

Finding clusters of similar events within clinical incident reports: a novel methodology combining case based reasoning and information retrieval

C Tsatsoulis, H A Amthauer

A novel methodological approach for identifying clusters of similar medical incidents by analyzing large databases of incident reports is described. The discovery of similar events allows the identification of patterns and trends, and makes possible the prediction of future events and the establishment of barriers and best practices. Two techniques from the fields of information science and artificial intelligence have been integrated—namely, case based reasoning and information retrieval—and very good clustering accuracies have been achieved on a test data set of incident reports from transfusion medicine. This work suggests that clustering should integrate the features of an incident captured in traditional form based records together with the detailed information found in the narrative included in event reports.

One of the goals of incident reporting systems is to allow their users to discover trends, identify patterns of organizational behavior, and predict future failures of the process. This is especially true for systems that collect reports across many organizations; such systems allow organizations to learn the shortcomings of others and to correct their own operating procedures before similar errors appear locally.

To achieve these objectives the users of an incident reporting system should be able to point to a specific report and then query the system for other incidents that are similar to it. In essence, users want to identify a cluster of event reports that are exemplified by the report in question. The identification of the cluster provides valuable information to the users of an incident reporting system:

- How many reports are in the cluster?
- What is their distribution in time (which, in turn, helps in establishing trends)?
- What are the exemplifying characteristics of the cluster?

Standard database retrieval cannot offer a measure of similarity; objects in a traditional database are accessed by exact matching of field values. While it is of some value to identify incident reports that have identical descriptions, it is a lot more probable that incident reports will only be similar—that is, will share some com-

mon features but will differ in others. In addition, even features that are different in two reports may share some common characteristics—for example, the incident time frame may be “4–8 am” in one report and “8–12 am” in another, but both times can be thought of as “morning”. When analyzing trends in medical incidents and when trying to identify best practices in response to incidents, medical personnel and quality assurance experts are interested in finding clusters of similar reports—that is, reports that share some important common characteristics—instead of looking for identical reports. Similarity requires both a *syntactic* and a *semantic* matching of the features describing an incident report. Syntactic matching compares two strings of characters—for example, “abc” and “Abc”—and determines if they are identical or, if not, how different they are (in this example they differ in one out of three letters). Semantic matching compares two concepts such as “male” and “man” and determines if they represent the same thing or idea and, if not, how close the two concepts are semantically. Case based reasoning (CBR) and information retrieval (IR), two techniques from the field of artificial intelligence, offer tools to identify similar incident reports.

In this paper we first describe CBR (box 1) and IR (box 2), and then describe our use of these techniques to identify clusters of similar documents in the Medical Event Reporting System—Transfusion Medicine (MERS-TM) event reporting system which is used to document incidents in transfusion services (box 3). When we applied CBR to the creation of clusters of similar reports, we first identified the features of a transfusion incident report that should be used as indexes (report descriptors useful in identifying similarity), assigned different weights to each index as an indicator of its importance in establishing similarity, and defined domain specific semantics to allow knowledge based matching of indexes. In addition, we used techniques from IR to analyze the textual description of the event attached to each report. We performed experiments on a set of incident reports collected through the MERS-TM transfusion medicine incident reporting system¹ using CBR retrieval, IR retrieval, and also integrating IR with CBR. The goal of our experiments was to determine whether the CBR and IR retrieval methodologies alone would identify as similar cases that experts in transfusion services would also consider as such, and whether a combination of CBR and IR

See end of article for authors' affiliations

Correspondence to:
Dr C Tsatsoulis, The
University of Kansas,
Lawrence, KS 66045,
USA; tsatsoul@ittc.ku.edu

retrieval would have superior retrieval performance to either technique alone. The results of each retrieval, clustering, and similarity assessment were evaluated with the help of experts in the area of quality assurance in transfusion medicine who calculated the number of false positives and negatives in the clusters of similar incident reports generated by our software.

Our results indicate that the integration of CBR with IR improves performance of the retrieval system and offers good recall and accuracy.

APPLICATION OF CBR AND IR TO INCIDENT REPORTS FROM TRANSFUSION SERVICES

The MERS-TM incident reports were analyzed by experts in the field of transfusion services who defined a subset of the

report features that should be used as indexes in our CBR system. These features include the discovery time, the discoverer's job description, the point in the process at which the event was discovered, where it first occurred, the causal and antecedent codes. The experts also assigned a weight of 1–5 to each index, where the higher weight indicated greater importance of a feature in matching and clustering. For example, where an event first occurred was weighted 5, the time an event was discovered was given 1, and the discoverer's job description 3.

For some attributes the experts gave conditional weights. For example, a causal code would receive a weight of 1 or 2 depending on whether it was based on a rough examination of the incident or on an in depth analysis.

The experts also defined hierarchies of attribute values that allowed us to define partial matches. For example:

Box 1 Case based reasoning (CBR)

Case based reasoning (CBR) is a problem solving paradigm based on psychological theories of human cognition which provides the foundations for a technology for intelligent systems.² It is based on the intuitive notion that human expertise is not based on rules or other formalized structures but on experiences. Human experts differ from novices in their ability to relate problems to previous ones, to reason based on analogies between current and old problems, and to use solutions from old experiences.

The process of reasoning using experiences or cases can be described by the following steps:

- **Retrieve:** Given a new problem, retrieve a similar past case from memory. The past case contains the prior solution.
- **Modify:** The old solution is modified to conform to the new situation, resulting in a proposed solution.
- **Test:** The proposed solution is tested for successful solution of the current problem.
- **Learn:** If the solution fails, explain the failure and learn it to avoid repeating it. If possible, repair the failure, generate a new proposed solution, and return to the *Test* step. If the solution succeeds, incorporate it into the case memory as a successful solution and stop.

Since our work concentrates on retrieval, this description will be constrained to this part of a CBR system.

A CBR system must select the best case or cases from memory. The question that must be answered is what constitutes an *appropriate* or *similar* case. What are the criteria of closeness or similarity between cases, and how should cases be indexed? Indexing a case is essential in establishing similarity since the indexes help to define the elements of a problem that are important.

During retrieval each case must be compared with the current problem and be assigned a degree of similarity. The retrieving program will then select the cases with the highest degree of similarity. We therefore need to define what we mean by "best match" or, as usually called in conceptual retrieval, what we mean by "similar(ity)". The simplest method would be to look at *structural* or *syntactic similarities* between the current problem and a case. This demands an exact match between index values in a manner identical to database retrieval. (Note that this is a simplification of structural matching.) One can demand a perfect syntactic match only of symbolic values—that is, non-numerical ones; the same is not true for numerical ones. For numbers a perfect match may be based on a formula: for example, "x is qualitatively equal to y if it is $y \pm 20\%$ ". If two values match structurally, we say that they match *perfectly* (or, if we wanted to assign a degree of match between 0 and 1 where 0 is absolute mismatch, a structural match would receive a value of 1.0). For example, we would say that "ABC" and "ABC" match perfectly (they are structurally identical—that is, they look the same), while "ABC" and "DEF" do not match since they do not look the same at all. On the other hand, we could define partial similarity, and say, for example, that "ABC" matches "XBC" with 67% match, since the two strings share two out of three letters.

Deciding whether two values match or not can also lead to a partial (or *semantic*) match. The concepts represented by the case indexes are placed on a hierarchy of classes and their subclasses. For example, one may say that "beef" and "chicken" are subclasses of "meat". Then, "beef" and "chicken" match partially since they are different concepts but they are both subclasses of the superclass "meat". We can assign a value to this partial match based on the level of the hierarchy where values match. For example, a complete match can be given 1.0 and, for moving up a level of the hierarchy, we may want to multiply the match by 0.7 (1, 0.7, 0.5, 0.35, ...). Creating a membership hierarchy is just one way to establish partial similarity of symbolic values. Some of the similarity can be rule based, where the rules are defined by experts. For example, an expert can give a rule that says that "the emotional state of anger is similar with degree 0.8 to the emotional state of rage".

Indexes can be assigned a weight (in an arbitrarily selected scale) that indicates the contribution of a particular index to establishing similarity. Usually, index weights are assigned by domain experts who are best suited to estimate which characteristics of a case are the most relevant ones.

After we determine which index values are qualitatively similar or equal, we compute a similarity value for the whole case. Usually this is done in a nearest neighbor method, which is a weighted average. For example, we can compute the degree of similarity as:

$$\text{similarity} = \frac{\sum w_i \times \text{sim}(f_c, f_p)}{\sum w_i}$$

where w_i is the weight for a matching feature, and *sim* is the degree of match between the old case f_c and the current problem f_p .

Box 2 Information retrieval (IR)

Information retrieval (IR) systems are used for indexing, searching, and recalling text or other unstructured forms of data. The primary basis of IR for text retrieval is through the use of weighted keywords. Since IR systems do not require any domain specific knowledge, IR systems can be applied in any domain where textual documents are available.

Traditionally, text documents are pre processed where common words (or “stop words”) such as “a”, “and”, “the”, etc are removed from the document. Next *stemming* is performed, where words are reduced to their stem so that, for example, “independence” and “independent” are represented by the common stem “independ”. The text tokens are then stored in a structure that allows quick comparison and retrieval.

One approach in IR for document retrieval is the vector space model (VSM)³ in which each document is represented by a list (vector) of terms. These terms have associated weights that describe a term’s value for a document. The weighting system for each term in the document uses a tf-idf scheme. (tf = term frequency; idf = inverse document frequency). In this term weighting scheme the tf and idf are calculated in the following manner:

tf = frequency of the term in the document/frequency of the most frequent word in the document

idf = $\log_{10}(\text{total number of documents in the collection} / \text{number of documents in the collection that contain the term})$

Thus, the weight of a term is calculated by:

Weight = tf * idf

Using the VSM makes it possible to compare two documents using vector algebra as, for example, the cosine measure of similarity.⁴ With this method the degree of similarity between two documents is determined by the cosine of the angle between the vectors that represent the two documents (the smaller the angle, the more similar), so that a document might be retrieved even if it shares only a few terms.

- For the attribute indicating when an incident was discovered, the values in the pairs (12–4 am, 4–8 am), (8–12 noon, 12–4 pm), and (4–8 pm, 8–12 midnight) were considered partially similar. So, for example, a report with value “8–12 noon” would have a partial match with a report with value “12–4 pm”.
- The values of the attribute indicating the job description of the person discovering the incident were organized into sets where the values were considered as matching partially—for example, supervisor, medical technician, quality assurance/quality control person, and registered nurse were all members of the same set of partial matching values.

In our CBR system we only had a single level of hierarchy of feature values, and every partial match was assigned a value of 0.7 (a perfect match received a value of 1.0 and a non-match a value of zero).

Our approach to the IR portion of this study uses the vector space model (VSM) and the cosine comparison measure, as described in box 2. In our case, a document is considered the free text of the report portion that describes what happened. The removal of noise from the text was difficult due to the domain specific abbreviations used. For example, “OR” was used mostly as an abbreviation for “operating room”, not as a conjunction. So as not to lose important abbreviations, no stop words were removed. Matching based on words that do not carry a lot of meaning due to their high frequencies is easy to identify, so the non-removal of stop words is easily handled.

Box 3 MERS-TM incident reporting system for transfusion medicine

The Medical Event Reporting System for Transfusion Medicine (MERS-TM) is an event reporting system developed for transfusion services and blood centers to collect, classify, and analyze events that could potentially compromise transfusion safety.¹ The incident reports of MERS-TM consist possibly of three parts. The first two parts are mandatory. One describes the incident with a set of surface features such as the time and date the incident was discovered, by whom it was discovered, when it occurred, location code of the point of occurrence, and so on. The other mandatory document is the quality assurance investigation report which includes codes describing the causal events (MERS-TM uses the Eindhoven classification system for causes of events⁵), any preventive actions taken, and the type of investigation conducted. In addition to the surface features and the causal event codes, the MERS-TM incident reports always include a brief (1–2 lines) textual description of the event. If the organization decides to perform a detailed investigation, it will generate the third optional part of the report which includes detailed information about the consequent and antecedent events.

Figure 1 shows the information that a user may enter to describe in MERS-TM the discovery of a transfusion medicine incident, and fig 2 shows part of a completed detailed investigation report entered in MERS-TM displaying the causal codes based on the Eindhoven classification system.⁵

We then performed a set of experiments to establish the efficacy of CBR and IR in clustering similar clinical incident reports. For the experiments we used an MERS-TM data set of approximately 600 reports collected by the transfusion services of two hospitals and made available to us by the MERS-TM group led by Dr Harold Kaplan of the Presbyterian Hospital of Columbia University, New York. The incident reports were indexed for CBR retrieval as indicated above and also preprocessed for IR retrieval. After the incident reports were indexed they were entered in a “case base”—that is, in a storage file that makes comparisons and similarity assessment possible through our software. Similarly, after IR preprocessing the incident reports were stored in a structure appropriate for IR retrieval.

The goals of our experiments were to determine:

- whether CBR retrieval would identify as similar cases that experts in transfusion services would also consider as such;
- whether IR retrieval would do the same;
- whether a combination of CBR and IR retrieval would have superior retrieval performance to either technique alone.

As a baseline test we performed retrieval using equal weights for all indexes; the goal was to establish whether the index weights given to us by the experts improved CBR retrieval or not.

To establish the usefulness of CBR for finding clusters of similar medical incident reports, we randomly selected 24 cases out of the approximately 600 incident reports in the data set (to avoid confusion with the cases in the case base the cases we used to match against will be called “reports” from now on), and for each of these reports we retrieved the 10 most similar cases from the case base created from the processed incident reports. An example of two matching transfusion incident reports is shown in fig 3 in which we show the case based

Section A–Discovery Information		
1.	Report date:	<input type="text"/> mm/dd/yyyy
2.	Discovery date:	<input type="text"/> mm/dd/yyyy
3.	Was this discovered on a weekend or weekday?	<input type="text"/>
4.	Discovery time:	<input type="text"/>
5.	Discoverer's job description:	<input type="checkbox"/> Clerk <input type="checkbox"/> MT <input type="checkbox"/> Supervisor <input type="checkbox"/> House staff <input type="checkbox"/> QA/QC <input type="checkbox"/> Other <input type="checkbox"/> MD/DO <input type="checkbox"/> RN <input type="checkbox"/> MLT <input type="checkbox"/> LVN/LPN
6.	Where discovered:	<input type="text"/>
	Location code (optional)	<input type="text"/>
7.	Describe briefly the event you discovered:	<input type="text"/>
8.	How did you discover this event?	<input type="text"/>
9.	This event was discovered:	<input type="text"/>
10.	Product/record action:	<input type="checkbox"/> Product retrieved <input type="checkbox"/> Additional testing <input type="checkbox"/> Product destroyed <input type="checkbox"/> Patient sample recollected <input type="checkbox"/> Record corrected <input type="checkbox"/> Other <input type="checkbox"/> Floor/Clinic notified

Figure 1 The “discovery information” section of MERS-TM. Here the user records how the transfusion medicine incident was discovered.

	Report accession number	100
1.	Consequent (discovery) code:	<input type="text"/> 1 <input type="text"/> AV <input type="text"/>
2.	Antecedent (1st occurrence) code:	<input type="text"/> US <input type="text"/>
3.	Significant antecedent (occurrence) code:	<input type="text"/> OE <input type="text"/>
4.	Additional description of event (optional)	<input type="text"/>
5.	Risk assessment:	QES .10 <input type="text"/> QEP .50 <input type="text"/> Final RAI 0.25
6.	Organizational risk?	<input type="text"/> None <input type="text"/>
7.	Follow up:	<input checked="" type="checkbox"/> Propose action <input type="checkbox"/> Consider action <input type="checkbox"/> Monitor <input type="checkbox"/> External report to other dept/org <input type="checkbox"/> FDA reportable
8.	If appropriate, describe the long term preventive action to be taken:	<input type="text"/>

Figure 2 Part of a completed detailed investigation report from MERS-TM. The quality assurance personnel performing the investigation have identified and recorded a number of causal and risk codes.

Feature	Case A	Case B
Report date	3/30/1999	4/2/1999
Discovery date	3/30/1999	4/2/1999
Discovery time **	4-8 PM	8-12 Midnight
Discoverer's job description *	MLT	MLT
Where discovered *	Trans. Serv.	Trans Serv.
What happened	PHLEBOTOMIST FAILED TO SIGN REQUISITION	PHLEBOTOMIST FAILED TO SIGN REQUISITION
How discovered	ON SAMPLE CHECKING	AT TIME OF SAMPLE CHECKING
Point in process discovered *	Before testing patient sample	Before testing patient sample
Product record action *	Patient sample recollected	Patient sample recollected
Date event occurred	3/30/1999	4/2/1999
Occurrence time	4-8 PM	4-8 AM
Person involved *	RN	RN
Where first occurred *	Sample collection	Sample collection
Consequent event type 1 *	3	3
Consequent event a *	SC	SC
Consequent event b *	099	099
Antecedent event a		
Antecedent event b		
Follow up *	Monitor	Monitor
Investigation type *	Routine investigation	Routine investigation
RL Cause code 1	HKK	HRM ***
RL Cause code 2	OK ***	HSS
RL Cause code 3	HRM ***	OK ***

- = No match (weight * 0)
- * = Exact match (weight * 1)
- ** = Partial match (0.7 *(weight * 1))
- *** = Group match (0.7 *(weight * 1)) (not listed in same order)

CBR matching score (partial and group matches are in **bold**; in each parenthesis the first value is the weight of the attribute and the second value is the matching value; attributes are listed in order from top to bottom): $((1 * 0) + (1 * 0) + (0.7 (1 * 1)) + (4 * 1) + (4 * 1) + 0 + 0 + (4 * 1) + (1 * 1) + (1 * 1) + (2 * 1) + (4 * 1) + (4 * 1) + (3 * 1) + (5 * 1) + (4 * 1) + (5 * 0) + (4 * 0) + (1 * 1) + (1 * 1) + (0.7 (3 * 1)) + (0.7 (3 * 1)) + (3 * 0))/59 = 0.73$

Figure 3 Example of matching of two incident reports using CBR. The features of each report are compared and, depending on the type of match (exact, partial, group, or no match), the weight of an attribute is multiplied by an appropriate value. The values are added up and normalized. The resulting “match value” for this example was 0.73.

matching of two reports with a degree of similarity of 0.73. The same two reports had a match value of 1.0 using IR (we compared the text under the attribute “what happened?”) since the two reports had almost identical textual descriptions.

To establish the usefulness of IR for finding clusters of similar medical incident reports we used the same 24 reports and identified similar ones using only an IR based keyword match of the text included with each case.

The results of the CBR and IR retrieval were then combined to establish whether the combination would offer superior performance for finding clusters of similar medical incident reports. We assigned to the matching percentage of each retrieval technique a weight between 0.9 and 0.1, in increments of 0.1, making sure that the sum of the two weights always equalled 1.0. In other words, the CBR match value was weighted by 0.9, 0.8, 0.7, ..., 0.2, 0.1, while the IR match value was weighted by 0.1, 0.2, 0.3, ..., 0.8, 0.9. For example, in fig 3 the CBR match value was 0.73 and the IR match value was 1.0. Combining the two values provided the similarity values shown in table 1.

The cases were then re-ranked based on the new combined matching value, resulting in nine new rankings. In our baseline test we performed CBR retrieval using equal weights for all indexes (the weights were set to 1 since the similarity value is normalized). The result of all these experiments was 12 sets of ranked cases which were similar to the original report (CBR only, IR only, CBR with no weights, and nine rankings with varying weights assigned to the CBR and IR similarity values).

In fig 4 we show a flow chart of operations, starting with the preprocessing of the MERS-TM incident reports and ending with the expert evaluation of the CBR and IR clustering.

CBR: IR	Matching score
100:0	0.73
90:10	0.75
80:20	0.78
70:30	0.81
60:40	0.84
50:50	0.86
40:60	0.89
30:70	0.92
20:80	0.95
10:90	0.97
0:100	1.0

The first column shows the contribution by each clustering technique and the second column shows the combined matching value. The first row represents 100% contribution by CBR and 0% from IR, so the resulting matching score is (CBR match * 1) + (IR match * 0) = 0.73 + 0 = 0.73. The second row represents 90% contribution by CBR and 10% from IR, so the resulting matching score is (CBR match * 0.9) + (IR match * 0.1) = 0.65 + 0.1 = 0.75, etc.

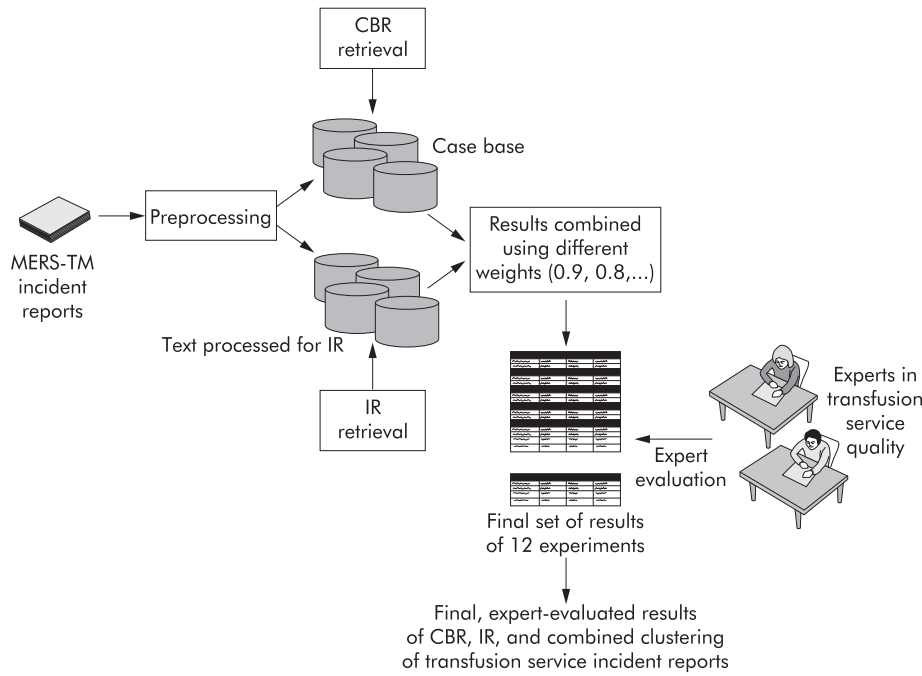


Figure 4 Flow chart showing operations of the system.

For evaluation, we collected the top five retrieved cases for each experiment for each report. Since many of the retrieved cases for the different experiments were the same, the result was a set of 10–20 cases for each report. To these cases we added one randomly selected case from the database to use as a control point for the evaluation. These cases were ordered randomly so as not to give any hint to the evaluators. The two experts who participated in the evaluation of our work were Ms Barbara Rabin Fastman of Columbia University’s New York Presbyterian Hospital and Ms Quay Mercer of the University of Texas Southwestern Medical Center at Dallas. Both are experts in transfusion services and quality assurance of medical and hospital processes. The experts were asked to evaluate whether the cases matched the report or not on a 4 point scale: “almost identical”, “similar”, “not very similar”, and “not similar at all”. This scale is clearly subjective and its intent is to give the experts the freedom to express their personal opinion of the quality of the performance of the similarity algorithm without having to understand how the algorithm works. Table 2 shows an example where the two incident reports (cases A and A2) were assessed by one of the experts to be “almost identical”, and in table 3 two incident reports are presented (cases R and R2) that were assessed by one of the experts to be “similar”.

The experiments and the evaluation of their results were performed during 2002 and early 2003. The incident reports were handled in an electronic format (transformed appropriately for CBR and IR as described above), and the CBR and IR clusterings were performed using software developed by us. The evaluation of the results by the experts was analyzed by statistical software to summarize it and to allow us to draw generalized conclusions.

ANALYSIS OF EXPERIMENTAL RESULTS

We analyzed the results of the system and the experts’ evaluation in the following manner:

- All cases ranked by the experts as “almost identical” and “similar” were classified as “retrievable”, while the other two rankings indicated cases that should be “non-retrievable”.

Table 2 Example of evaluation of two incident reports assessed by one of the experts to be “almost identical”

Attribute	Case A	Case A2
Report date	3/30/1999	7/15/1999
Discovery date	3/30/1999	7/13/1999
Discovery time	4–8 am	12–4 pm
Discoverer’s job description	Medical laboratory technician	Medical laboratory technician
Where discovered	Transfusion service	Transfusion service
What happened	Wrong requisition used for crossmatch	Wrong requisition used for group and screen
How discovered	On sample check in	At requisition check in
Point in process discovered	Before testing patient sample	After component process, before issue
Product record action	Patient sample recollected	Patient record corrected
Date event occurred	3/30/1999	7/13/1999
Occurrence time	4–8 am	12–4 pm
Person involved	Registered nurse	Medical laboratory technician
Where first occurred	Sample collection	
Consequent event type	3	3
Consequent event a	SC	SH
Consequent event b	099	099
Antecedent event a		
Antecedent event b		
Follow up	Monitor	Monitor
Investigation type	Routine investigation	Routine investigation
RL cause code 1	OK	HRM
RL cause code 2	HKK	HKK
RL cause code 3	OM	OK

- The results of the 12 experiments (CBR only, IR only, CBR with no weights, and nine rankings with varying weights assigned to the CBR and IR similarity values) were studied for different similarity matching thresholds (ranging from 0.1 to 1.0 in 0.1 increments). These thresholds indicate what cases should be added to the cluster of similar ones. For example, a 0.4 threshold would include in the cluster cases which match with a similarity value of 0.4 and above.
- Our two quality criteria were *recall* and *accuracy*, which is the percentage of retrievable cases retrieved and the

Table 3 Example of evaluation of two incident reports assessed by one of the experts to be "similar"

Attribute	Case R	Case R2
Report date	9/9/1999	4/11/1999
Discovery date	9/9/1999	4/11/1999
Discovery time	12-4 am	12-4 am
Discoverer's job description	Medical laboratory technician	Medical laboratory technician
Where discovered	Transfusion service	Transfusion service
What happened	Wrong RBC expiration entered in Hemocare	Wrong expiration date entered when RBC modified to irradiated
How discovered	Upon transfusion reaction investigation	MLT discovered immediately issue record
Point in process discovered	After issue, before infusion	After issue, before infusion
Product record action	Unit destroyed	Unit destroyed
Date event occurred	8/9/1999	4/11/1999
Occurrence time	12-4 pm	12-4 am
Person involved	Registered nurse	Registered nurse
Where first occurred		
Consequent event type 3		3
Consequent event a	PC	UM
Consequent event b	002	001
Antecedent event a		
Antecedent event b		
Follow up	Monitor	Monitor
Investigation type	Routine investigation	Routine investigation
RL cause code 1	HSS	HSS
RL cause code 2	TEX	TD
RL cause code 3		TEX

percentage of non-retrievable cases not retrieved.* In other words, recall tells us how many of the appropriate reports we are finding, while accuracy tells us how many of the inappropriate reports we are avoiding (one minus accuracy would give us the percentage of the incorrect reports we are including in our similar cluster, indicating false positives). Clearly, we want high recall and accuracy.

*These are also known as "true positives" and "true negatives", respectively.

We expected that, as the similarity threshold was raised, recall would be lower and accuracy would improve: a lower similarity threshold would assume that most cases were similar and, as a result, would include all the retrievable cases but also many non-retrievable cases; as the threshold is increased, fewer cases are considered similar, excluding some retrievable ones, but, hopefully, also excluding most non-retrievable ones. We also expected that our CBR system would do better than the CBR with no weights since the weights were assigned by experts specifically to assist in matching and similarity assessment. We had no expectations about the performance of the integrated CBR and IR retrieval since no similar experiments had been performed previously.

Some of the results of our experiments are shown in tables 4-8. As expected, as the matching threshold is increased, recall is lowered but accuracy increases greatly (table 4). The asterisks indicate that no cases were retrieved which were above the listed matching thresholds. Clearly, the CBR only retrieval does very well with recall but poorly with accuracy. This may be an indication that the report fields used as indexes are superficial descriptors of an event and, as such, do not offer the detail necessary to distinguish between dissimilar reports.

We next compared the recall and accuracy of the CBR system using expert assigned weights versus the CBR system using equal weights. In table 5 we list the difference in the quality of recall and accuracy as a function of the matching threshold. The asterisks indicate that no cases were retrieved that were above the listed matching thresholds. As expected, the recall of the CBR system with weights is substantially better than that of the CBR system with equal weights. Table 5 would seem to indicate that CBR with equal weights has a better accuracy, but closer inspection of the results showed this not to be the case since CBR with equal weights classified almost *all* cases as not similar and, thus, would trivially exclude non-retrievable ones.

We also examined the accuracy and recall of the IR retrieval as a function of the matching threshold. As expected, as the matching threshold is increased, recall is lowered but accuracy increases greatly. The problem with IR retrieval is that the fall off in recall is extremely steep. Our hypothesis is that the text in the MERS-TM reports stresses case specific details that allow differentiation between

Table 4 Analysis of the clustering quality of CBR by examining recall and accuracy over different matching thresholds

	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
Recall	1.00	1.00	1.00	1.00	1.00	0.95	0.76	*	*	*
Accuracy	0.00	0.00	0.00	0.00	0.00	0.05	0.57	*	*	*

The matching thresholds indicate the value over which two reports were considered similar. "Recall" refers to the percentage of reports that the experts deemed as similar which we did retrieve. "Accuracy" refers to the percentage of reports that the experts deemed as not similar which we did not retrieve. The asterisks indicate that no cases were retrieved which were above the listed matching thresholds.

Table 5 Comparison of the clustering quality between CBR with expert supplied weights and CBR with equal weights

	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
Recall _{CBR} -	0.05	0.15	0.21	0.30	0.45	0.57	0.50	*	*	*
Recall _{CBRequal}										
Accuracy _{CBR} -	-0.35	-0.53	-0.73	-0.81	-0.88	-0.91	-0.40	*	*	*
Accuracy _{CBRequal}										

Quality is evaluated as the difference in recall and accuracy of the two techniques where positive numbers indicate better performance. The results indicate that CBR performs better when attributes used in matching are weighted by experts to signify the contribution of an attribute to the clustering decision. The asterisks indicate that no cases were retrieved that were above the listed matching thresholds.

Table 6 Analysis of the clustering quality of IR experiments examining recall and accuracy over different matching thresholds

	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
Recall	1.00	1.00	1.00	0.96	0.79	0.52	0.38	0.25	0.22	0.22
Accuracy	0.00	0.00	0.11	0.46	0.64	0.89	0.93	0.96	1.00	1.00

As for table 4, the matching thresholds indicate the value over which two reports were considered similar. "Recall" refers to the percentage of reports that the experts deemed as similar we did retrieve. "Accuracy" refers to the percentage of reports that the experts deemed as not similar we did not retrieve.

Table 7 Best recall and accuracy results achieved by combining CBR with IR

	CBR+IR (10:90) Threshold 0.3	CBR+IR (40:60) Threshold 0.4	CBR+IR (30:70) Threshold 0.4	CBR+IR (20:80) Threshold 0.4	CBR+IR (10:90) Threshold 0.4	CBR+IR (60:40) Threshold 0.5
Recall	0.79	0.80	0.75	0.73	0.79	0.75
Accuracy	0.73	0.75	0.80	0.78	0.73	0.80

The table shows the combination of the CBR matching value with the IR matching value at different percentages (e.g. 10% weight to the CBR value and 90% weight to the IR value). Recall and accuracy are as previously defined.

dissimilar ones, but also precludes the identification of similar ones that may share more general characteristics (table 6).

We next performed similarity retrieval using a weighted combination of IR and CBR and the results are shown in table 7 (note that we only list results where recall and accuracy are above 70%). Interestingly, the best results occurred in lower matching thresholds and when, in general, the contribution of IR retrieval is greater or even dominant.

In general, the CBR system had better recall but worse accuracy than the IR system. The integration of CBR with IR produced the best results, since it combined the strengths

of both techniques. Figure 5 shows the sum of accuracy and recall plotted against the matching threshold for 11 experiments (the CBR with equal weights is not included since it was used only as a baseline test). The best combined recall and accuracy values were obtained for matching thresholds of 0.40 and 0.50 and for combined weighted CBR and IR retrieval.

Since CBR seemed to identify retrievable cases well, and IR seemed to identify non-retrievable cases with over 90% accuracy, we performed one more experiment to examine an integrated CBR and IR system where each technique is used independently and then their results combined to exploit the strength of each method. The additional experiments were conducted as follows: CBR retrieval using the thresholds where CBR gave the best recall result (at 0.7, 0.6 and 0.5), IR retrieval using the thresholds where IR had the best accuracy (0.5 and 0.4), then the six possible intersections of the three CBR and two IR sets were created. Table 8 summarizes the results. As can be seen, these results are not substantially better than those achieved in the previous experiments, although they have outstanding recall and, in one case (CBR 0.6, IR 0.5), very good accuracy.

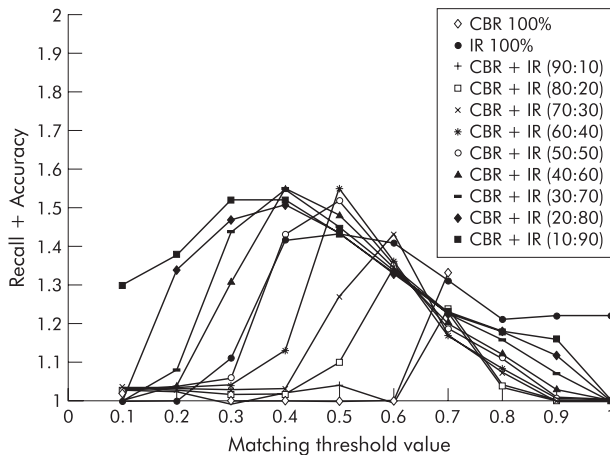


Figure 5 Plot of the sum of recall plus accuracy versus the matching threshold for CBR, IR, and combined weighted CBR+IR.

DISCUSSION

Our experiments showed that CBR is useful in identifying similar medical incident reports but its accuracy is poor. It seems that a lot of the detail of a case is contained in the textual description provided by the reporters of the event, and this is indicated by the retrieval accuracy of IR. On the other hand, the text in the reports is too detailed to provide sufficient abstract descriptions, leading to good accuracy but poor recall for IR only retrieval. The combination of CBR and IR techniques, either as a weighted sum of similarity values

Table 8 Results of integrating CBR with IR retrieval for the best retrieval thresholds of each technique

	CBR 0.7 IR 0.5	CBR 0.7 IR 0.4	CBR 0.6 IR 0.5	CBR 0.6 IR 0.4	CBR 0.5 IR 0.5	CBR 0.5 IR 0.4
Recall	1.00	0.99	0.92	0.93	0.93	0.88
Accuracy	0.13	0.17	0.66	0.25	0.40	0.52

See tables 4–6 for the clustering results for CBR and IR. The matching values at each threshold (CBR 0.7, 0.6 and 0.5 and IR at 0.4 and 0.5) were weighted equally to generate the combined clustering value. Recall and accuracy are as previously defined.

Pointers for future research

- Analyze a larger (>10 000) corpus of incident reports to determine the effects of the data size on the results of similarity clustering.
- Study the optimal combination of case based reasoning and information retrieval in the creation of clusters of similar incident reports.

or as the intersection of separate trials, greatly improved accuracy and recall. Based on our results, we strongly recommend that future systems developed to cluster medical event reports integrate both the field values and the text of the reports in their methodological approach.

Only a small amount of research has been done on the combination of CBR and IR for clustering, and it tends to support our findings. Specifically, the DRAMA system used text to enhance its CBR process by analyzing free form text that was part of aircraft design documents to capture rationale and interrelationships of design choices. In an example presented by Wilson and Bradshaw,⁶ the information in the text associated with a case improved retrieval, but the authors did not provide a systematic evaluation of the integration of the technologies.

There are two directions that our future work could take: the analysis of a large corpus of incident reports and the theoretical and experimental analysis of the best combination of CBR and IR. Our sample of incident reports (approximately 600 in total) is small compared with the size of databases of medical reports being created globally. It would be interesting to study how the size of the underlying database of reports affects the performance of clustering. Also, our work has provided some indications that CBR and IR work best when combined; future work should examine the circumstances under which each technique offers the best benefit—for example, more detailed versus more abbreviated text—and how the two techniques can be best combined to provide optimal clustering and retrieval results.

Key messages

- Databases of medical incident reports need to become active and provide answers instead of simply history. One way to do so is to identify clusters of similar incident reports that help in determining patterns, trends, and best practices.
- Case based reasoning offers a useful methodology for identifying similar medical incident reports, but has accuracy problems.
- The integration of case based reasoning and information retrieval greatly improves the recall and accuracy in clusters of similar medical incident reports.

ACKNOWLEDGEMENTS

We would like to thank Barbara Rabin Fastman from Columbia University's New York Presbyterian Hospital and Quay Mercer currently with University of Texas Southwestern Medical Center at Dallas for assisting in the evaluation of our system. We would also like to thank the anonymous reviewers for improving the exposition of our work. The work described in this paper was supported in part by grant RO1 HL53772 from the National Heart, Lung and Blood Institute of the National Institutes of Health, through a subcontract from Columbia University.

.....

Authors' affiliations

C Tsatsoulis, H A Amthauer, Department of Electrical Engineering and Computer Science, The University of Kansas, Lawrence, KS 66045, USA

REFERENCES

- 1 **Kaplan HS**, Callum JL, Rabin Fastman B, *et al*. The medical event reporting system for transfusion medicine (MERS-TM): will it help us get the right blood to the right patient? *Transfusion Med Rev* 2000;**16**:86–102.
- 2 **Kolodner J**. *Case-based reasoning*. San Mateo, CA: Morgan Kaufmann, 1993.
- 3 **Salton G**. *The SMART retrieval system—experiments in automatic document processing*. Englewood Cliffs, NJ: Prentice Hall, 1971.
- 4 **Salton G**, McGill MJ. *Introduction to modern information retrieval*. New York: McGraw-Hill, 1983.
- 5 **van der Schaaf TW**. *Near miss reporting in the chemical process industry*. Eindhoven, Netherlands: Technical University of Eindhoven, 1992.
- 6 **Wilson DC**, Bradshaw S. CBR textuality. *Expert Update* 2000;**3**:28–37.