

Incremental Discovery of Hidden Structure: Applications in Theory of Elementary Particles

Jan M. Żytkow† and Paul J. Fischer‡

†Computer Science Department, Wichita State University, Wichita, Kansas 67260-0083
and Institute of Computer Science, Polish Academy of Sciences, Warsaw;
‡Sterling Commerce, 15301 Dallas Parkway, Suite 400, Dallas TX. 75248;
zytkow@cs.twsu.edu, pfischer@gte.net

Abstract. Discovering hidden structure is a challenging, universal research task in Physics, Chemistry, Biology, and other disciplines. Not only must the elements of hidden structure be postulated by the discoverer, but they can only be verified by indirect evidence, at the level of observable objects. In this paper we describe a framework for hidden structure discovery, built on a constructive definition of hidden structure. This definition leads to operators that build models of hidden structure step by step, postulating hidden objects, their combinations and properties, reactions described in terms of hidden objects, and mapping between the hidden and the observed structure. We introduce the operator dependency diagram, which shows the order of operator application and model evaluation. Different observational knowledge supports different evaluation criteria, which lead to different search systems with verifiable sequences of operator applications. Isomorphism-free structure generation is another issue critical for efficiency of search. We apply our framework in the system GELL-MANN, that hypothesizes hidden structure for elementary particles and we present the results of a large scale search for quark models.

Introduction

Intense research in physics during the 1950s and early 1960s centered on the discovery of elementary particles. After more than one hundred elementary particles were known, many arranged into groups with internal symmetries (e.g., hadron octet shown in Figure 1a and meson octet in Figure 3), physicists in the 1960s started to postulate theories of smaller particles, called quarks, proposing their number, properties, and structures they form. Eventually, the standard quark model became one of the foundations of physics.

The discovery problem has been: “Given a set of observed objects and observational knowledge about them, postulate a hidden layer of objects and their structure that explains observed objects”. This problem has been considered many times in the history of science. Examples of successful discoveries include elements, atoms, ions, genes, and quarks. Today the same problem is being re-visited in particle physics, where

scientists search for the next layer of structure beneath quarks.

Discovery of hidden structure has been the subject of various case studies, that led to a number of discovery systems. STAHL (Langley, Simon, Bradshaw, & Zytkow 1987) and STAHLp (Rose & Langley 1986) discover componential models, while DALTON (Langley et. al 1987) discovers atomic models. REVOLVER (Rose 1989) deals with revision of beliefs about hidden structure, MECHEM (Valdes-Perez 1992) infers plausible intermediate structure in chemical reactions, and BR-3 (Kocabas 1991) demonstrates how hidden properties can be postulated for observable objects. Sleeman, Stacey, Edwards, and Gray (1989) have suggested a search space for hidden qualitative models of chemical structure. Valdes-Perez, Simon, and Zytkow (1993) introduced a matrix representation of structure that works for many discovery systems.

This paper presents a general framework which can be used to design various systems that search for hidden structure in different domains. We discuss representation of hidden structure, operators which construct tentative solutions, and the solution evaluation. All these elements are combined in the operator dependency chart that is instrumental in construction of discovery systems. We then discuss GELL-MANN, a system which can postulate hidden structure of elementary particles. It has produced the standard quark model, various alternatives to that model, and many other models of hidden structure.

GELL-MANN is a case study in automated discovery. Instead of speculating on the nature of general purpose automated discoverers, we follow the program of scientific research which relies on an accumulation of case studies that can be used as facts of high order. Experience of empirical science shows that accumulation of many such cases eventually leads to striking theories.

Hidden structure

Hidden structure can be described in the same way as visible structure. However, since hidden objects are not observed, a description of hidden structure must also include a mapping to the level of observation. We

define hidden structure by the following components:

1. A set $T = \{t_1, \dots, t_N\}$ of different types of hidden objects (elements). The number of objects within each type is not limited.
2. A set $C = \{c = (t_{j_1}, \dots, t_{j_s}) \& \varphi(c)\}$ of admissible microstructures (microcompounds), each defined as a bag (multiset) of objects t_{j_p} from T . $\varphi(c)$ is a constraint on admissible structures. For instance, the constraint used by GELL-MANN requires that each bag contains the same number M of hidden objects.
3. A set of attributes $P = \{P_1, \dots, P_k\}$ for objects in T and in C .
4. A set V_i of admissible values for each attribute in P .
5. Specific attribute values for each object type and each attribute, $P_i : T \rightarrow V_i, i = 1, \dots, k$.
The properties of admissible combinations in C are related to the attribute values of the components by additivity:
For each property P_i , each object c , and all components c_1, \dots, c_M of c in C : $P_i(c) = \sum_{j=1}^M P_i(c_j)$.
6. A set R of reaction schemes $(c_{i_1}, \dots, c_{i_p}) \mapsto (c_{o_1}, \dots, c_{o_r})$, in terms of inputs and outputs on the microlevel.
7. Partial mapping $\psi : C \rightarrow \Omega$ between microstructures and the set Ω of observable objects.

Not all components 1–7 must be present in each model. On the other hand, this definition can be augmented by further relationships between hidden objects, such as chemical bonds and spatial proximity. A computer model of hidden structure is a data structure that fits our definition.

Discovery of hidden structure

We will concentrate on the architecture of GELL-MANN, but to illustrate the generality of our framework, we will also use examples from DALTON, a system that simulates discovery of atoms and molecules (Langley et al. 1987). DALTON uses components 1–2 and 6–7 of our definition; 1 and 7 are trivial: one type of atom corresponds to each chemical element, and one type of molecule to each substance. GELL-MANN has been developed to explore quark models in the domain of elementary particles. GELL-MANN uses components 1–5 and 7 of our definition.

The search for hidden structure should propose as much of the data structure as fits our definition, as can be tested by the evidence at the observable level.

Input and output of model construction. The input to GELL-MANN is a family of elementary particles, their properties, and values for each property of each particle. Figure 1b gives an example of input (particle family of hadrons), from which GELL-MANN infers, as output, two models of underlying structure shown in Figure 2. These two quark models postulate three types of hidden objects (a, b, c) occurring in triplets, which are mapped to the input particles. For instance,

in Model 1 the quark a has a charge of $2/3$, a third component of isospin of $1/2$ and a strangeness of 0. In this model, proton p consists of two quarks a and one quark b . Figure 3 provides another example of GELL-MANN's input and output: the meson family.

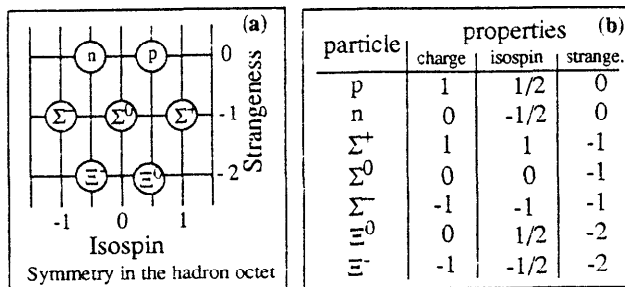


Figure 1: (a) symmetry in the hadron octet family; (b) hadron octet as input given to GELL-MANN

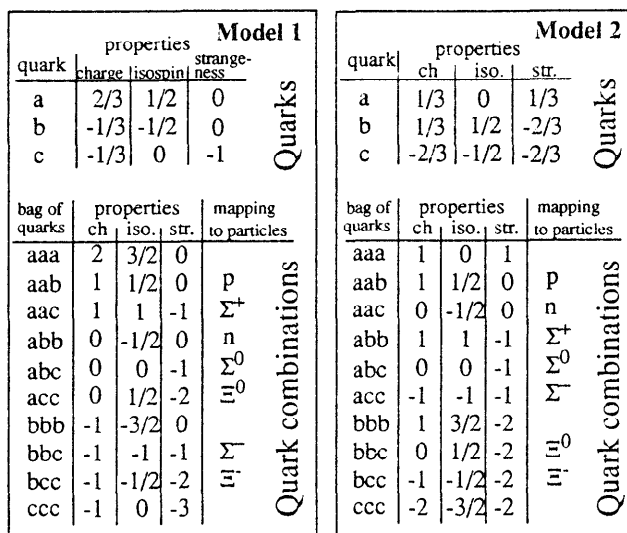


Figure 2: Output of GELL-MANN for the hadron octet. Both models are equally simple. Model 1 is the standard quark model in physics. Quarks u,d,s are GELL-MANN's quarks a,b,c.

Model evaluation. Philosophers since Democritus have speculated about the makeup of atomic structure, but they could justify neither concrete properties of atoms nor concrete atom combinations, because it was difficult to find observational consequences of specific claims about hidden structure. At certain times, however, the knowledge about hidden structure has progressed remarkably. Historically, such progress has occurred when simple symmetries or combination laws expressed in terms of small integers have been detected at the observable level. At the beginning of the 19th century, the law of constant proportions in chemical reactions and Gay-Lussac's law of combining volumes created such an opportunity. Later, the Prout's hypothesis on atomic numbers of elements and the periodic table stimulated models of the atom and its nucleus. After Mendel discovered simple combination rules for properties of the pea, he postulated the

gene model. Similarly, Murray Gell-Mann proposed the quark model after grouping elementary particles into small families.

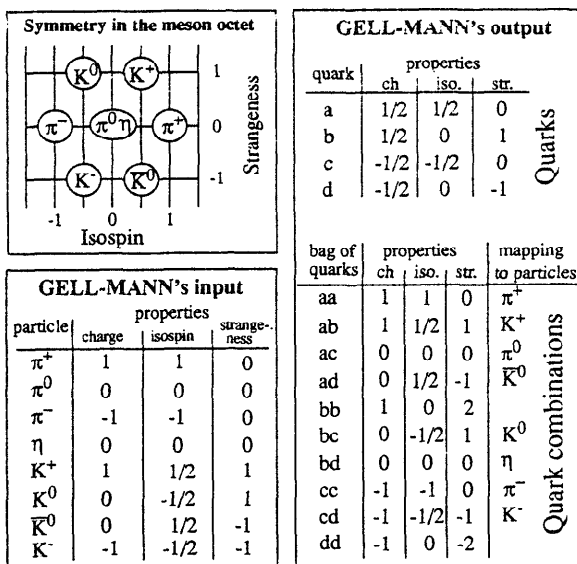


Figure 3: When given the meson octet as input, GELL-MANN finds one model which consists of four quarks in combinations of two.

Facts useful for evaluation come from two basic sources: attributes of the observed objects, and descriptions of observed reactions. In GELL-MANN, evaluation is based on properties of elementary particles, whereas DALTON tests its models against knowledge of macro-reactions. GELL-MANN uses each property P of each observed object O in the input to verify a model, after it proposes a mapping between O and a microstructure c , and applies the *additivity principle* (cf. component 5 of hidden structure definition) to hidden components c_1, \dots, c_M of c . The model is confirmed if $P(c)$ is equal to the observed value $P(O)$.

Reactions can be used in a similar manner. DALTON uses knowledge of combining volumes in a reaction on the observable level and the postulated microstructure of each substance in the input of the reaction. Then the micro-output is computed based on the conservation of elementary parts, so that the number of atoms of each type is equal before and after reaction. Finally the output is interpreted in terms of volumes on the observable level, and compared with the observed output.

It is not sufficient to confirm a model by observational consequences. If there are many models, all justified by their observational consequences, what are the reasons to claim that one of them is true. Each model is questionable because they make mutually inconsistent claims about the hidden level. We cannot require absolute uniqueness, because for each model there are many models which are more complex and observationally equivalent. We can accept model M when all other models are more complex, that is, when M is *unique in the simplest class in which a model exists*.

When the search is arranged in the order of growing complexity, if it is successful, it finds the simplest model. Complexity is measured by model parameters, such as the number of elements postulated and the number of elements in each microcompound.

Operator Dependency Chart

Each model can be constructed gradually in steps that correspond to items 1-7 in our hidden structure definition in Section . Each item in the definition can be represented by operators that build the corresponding parts of the model: add objects to T , postulate their properties, their structures, and so forth.

Operator selection. Not all components of hidden structure are discovered by every system. It does not make sense to propose components which cannot be verified. For instance, the observational data for DALTON do not include properties of molecules and therefore, do not permit verification of hypotheses about properties of atoms. No data on reactions can be used by GELL-MANN, so the inference of hidden structure of reactions is not feasible.

Dependencies among operator application. The order of model construction must satisfy the preconditions at each step. The preconditions can be inferred from the definition. For instance, assigning an attribute value to an object requires an object, an attribute, and a candidate value. Similarly, one cannot create structures without having postulated objects. Activities which lead to model generation and their preconditions can be arranged in a chart, depicted in figure 4a.

Operator dependency chart and search. Different components of hidden structure are postulated by operators. Alternative models are constructed by following alternative paths; that is, by different operator instantiations. Figure 4b shows the subset of all operators used in GELL-MANN, while figure 4c shows the subsets of operators used by DALTON. Only the operators that lead to observational evaluation have been left in Figures 4b and 4c.

The constraints imposed by preconditions leave a great deal of freedom in arranging the search control, so that additional requirements of efficiency can be satisfied. To construct efficient search for hidden structure in a given domain we can use the operator dependency chart and several principles:

1. All operators which do not contribute components of structure which can be evaluated by the existing evidence should be removed.
2. The evaluators must be used as early as possible.
3. Consider each possible model exactly once. The search should be systematic, so that no model is overlooked, but also isomorph-free (non-redundant).
4. Try models in the order of increasing complexity.
5. Use depth first search within each simplicity class.

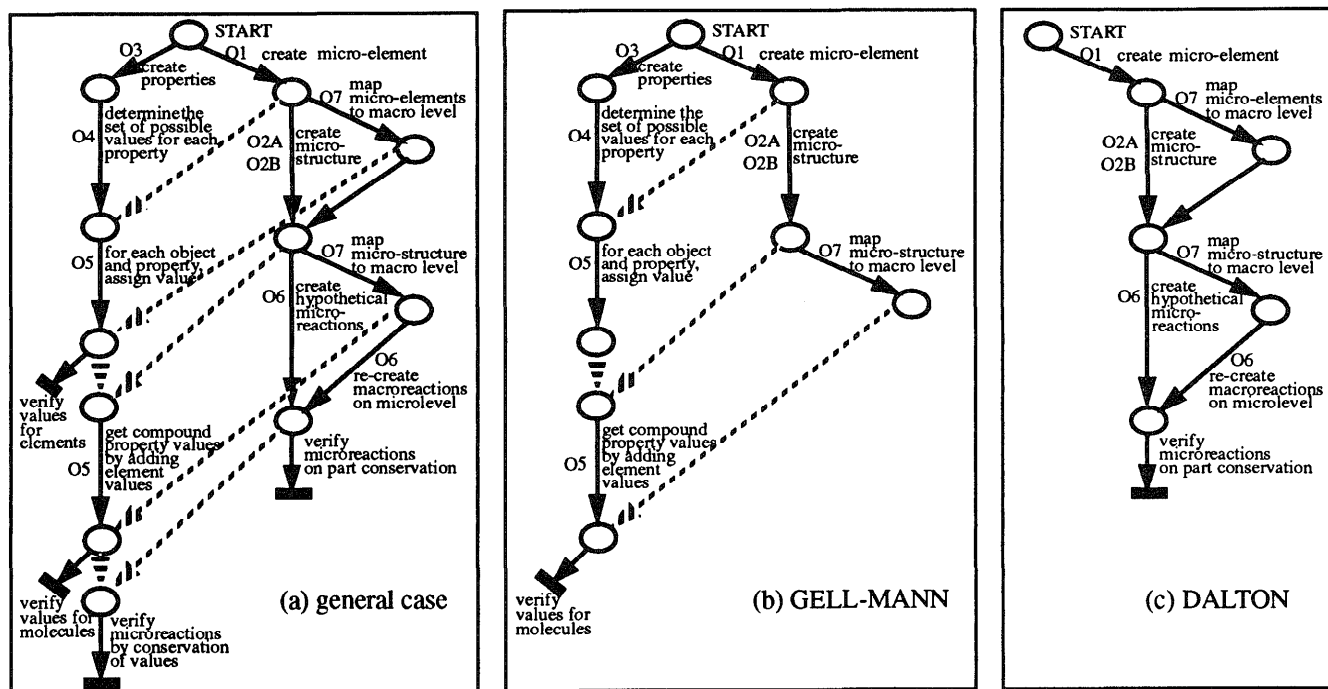


Figure 4: Operator dependency chart for discovery of hidden structure. It shows preconditions for each operator. Labeled, solid arcs represent operators. Unlabeled, dashed arcs show preconditions which must be satisfied before the subsequent operator can apply. (a) general case, (b) GELL-MANN, (c) DALTON. Only the operators that lead to available tests have been retained in (b) and (c).

6. The operators should be used in the most efficient order that satisfies all other principles.

Non-redundancy and exhaustiveness (cf. 3, above) have been used by DENDRAL developers, Lindsay, Buchanan, Feigenbaum, and Lederberg (1980) as requirements for their structure generator.

GELL-MANN's search

According to these principles, GELL-MANN has been arranged in a three phase search (Figure 5abc), implemented in common lisp. Each phase generates a part of hidden structure and verifies it by specific evaluators.

In Phase I (Figure 5a) GELL-MANN searches for admissible classes of microcompounds. It postulates hidden object types in T (operator O_1 , "Create Micro Elements"), then the number of objects in a microcompound (Operator O_{2A} "Create Micro-Compounds"). For a given set of hidden types and a given number of elements in a microcombination, GELL-MANN creates all their combinations (Operator O_{2B} "Create Micro-Compounds"). Those combinations are often called bags. The same object can occur several times in a bag, but the order in a bag does not matter.

The isomorph-free bag generator uses the order of elements in T and creates each bag in that order, with possible repetition or omission of some elements.

GELL-MANN starts its search from one hidden object and keeps adding objects until a solution is found,

or the number N of objects in T reaches the number of objects in the input family, so that the model fails to simplify the input. Operator O_{2A} starts from two elements per bag, as one element would make the structure identical with a single part. O_{2A} increments the number M of elements in a bag by one. For a given N and M , a set of bags is admitted to the next phase when the number of bags is not smaller than the number of particles in the input family, but no greater than three times the number of observed objects. Indeed, we want different observable particles explained by different quark combinations. We also do not want to postulate a much larger number of hidden structures than is supported by the number of observed objects. The limits on quark combinations will cause the search to stop even if no model has been found. This makes sense, not because more complex models are impossible, but because the available data do not support speculations about them.

The next tasks, according to the dependency chart are to determine: (1) what attributes will occur in the model, (2) what attribute values are admissible. GELL-MANN uses all attributes provided in the input (operator O_3 "Create Properties"). The more attributes used, the more demanding is the evaluation. GELL-MANN starts Phase II by postulating candidate values for each attribute (O_4 , "Determine Possible Values"). Too many values would result in huge search

spaces. Too few values may exclude valid solutions. We turned to the observed objects for guidance. Let v_i be the largest absolute value exhibited at the observed layer for attribute P_i . GELL-MANN takes as admissible values V_i of P_i all positive and negative integers between $-v_i$ and v_i . In addition, GELL-MANN postulates rationals compatible with those in the input, and rationals with denominators equal to the number of hidden objects currently postulated per bag. For example, if three hidden objects are postulated per bag, values down to thirds are used.

The exhaustive search must try all assignments of values to hidden objects. However, such a search is typically too complex. For 3 quarks, 3 attributes, and 10 admissible values per attribute, a straightforward search would try 10^9 models. Can we eliminate invalid partial solutions? GELL-MANN's Phase II of search is an answer. GELL-MANN considers each attribute P_i separately trying all assignments of values in V_i to N elements. This is another application of isomorph-free bag generator, this time generating bags of size N . For each assignment of candidate values, GELL-MANN uses the additivity principle to compute the value of P_i for each combination generated by the first phase of the search. An assignment is admissible if a mapping *exists* from each observed particle to a microcompound with the same P_i value. Figure 5b depicts that phase. The output is typically a small set of admissible N -tuples of values for each attribute.

In Phase III, these separate solutions for each attribute are combined to form a solution which works across all attributes and all objects in T . At the same time, to enable evaluation, concrete mappings ψ are tried between particles and quark combinations (operator $O7$ "Map Micro Objects to Macro Level"). Phase III is depicted in figure 5c. If GELL-MANN finds a solution valid for the given N and M , it continues, searching for all solutions for the same N and M and halts. If it cannot find a solution, the search returns to phase one and increments N or M .

In Phase III, the isomorph-free generation faces the biggest challenge. Consider different mappings for proton p , the particle which opens the search depicted in Figure 5c. It could be assigned ten different bags for $N = 3$ (let us call the three quarks a, b, c) and $M = 3$. But, for instance, the bags $(a a a)$, $(b b b)$, and $(c c c)$ lead to isomorphic solutions. Only one of them should be considered. The situation becomes more complicated after partial solutions have been proposed, but some quarks are still indistinguishable from others by their attribute values. Here GELL-MANN's generator uses three principles: (1) use the order in which the elements in T are listed; (2) do not skip any elements; (3) each next element can be listed no more times than the previous one. For $N = 3$ and $M = 3$, only the bags $(a a a)$, $(a a b)$, $(a b c)$ will be generated.

GELL-MANN can search for a model of a single input family. However, many families of particles exist.

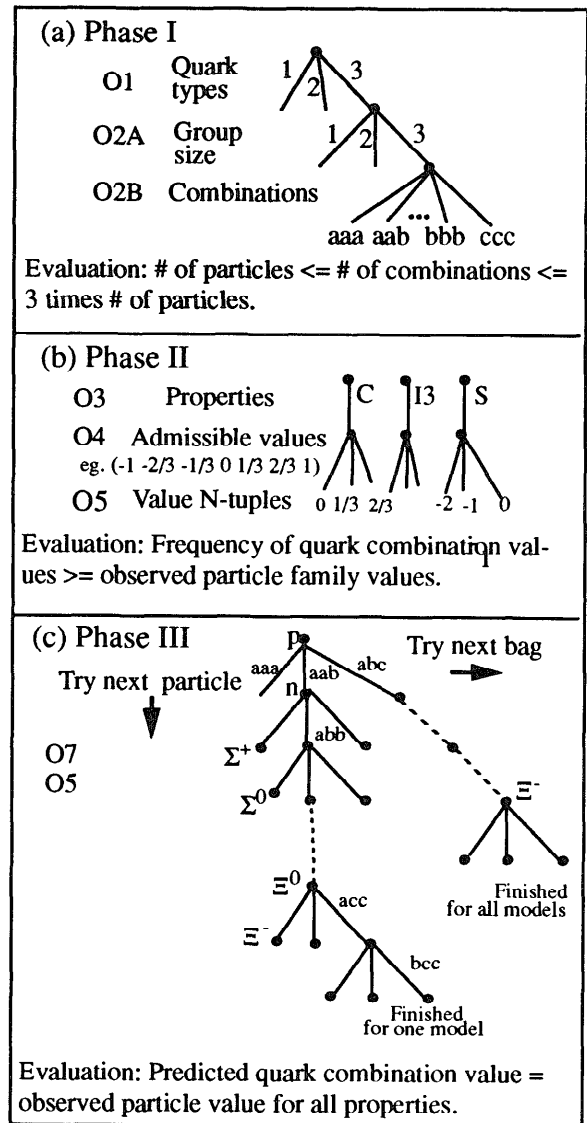


Figure 5: The three phases of search in GELL-MANN. Phase I generates quark types, bag sizes and quark combinations. Phase II generates attributes of quarks and values for each attribute. Phase III maps attribute values to quarks and quark combinations to particles.

Is a joint solution possible? Can it be reached incrementally? GELL-MANN handles multiple families by working with each in succession. For each next family it tries a solution based on the quarks used in the models that worked for all previous families, and adding new quarks only if necessary. Although the same set of quarks is used for all families, the bag size for each family can be different.

Incremental search. GELL-MANN proceeds to the second family, using the same three search phases. It first tries the known quarks, adding new ones only if

no solution has been reached for the existing ones.

The order of input families can be historical, but we can also try different orders of processing the same set of input families. In that case, we want GELL-MANN to seek the simplest global solution for all families.

Generalizations. Attributes such as spin can be combined by vector addition of quantum numbers. The rule of vector addition and other combinatorial rules can be plugged into GELL-MANN. Phase II can be eliminated altogether when attribute values are found by solving matrix equations (Valdes-Perez, Zytlow & Simon, 1993). This approach has been implemented by Valdes-Perez in YUVAL (Valdes-Perez & Zytlow, 1996). For small sets of values it turns out that GELL-MANN's generating and testing value combinations works faster than solving equations.

Results of Experiments

Early historical data. Initially, three families of elementary particles, postulated by Murray Gell-Mann, formed the theoretical basis for the quark model: hadron octet, meson octet, and baryon decuplet. Figures 2 and 3 present results of GELL-MANN's non-incremental search on the inputs of hadrons and mesons. Later, we applied our incremental search for the joint quark model to these three families, in all six permutations. Figure 6 depicts our experiment. Each solid arc is labeled with the particle family given as input to GELL-MANN. Each node in the tree, except for START, represents the output of GELL-MANN: the number of quark models found and the complexity of the quark models ($N * M$) in terms of the number N of quarks postulated and the number M of quarks per bag. GELL-MANN incrementally builds on the previous quark models on the direct path from the root.

The initial hadron octet run produced two models of complexity $3*3$ (three quarks in groups of three; Figure 2). Of these two models, one was postulated by the physicist Gell-Mann. The other was a new model. The meson octet run yielded one model of complexity $4*2$ (Figure 3), much simpler than the $6*2$ model accepted in physics. The baryon decuplet family yielded one model $3*3$, the model postulated by Gell-Mann. All these models are shown in Figure 6 as direct descendants of the START node.

The remaining two families of particles were added, one at a time, to the initial models. Building on the meson octet unconventional $4*2$ model, 5 models have been found for the baryon decuplet in the $7*3$ category, which is more complex than the standard $6*3$ model. The hadron octet family led to a $6*2$ model, but adding the baryon decuplet produced no solution, so this path was discarded.

Building on the unconventional hadron octet model (Model 2 in Figure 2), the baryon decuplet family produced no solution. For the standard model (Model 1 in Figure 2), a $3*3$ solution was found. When the meson

octet family was added to the remaining model, two solutions of $6*2$ were produced.

Building on the single baryon decuplet model depicted as the rightmost child of START in Figure 6, no new quarks were needed to explain the hadron octet. Then adding the meson octet produced the same two $6*2$ solutions found following the hadron octet path from the root.

Additional results. Expanding the models built for each of the three particle families, we added two additional families, the charmed mesons and charmed hadrons, also depicted in Figure 6. Each $6*2$ model has been expanded to a $8*2$ model for charmed mesons. When the charmed hadrons were added, only one $8*2$ model could be adapted to handle the new class. The output was one $8*3$ model, known by physicists as the standard model. We could not further expand our search because too few particles are known to contain the bottom quark, so the search is not constrained enough.

The other paths in the tree end with either no solution, or solutions of greater complexity than the standard model. This implies that in the space examined by our incremental search without backtracking, the standard model is indeed unique in its simplicity class.

As the complexity of the quark model grows, so does the size of the search space and the program execution. For instance, the initial run on the hadron octet found two $3*3$ solutions in approximately 4 sec of CPU time at 6MIPS. But more complex searches took days and even weeks, as indicated in Figure 6.

Many other results have been reported by Fischer & Zytlow (1990) and by Valdes-Perez & Zytlow (1996).

Conclusions and future work

We have presented a theoretical framework for the discovery of hidden structure and have demonstrated how the discovery system GELL-MANN fits that framework. We presented results of a large-scale exploration in the space of quark models. The same operators, similar knowledge representation, and the same elements of search are used by many systems that discover hidden structure. This unification suggests a unified computer system which would be able to discover hidden structure in different domains.

Two problems must be solved before we can build an autonomous system capable of discovering hidden structure in various domains. First, we lack a general algorithm that would use domain knowledge on the observational level to set up the search for hidden structure. Using the experience accumulated by construction of several systems, it is relatively easy to *manually* generate all elements of search: to select operators, to define operator instantiation and evaluation criteria, and to organize them in a simple, "bottom-line" system: simple in structure and able to search exhaustively, yet entirely inefficient. Thus, the second problem is concerned with the automated generation of

