

## Iterative Design of Case Retrieval Systems

Eric K. Jones and Aaron Roydhouse

Department of Computer Science  
Victoria University of Wellington  
P.O. Box 600, Wellington, New Zealand  
{Eric.Jones,Aaron.Roydhouse}@comp.vuw.ac.nz

### Abstract

We argue that the process of designing effective systems for case retrieval should not be divorced from the process of evaluation. Developing a suitable representation vocabulary for indices, and suitable mechanisms for case retrieval should proceed in an iterative fashion, guided by the joint application of sound AI principles, domain expertise, and empirical evaluation.

We lay out the key components of a case retrieval system and suggest procedures for evaluating each component. We illustrate our claims using the MetVUW Workbench, a system for intelligent retrieval and display of historical meteorological data.

### Introduction

Case retrieval is an important component of case-based reasoning. Design of a case retriever requires addressing the *indexing problem* (Domeshek 1991; Schank 1982; Schank & Fano 1992): the problem of accessing appropriate cases using partial descriptions of their content. In this paper, we argue that solutions to the indexing problem should be developed by an iterative process of design and evaluation, and that this process is to a significant extent an experimental and empirical enterprise. Evaluation should not be carried out only after completion of a system, but should serve as an integral part of the development process.

We draw on our experience with the design and implementation of MetVUW Workbench, a system for intelligent retrieval and display of historical meteorological data (Roydhouse *et al.* 1993; Jones & Roydhouse 1994). The system is intended to serve as a “memory amplifier” for meteorologists that allows them to rapidly locate historical situations of interest. In particular, weather forecasters should be able to use the system to quickly retrieve past cases that are similar to the current situation in meteorologically significant respects, providing them with an additional source of information to supplement the output of numerical models.

Each case is a slice of time for which meteorological data is available. The data available to us includes satellite imagery stored both in digital form and on laser disc, a document archive, and numeric fields from the European Center for Medium-range Weather Forecasting (ECMWF). Examples of numeric fields include pressure, temperature, relative humidity, wind speed, and relative vorticity, all of which are available at 14 different levels of the atmosphere.

We are working with a portion of this numeric data that covers the Australasian region.

Cases are indexed and retrieved in terms of high-level descriptions of the content of their numeric fields. Numeric fields are available at 12-hour intervals. We currently possess 3.5 years (1.6Gb) of data, which allows us to build a case base containing some 2500 cases. The experiments described in this paper, however, were carried out on a 550-case subset of the data. In the near future, we plan to obtain an additional 10 years of data, allowing us to expand the case base to 10,000 cases (6Gb).

MetVUW Workbench is intended to serve as an intelligent assistant that empowers a human expert. For example, forecasters at the Meteorological Service of New Zealand currently engage in a “map session” every day at 10am, at which they compare competing prognoses to arrive at a consensus forecast. MetVUW Workbench is designed to provide valuable input to discussions such as these. We envisage that forecasters will follow three steps when using the system:

1. Determine the meteorological processes and systems that appear to be driving the current weather situation.
2. Formulate a query to MetVUW Workbench in terms of these processes and systems.
3. Analyze and summarize the cases that the system returns, or display any of the data associated with particular retrieved cases.

MetVUW Workbench is a case retrieval system as opposed to a full case-based reasoner. Case retrieval is viewed as a process of comparing queries (sometimes called *index probes* (Domeshek 1991)) against explicitly represented *index labels*. Queries identify particular features of the current situation that appear to be meteorologically significant. Index labels are representations of high-level features of the weather situation in the case. As far as possible, index labels are extracted automatically or semi-automatically from the raw data. We are currently focusing on features such as local minima and maxima that are easy to derive automatically from ECMWF fields. In the medium term, we plan to address the much harder problem of extracting high-level features from satellite imagery.

The remainder of this paper is structured as follows. We first describe the architecture of MetVUW Workbench and identify four kinds of evaluation for MetVUW Workbench

and other case retrieval systems with this architecture. Next, we sketch MetVUW Workbench's representations for index labels and describe similarity assessment. We then introduce important issues in case selection in MetVUW Workbench and report some preliminary evaluations of algorithms for case selection, explaining how they have contributed to a process of iterative design and evaluation. We conclude with a brief discussion of future work.

## System Architecture and Evaluation

The appropriate role of evaluation in the design of a case retrieval system is in part a function of the architecture of the retrieval engine. In this section, we describe the architecture of MetVUW Workbench, and then use this description to identify four kinds of evaluation that are appropriate for this architecture.

### The architecture of MetVUW Workbench

Figure 1 depicts the architecture of MetVUW Workbench. The user interacts with the system via a *query constructor*, which allows meteorologists to formulate queries in terms of high-level descriptions of the meteorological systems and processes that appear to be relevant in the current situation. Users can describe high-level features of past weather situations. These features include low and high pressure systems, ridges, troughs, and jet streams. Eventually, it should also be possible to represent properties of these objects as they develop over time, such as the track of the pressure minimum of a low-pressure system.

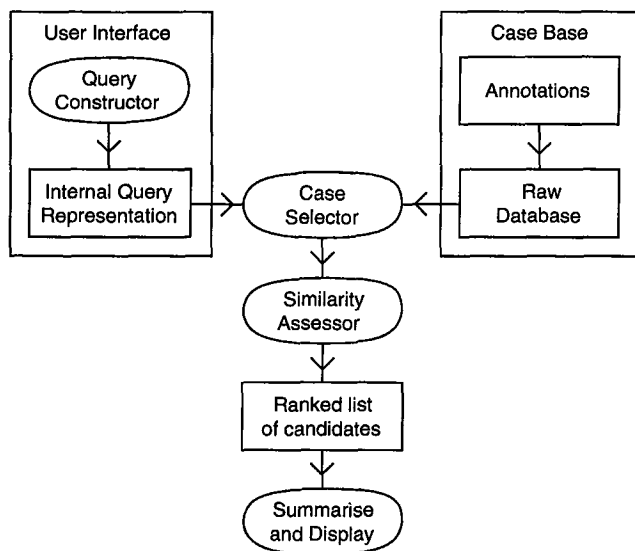


Figure 1: Architecture for case retrieval

Users construct queries by drawing representations of meteorological phenomena on a rectangular canvas that depicts a region of the globe. These phenomena are represented by points, lines, and regions, using a notation similar to an abstract weather map. This kind of interface was chosen so that the components of queries can be entered quickly and easily using a notation that is familiar to forecasters.

A sample query is shown in figure 2. The query requests a low-pressure system with an upper-level trough situated to the northwest at a distance of 400–600km. The “+30mb” marking indicates that there is a 30mb pressure difference between the contour so labeled and the lowest pressure at any point within it.

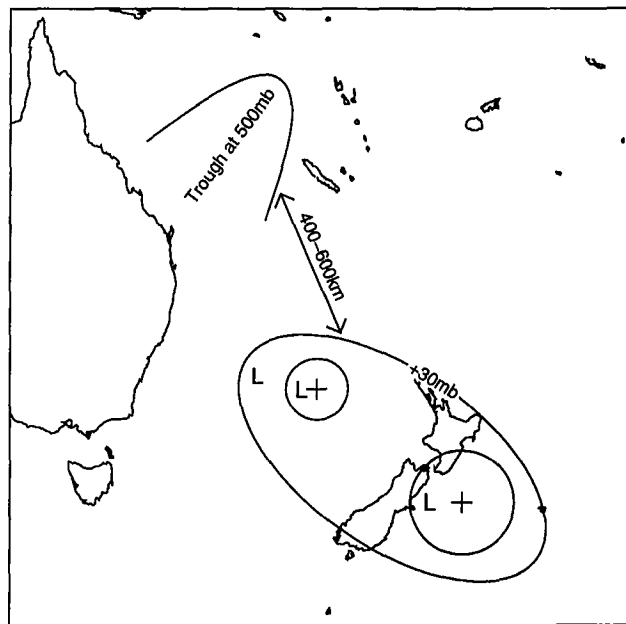


Figure 2: A query

After a query is constructed, a symbolic representation of it is passed to the *case selector*, which uses key features of the query as indices for retrieving past situations or cases. The purpose of case selection is to allow the time it takes to retrieve a given number of cases to scale sub-linearly with the size of the case base. In particular, case selection cannot simply scan the entire case base.<sup>1</sup>

Case selection presupposes a suitable memory organization for index labels, and a retrieval mechanism that exploits this memory organization to efficiently access cases whose labels match a query. As we discuss in more detail below, MetVUW Workbench uses an advanced relational database called Postgres to implement case selection (Stonebraker & Rowe 1986). This “heavy duty” machinery was chosen so that the system can readily scale up to a very large case base. Case selection is efficient because it employs a number of sophisticated indexing strategies provided by Postgres, including R-tree indexing of spatial data (Guttman 1984).

Once a collection of past cases has been retrieved, they are passed to the *similarity assessor*, which uses a knowledge-intensive partial matching process to rank them according to how well they match the query. Those cases whose match quality falls below some threshold are discarded.

Case retrieval thus proceeds in two stages. First, case selection uses efficient retrieval machinery—in MetVUW Workbench, a sophisticated relational database—to quickly

<sup>1</sup>We assume serial computer hardware.

identify a manageable number of candidate cases that need to be explicitly considered. Second, these cases are further filtered and ranked using a more costly but more accurate process of similarity assessment. This combination of case selection and similarity assessment facilitates fast, high quality retrieval.

It is important to distinguish two notions of retrieval that arise in this architecture: case retrieval itself, and what we are calling case selection. Each has its own representation vocabulary for index labels. Case selection employs simple labels that meet the often strict constraints of its efficient retrieval mechanisms. In contrast, case retrieval employs a richer representation vocabulary that also supports similarity assessment.

For example, Postgres' R-tree indexing mechanism allows efficient retrieval of rectangular regions oriented along lines of latitude and longitude that are spatially related in some way to a similarly oriented rectangular region in a query. Spatial relations include contains, contained-in, and overlaps. Because R-tree indexing is available, queries in MetVUW Workbench employ rectangular bounding boxes of low and high-pressure systems as index labels for case selection. Moreover, queries request cases containing systems that bear one of Postgres' build-in spatial relations to the corresponding system in the query. In contrast, similarity assessment uses arbitrary polygons to represent low and high pressure systems, and computes a more complex (and more informative) relation of graded match between corresponding systems in the query and the case.

#### Four kinds of evaluation

The architecture of MetVUW Workbench suggests evaluating the system along two dimensions. Along the first dimension, it is possible to separately evaluate either case selection or similarity assessment components. Along the second dimension, one can evaluate either the effectiveness of each component's representations for index labels or the utility of its algorithms. We now examine each of the four kinds of evaluation that these possibilities give rise to.

Similarity assessment must be evaluated with respect to the requirements of end users of the system. Recall that similarity assessment produces a ranked list of cases for further processing. Evaluation of similarity assessment focuses on the following questions:

1. *Representation vocabulary.* Do the index labels encode representational distinctions that matter to the user?
2. *Algorithms.* Do computations of graded match between queries and cases produce an appropriate ranking of the cases?

Answering these questions typically requires detailed consultations with users and domain experts. We do not consider evaluation of similarity assessment further in this paper.

Case selection, in contrast, can often be evaluated without consulting users or experts, assuming that appropriate procedures for similarity assessment are already in place. Case selection exists for efficiency reasons only: if similarity assessment were adequately efficient, or scaling considerations unimportant, then case retrieval could be implemented

simply by applying similarity assessment to the entire case base and returning those cases that well match the query. Let us call the list of cases produced in this way the *benchmark list* (for a given query).

Benchmark lists of cases produced by similarity assessment provide a convenient standard for evaluating case selection. A procedure for case selection cannot be faulted if when using this procedure, case retrieval successfully returns exactly the benchmark list.

Like similarity assessment, case selection can be evaluated with respect to both index representations and algorithms:

3. *Representation vocabulary.* Are the index labels sufficiently expressive? Are they liberal enough to permit retrieval of most cases judged to be a good match by similarity assessment? Are they conservative enough to filter out most of the irrelevant cases?
4. *Algorithms.* Does the time complexity of case selection scale well with the size of the database?

In contrast to similarity assessment, however, questions of efficiency are more central to case selection. This is not to imply that efficiency is not also important for similarity assessment. However, the quality of case selection directly impacts the efficiency of the entire retrieval system: if case selection retrieves many irrelevant cases, then similarity assessment will take much longer, because it must scan all of the cases that case selection returns.

Two standard evaluation metrics from the information retrieval literature are useful in evaluation of case selection (Jacobs 1992):

- The *precision* of an information retrieval system is the percentage of items retrieved in response to a query that are good answers to the query.
- The *recall* of an information retrieval system is the proportion of good answers to a query that the system retrieved, expressed as a percentage of the total number of good answers that exist in the entire database.

In the context of case selection, precision measures the percentage of candidates returned by case selection that are contained in a query's benchmark list. In contrast, recall measures the proportion of candidates returned by case selection that are in the benchmark list.

For a wide variety of retrieval problems, there is a tradeoff between recall and precision. A given algorithm for case selection may be able to achieve a high level of recall, but in general this requires sacrificing precision. The extreme case is to retrieve the entire case base, which achieves 100% recall at the price of very low precision.

It follows that a good measure of the overall quality of a case selection algorithm is the precision it achieves for a given (acceptable) level of recall. This measure is especially useful for our purposes because it correlates closely with the time efficiency of the overall case retrieval system.

In the next section, we outline some of MetVUW Workbench's representations for index labels and algorithms for case selection and similarity assessment. We then report on some initial experiments concerning the evaluation of case selection in MetVUW Workbench.

## Representing Queries and Index Labels

Queries are constructed using graphical objects such as points, vectors, and regions, which denote high-level meteorological features. There are three kinds of features: static, dynamic, and relational. Static features describe meteorological phenomena at particular points in time. Examples include the center and extent of a low pressure system, the orientation of a ridge or trough, and wind direction. Dynamic features encode properties of these phenomena as they vary over time, such as the path followed by a low pressure system or whether a trough is intensifying. Relational features encode spatial constraints between features.

Corresponding to each graphical object in a query there is an underlying symbolic representation that is used in case retrieval. Index labels are also constructed using the same representation vocabulary. In this paper, rather than attempting a survey of the full range of available representations, we restrict our attention to representation of low pressure systems.

Our representation for low pressure systems encodes key aspects of their shape, geographical extent, and intensity that experts consider to have important meteorological implications. Simple low pressure systems—that is, systems with a single pressure minimum—have few interesting shape properties; they are generally close to circular or elliptical. Only the geographical extent of these systems and the location of their pressure minima are considered meteorologically significant by forecasters.<sup>2</sup>

Although simple low pressure systems have uninteresting shapes, they are often clustered together into groups of two or more that overlap; see figure 3. The entire region covered by a cluster is likely to experience adverse weather, especially near centers of low pressure. It is therefore important to encode the shape of clusters in a way that makes clear the location of regions of low pressure.

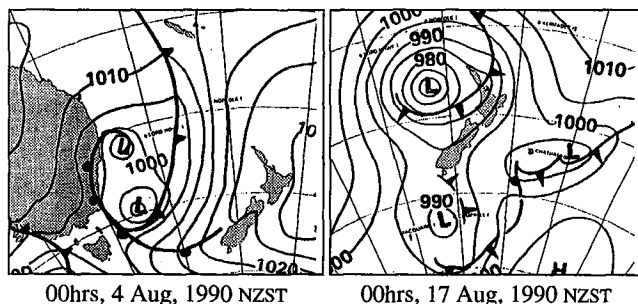


Figure 3: Two complex low pressure systems.

A cluster is represented as a tree whose leaves are simple low pressure systems and whose internal nodes are smaller clusters. One cluster is an ancestor of another in the tree if and only if there is a pressure contour associated with the ancestor that completely encloses the contours of its descendant. This tree structure is used during similarity assessment to compute a measure of the degree of structural match between a query and a case.

<sup>2</sup>High pressure systems are another matter, as we discuss in (Jones & Roydhouse 1994).

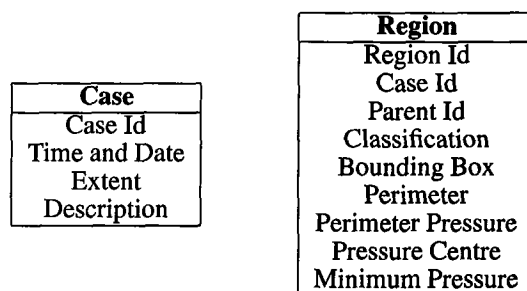


Figure 4: Relation scheme for low-pressure systems.

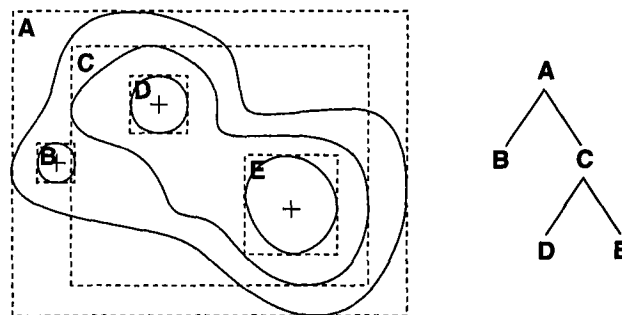


Figure 5: Representation of a typical cluster.

All of this information is automatically extracted from the numerical ECMWF data, and stored in the Postgres relational database using the relation scheme depicted in figure 4. Each case is allocated a unique identifier, as is each region of low pressure. Clusters are built up out of individual regions of low pressure; the tree structure of a cluster is encoded using the "Parent Id" field of the "Region" relation. Other high-level features of cases such as ridges, troughs, and high-pressure systems are encoded using further relations similar to "Region."

Regions of low pressure are used to represent both the simple low pressure systems at cluster leaves and the enclosing contours that make up the internal nodes of the tree. The algorithm that automatically extracts these contours from the ECMWF data selects closed contours that are as large as possible and yet preserve cluster structure. For each region, we store the location and mean sea-level pressure of the pressure minimum, together with a high-resolution polygon that represents its outer boundary. Case selection works with a rectangular bounding box that approximates this polygon, but similarity assessment can access the full polygon to compute such quantities as degree of overlap with a region in the query. The mean sea-level pressure at the polygonal boundary is also stored. The representation of a typical cluster is depicted in figure 5.

## Similarity Assessment

The similarity assessor is handed a query and a set of cases retrieved by the case selector, and sorts the cases by their degree of partial match to the query. Cases that match too poorly are discarded.

Developing appropriate representations and algorithms for similarity assessment requires exploring a space of de-

sign decisions. It is useful to divide similarity assessment into two subtasks: assessing similarity between individual low pressure regions in a query and a case, and using cluster structure to combine those assessments into an overall assessment of goodness of match. We now consider each of these subtasks in turn.

### Assessing similarity of individual low-pressure regions

For purposes of pairwise matching, similarity assessment makes no distinction between regions internal to a cluster and the simple low pressure systems at the cluster's leaves. To compute the similarity between a low pressure region in a query and a corresponding region in a case, MetVUW Workbench computes the following five parameters.

1. *Relative location of pressure minimums.* A measure of the displacement between the pressure minimums of a region in the query and a region in the case is computed. Longitudinal (east-west) displacements are usually of less concern than latitudinal (north-south) displacements, because weather systems tend to naturally progress from west to east, and small displacements in this dimension can be attributed to infrequency of sampling. Latitudinal displacements, on the other hand, are penalized more heavily because the observed weather at a given location is considerably affected by such deviations, and because the physical behaviour of weather systems is strongly conditioned by their latitude.
2. *Aspect ratio.* The aspect ratio of a low pressure system is the ratio of its lengths of the sides of its bounding box. Aspect ratio gives a cheap measure of the shape of a low pressure region.
3. *Density.* We define the density of a low pressure region to be the ratio of the area of its high-resolution polygonal approximation to the area of its bounding box. Density gives a second inexpensive measure of region shape.
4. *Aligned overlap.* The pressure minima of two low pressure systems to be compared are aligned, then the ratio of the area of the overlap between their bounding boxes to the combined area they cover is computed. This ratio gives a third way of computing similarity of shape.
5. *Area.* The three previous parameters all facilitate comparisons of shape independent of a region's size. However, the relative size of low pressure regions to be compared is also important, so the area of the high-resolution polygonal approximation to each region is also considered as a separate parameter (even though it is also used to compute density).

These parameters are combined to compute a goodness of match between a pair of low pressure regions. Eight numbers are computed in all: the degree of match on relative location of pressure minima, the degree of aligned overlap, and the aspect ratio, density, and area of each of the regions to be compared. The goodness of match is computed as the sum of the first two parameters together with the ratios of the smaller to the larger value for each of the other parameters, using appropriate weights for each summand. Similarity assessment does not currently consider the intensity of low pressure systems; this needs to be fixed.

### Assessing similarity of clusters

MetVUW Workbench computes an overall goodness of match between clusters in a query and a case by assessing both the goodness of match between individual low pressure regions and the degree of structural match between the clusters. Exactly how these assessments should trade off against one another is a difficult issue, and we are still experimenting with different alternatives. MetVUW Workbench currently employs a form of breadth-bounded depth-first search with successor ordering (Pearl 1984) that appears to strike an appropriate balance.

We enforce the following structural constraints on the possible sets of correspondences between low pressure regions of a query and a case. First, the relation that the correspondences define must be one-to-one. That is, no region in a query can correspond to more than one region in the case, and no region in the case can correspond to more than one region in the query. Second, if a region  $R_q$  in the query corresponds to a region  $R_c$  in a case, then no child of  $R_q$  may correspond to a parent or sibling of  $R_c$ , and no child of  $R_c$  may correspond to a parent or sibling of  $R_q$ . A correspondence that violates these constraints is said to be *structurally inconsistent* with the correspondence  $(R_q, R_c)$ .

An informal description of the matching algorithm is as follows. First, the goodness of match between each pair of low pressure regions in the query and the case is computed. Pairs whose similarity falls beneath a certain threshold are not considered further in the search. The remaining pairs are sorted according to their goodness of match and stored on an ordered list  $L$ .

Next, a breadth-bounded backtracking search is performed through the space of possible correspondences between regions in the query and regions in the case. At each stage of the search, the best remaining pair  $(R_q, R_c)$  on  $L$  is selected and used to extend the current partial match. All pairs that are structurally inconsistent with  $(R_q, R_c)$  are removed from  $L$ . The search proceeds in a depth-first manner, bottoming out when  $L$  is empty. Each terminal node of the search represents a possible best match (set of pairs) between the query and the case. Upon backtracking, the previous value of  $L$  is restored (except that the pair most recently used to extend that node of the search is discarded), and the next best remaining pair on  $L$  is used to extend the partial match. At each level of the search, only the  $k$  best pairs are considered, for some small integer  $k$ . The parameter  $k$  constrains the breadth of the search, and can be tuned as needed. A small  $k$  produces very fast search, but risks not finding a good match.

### Case Selection

High quality case selection requires a choice of indices that achieve good recall without unduly sacrificing either precision or speed of case selection itself. This requires finding indices that closely mimic the behavior of similarity assessment. However, the available indices are constrained by the kinds of index mechanisms that Postgres provides.

Two kinds of design decision must be made to arrive at a procedure for case selection:

1. *Representation choice.* Which features of the full index label should case selection be sensitive to? For example,

Method	Description	Sample Postquel query
1	The bounding boxes overlap	retrieve (r.situation) from r in region where r.classification = "low"::char16 and r.bounds && "(162.8,-53.8,183.6,-29.3)"::box
2	The pressure minimum of one system is contained in the bounding box of the other	retrieve (r.situation) from r in region where r.classification = "low"::char16 and (r.centre → "(162.8,-53.8,183.6,-29.3)"::box or "(173.1,-41.4)"::point → r.bounds)
3	Both the pressure minima are within the overlap of the bounding boxes	retrieve (r.situation) from r in region where r.classification = "low"::char16 and r.centre → "(162.8,-53.8,183.6,-29.3)"::box and r.bounds → "(173.1,-41.4)"::point

Figure 6: Index representations

should case selection be sensitive to cluster structure or not? Should case selection consider the pressure minimum of a low pressure region, or the region's geographic extent, or both?

2. *Adjustable parameters.* If geographic extent is considered, then Postgres' R-tree indexing mechanism must be employed to achieve adequate efficiency, with its restriction to rectangular regions. How should suitable dimensions for rectangular boxes in a Postgres query be computed from the low-pressure regions present in a user's query? For example, should bounding boxes for these regions be employed, or larger rectangles? As we will see below, larger rectangles turn out to yield better performance. A parameter describing exactly how much larger to make these rectangles needs to be adjusted on the basis of empirical data.

In the next section, we describe several alternative approaches to case selection in more detail.

## Evaluation

It should be clear from the discussion so far that developing representations and algorithms for case selection and similarity assessment necessarily involves a large number of design decisions, many of which should only be made on the basis of empirical evaluation.

There is an important sense in which designing a case retrieval system requires a kind of hill-climbing search, guided by empirical considerations. At each stage in the search, a range of possible choices is considered and evaluated, and the best one selected. To illustrate this process, we next briefly describe two preliminary experiments regarding retrieval of low pressure systems that we carried out as part of developing MetVUW Workbench's current algorithm for case selection.

### Experiment 1: Representation choice for case selection

We evaluated three possible index representations for case selection involving queries that mention only low-pressure systems. These index representations are summarized in figure 6. The third column exemplifies actual queries to Postgres that case selection constructs in response to a query from the user involving a single low pressure system. Queries to Postgres are represented in Postquel, Postgres' query language.

Five representative queries were chosen and tested using each index representation. The average precision and recall of each was computed with respect to the "benchmark set" obtained by applying similarity assessment to the entire case base. The results of our analysis are presented in figure 7.

Method	Precision	Recall
1	0.14	0.80
2	0.22	0.60
3	0.58	0.28

Figure 7: Results of experiment 1.

It is instructive to recount the history of design decisions that led us to consider these choices of index representation for case selection. Throughout the design process, our aim was to maximize precision without unduly sacrificing recall.

We started with the most simple-minded approach possible: retrieve a case in response to a query if each low pressure system in the query overlaps a low pressure system in the case (method 1 in figure 7). We quickly discovered that this strategy produces quite low precision. We therefore tried incorporating pressure minima into case selection in an effort to improve precision. As can be seen from figure 7, quite high precision was obtained for method 3 but at the cost of unacceptably low recall.

The results of this first experiment led us to explore new ways to improve recall while maintaining precision. Our next hypothesis was to modify methods 2 and 3 to retain sensitivity to pressure minima, but allow a degree of latitudinal and longitudinal displacement between systems in a query and corresponding systems in a case. Our rationale was as follows: similarity assessment also permits a certain amount of displacement, so if we also allow a little displacement in queries, case selection should mimic similarity assessment more closely. It is therefore reasonable to expect an improvement in recall without incurring a huge penalty for precision.

### Experiment 2: Parameter choice—geographic displacement

In this experiment, we made two changes to the Postgres queries constructed by case selection to permit geographic displacements between low pressure systems in queries and cases. First, we introduced two new adjustable parameters *Lat* and *Long* that specify amounts to expand rectangular bounding boxes of low pressure systems in queries in

the north-south and east-west directions respectively. The new queries permit a degree of geographical displacement between regions in queries and corresponding regions in cases. Second, we replaced each pressure minimum in the query with a rectangular region whose dimensions are *Lat* and *Long*, and we replaced each test of point containment involving this minimum with a corresponding test of overlap with the new region. It is easy to see that this partially relaxes the constraint that pressure minima in the query must be contained in a low pressure system in the case.

We tried three different settings for *Long* and two different settings for *Lat*. The results are summarized in figure 8. We did not carry out tests for cells whose expected precision was less than 10%. The results for method 2 are largely negative: there appears to be a straightforward and relatively uninteresting tradeoff between precision and recall, and no choice of values for *Lat* and *Long* seems to yield results clearly superior to the methods of experiment 1. It appears that the requirement that only one pressure minimum be approximately contained in the bounding box of the region it is matched against does not sufficiently constrain case selection to produce good results.

The results for method 3 are more encouraging. Relatively high levels of both precision and recall are achieved for a number of settings for *Lat* and *Long*. For example, the setting  $\{Lat = 4^\circ, Long = 10^\circ\}$  achieves the same recall as method 2 of experiment 1, but with significantly better precision.

Lat Long	Method 2			Method 3		
	0°	4°	8°	0°	4°	8°
0°	0.60 0.22	0.68 0.18	0.75 0.14	0.28 0.58	0.32 0.31	0.35 0.27
5°	0.73 0.20	0.82 0.13	0.87 0.10	0.45 0.56	0.53 0.37	0.57 0.27
10°	0.82 0.17	0.90 0.10	— —	0.48 0.48	0.60 0.34	0.65 0.21
15°	0.82 0.14	— —	— —	0.53 0.41	0.68 0.28	0.73 0.18

Lat = Permitted latitudinal displacement  
Long = Permitted longitudinal displacement

Figure 8: Results of experiment 2 (recall/precision).

## Conclusions and Future Work

The two experiments describe above only consider a small region of the space of design decisions that have gone into MetVUW Workbench's case retrieval system. For example, we have not discussed evaluation of similarity assessment. Nevertheless, these experiments well illustrate the utility of a process of iterative design. More generally, we believe that the large number of decisions that go into the design of a case retrieval system, and the difficulty of assessing their impact *a priori*, make the design process an enterprise with a necessarily empirical component.

Earlier we used a metaphor of hill-climbing search to describe the design process. We are actively considering

whether it might be possible to take this metaphor literally, and automate part of the design-evaluation cycle. Design decisions can be divided into two categories: choices of appropriate representational distinctions and "parameter tweaking," in which a good value for a numeric parameter is selected. Decisions as to representational distinctions cannot be easily automated, as formulating the available options appears to require both human intuition and AI expertise. Choosing appropriate values for numerical parameters, however, may be easier to automate.

In this paper, we have seen a number of adjustable numerical parameters: weights for matching of latitudinal and longitudinal displacement, aspect ratios, density, area, and the like when matching regions of low pressure; width of search in structure matching; and degree of allowed latitudinal and longitudinal displacement in case selection. Perhaps hill-climbing techniques from machine learning can be used to explore this parameter space in an attempt to optimize the choice of parameter values? We intend to explore this idea in our future research.

## Acknowledgements

This research was carried out in collaboration with James McGregor of the Institute of Geophysics at Victoria University of Wellington. Thanks to Linton Miller for commenting on a draft of this paper, and to Erik Brenstrum of the Meteorological Service of New Zealand for useful discussions regarding the meteorologically significant features of low and high pressure systems. The Meteorological Service of New Zealand provided the situation maps in figure 3.

## References

- Domeshek, E. A. 1991. Indexing stories as social advice. In *Proceedings, Ninth National Conference of Artificial Intelligence*, volume 1, 16–21. AAAI Press.
- Guttman, A. 1984. R-trees: A dynamic index structure for spatial searching. In *Proceedings of the 1984 ACM-SIGMOD Conference on Management of Data*. Boston, Massachusetts: ACM-SIGMOD.
- Jacobs, P. S., ed. 1992. *Text-Based Intelligent Systems*. Lawrence Erlbaum Associates.
- Jones, E. K., and Roydhouse, A. 1994. Intelligent retrieval of historical meteorological data. To appear in *AI Applications*. Also available as Technical report CS-TR-93/8, Victoria University of Wellington. URL: <ftp://ftp.host.comp.vuw.ac.nz/doc/vuw-publications/CS-TR-93/CS-TR-93-8.ps.gz>.
- Pearl, J. 1984. *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. Reading, MA: Addison-Wesley.
- Roydhouse, A.; Miller, L.; Jones, E.; and McGregor, J. 1993. The design and implementation of MetVUW Workbench version 1.0. Technical Report CS-TR-93/7, Victoria University of Wellington. URL: <ftp://ftp.host.comp.vuw.ac.nz/doc/vuw-publications/CS-TR-93/CS-TR-93-7.ps.gz>.
- Schank, R. C., and Fano, A. 1992. A thematic hierarchy for indexing stories in social domains. Technical report, The Institute for the Learning Sciences, Northwestern University.
- Schank, R. C. 1982. *Dynamic Memory*. Cambridge, England: Cambridge University Press.
- Stonebraker, M., and Rowe, L. 1986. The design of Postgres. In *Proceedings of the 1986 SIGMOD Conference on Management of Data*. ACM Press.