

## Clustered Disease Findings: Aspects of Expert Systems

*Charles Turner, PhD, Michael Lincoln, M.D., Peter Haug, M.D., Homer Warner, M.D.,*

*John Williamson, M.D., Neal Whitman, PhD and James Buchanan, M.D.*

*Departments of Psychology, Medical Informatics, Internal Medicine,  
and Family and Preventive Medicine*

*University of Utah*

*Salt Lake City, Utah, USA*

A variety of expert systems have been developed to provide assistance in making medical diagnoses. These systems combine various medical findings from a patient to identify a plausible diagnosis to account for these findings. These findings typically occur in natural groupings or "clusters." That is, several related findings almost always occur together for some diseases; they are almost always absent for other diseases. Efforts to accomodate these "clustered" findings in expert systems lead to two types of difficulties. Systems that use all of the related findings can be unreliable since they produce overconfident decisions. However, systems that use only part of the related findings become too sparse, and they are unable to explain all of the data for real patients.

Clustered knowledge representations can solve both of these problems. An expert system can manage natural groupings by forming hierarchical "clusters" of findings; these clusters describe higher medical concepts such as pathophysiological processes and diseases. Clustered representations can be useful in knowledge engineering of expert systems. However, when medical experts attempt to create clustered knowledge representations, they can be subject to heuristic biases. In this paper, we examine some alternative methods that can be used to identify hierarchical clusters of findings appropriate for knowledge engineering tasks. The present paper also compares the use of hierarchical clusters in knowledge engineering to recent developments of new models in the cognitive sciences.

### 1. INTRODUCTION

Medical expert systems such as HELP (Health Evaluation through Logical Processes, Iliad, QMR (Quick Medical Reference) and MYCIN have provided a technology to support the medical decisions which are made in caring for patients [1]. Until recently, most medical expert systems made little or no use of clustered knowledge structures. Bayesian systems, such as HELP, process individual findings sequentially. Such systems traditionally avoid findings that are conditionally dependent because conditional dependence causes the system to produce overconfident diagnoses [1,2]. Yet, much medical information has a high level of conditional dependency. For example, the egophony heard on chest examination over an area of lung consolidation is not conditionally independent from rales heard in the same area [3]. That is, a patient is likely to have either both or neither of these symptoms.

Writers of Bayesian HELP frames traditionally have deleted many such conditionally non-independent findings, retaining only "key" findings. The resulting frames contain a "sparse" set of findings; the frames can have good diagnostic accuracy in some settings, but may not accurately represent the richness of expert knowledge. Rule based systems, such as MYCIN, can create rules encompassing simple clusters [1]. But complex concepts require increasingly numerous, complex, and unnatural sets of rules. Complex knowledge representation is difficult with rules, and such rules do not conform with the way that clinicians organize diagnostic knowledge. By definition, tree-based expert algorithms process one test or

decision node at a time and avoid clustered knowledge representation.

Recently we proposed that clustered concepts represent the fundamental basis of organizing medical knowledge [3,4]. Iliad (not an acronym) is a computerized expert system that explicitly represents this approach to knowledge representation. The system contains complex, clustered pathophysiologic representations as a central part of the knowledge model. Iliad is a derivative version of the HELP system. This version is a microcomputer based expert system designed to assist medical students to develop diagnostic problem solving skills. For this teaching tool, we needed to create an expert system that was able to handle the multiple manifestations of pathophysiologic processes. We are currently in the process of evaluating the effectiveness of this teaching tool on the diagnostic skills of third year medical students.

As part of our evaluation of the clustering model, we have examined the reliability of clustered and unclustered versions of pulmonary knowledge frames [2]. As with the HELP system, Iliad uses a frame-based, sequential Bayesian expert system. Sequential Bayesian systems require that diagnostic findings used in the frames be conditionally independent. In other words, the findings must not tend to co-occur in the same patient. However, many medical findings are conspicuously non-independent. As an example, the findings of rales, bronchial breath sounds, and dullness to percussion tend to co-occur in pneumonia. When this assumption of conditional

---

**Title of Cluster: Lung Consolidation**
**Findings And Decision Logic**

For "Lung Consolidation" to be present, the following combination of findings must occur:

**Either**

Chest x-ray:

Lung infiltrate with alveolar pattern and air bronchograms

**OR**both of

Coarse crackles

Bronchial breath sounds

**And**one of

Egophony

Increased vocal fremitus

Dullness to percussion

Whispered pectoriloquy

---

Figure 1: A sample clustered frame using Boolean logic

---

**Title of Frame: Pneumonia**

The a priori (prevalence) of "Pneumonia" in hospitalized settings is .025.

<u>FINDINGS:</u>	<u>Probability in Patients with the Disease</u>	<u>Probability in Patients without the Disease</u>
Lung consolidation	.99	.07
Signs of systemic infection	.90	.20
Pleuritic chest pain	.25	.02
else		
Pleural effusion	.05	.01
Hypoxemia	.40	.10
Dyspnea	.40	.10
Cough with sputum	.90	.15
else		
Cough	.90	.25

---

Figure 2: A sample frame using Bayesian logic

independence is violated, Bayesian systems suffer from diagnostic overconfidence [1,2].

Despite our best efforts to minimize conditional non-independence, overconfidence has been a perennial problem in HELP and Iliad. Before using clusters, we attempted to prune out findings of secondary importance that tended to co-occur with "key" findings. This strategy resulted in disease frames containing a "sparse" set of findings. Sparse

frames were unsatisfactory for several reasons. First, users often attempted to enter seemingly important findings that had been pruned from the frame. They had no way of knowing which "key" findings had been retained. When the diagnostic probabilities did not change in response to these "pruned" findings coming true, users assumed the frames were faulty. Second, sparse frames had limited usefulness as teaching tools. These frames did not reflect the rich patterns of disease presentation that students must learn to recognize.

Finally, sparse frames did not completely eliminate conditional non-independence, so Iliad was still overconfident, substantially impairing user confidence in Iliad's diagnoses. To solve these problems, we adopted a new model of disease frame knowledge representation called "clustering" [2,3,4].

A cluster is a medical concept described by a subset of conditionally dependent patient findings. These clustered concepts usually describe pathophysiologic entities and often embody useful teaching paradigms. A typical, rich cluster, lung consolidation, is reproduced in Figure 1. The individual findings in clusters can be patient history items, physical exam findings, or test results. By definition, findings in a cluster often co-occur in patients.

When students master the pattern recognition skills needed to detect clusters, they can make decisions about many diseases. For example, lung consolidation is a pathophysiologic process found in many types of pneumonias. The lung consolidation cluster can be used in combination with other clusters. A cluster denotes the focus of a number of causal connections. These causal connections relate disparate findings such as rales, and pectoriloquy and chest x-ray with lung infiltrates into recognizable patterns.

We have developed other clustered frames to represent other pathophysiologic processes such as "signs of systemic infection," "pleuritic chest pain," and "hypoxemia". The decisions made from each of these clusters can be combined within a Bayesian frame into a hierarchical decision to estimate the probability of a disease process such as pneumonia. Thus, the Iliad system contains both Boolean and Bayesian frames. Decisions involving conditionally dependent findings are made using Boolean principles while decisions involving independent processes are made using Bayesian logic. Psychological research indicates that facts organized into causal relationships (clusters) are easier to understand than isolated facts [5,6,7,8]. This research suggests that clusters of medical facts organized by causal relationships should be easier to understand than unorganized facts [5,6].

Clusters not only assist understanding of complicated medical concepts, but they offer ways to improve the diagnostic accuracy of Bayesian expert systems. When findings in a disease are causally related they usually co-occur in the same patient. Stated in probabilistic terms, the findings are conditionally non-independent. Because conditionally non-independent findings make Bayesian expert systems overconfident, a solution is required to overcome the problem of non-independent findings. To the extent that causal relationships are common in medicine, the problem of overconfidence is pervasive.

Although clusters can solve the overconfidence problem in Bayesian expert systems, cluster development is difficult for three reasons. First, cluster development is time consuming because experts must be used both for development and validation of the clusters. Second, a procedure must be developed to establish the validity of the clusters. Because clusters are not often explicitly stated in the traditional medical

literature, this literature is not a source to discover or validate clusters. Alternative sources of knowledge, such as the knowledge base in QMR, may contain mathematical relationships among findings that allow us to discover and validate clusters. Given this situation, one way to validate clusters would be to tap the implicit knowledge of experts in an expert review of the clusters. We have shown that experts can create clusters [4] but they are frequently vulnerable to heuristic biases [1].

We reasoned that expert systems that are rich in findings would also contain implicitly structured knowledge [3]. We developed a procedure to identify these implicit clusters within expert systems. As part of this research we compared clusters developed explicitly for Iliad to "implicit" clusters found by examining patterns of evoking strengths in another system QMR (Quick Medical Reference). The process of detecting implicit clusters in QMR required factor and cluster analysis of evoking strengths associated with the clustered findings. Our analysis focused on pulmonary diseases because both Iliad and QMR have comprehensive pulmonary knowledge bases. We continue our analysis of the implicit clusters in QMR to those explicitly developed previously for Iliad. Our initial findings indicated that clustered knowledge models do appear in QMR even though the system was not explicitly developed using clustered knowledge representations system. Since the natural structure of human knowledge appears to be clustered, the use of clusters may save time building knowledge bases, increase the accuracy of diagnosis, and provide a better model for expert systems such as Iliad, QMR, and HELP (Health Evaluation through Logical Processes).

In contrast to the sparse frames in Bayesian expert systems, QMR disease profiles are rich in conditionally dependent findings. The decision model in QMR is not based on Bayes theorem, and the system does not require an assumption of conditional independence of findings. Consequently, the dependencies among findings in QMR do not necessarily result in overconfident diagnoses. These disease profiles, rich in conditionally dependent findings, are a robust knowledge base containing information that can be examined mathematically for the presence of clusters (i.e., conditionally dependent findings).

The original Iliad pulmonary medicine clusters were produced before we completed the mathematical analysis of QMR. An expert in pulmonary medicine produced these Iliad clusters spontaneously. He relied only on his knowledge of pulmonary diseases and on the twin goals of reducing conditional non-independence and increasing the teaching value of Iliad frames. Another paper in this series describes the important enhancement in diagnostic accuracy that resulted from the introduction of clustering. In this paper, we compare the clusters found in QMR to those developed for Iliad and explore the implications of the common ground between them. We discuss how an analysis of other disease areas in QMR could be used to predict and guide further cluster development in

Iliad. Finally, we discuss the teaching advantages offered by clustered knowledge representation in Iliad.

We hypothesized that the inherently clustered nature of medical knowledge was implicitly reflected in QMR's knowledge base [3]. The purpose of our earlier research was to determine whether we could extract these clusters by analyzing the mathematical relationships between findings and disease profiles. The results of the cluster analysis included neoplasm, pneumonias, obstruction, restriction, and pleural diseases. We thought it likely that these groups each contained several clusters. For instance, lung consolidation might be a cluster we would expect to find among the pneumonias disease group. Some clusters, such as lung consolidation, occur across multiple diseases within a disease group as well as in diseases from different disease groups. Lung consolidation can occur in pulmonary neoplasms, as a consequence of endobronchial obstruction and distal collapse and filling of lung parenchyma. Hypoxemia is another example of a cluster seen in many diseases. Other clusters, such as Solitary pulmonary nodule, occur in a single disease group [3].

Clusters have important implications for representation of human knowledge in expert systems. The finding of clustered knowledge representations in ostensibly non-clustered systems like QMR is evidence for the universality of the cluster concept. Humans naturally cluster knowledge, in part because of the principle of "bounded rationality." This means human beings are inherently limited in their capacity to simultaneously hold in mind and process multiple data items. Simplifying heuristics are essential if we are to deal with complex problems. Teachers of medicine have perhaps instinctively realized this and focused much of their efforts on teaching students to recognize patterns of disease findings, especially where such patterns were mechanistically related. Professors call these patterns "pathophysiologic mechanisms." We call them clusters. Few expert systems use clusters, yet many (including early versions of HELP and Iliad) go through great trouble to consciously work around them. It is not only simpler to include these natural human knowledge constructs in a computerized expert system, but it improves the accuracy of the results.

These new models of knowledge representation in Iliad are consistent with new developments in the cognitive sciences. These models characterize memory representation describe knowledge as network models or parallel distributed process (PDP) models.[7,8] The PDP models provide a useful framework for representing the hierarchically clustered knowledge structures used in medical knowledge engineering.

The development of computer algorithms to solve medical problems can provide important contributions to our understanding of human cognitive processes. Great breakthroughs have occurred in our understanding of cognition as a result of development of high level computer languages and powerful algorithms. However, earlier versions of human memory relied heavily upon the computer metaphor described by von Neumann's model. This metaphor

characterized human thought as occurring in a sequence of limited capacity, information processing steps. This sequential model seems to be very inadequate in describing much of human cognition. Advocates of the Parallel Distributed Processing models note that humans sometimes seem so much smarter than machines. We are not quicker or more precise in the decisions that we make. A computer can perform 1000's of successive computations in the same time required for a human to make a single decision (or computation). However, people seem to be far better than computers in perceiving relationships (e.g., objects in a natural context), at understanding language, and in making contextually appropriate actions.

That is, most humans can quickly handle large arrays of information (e.g., in perceptions of visual fields). Suppose that the individual processed all of the elements in the array in a sequential manner. Physiological limitations in the human nervous system greatly restrict the rate at which a single decision can be made (i.e., neuronal transmission can occur). Because each successive step in the decision process would be delayed until previous decisions were made, a sequential string of decisions would be very slow. Recent views characterize human thought as occurring in parallel. That is, the computational framework for modeling human cognitive processes can be more appropriately characterized as a set of the simultaneously occurring processes. From this perspective, one of the most important questions to answer in a description of the cognitive process is the nature of the relationships among the processing units.

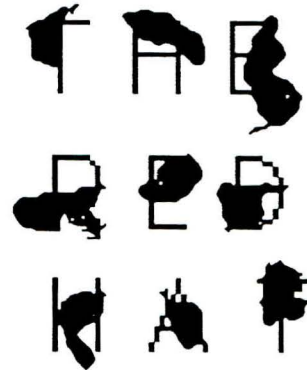


Figure 2. Ambiguous letters (see rum, mcll).

The problem presented in Figure 2 is a demonstration of one of the features of parallel processing. The reader is encouraged to try to read the material in the figure. Nine ambiguous letters are presented together in the figure with each letter partially obstructed. However, the letters can be easily recognized as the three words: THE, RED, HAT. A paradox arises in this example. Few of the letters can be identified easily by themselves. Yet, the combination of the letters makes it possible to identify all of the letters with relatively little effort. That is, the presence of other ambiguous letters makes it possible to disambiguate the remaining letters. Still, all of the individual letters are partially obscured. How

can unrecognizable letters help us to recognize other unrecognizable letters? Apparently, the presence of each letter fragment imposes constraints upon the possible values that can be taken by the other letters. Thus, the restrictions on the identity of each letter imposes a constraint or a context for the processing of the other letters. This example demonstrates how perceptual processes reflect patterns of interdependent processes. If all of the mutual constraints were processed in sequence, the reader might take a very long time to recognize the three words. However, if the process occurs simultaneous manner, then the process can procede rapidly.

We store knowledge in terms of structures variously known as frames, scripts, or schemata. Such knowledge structures are assumed to be the basis for comprehension. However, most typical everyday experiences can not be assigned to a single knowledge frame. Rather understanding occurs through the combination and interaction of several knowledge frames. Thus, knowledge representation involves the distribution of knowledge across a set of frames which are concurrently acting to influence comperehension of events. Similarly, expert system such as Iliad and QMR contain representations of medical knowledge with a mixture of implicit and explicit relationships distributed across a number of elements in the system. The model of human cognition proposed in the Parallel Distributed Processing models is generally consistent with the knowledge representation contained in expert systems such as Iliad and QMR.

#### REFERENCES

- [1] Weinstein MC, Fineberg HV. Clinical decision analysis. 1st ed. W. B. Saunder Company, 1980; 156-158.
- [2] Yu H, Haug PJ, Lincoln MJ, Turner CW, Warner HR. Clustered knowledge representation: Increasing the reliability of computerized expert systems. Symposium on Computer Applications in Medical Care 1988; 12: in press.
- [3] Lincoln MJ, Turner C, Hesse B, Warner H, Miller R. Discovering clustered disease findings: Prospects for enhancing expert systems. Symposium on Computer Applications in Medical Care 1988; 12: in press.
- [4] Lincoln MJ, Haug PJ, Yu H, Turner CW, Warner HR. Expert biases prevent accurate estimation of population statistics for clustered disease frames. International Symposium of Medical Informatics, 1989, in press.
- [5] Estes WK. Toward a framework for combining connectionist and symbol-processing models. Journal of Memory and Language 1988; 27: 196-212.
- [6] Gluck MA, Bower GH. Evaluating an adaptive network model of human learning. Journal of Memory and Language 1988; 27: 166-195.
- [7] Rumelhart DE, McClelland JL and the PDP Research Group. Parallel distributed processesing: Explorations in the microstructure of cognition. Volume 1: Foundations. Cambridge, MA: Bradford, 1986.
- [8] McClelland JL, Rumelhart DE and the PDP Research Group. Parallel districuted processesing: Explorations in the microstructure of cognition. Volume 2: Psychological and Biological Models. Cambridge, MA: Bradford, 1986.