

A Comparative Evaluation of Name-Matching Algorithms

L. Karl Branting

LiveWire Logic, Inc.

2700 Gateway Centre Blvd., Ste. 900

Morrisville, NC 27650, USA

karl.branting@livewirelogic.com

ABSTRACT

Name matching—recognizing when two different strings are likely to denote the same entity—is an important task in many legal information systems, such as case-management systems. The naming conventions peculiar to legal cases limit the effectiveness of generic approximate string-matching algorithms in this task. This paper proposes a three-stage framework for name matching, identifies how each stage in the framework addresses the naming variations that typically arise in legal cases, describes several alternative approaches to each stage, and evaluates the performance of various combinations of the alternatives on a representative collection of names drawn from a United States District Court case management system. The best tradeoff between accuracy and efficiency in this collection was achieved by algorithms that standardize capitalization, spacing, and punctuation; filter redundant terms; index using an abstraction function that is both order-insensitive and tolerant of small numbers of omissions or additions; and compare names in a symmetrical, word-by-word fashion.

1. INTRODUCTION

An important task in many legal information systems is *name matching*, recognizing when two different strings are intended to denote the same entity. Many government agencies and private-sector companies are required by law to compare names on documents, such as passports and credit cards, to watch lists of suspected terrorists and their supporters [1]. In legal case management systems (LCMSs), name matching is important in performing searches and in detecting redundant entries. In judicial case management systems, such as the United States federal court system's CM/ECF system [2], name matching is needed in detecting potential conflicts of interest. Conflicts of interest can arise when an attorney has a personal stake in the outcome of a case or if a judge has any connection to a case that might affect the judge's objectivity. For example, United States federal judges are required to disqualify, or "recuse," themselves "in any proceeding in which [their] impartiality might reasonably be questioned," such as when a judge "has a personal bias or prejudice concerning a party, or personal knowledge of disputed

evidentiary facts concerning the proceeding ..." or if a relative of the judge has a financial interest in the controversy [3].

Potential conflicts of interest can be detected by comparing the contents of a *conflict file*—which contains names of entities that could give rise to potential conflicts—to the names of the parties and attorneys in cases that are candidates for assignment to the attorney or judge. If there is a match between an entry in the conflict file and a party or attorney in a case, the potential conflict can be flagged and the case reassigned. The effectiveness of this scheme depends on the accuracy with which names in a conflict file can be recognized as denoting the same entity as a name occurring in the party or attorney field of a case record.

Unfortunately, errors and stylistic inconsistencies can lead a single legal entity to be designated by multiple distinct expressions. For example, expressions denoting a single entity may differ in word order, spelling, spacing, punctuation, and use of abbreviations or organizational terms (such as "LLP" or "Ltd"). A direct comparison of names occurring in conflict files with names in case files may therefore fail to detect a significant proportion of potential matches. An effective name-matching algorithm for legal applications must be capable of recognizing that two expressions potentially designate the same entity notwithstanding naming variations typical of the application.

In view of the importance of name matching, one might expect there to be a substantial literature concerning appropriate algorithms for this task. Surprisingly, however, name-matching algorithms are generally proprietary and therefore not available for independent evaluation, comparison, or improvement. There is an extensive literature on the general problem of sequence matching, much of it directed to text retrieval, spelling correction, and computational molecular biology [4,5]. However, these general-purpose sequence-matching algorithms were designed to overcome typographical or genetic transcription errors rather than the naming variations peculiar to legal cases. Intuitively, one would expect that algorithms capable of exploiting knowledge of typical naming variations could achieve higher efficiency and accuracy than general-purpose matching algorithms.

This paper analyzes the name-matching task as it arises in LCMSs, identifies the naming variations characteristic of LCMSs, proposes a three-stage framework for performing the name-matching task that exploits knowledge of these characteristic variations, describes several alternative approaches to each stage, and evaluates the performance of various combinations of the alternatives on a collection of names drawn from a United States District Court LCMS.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICAIL '03, June 24-28, 2003, Edinburgh, Scotland, UK.

Copyright 2003 ACM 1-58113-747-8. \$5.00.

2. THE NAME-MATCHING TASK

The information-processing requirements of the name-matching task are as follows:

Given:

- A pattern string
- A collection of target strings

Do:

- Find each target that matches the pattern well enough that there is significant likelihood that the pattern and target denote the same entity.

In conflict detection, each pattern is an entry in a conflict file, and the target strings consist of the names of parties, attorneys, and law firms associated with a case.

A target string returned by a name-matching algorithm is termed a *match* or a *positive*. If the match does in fact denote the same entity as the pattern, the match is a *true positive*, whereas if the match does not denote the same entity as the pattern it is a *false positive*. A target string that is not a match is a *negative*. A *false negative* is an unmatched target string that denotes the same entity as the pattern, and a *true negative* is an unmatched target string that does not.

Two distinct types of evaluation criteria for name matching can be distinguished: match accuracy and computational efficiency. Match accuracy is a function of two complementary measures of performance:

- Precision, the proportion of matching target strings that denote the same entity as the pattern, and
- Recall, the proportion of entities denoting the same entity as the pattern that were matched.

Precision is equal to $|\text{true positives}| / (|\text{true positives}| + |\text{false positives}|)$, that is, the proportion of actual matches that should in fact have been matched. Recall is equal to $|\text{true positives}| / (|\text{true positives}| + |\text{false negatives}|)$, that is, the proportion of targets that should have been matched that were in fact matched.

There is a tradeoff between recall and precision. If every target is matched, recall will be 100%, but precision may be very low. Conversely, if only identical strings are matched, precision will be 100%, but recall may be very low. The most desirable algorithm is one that optimizes this tradeoff, i.e., making recall as high as possible without sacrificing precision. To express this optimization, recall and precision can be combined into a single measure of overall performance, such as the F-measure [6]. If recall and precision are weighted equally, the F-measure is the harmonic mean of recall and precision:

$$F = \frac{2PR}{P + R}$$

where P is precision and R is recall.

The potential size of LCMS databases places a practical upper bound on the computational cost of name-matching algorithms. A single, incremental change to a conflict file (or case database)

requires a number of comparisons proportionate to the size of the database (or, respectively, conflict file). However, judicial LCMSs are typically required to perform a conflict-detection screening for the entire conflict file and case database. This can potentially entail a very large number of comparisons. For example, one U.S. Federal District Court LCMS studied by the author had 12,890 conflict-file entries and 268,104 strings occurring in cases. A usable name-matching algorithm must be capable of comparing files of this size in (at most) minutes; other demands on LCMS hardware make day-long computations unacceptable.

3. NAME MATCHING AS METRIC-SPACE PROXIMITY QUERYING

The name-matching task is an instance of the more general task of proximity querying in metric spaces. A metric space consists of a universe U of objects on which is defined a distance function

$$d: U \times U \rightarrow \mathbb{R}$$

that satisfies the triangle inequality. Fulfilling a proximity query with query q and range r requires retrieving every element e of a database for which the distance $d(q,e) < r$. Proximity queries in metric spaces arise in a wide range of domains, including molecular biology, multimedia, pattern recognition, chemical databases, and spell correction [7].

Many metric spaces have straightforward representations as vector spaces in which geometric distance corresponds to similarity under the metric. In such cases, spatial access methods, such as k - d trees, are well-suited for proximity queries (provided that the number of dimensions is not too high) [8]. However, in many domains there is no obvious vector space representation that is proximity-preserving, that is, that preserves the similarity ordering in the original metric space. Development of proximity query methods for such metric spaces is an active research area.

In general, metric-space proximity query methods involve creation of an indexing structure based on a set of k elements selected as pivots [7]. The elements of the database are partitioned into equivalence classes based on distance to each of the pivots. A query q is satisfied by determining the distance from q to each pivot, then using the triangle inequality to select equivalence sets that are guaranteed to include all elements within distance r of q (as well as potentially other elements as well). The distance between q and each element of the equivalence sets is then calculated to determine which of the candidate elements is actually within distance r [7]. The computational expense of constructing the index is typically proportional to the cost of calculating the distance between a pair of elements times kn , where n is the number of elements in the database

The name-matching task as it occurs in many LCMSs is characterized by a relatively expensive distance calculation. As described in Section 4, the equivalence criteria for pairs of names includes a number of complex, domain-specific rules. Moreover, queries can involve large numbers of both texts and patterns and are sometimes limited to a very short execution time. As a result, metric-space indexing schemes that require applying the distance calculation kn times to build an index can be prohibitively expensive in such systems.

When pivot-based indices are too costly to construct, an alternative approach is to partition the database into equivalence classes using an abstraction method that is much cheaper than calculating the distance to k pivots. For example, in the two-stage retrieval methods developed for analogical reasoning systems, such as MAC/FAC, labeled graphs are abstracted into bags of labels. The expensive subgraph-isomorphism similarity metric is calculated only on those candidates sharing the highest proportion of labels with the query [9]. This approach, in which the candidates for match with a query q are those elements in the same equivalence class as q , will be referred to as *non-pivot indexing*.

The next section describes the categories of orthographic variations characteristic of names occurring in legal case management systems. Section 5 proposes a three-stage model that performs non-pivot indexing using computationally inexpensive methods to abstract strings into equivalence classes of strings that are likely to be equivalent with respect to the variations described in Section 4.

4. A TAXONOMY OF NAME VARIATIONS

The difficulty of the name-matching task, and the requirements for an effective algorithm to perform this task, depend on the type and degree of name variations typical in the collection to which the name-matching algorithm is applied. To determine the name variations typical of LCMSs, an informal analysis was performed of a United States federal district court database containing 41,711 name occurrences.

Nine primary categories of variations were apparent within this database:

1. Punctuation, e.g., “Owens Corning” vs. “Owens-Corning”; “IBM” vs. “I.B.M.”
2. Capitalization, e.g., “citibank” vs. “Citibank”; “SMITH” vs. “Smith”
3. Spacing, e.g., “J.C. Penny” vs. “J. C. Penny”
4. Qualifiers, e.g., “Jim Jones” vs. “Jim Jones d.b.a. Jones Enterprises”
5. Organizational terms, e.g., “corporation” vs. “incorporated”
6. Abbreviations, e.g., “cooperative” vs. “coop”; “General Motors” vs. “GM”
7. Misspellings:
 - a. Omissions, e.g. “Collin” vs. “Colin”
 - b. Additions, e.g., “McDonald” vs. “MacDonald”
 - c. Substitutions, e.g., “Smyth” vs. “Smith”
 - d. Letter reversals, e.g., “Peirce” vs. “Pierce”
8. Word omissions, e.g., “National Electrical Benefit Fund and its Trustees” vs. “National Electrical Benefit Fund”
9. Word permutations, e.g., “State of Missouri District Attorney” vs. “District Attorney, State of Missouri”

While it is impossible to determine precisely the relative frequency of these variations without a systematic analysis of errors occurring in a representative collection of LCMSs, informal inspection of the district court database suggests that word omissions and variations in capitalization, punctuation, spacing, abbreviations, and organizational terms are relatively common. Word permutations are somewhat less common, and misspellings and qualifier variations are relatively rare.

Recognizing the similarity between pairs of expressions is computationally straightforward for some of the variations. For example, the similarity between a string representing a correctly spelled word and a string with minor spelling errors can be recognized using standard dynamic programming techniques to determine the minimum edit distance between the strings [10]. However, these techniques are ill-suited to variations in organizational terms, qualifiers, word omissions, and differences in word order.

5. ALGORITHMS FOR NAME MATCHING

This section distinguishes three stages in the name-matching process, identifies the stages at which each of the nine naming variations described above can be addressed, and distinguishes alternative design options for each of the stages.

5.1 Stages of Name Matching

Three distinct stages can be distinguished in the name-matching task: normalization, indexing, and similarity assessment.

Normalization is the process of transforming pattern and target strings into a standard form by eliminating inessential textual variations that can prevent matching. Normalization operations include adopting a standard convention for capitalization (e.g., all uppercase or all lowercase), punctuation (e.g., removing all punctuation), and stop-word filtering (e.g., removing uninformative, common words, such as “the” and “LLC”).

Indexing is the process of selecting a set of candidates from the targets for comparison with the pattern. The simplest indexing method is exhaustive retrieval, that is, selection of the entire set of targets for comparison with the pattern. Alternatively, each string can be abstracted into a simplified representation that can be used to index strings through a hash table, decision tree, or other retrieval mechanism. The motivation for abstraction is that multiple similar strings may have the same abstraction. If the abstraction of a pattern is used to index target strings with the same abstraction, only a small number of comparisons will be needed for each pattern. This approach can reduce the number of comparisons without compromising accuracy if pairs of strings intended to denote the same entity usually have the same abstraction.

Exhaustive retrieval is too slow for any but the smallest target sets. In the example mentioned above, a direct comparison between each of 12,890 conflict-file entries and 268,104 case strings would require 3,455,860,560 comparisons. Even if each individual comparison were very fast, the entire process would be unacceptably slow.

Similarity Assessment is the process of determining whether there is sufficient similarity between a normalized pattern and a normalized target to indicate a significant probability that the target designates the same entity as the pattern. Approaches to comparison can differ in granularity—whether the comparison is word-by-word or on the entire string—and match criterion. Possible match criteria include string equality, a sub-string relationship between the pattern and the target, and approximate matching, which consists of determining whether the edit distance between the pattern and target is less than a given mismatch threshold. The edit distance between two strings consists of the number of insertions, deletions, or substitutions required to make the pattern equal to the target. Testing for string equality or sub-string matching is relatively efficient. The approximate matching task is inherently more computationally expensive than exact matching, with the time complexity of the dynamic programming algorithm proportional to the product of the lengths of the strings being compared [4].

5.2 Name Variations Addressed at Each Stage

The first five sources of variation enumerated in Section 3—punctuation, capitalization, spacing, qualifiers, and organizational terms—can all be addressed by normalization. No normalization scheme is likely to be entirely infallible, however, for two reasons. First, corporate names sometimes consist entirely of organizational terms or stop-words that in other contexts can cause mismatches. For example, if a company were named “U.S. Association of Corporations,” all the words in the company’s name would be filtered if “U.S.,” “Association,” “of,” and “Corporation” were all stop words, resulting in an unmatchable null string. If one or more of the terms were excluded from the stop list, however, some matches might not be detected, e.g., “Smith Corporation” might fail to match “Smith Incorporated” if “Corporation” or “Incorporated” were not on the stop list. Second, irregular spacing in some acronyms can make them indistinguishable from stop words, e.g., if “American Tomato Originators Network” were abbreviated “ATON,” extra spaces could generate “A TO N” and “AT ON,” which appear to contain stop words (“A,” “TO,” and “ON”). Tables of common corporate names (e.g., “AT&T”) can reduce, but not eliminate, this problem.

The 6th source of variations—abbreviations—can be addressed during normalization or during similarity assessment by using a table of abbreviations to recognize the equivalence between abbreviated and unabbreviated forms.

The 7th variation source—misspellings—can be addressed by approximate matching, a technique specifically developed to recognize omissions, additions, and substitutions. As discussed below, approximate matching can easily be extended to transpositions of adjacent letters.

Word omissions, the 8th source of variations, can be addressed by a word-by-word similarity-assessment procedure under which a pattern matches a target if each word of the pattern matches a unique word in the target. As discussed below, this similarity assessment can be either symmetrical, meaning that every word in the string with the fewest words must match a word in the other string, or asymmetrical, meaning that every word in the pattern

must match a unique word in the target (but not necessarily the converse).

Symmetrical matching seems desirable, but might give rise to the danger of false positives from extraneous text in the party or attorney fields of cases. For example, if “Attorney” appeared as an attorney field in a case and matching were symmetrical, every conflict record containing the words “Attorney” would be matched (unless “Attorney” were a stop word). One could argue that quality control is, in general, much easier in conflict files than in case files and that it is therefore more important to match every part of a conflict record than to match every part of a case record. This argument suggests that an asymmetrical match policy would be preferable. On the other hand, the possibility that attorneys and judges might include nonessential text in conflict-file entries suggests that asymmetrical matching risks false negatives. In a system that can filter recurring false-positives, it may be better to err on the side of false positives rather than false negatives, that is, to weigh recall more heavily than precision. The experiments below evaluate the relative performance of symmetrical and asymmetrical matching.

Word-by-word similarity assessment can also address the 9th source of variations—word permutations—if the similarity-assessment procedure does not constrain words to appear in the same position in pattern and target.

5.3 Design Options for Each Stage

A variety of algorithms representing different combinations of design options for normalization, indexing, and similarity assessment are possible. The following algorithms were implemented in the evaluation described in Section 6:

Exact-Match. The exact-match algorithm is intended as a benchmark for name-matching speed. The only normalization is conversion to upper case, removal of all punctuation, and normalization of spaces, which consists of trimming beginning and ending white space, replacing multiple successive spaces, and removing spaces in abbreviations, e.g., replacement of “I.B. M.” with “I.B.M.”. Candidates are indexed by hashing on the normalized string, and similarity assessment consists of testing for string equality. The processing time for exact-match represents a lower bound on the time required for a reasonable job of matching.

Palmer. Doug Palmer, a U.S. District Court system administrator, implemented a modification of the exact-match algorithm intended to improve efficiency. In Palmer’s modification, normalization consists of capitalization, punctuation removal, and removal of stop-words. Indexing is by hashing on an abstraction formed by removing vowels, double letters, and terminal “s”s. There is no further similarity assessment after retrieval, i.e., every retrieved candidate is assumed to be a match.

Palmer’s normalization and abstraction often yields an empty string. A policy issue concerns how strings with an empty abstraction should be treated. One approach is to treat every target with an empty abstraction as a candidate for matching with every pattern with an empty abstraction. Alternatively, strings with empty abstractions can be treated as matching nothing. Empirical evaluation indicated that the former approach leads to large

numbers of spurious candidates. As a result, the latter policy was used in the experiment described below.

All the remaining algorithms use identical normalization, consisting of capitalization, removal of all punctuation, space normalization, and removal of stop-words. Each of the algorithms uses a different combination of choices for abstraction, granularity, symmetry, and similarity assessment.

Abstraction. Four options for abstraction were implemented. In each abstraction method, target strings were stored in a hash-table entry indexed by each string's abstraction. Since multiple strings can have the same abstraction, the hash table entries consisted of lists of target strings.¹

Soundex is a phonetic encoding developed by the U.S. Bureau of Census and used to index all individuals listed in the U.S. census records starting in 1880 [12]. Soundex encodes each string as the first letter of the string followed by 3 numbers representing the phonetic categories of the next 3 consonants, if any, in the string.²

Unordered-sounds. A limitation of Soundex is that the abstraction it produces is dependent on word order. As a result, permutations of identical words have different Soundex encodings. For example, "Social Services Dept., State of Alaska" has a Soundex encoding of S243, whereas the encoding of "State of Alaska Social Services Dept." is S331. Unordered-sounds is a variant of Soundex whose encoding is independent of word order. Specifically, Unordered-sounds encodes a multiple-word string in 19 bits that indicate the category of sounds that occur in the 1st, 2nd, or 3rd positions of any words. The first 7 bits indicate whether any word in the string starts with the corresponding one of Soundex's 6 categories of letters or with a letter disregarded by Soundex (A, E, I, O, U, H, W, Y). The next 12 flags indicate whether a letter in any of the 6

categories occurs in the second or third position of any word.

Nsoundex. Unordered-sounds has the disadvantage that an omission or addition of a single word can cause two strings to have different encodings if any of the first 3 sounds of the word occur in a different position in some other word. Nsoundex is a variant on Soundex intended to address Soundex's order sensitivity without introducing Unordered-sounds' sensitivity to extra words. Nsoundex removes stop words and sorts the remaining words alphabetically before applying Soundex. An extra target word will prevent indexing only if the extra word starts with a letter earlier in the alphabet than the first word in the pattern.

Redundant is similar to Nsoundex except that each string is redundantly indexed by both the first and last words in the sorted, normalized, stop-word-free string. If the similarity-assessment procedure uses approximate matching, the Soundex of the first and last words are used as indices for the string; if string equality is used for similarity assessment, the words themselves are used as indices. Redundant is less sensitive to omitted or extra words than Nsoundex but incurs the added cost of indexing every string twice.

Granularity and Symmetry. Three approaches to granularity and symmetry were implemented:

Entire-string consists of similarity assessment of the entire pattern with the entire target, after stop words have been filtered from both.

Word-by-word consists of splitting the normalized pattern and target strings into individual words. Word-by-word comparison can be symmetrical or asymmetrical.

Asymmetrical. After stop-words are removed from both lists of words, each pattern word is compared to every target word in turn until a word is found that satisfies the applicable similarity assessment criterion (discussed *infra*) or which exactly matches the standard abbreviation of the pattern word. The abbreviation table is based on the abbreviations found in *The Bluebook: A Uniform System of Citation* [13] and *The Chicago Manual of Style* [14]. Each target word is permitted to match only a single pattern word, and isolated letters (such as initials in names) are required to match exactly. Under this approach, the pattern "John Jones" would match target "John Q. Jones", but pattern "John Q. Jones" would match neither "John Jones" nor "John A. Jones".

Symmetrical matching succeeds if every string in the shorter name matches a string in the longer name, regardless of order. Under

¹ Not tested in this experiment, but worth considering, is lexical vector-space indexing [11] applied to names represented as a bags of string abstractions. For example, each name could be indexed as a bag of tokens produced by applying Soundex to each word in the name. This approach might increase recall, although it would significantly increase retrieval time.

² The categories of consonants are:

- 1) B,P,F,V
- 2) C,S,K,G,J,Q,X,Z
- 3) D,T
- 4) L
- 5) M,N
- 6) R

Vowels are ignored and adjacent letters from the same category are represented with a single digit. For example, "Washington" would be encoded as "W252": W is the first letter, 2 for the S, 5 for the N, 2 for the G, and the remaining letters disregarded.

this approach, the pattern “John Jones” would match target “John Q. Jones”, and pattern “John Q. Jones” would match “John Jones”. However, “John A. Jones” would not match “John Q. Jones”.

Similarity-Assessment Criterion. Two similarity-assessment criteria were used in the test described below:

String equality.

Approximate match. Dynamic programming was used for approximate matching with a modification so that a separate penalty could be assigned for reversals of pairs of adjacent letters. The motivation for this modification is that letter reversals are a common typing mistake. In all experiments, a penalty of 1.0 was assigned to insertions, deletions, and substitutions, and a penalty of 0.6 was assigned to letter reversals. The mismatch threshold was set at 15% of the number of letters in the pattern when used in word-by-word matching, meaning that a match would succeed in a word of at least 7 letters if there were a single insertion, deletion, or substitution, and a word of 5 letters or more would match if there were a single reversal of a pair of adjacent letters. In entire-string similarity assessment, however, the mismatch threshold was set to 10% of the number of letters in the pattern.

Each combination of choices for abstraction, granularity, symmetry, and similarity assessment constitutes a distinct name-matching algorithm. Each algorithm is identified below by an acronym consisting of a concatenation of the first letter of the name of the algorithm’s abstraction method, granularity, symmetry, and similarity assessment criterion. For example, Soundex abstraction combined with Word-by-word granularity, Asymmetrical matching, and Approximate matching is denoted “SWAA.”

Suppose that SWAA were called with pattern “Jones Environmental Systems and Service Corporation” and a set of targets that includes “Jones Env. Services Systems, Inc.” The corresponding normalized strings would be “JONES ENVIRONMENTAL SYSTEMS SERVICE CORPORATION” and “JONES ENV SERVICES SYSTEMS INC”, respectively. The target string “JONES ENV SERVICES SYSTEMS INC” would be indexed by its Soundex encoding of J525. The pattern also has a Soundex encoding of J525, so the target would be retrieved for matching. The first pattern word, “JONES”, matches the first target word perfectly. The second pattern word, “ENVIRONMENTAL” doesn’t match any word in the target, but its abbreviation, “ENV”, matches the second word of the target. The third pattern word, “SYSTEMS”, matches the fourth target word. The fourth pattern word, “SERVICE”, is an approximate match to “SERVICES”, with an edit distance of 1 deletion and 1 addition. Finally, “CORPORATION” is a stop word that does not need to be matched. The pattern therefore matches the target. Note that the difference in word order is irrelevant for the matching performed by SWAA and that words in the target but not in the pattern are simply ignored.

6. EXPERIMENTAL EVALUATION

6.1 Procedure

To identify the best combination of design options, the performance of Exact-Match, Palmer, NEA (Nsoundex, entire-word comparison, approximate-match), RWSE (Redundant, word-by-word symmetrical, exact), RWSA (Redundant, word-by-word symmetrical, approximate), and every combination of {Nsoundex, Soundex, Unordered} \times {word-by-word symmetrical, word-by-word asymmetrical} \times {exact-match, approximate} was evaluated on a static copy of a U.S. District Court database containing 41,711 records for cases assigned to 20 judges. Unfortunately, this database did not include actual conflict files. An artificial conflict file was therefore created for each of 20 judges by randomly selecting 700-800 entries from the case records for each judge. This resulted in a total of 15,478 conflict file entries.

The testing procedure copied each judge’s conflict records into one array and the party and attorney fields of all cases assigned to that judge into a second array. Each algorithm in turn was called to determine the matches between the two arrays. The execution of the algorithm was timed, the total number of matches counted, and matches between non-identical strings stored in a match file for that algorithm (since matching identical strings is trivial and is performed equally well by all algorithms, these matches were not included in the calculation of precision, recall, and F-measure).

After every algorithm was tested, the match files for all algorithms were merged and recorded in a file called approx-matches, containing all non-identical strings returned as a match by any algorithm. For each algorithm, the elements of approx-matches not found by that algorithm were written into that algorithm’s miss-file.

To estimate true and false positives and false negatives, the approx-matches file was manually edited to tag its entries as true or false positives. The contents of each algorithm’s miss-file were compared to approx-matches to determine that algorithm’s true and false positives and apparent false negatives (i.e., matches not found by the algorithm that were found by some other algorithm). Only the apparent false negatives could be determined under this procedure because there was no oracle to determine whether there were any targets that should have been matched but were missed by all of the algorithms.

6.2 Results

Figure 1 sets forth the F-measure of each algorithm. The highest F-measure was obtained by the two algorithms that used redundant indexing: RWSA and RWSE (i.e., redundant, word-by-word symmetric, and approximate or exact, respectively). RWSA had higher recall, but lower precision, than RWSE. In general, algorithms that used word-by-word, symmetrical similarity assessment outperformed equivalent algorithms that used asymmetrical or entire-word similarity assessment. Approximate matching yielded much higher recall, but lower precision, than string equality, leading to little difference in F-measure between approximate and exact matching. Exact-match had the lowest F-measure because of its low recall.

Figure 1. F-measure of match algorithms.

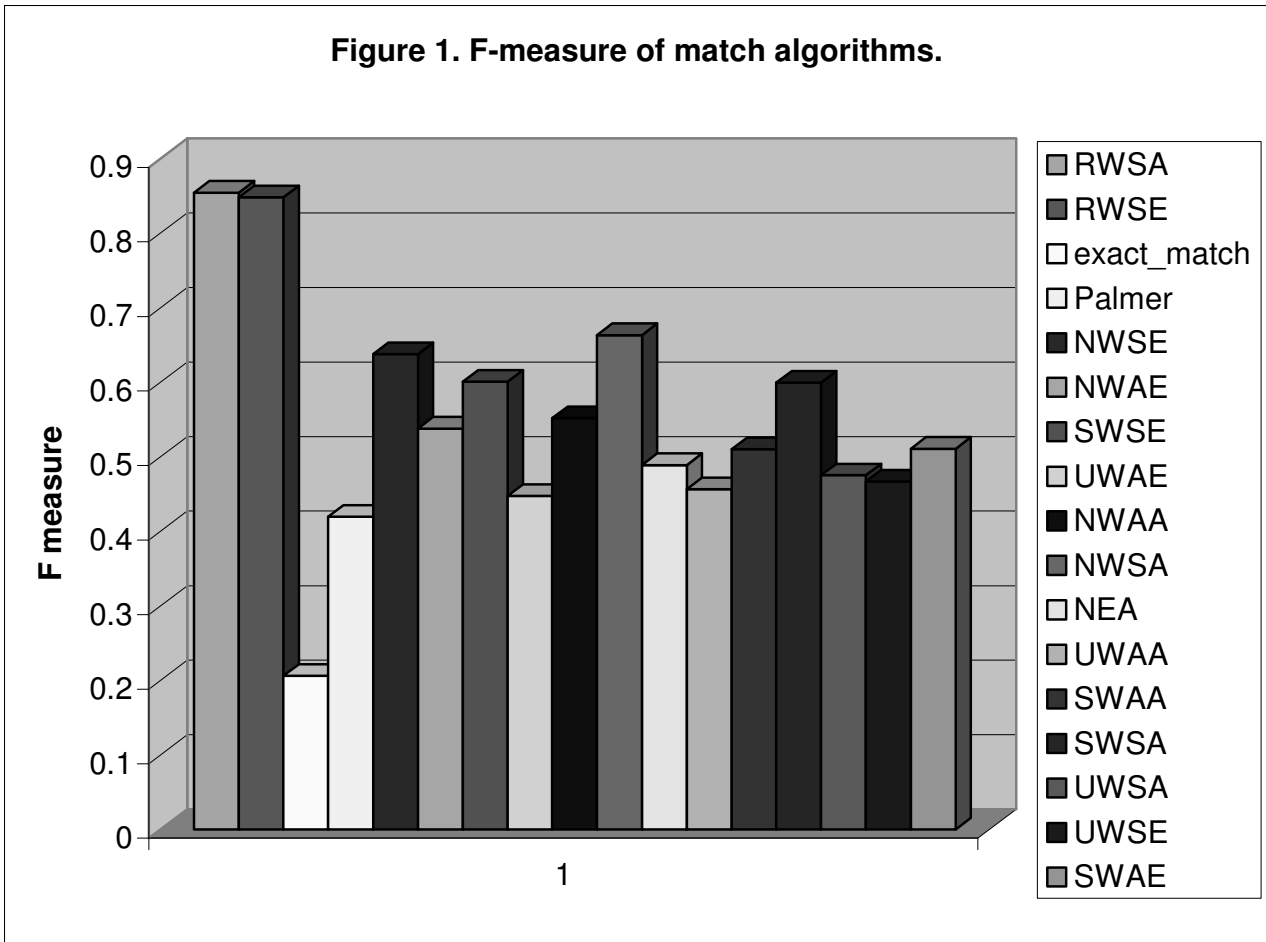


Figure 2 shows the computation time of the same set of algorithms, normalized by the computation time of exact-match (i.e., the computation time of each algorithm was divided by the computation time of exact match).³ The slowest algorithm was RWSA, because it performs many more similarity assessments than the algorithms with non-redundant indexing. The next slowest algorithm was NEA, illustrating the high computational cost of entire-word approximate matching.⁴ RWSE had almost the same accuracy as RWSA but was more than four times as fast.

6.3 Discussion

The most striking aspect of the results set forth in Figures 1 and 2 is that the algorithms with redundant indexing had much higher recall—and therefore much higher F-measure—than any algorithms with non-redundant indexing. The likely explanation for this phenomenon is that every partition of a metric space into two or more equivalence classes places some pairs of points separated by less than any arbitrary distance r into separate classes. As a result, no non-pivot indexing method that uses only

one abstraction method—and therefore only one partition of the metric space—can have perfect recall.

It is entirely possible, however, for there to be a pair of partitions of a given metric space such that every pair of elements closer than some threshold r share some equivalence class. For example, Figure 3 represents schematically universe U containing query q and elements $e1$ and $e2$, both of which are within distance r of q . Abstraction 1 generates a partition under which q and $e1$ are in the same equivalence class but q and $e2$ are not. Abstraction 2 generates a partition in which the converse is the case. Neither abstraction is sufficient in itself for perfect recall, but the union of both equivalence classes containing q also contains both $e1$ and $e2$.

The redundant indexing methods, RWSE and RWSA, each use two independent abstraction methods that produce two independent partitions of the target strings. Taking candidates from both equivalence classes of which a pattern is a member—that is, selecting targets that are identical to the pattern under either abstraction—greatly reduces the probability of missing a target that matches the pattern. These results strongly suggest that high recall in non-pivot metric spaces requires multiple, independent abstractions.

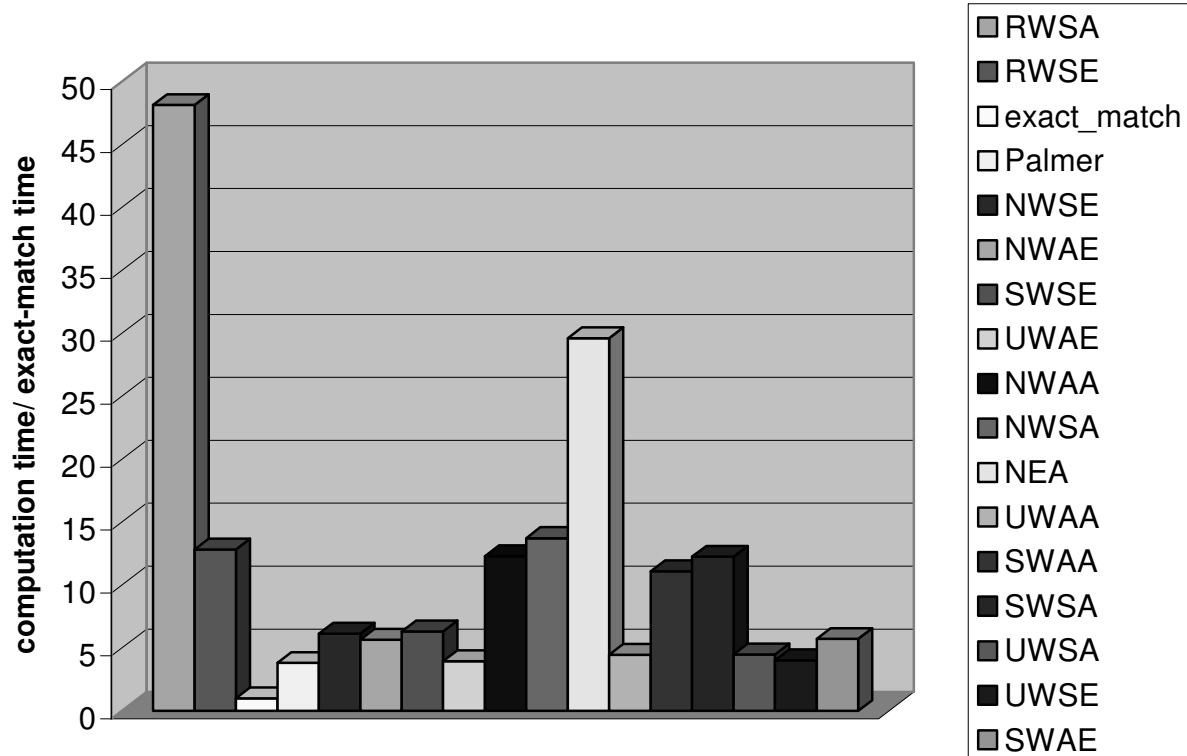
7. SUMMARY

The most accurate name-matching algorithm for name matching in a given LCMS depends on the relationship between the choices attorneys and judges make in expressing their potential conflicts and the conventions governing party and attorney names in case captions in that LCMS. If there is a high degree of consistency

³ The actual run time for exact-match, running in interpreted Perl on an elderly, multi-user Sun, was approximately 5 seconds. On more modern machines, the run time for files of comparable size should be considerably lower.

⁴ As mentioned above, the computational complexity of the dynamic programming algorithm for approximate matching through is proportionate to the product of the lengths of the two strings being compared [4]. Pair-wise approximate matching between partitions of two strings is therefore much faster than approximate matching between the entire strings.

Figure 2. Computation time of match algorithms.



between conflict records and case records, matching is straightforward. To the extent that there are variations, matching algorithms can be expected to achieve satisfactory precision and recall only to the extent that they embody matching techniques that compensate for those variations.

The results of the empirical evaluation suggest that, in databases with name variations similar to those occurring in the U.S. District Court database, name-matching accuracy is maximized by matching algorithms with the following characteristics:

- Normalization by capitalization, removal of all punctuation, space-normalization, abbreviation-replacement, and stop-word removal.
- Indexing using an abstraction function that is order-insensitive and tolerant of small numbers of omissions or additions in the strings being matched. Redundant indexing appears to achieve these goals better than Nsoundex, Soundex, or Unordered-sounds.
- Symmetrical, word-by-word similarity assessment.
- If time is not critical and recall is much more important than precision, approximate matching should be used. If time is critical or recall is no more important than precision, string equality should be used instead.

RWSE and RWSA had the highest F-measures on the U.S. District Court database, and RWSE (unlike RWSA) was relatively fast, *i.e.*, only about 13 times slower than exact-match. In view of these experimental results, RWSE was adopted as the name-matching algorithm in CM/ECF in early 2002.

The experimental results are tentative because of two factors: (1) uncertainty concerning the typicality of the name variations occurring in the U.S. district court database that was the source of the data used in the evaluation and (2) the absence of a definitive list of false negatives, *i.e.*, target strings that should have matched the pattern but which were matched by no algorithm. A more conclusive evaluation of the relative accuracy of alternative name-matching algorithms must await the collection of more data on name variations from a representative sampling of LCMSs. Creation of publicly available datasets, in the spirit of the UCI machine-learning data repository [15] would significantly advance the development of name-matching algorithms by permitting replicable evaluation of alternative algorithms.

ACKNOWLEDGMENTS

The research described in this paper was performed while the author was a United States Supreme Court Judicial Fellow at the Administrative Office of U.S. Courts. However, the views expressed in this paper are those of the author and do not necessarily reflect the views of the Judicial Conference or the Administrative Office of U.S. Courts.

8. REFERENCES

- [1] S. Milstein, *Taming the Task of Checking for Terrorists' Names*, The New York Times, C4, Monday, December 30, 2002.
- [2] Judicial Conference of the United States, *Electronic Case Files in the Federal Courts: A Preliminary Examination of Goals, Issues, and the Road Ahead*, Administrative Office of the U.S. Courts (March 1997).
- [3] 28 USCS § 455. Canon 3 of the American Bar Association's Model Code of Judicial Conduct (1990) sets forth a similar standard.
- [4] D. Gusfield, *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*, Cambridge University Press (1997).
- [5] J. Setubal and J. Meidanis, *Introduction to Computational Molecular Biology*, PWS Publishing Co. (1997).
- [6] C. van Rijsbergen, *Information Retrieval*. London: Butterworths, 2nd Edition (1979).
- [7] E. Chávez, G. Navarro, R. A. Baeza-Yates, and J. Marroquín, Searching in metric spaces, *ACM Computing Surveys*, 33(3):273-321 (2001).
- [8] V. Faede and O. Guenther, Multidimensional access methods. *ACM Computing Surveys*, 30(2):170-231 (1992).
- [9] K. Forbus, D. Gentner, and K. Law, MAC/FAC: A model of Similarity-based Retrieval. *Cognitive Science*, 19(2):141-205 (1995).
- [10] P. Sellers, The theory and computation of evolutionary distances: pattern recognition. *Journal of Algorithms*, 1:359-373 (1980).
- [11] C. Manning & H. Schütze, *Foundations of Statistical Natural Language Processing*. The MIT Press (1999).
- [12] For an introduction to Soundex, see *The Soundex Indexing System*, National Archives and Records Administration, <http://www.nara.gov/genealogy/coding.html>.
- [13] Harvard Law Review Association, 17th Edition (2000).
- [14] University of Chicago Press, 13th Edition (1982).
- [15] C. Blake and C. Merz, UCI Repository of machine learning databases [<http://www.ics.uci.edu/~mllearn/MLRepository.html>]. Irvine, CA: University of California, Department of Information and Computer Science.