests an important goal for research in machine learning: if people have to understand and validate machine–generated knowledge, then machine learning systems should be equipped with adequate *explanation capabilities*. Furthermore, knowledge created by machines should be expressed in forms closely corresponding to human descriptions and mental models of this knowledge; that is, such knowledge should satisfy what this author calls the *comprehensibility principle*."

However, just plain understanding of the automatically generated knowledge is not enough if one wants to verify it, or possibly learn some general lessons from the acquisition process itself. In such cases, understanding of the methods that led to such results is critical. Therefore, these methods should be explainable at the problem level; The same author wrote in [7]:

> "...one should strive to facilitate human understanding not only of the surface results but also of the underlying principles, assumptions, and theories that lead to these results."

Because the framework for inductive learning, along with clearly understood mechanisms manipulating knowledge, has been described in [6], it would seem natural to apply it to the restricted attribute–based spaces as well. This would make both the processing and its outcome more understandable and meaningful. However, the most known systems for concept learning are offspring of two dominating methodologies: AQ and ID, which both use different algorithmic mechanisms.

The AQ family uses a rule–based framework. It simultaneously generalizes and specializes the current knowledge, but it does so by logic–based operators of negation, union, and intersection. It uses two kinds of heuristics: explicit measures restricting the number of plausible disjunctions, and implicit heuristics implemented in the iterative control. The ID family uses decision trees constructed by iterative partitioning of the event space. This approach is equivalent to applying only specialization. Moreover, such specializations are performed based on the assumption of attribute independence.

The reasons for such a diversion between theory and practice are the immense search spaces of attribute–based descriptions. For example, using ten descriptive attributes with three values per domain we are faced with $7^{10}$ number of different rules. Then, the number

of possible concept descriptions is $2^{7^{10}}$, or almost $10^{100000000}$. One solution to this problem would be to use hill-climbing techniques. However, the available heuristics, relying mostly on partial measures of completeness and consistency, could easily lead to local traps. Therefore, irrevocable strategies are quite inapplicable. Another solution would be to use some tentative techniques. However, the extremely large search spaces require any such method to be extremely well informed, or otherwise the explored database would grow unmanageably fast. Again, the available heuristics are not strong enough to provide such qualities.

To solve this problem we use a search mechanism that keeps the database size under control by retaining only a fixed number of states. This is achieved by abandoning states that show little potentials, and fully exploring those most promising. To implement such control we use the mechanisms of genetic algorithms, which exhibit exactly such behavior and was shown to be very robust.

In the rest of this paper we first briefly describe inductive learning in attribute–based spaces and the ideas of genetic algorithms. Then, we describe our new approach followed by some experimental results. Finally, we close by drawing some conclusions.

# 2  Concept Learning in Attribute–Based Spaces

Concept learning is a fundamental cognitive process that involves learning descriptions of some categories (*i.e.* sets drawn from a common class) of objects. Such descriptions can be drawn from different universes. We are considering here the restricted attribute–based spaces, spanned by a number of variables, each of which has a finite set of allowed values (in essence, this is a generalization of boolean functions from two– to multi–valued domains). A priori knowledge consists of a set of events, *i.e.* examples of the space. Moreover, we are dealing here with the case of supervised learning, *i.e.* learning assuming that each a priori event is preclassified as an example of one of the concepts to be learned. The task is to generalize the a priori knowledge in order to produce descriptions of such concepts.

An important issue is that of defining both the input and output language. While the most natural input language allows for specification of objects by conjunctions of their features related to different attributes, the choice of the output language is more elaborate: the most known are decision trees ([10]) and rules ([6]). However, it is clear that the same language on both ends of the learning process facilities some desirable learning characteristics *e.g.* incremental, closed–loop and multistrategy task–adaptive learning. In this sense, rules are advantageous.

One widely used language, which is not only closely associated with rules but also normally used to represent input examples for majority of attribute–based programs, is $VL_1$ ([7]). Variables (attributes) are basic representational units. Each variable has a multi–valued domain. These domains might be of different types, depending on the relations among their members (*e.g.* *linear* for linearly ordered values, *nominal* for unrelated values, and *structured* for partially ordered sets). Relations associate variables with their values (*e.g.* $=, \neq, \geq$). A selector has the form [*variable relation value*], where the *value* might be a set of values from the corresponding domain (called an internal disjunction). A conjunction of selectors forms a complex. The selectors can be used to express con-

ditions of a rule–based paradigm; Consequently, complexes can be used to express rules (*e.g. complex* ::> *decision*), and their sets can be associated with sets of rules constituting a system's decisions. As an input language, $VL_1$ expresses each category, or a set of events of a given class, as a set of complexes such that each one is a conjunction of all single–valued attributes. On the output, a disjunction of any complexes specifies the set of rules describing the given category. In practice, such formulas are quasi–optimized with respect to some criteria, as there might be many different formulas that are complete and consistent with the given input set. In order to generate the concept descriptions, induction must be applied to the a priori knowledge.

According to Holland ([2]), an inductive process in a rule–based framework must accommodate for both revision of rules and generation of new ones, where generation of new ones is connected to two sources: the environment and the existing knowledge. According to Michalski ([6]), an inductive learning system must accommodate for production of new rules, as well as generalization and specialization of existing ones. The same paper ([6]) provides a detailed description of various inductive operators that constitute the process of inductive inference. In the somehow restricted language $VL_1$, the most important are: *condition dropping* — *i.e.* dropping a selector from a concept description; *Adding alternative rule* and *dropping a rule* — adding/removing one rule from a description; *Extending a reference* — extending an internal disjunction; *Closing an interval* — for *linear* domains filling up missing values between two present values in a condition; *Climbing generalization* — for *structured* domains climbing the generalization tree; *Turning a conjunction to a disjunction*; *Inductive resolution* — analogous to the resolution principle. These operators are either generalizing or specializing existing knowledge.

Despite the availability of so clearly defined methods, the two most known symbolic systems use somehow different algorithmic approaches. AQ–based systems iteratively construct covers for all still uncovered positive events, and then accumulates them by means of logical unions. Each cover is generated by iteratively intersecting partial covers consistent with exactly one negative event. To keep the complexity manageable, only a number of best partial descriptions is retained at any time. The result is always a complete and consistent description, assuming consistency in data. ID–based systems iteratively partition the event space, starting with the whole space and working on all partitions until they are consistent. Naturally, such partitions are always complete as well. At any given moment, the partition is based on the attribute giving the highest increase in the information contents. Therefore, this algorithm treats all attributes independently.

# 3   Genetic Algorithms

Genetic algorithms (GAs) are adaptive methods of searching a solution space by applying operators modeled after the natural genetic inheritance and Darwinian struggle for survival. They belong to the class of probabilistic algorithms, yet are distinguished by their different search method and relative robustness to local traps.

In general, a GA performs a multi–directional search, and it encourages information formation and exchange between such directions. It does so by maintaining a population of proposed solutions (chromosomes) for a given problem. This population undergoes a simulated evolution: relatively "good" solutions reproduce producing offspring, which

subsequently replace the "worse" ones. The quality of a solution is estimated based on an evaluation function, which plays the role of an environment. The existence of this population provides for the superiority of genetic algorithms over pure hill–climbing methods, for at any time the GA provides for both local exploitation of the most promising solutions and global exploration of the whole search space. Moreover, the restriction on the population size provides for the superirity over deterministic search methods in huge spaces.

The algorithm works iteratively until a termination condition is satisfied. Each iteration, called a reproduction cycle, is performed in three steps. During the selection phase a new population is formed from stochastically best samples (with replacement). Then, during the recombination phase some of the members of the newly selected population are altered. Finally, all such altered individuals are evaluated.

The alteration process (recombination) is based on the application of two operators: mutation and crossover. Mutation introduces random variability into the population, and crossover exchanges random pieces of two solutions in the hope of propagating partial solutions. Because both these operators are defined on syntactical pieces of some underlying representation, the search has domain–independent properties. The robustness is, on the other hand, a direct result of the existence of the population, which allows the search to be performed simultaneously in many different directions.

# 4   The New Approach

We want to organize the available inference mechanisms in a rule–based production system with a non–deterministic control borrowed from genetic algorithms. Such control restricts the number of currently open states for expansion, but does it in a very intelligent, stochastic way: states having stochastically better outlook are explored with a proportionally increasing effort. This also means that states that show little promise slowly die out.

Production systems were originally inspired by attempts to model the human cognitive process, and were first proposed by Post in 1943. Their idea was to use production rules, which change states in a way of firing neurons. What artificial intelligence found interesting in production systems was the clear separation of various elements of the paradigm: control, current situation, and production rules. This, in turn, allowed for transparency, modular design, and ease of both maintenance and knowledge refinement. Many different versions and generalizations of such a design were proposed by the artificial intelligence community: some specific for certain domains. However, they all share the ideas of the high level modularity, application of rule–like operators, and the name of production systems, AI production systems, production rule systems, *etc.* ([9]).

We describe the system in terms of the same three top level components: control, database, and inference rules. For simplification of discussion, we assume that we are learning single concepts only. This scenario lets us deal only with one concept description, and lets us assume that the uncovered portion of the event space represents the negation of the concept. This, in turn, allows us to learn sets of $VL_1$ complexes in place of actual rules — the decision is implicit and associated with the concept being learned. This system uses the same full memory assumption as previously mentioned AQ and ID systems.

We use the rule–based framework of the $VL_1$ language. Consequently, the search space becames the space of all $VL_1$ descriptions. The goal is to find the best state that fits some criteria (completeness, consistency, and possibly some learning bias), and not the path leading to such a state. Then, each state of the current database is a sets of $VL_1$ complexes: a potential solution. Each complex of these sets is a conjunction of conditions that must be satisfied. Finally, each condition is related to exactly one attribute, and is equivalent to the $VL_1$ selector.

The operators transform states of the database to new (possibly better) states in the search space. Since the system operates in the space of $VL_1$–based descriptions, the operators model those of the inductive learning methodology when restricted to the attribute–based spaces. To fully use the idea of the population, and to provide similar qualities as those of the genetic algorithm search (*i.e.* robustness), we introduce an additional operator which exchanges pieces of information between different states (rules). We also define few additional operators utilizing some additional ideas of inductive learning, as replacing a number of rules by their most specific generalization or their most general specialization. Having the advantage of knowing the AQ and ID algorithms, we define two operators simulating basic steps of those two systems. Thus, we have an operator which, given a rule and a negative inconsistent event, replaces this rule by all maximal rules consistent with the old rule and that event (same as AQ's kernel action). Following the ideas of the ID algorithm, we have an operator which partitions a rule using exactly one attribute. Finally, we also have a rule which, given a state (*i.e.* a rule set) and a positive uncovered event, adds this event as a new rule to this current description.

Each operator is given an initial probability of application to its type of structure (based on the level of its definition: set of complexes, complex, and selector levels). These probabilities have a dynamic character with respect to the current context, *i.e.* to the current completeness and consistency. Priori probabilities of generalizing operators are increased for applications to structures (rule sets, rules, conditions) that are incomplete. Priori probabilities of specializing operators are increased for applications to structures that are inconsistent. Moreover, the levels of probability increase/decrease are based on the levels of inconsistency/incompleteness. In other words, these two measures serve as heuristics guiding the selection of appropriate operators.

The control follows that of genetic algorithms. First, the initial database is filled with both random rules and positive training events. Then, all the rule sets are evaluated to provide some heuristic measures of accomplishments. Such measures use the idea of normalized correctness based on completeness and consistency taken with equal weights. To prevent redundancy in the descriptions, the measure is slightly adjusted to accommodate for the number of complexes in a description (however, we do not use any redundancy removing axioms). In addition, one may want to include an additional bias based on some learning criteria. For example, in the experiments shown here we also accommodated a complexity measure which reflected the number of conditions (as in [13]). Following these preliminaries, the algorithm enters a stage in which the three iterative steps are performed (selection, reproduction, evaluation) until a desired state is found. Since in general we do not know the sought description a priori, the termination condition may be based on the amount of resources available or some criteria for the solution.

During selection a new database is generated by choosing stochastically better states (with repetitions). In other words, the most promising states may appear multiple number

of times in the new database, while the weak ones may be abandoned. During reproduction the set of inference rules hypothetically applies to all structures of the newly selected database. This action generates a set of new states, which are offspring of the higher–valued previous states. Some of the operators require a simple condition to be satisfied (*e.g.* the operator generating a new rule from an uncovered positive training event) to candidate for firing, others do not have any preconditions. However, the actual firing is non–deterministic, with probabilities based on priori values and the context in which each operator applies (completeness and consistency of current structures). All operators found to fire perform their actions: there is no conflict resolution. Finally, all new states are reevaluated using the available heuristics.

Our operators are defined as simple atomic actions, which can be performed very efficiently on the $VL_1$ descriptions if a special representation is used (binary representations of selectors, see [4]). However, some of the operators require finding a proper positive or negative event. Moreover, the reevaluation of a state requires an extensive pattern matching against all training events in order to estimate its completeness and consistency. This may be uneconomical and highly inefficient. To deal with this problem we precompile the training data into binary coverage vectors, and we subsequently operate on such structures. The idea is as follows: rather than storing data in terms of features, store features in terms of data coverage assuming some artificial data enumeration. Then, finding an uncovered positive event, or an inconsistent event, of a rule can be easily accomplished by selecting a binary one from an appropriate vector. Also, evaluating the completeness and consistency can be accomplished by counting the number of such binary ones. Moreover, such coverage vectors can be easily updated incrementally after an operator fires, based on understanding of the action of each operator. For example, having such binary vectors $b_l^i$ and $b_m^i$ for two rules $l$ and $k$ of a state $i$ requires only a simple bitwise AND to construct the coverage of their most general specification $b_{l,m}^i = b_l^i$ AND $b_m^i$.

# 5    Some Experiments

| | Learning Scenario (Positive%/Negative%) | | | | |
|---|---|---|---|---|---|
| System | 6%/3% | 10%/10% | 15%/10% | 25%/10% | 100%/10% |
| AQ15 | 22.8% | 5.0 % | 4.8 % | 1.2% | 0.0% |
| BpNet | 9.7% | 6.3% | 4.7% | 7.8% | 4.8% |
| C4.5 | 9.7% | 8.3% | 11.3% | 2.5% | 1.6% |
| CFS | 21.3% | 20.3% | 21.5 % | 19.7% | 23.0% |
| GIL | 4.3% | 1.1% | 0.0% | 0.0% | 0.0% |

Table 1: Quantitative results: average error as function of the percentage of seen positive and negative examples.

To obtain some experimental results we implemented a simple version of our algorithm (called GIL) and ran it on two standard data sets: robots of the Emerald system ([5]) and the three–address multiplexer $f_{11}$.

In the robots case we exactly repeated the experiments described and reported in [13]. In there, four learning systems were used: rule–based AQ15, decision tree constructing C4.5 (with a special module to convert trees to rules), a connectionist implementation of learning by back propagation of a failure BpNet, and a genetic algorithm in a classifier system CFS. This allows us use those results for comparisons. The robots were described by the following six attributes:

| Attribute | Values |
|-----------|--------|
| *Head* | *Round, Square, Octagon* |
| *Body* | *Round, Square, Octagon* |
| *Smiling* | *Yes, No* |
| *Holding* | *Sword, Balloon, Blue* |
| *Color* | *Red, Yellow, Green, Blue* |
| *Tie* | *Yes, No* |

and were classified into the following five categories (created by a task–unaware human):

| Concept | Description |
|---------|-------------|
| $C_1$ | *Head* is *Round* and *Color* is *Red* or *Head* is *Square* and *Holding* a *Balloon* |
| $C_2$ | *Smiling* and *Holding* a *Balloon* or *Head* is *Round* |
| $C_3$ | *Smiling* and not *Holding* a *Sword* |
| $C_4$ | *Jacket* is *Red* and no Tie or *Head* is *Round* and is *Smiling* |
| $C_5$ | *Smiling* and *Holding* either a Balloon or a *Sword* |

The task was to learn a description of each concept while seeing only a percentage of the positive and the negative examples. There were a total of 432 different robots in this world. Two different measures of performance were described in [13]: error rate in recognizing both training and testing events after a learning session, and complexity of the acquired knowledge measured by the number of rules and conditions when viewed in a rule–based framework.

| | Learning Scenario (Positive%/Negative%) | | | | |
|--------|---------|----------|----------|----------|-----------|
| System | 6%/3% | 25%/10% | 50%/10% | 75%/10% | 100%/10% |
| AQ15 | 2.6/4 | 1.6/3 | 1.6/3 | 1.6/3 | 1.6/3 |
| BpNet | NR | 18/29 | NR | NR | 32/54 |
| C4.5 | 6.8/12.2 | 4.4/9.2 | 4.8/9.2 | 4.8/9.2 | 3.8/7.3 |
| CFS | NR | NR | NR | NR | NR |
| GIL | 1.4/2.6 | 1.6/3 | 1.6/3 | 1.6/3 | 1.6/3 |

Table 2: Qualitative results: average #rules/#conditions.

Table 1 reports the average error rate for the five learned concepts for all five systems. Surprisingly, GIL produced the highest recognition rate, especially in cases of small training sets. Table 2 reports the average acquired knowledge complexity as learned by all five systems, for different learning scenarios; The NR entry indicates that the complexity was large and not reported in our reference paper. The reason for the quite higher complexity

of the connectionist approach is that this is a subsymbolic system operating on numerical weights rather than on the problem symbols. On the other hand, the high complexity of the genetic approach can be attributed to the fact that the symbolic processing was being done in the representation rather than the problem space. This result is rather common for genetic algorithm approaches; Therefore, it was a big surprise to find out that GIL's knowledge was at the same complexity level as that of the highly acclaimed AQ15. An interesting fact to observe is that GIL was the only one to overgeneralize one of the descriptions — a result of the strong bias toward simple descriptions used in these tests. We believe that this bias was the main reason the superior results in this case, since the sought descriptions were all fairly simple. For this experiment, a single run took an average of about 10 CPU seconds on a DEC3100 station (for 100 iterations).

| % training events | Accuracy |
|---|---|
| 5% | 77% |
| 10% | 88% |
| 20% | 94% |

Table 3: GIL's accuracy on multiplexer $f_{11}$.

The other data set that we used was the three-address multiplexer $f_{11}$. For each integer $k = 1, 2, \ldots$ there is a multiplexer boolean function defined in the following way: the function's inputs are the $k$ bits (called addresses), and there are exactly $2^k$ outputs (called data bits). The function of a multiplexer is to activate the data bit whose address (in binary, assuming some ordering of the data bits starting at 0) is specified by the address bits. In our case we have three address bits and eight data bits, for a total of eleven variables. Therefore, the size of the event space is 2048. Table 3 reports an average accuracy while learning with varying sizes of the training set. These results are similar to those from other systems (*e.g.* [11]), but a different experimental methodology doesn't allow for a direct comparison. An interesting fact observed here was that this learning was quite slower than for the robots: it took an average of 20 CPU minutes on the same DEC3100 station (for 2000 iterations). There seem to be two major reasons for this increase. Firstly, the complexity of the sought here description was 48 vs. an average of 6.2 in the previous case. Secondly, the number of possible complexes is much larger in the second case, giving a much larger search space: approximately $10^{60000}$ vs. $10^{15000}$. Indirectly, the larger search space required more iterations for the learning, and the larger average state size caused a longer processing on each iteration. However, this complexity still compares very favorably with other GA approaches.

Finally, for an illustration, figure 1 traces a sample run while training with 20% of the available events. During this learning session, the exact concept was learned after 1700 iterations. A consistent and complete description of the training events was found shortly after 1000 iterations, and the remaining 700 cycles were required to simplify the generated description. The first of these two graphs traces completeness and consistency of the currently best database individual at 17–iteration intervals (100 data points). The other graph traces the complexity of such best states. It is interesting to note that the
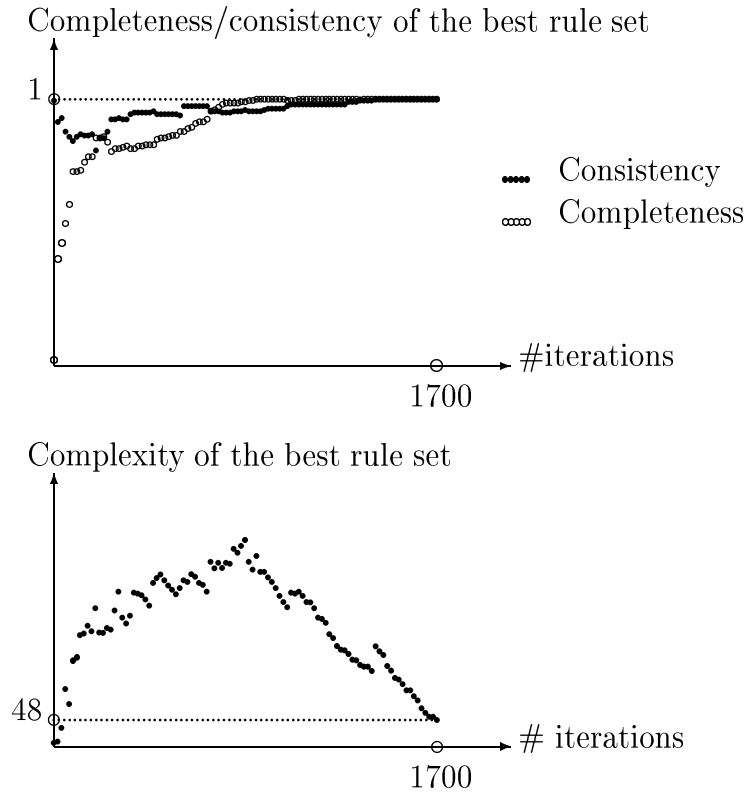
Completeness/consistency of the best rule set



Complexity of the best rule set



Figure 1: A sample behavior on multiplexer $f_{11}$.

complexity rises during the learning, as a result of not finding simple enough complete and consistent descriptions, and then decreases — forced down by an increasing cost influence and formation of such better descriptions.

# 6    Summary

We described here a new system for learning concept descriptions from examples of attribute–based spaces. It combines two well known ideas: inference rules of the inductive methodology applied in a production system like formalism, and search mechanisms of genetic algorithms. The task–specific operators implement deep knowledge of the problem solving methodology, and the population–oriented search provides a robust way of dealing with exponential explosion. This method shows big potentials, as it can be seen as simultaneous multi–directional search, with means for information sharing among these different states. Moreover, the transparency and independence of the inference mechanisms allow us to use other known mechanisms as means of finding better states. This way, the new system can be seen as a framework for a hybrid system providing both competition and cooperation among other algorithmic methods.

A special implementation allows farly efficient runs. We hope that this current complexity can be even further reduced. Other possible improvements involve abstracting the huge parameter space of the current implementation, dealing with multiple concepts, and extending the language to extensions of $VL_1$.

This research was partially supported by a computational grant from the North Car-

olina Supercomputer Center.

# References

[1] DeJong, K.A., Spears, W., *"Using Genetic Algorithms for Supervised Concept Learning"*, *Proceedings of the Tools for AI International Conference, 1990.*

[2] Holland, J.H., Holyoak, K.J., Nisbett, R.E., and Thagard, P.R., *Induction*, The MIT Press, 1986.

[3] Holland, J.H., *"Escaping Brittleness"*, *Machine Learning*, Vol. 2, R. Michalski, J. Carbonell, and T. Mitchell (eds.), Morgan Kaufmann Publishers, 1986.

[4] Janikow, C.Z., *"Inductive Learning from Attribute–Based Examples: A Knowledge–Intensive Genetic Algorithm Approach"*, Doctoral Dissertation, University of North Carolina at Chapel Hill, 1991.

[5] Kaufman, K.A., Michalski, R.S., Schultz, A.C., *"EMERALD 1: An Integrated System of Machine Learning and Discovery Programs for Education and Research"*, Center for AI, GMU, User's Guide. No. MLI–89–12, 1989.

[6] Michalski, R., *"A Theory and Methodology of Inductive Learning"*, *Machine Learning: An Artificial Intelligence Approach* Vol. 1, R. Michalski, J. Carbonell, and T. Mitchell (eds.), Tioga Publishing Co., Palo Alto, CA.

[7] Michalski, R., *et al.* , *"The AQ15 Inductive Learning System: an Overview and Experiments"*, Dept. of Computer Science, Univ. of Illinois, 1986.

[8] Michalski, R., *"Understanding the Nature of Learning"*, *Machine Learning II*, Morgan Kaufmann, 1986.

[9] Nilsson, N., *Principles of Artificial Intelligence*, Morgan Kaufmann, 1980.

[10] Quinlann, J.R., *"Induction of Decision Trees"*, *Machine Learning*, Vol. 1, No. 1.

[11] Quinlan, J.R., *"An Empirical Comparison of Genetic and Decision–tree Classifiers"*, *Proceedings of the Fifth International Conference on Machine Learning*, 1988.

[12] Reinke, R.E. & Michalski, R., *"Incremental Learning of Concept Descriptions"*, *Machine Intelligence XI*, D. Michie (ed.), 1985.

[13] Wnek, J., Sarma, J., Wahab, A., Michalski, R., *"Comparing Learning Paradigms via Diagramatic Visualization"*, *Methodologies for Intelligent Systems 5*, M. Emrich, Z. Ras, M. Zemankowa (eds), 1990.