

1           A New Iterative Algorithm for Computing a Quality  
2           Approximate Median of Strings based on Edit Operations.

3   J. Abreu<sup>\*a</sup>, J. R. Rico-Juan<sup>b</sup>

4           <sup>a</sup>*Dpto Informática, Universidad de Matanzas, Carretera a Varadero Km. 3 1/2, Matanzas, Cuba*

5           <sup>b</sup>*Dpto Lenguajes y Sistemas Informáticos, Universidad de Alicante, San Vicente del Raspeig, Alicante,*  
6   *Spain*

---

7           **Abstract**

This paper presents a new algorithm which can be used to compute an approximation to the median of a set of strings. The approximate median is obtained through the successive improvements of a partial solution. The edit distance from the partial solution to all the strings in the set is computed in each iteration, thus accounting for the frequency of each of the edit operations in all the positions of the approximate median. A goodness index for edit operations is later computed by multiplying their frequency by the cost. Each operation is tested, starting from that with the highest index, in order to verify whether applying it to the partial solution leads to an improvement. If successful, a new iteration begins from the new approximate median. The algorithm finishes after all the operations have been examined without a better solution being found. Comparative experiments involving Freeman chain codes encoding 2D shapes and the Copenhagen chromosome database show that the quality of the approximate median string is similar to benchmark approaches but achieves a much faster convergence.

8           *Key words:* approximate median string, edit distance, edit operations

---

9           **1. Introduction**

10           Extending the concept of “median” to structural representations such as strings  
11           has been a challenging issue in Pattern Recognition for some time, as is shown in the  
12           review presented in Jiang et al. (2004). This problem arises in many applications such

---

\*Corresponding author. Fax:965909326

13 as 2D shape representation and prototype construction (Jiang et al., 2000; Bunke et al.,  
 14 2002), the clustering of strings (Lourenço and Fred, 2005), Self-Organized Maps of  
 15 strings (Kohonen, 1998; Fischer and Zell, 2000) or the combination of multiple source  
 16 translations (González-Rubio and Casacuberta, 2010).

17 Formally, given a set  $S = \{S_1, S_2, \dots, S_n\}$  of strings over the alphabet  $\Sigma$  and a  
 18 distance function  $D(S_i, S_j)$  which measures the dissimilarity between strings  $S_i$  and  
 19  $S_j$ , the distance from a string  $S'$  to all strings in  $S$  can be computed by the expression  
 20 (1).

$$21 \quad SOD(S') = \sum_{S_i \in S} D(S', S_i) \quad (1)$$

22 The *median string* is the string  $\hat{S} \in \Sigma^*$  that minimizes (1). This string is also  
 23 denoted as the *generalized median string*. A common approximation to the true median  
 24 string is the *set median*, a string in  $S$  which minimizes (1). It is not necessary for either  
 25 the median string or the set median to be unique.

26 An exact algorithm with which to compute the median of a set of strings was pro-  
 27 posed by Kruskal (1983). However, in most practical applications this is not a suitable  
 28 approach due to the high computational time requirements. As Casacuberta and An-  
 29 tonio (1997) and Nicolas and Rivals (2005) point out, there are various formulations  
 30 of this problem within the NP-Complete class. Several approximations have therefore  
 31 been proposed. One approach that has been studied by several authors is that of build-  
 32 ing the approximate median by using the successive changes of an initial string. Per-  
 33 turbations can be applied one or more at a time, as in the works of Martínez-Hinarejos  
 34 et al. (2003) and Fischer and Zell (2000), respectively. The results of empirical testing  
 35 show that the first approach leads to high quality approximations but requires more  
 36 computational time. The principal motivation of this work is to describe a new algo-  
 37 rithm that is able to compute a quality approximation to the median string like that  
 38 of Martínez-Hinarejos et al. (2003), but requires significantly less computational ef-  
 39 fort. In the Section 2 some related works are examined. Section 3 describes the pro-

40 posed approach and provides an analysis of the computational cost bounds for the algo-  
41 rithm. Various comparative experiments are described in Section 4, and finally Section  
42 5 shows our conclusions and some lines for further research.

## 43 **2. Related works**

44 Many approximated solutions have been described since Kruskal (1983) proposed  
45 an exact algorithm that could be used to compute the median string for a given set  
46  $S$  of  $N$  strings of a length of  $l$  and the Levenshtein (Levenshtein, 1966) metric. This  
47 algorithm runs in  $O(l^N)$  proportional time. A number of heuristics therefore address  
48 this difficulty by reducing the size of the search space. Some authors, such as Oli-  
49 vares and Oncina (2008), have studied the approximation to the median string not only  
50 under the Levenshtein edit distance but also under the stochastic edit distance (Ristad  
51 and Yianilos, 1998). In other works, the search for the approximate median is not per-  
52 formed directly in the string space but rather in a vectorial space in which the strings  
53 are embedded; this is the approach studied in Jiang et al. (2012) which also relies on  
54 the weighted median concept described by Bunke et al. (2002).

55 One general strategy is to construct the approximate median letter by letter from an  
56 initial empty string. In order to decide which symbol is the next to be appended it is  
57 necessary to define a goodness function. The greedy procedure described in Casacu-  
58 berta and Antonio (1997) implements this approach. An improvement to the aforemen-  
59 tioned method is described in Kruzslicz (1999) through the use of a refined criterion  
60 which allows the next letter to be selected. Another approach that has been studied by  
61 several authors is that of building the approximate median by using successive pertur-  
62 bations of an initial string. Two important issues regarding this kind of method are how  
63 to select a perturbation leading to an improvement and how to make the algorithm con-  
64 verge faster without spoiling the results. Another interesting topic is that of studying  
65 the effect of performing modifications one by one or simultaneously. Kohonen (1985)

66 starts from the set median and systematically changes the guess string by applying  
67 insertions, deletions and substitutions in every position. In Martínez-Hinarejos et al.  
68 (2003) the authors proposes to improve a partial solution  $\hat{S}$  generating new candidates  
69 by applying all possible substitutions, insertions or deleting the symbol at a position  $i$ .  
70 The new partial solution is the string, selected from all the new candidates and  $\hat{S}$ , which  
71 minimizes (1). This procedure is repeated for every position  $i$ . The effect of choosing  
72 a different initial string as the set median or a greedy approximation is also studied.  
73 Theoretical and empirical results show that this method is capable of achieving very  
74 good approximations to the true median string. Note that these methods do not define  
75 a criterion to compare the operations in order to select which one can lead to better  
76 results in each case. In Martínez-Hinarejos et al. (2002) authors describe alternatives  
77 to speed up the computation of the approximated median string. Based on information  
78 provided by the weight matrix used to compute the edit distance certain operations are  
79 preferred instead others. For example, no all possible substitutions are tried but only  
80 the two closest symbols to the one in the analysed position.

81 Some heuristic knowledge which can help to assess how promising a modification  
82 will be are included in Fischer and Zell (2000) and Mollineda (2004). The quality  
83 of a partial solution  $\hat{S}$  is evaluated by computing its distance from every string in the  
84 set, and it is thus also possible to discover the sequences of edit operations. In an  
85 attempt to speed up the convergence of the search procedure, these authors propose the  
86 simultaneous performance of several modifications by applying the most frequent edit  
87 operation, including “do nothing” in each position of the partial solution. This process  
88 is repeated while modifications increase the quality of the partial solution.

89 This approach has two potential drawbacks. Applying the most common operation  
90 in every position does not guarantee the best results, and a further issue is that although  
91 it might be relatively simple to figure out how applying just one operation will affect  
92  $SOD(\hat{S})$ , this does not hold when several changes are made at the same time. For

93 example, let  $\hat{S}$  be a partial solution and  $op_i$  be an edit operation which occurs several  
 94 times when computing the distance from a partial solution to strings in  $S$ .  $Op_i$  thus  
 95 determines a subset  $S^{YES} \subseteq S$  of those strings in which  $op_i$  occurs when computing  
 96 the distance from  $\hat{S}$ . There is also another set  $S^{NO} = S - S^{YES}$ . Let  $\hat{S}'$  be a new  
 97 solution after applying  $op_i$  to  $\hat{S}$ . Intuitively, it may be expected that the distance from  
 98  $\hat{S}'$  to strings in  $S^{YES}$  decreases as regards  $\hat{S}$ . A formal discussion of this result can be  
 99 found in Bunke et al. (2002). The effect on the strings in  $S^{NO}$  clearly needs to be taken  
 100 into account. Since sets induced by each operation may be different when applying  
 101 multiple operations, it might be very difficult to characterize the effect on  $SOD(\hat{S})$ .  
 102 Empirical results, which will be discussed later, suggest that methods which apply  
 103 multiple perturbations at the same time are able to find a better approximation than  
 104 the set median very quickly. However, approaches which perform modifications one  
 105 by one, such as Martínez-Hinarejos et al. (2003), significantly outperform the former  
 106 methods with respect to the average distance to the set of the approximate median  
 107 computed.

### 108 **3. A new algorithm for computing a quality approximate median string**

109 As was noted earlier, a general scheme that can be used to search for an approximate  
 110 median string is:

- 111 - select an initial coarse approximation to the median, as the set median.
- 112 - generate a new solution by performing some modifications to the current solu-  
 113 tion.
- 114 - repeat while a particular modification leads to an improvement or another stop  
 115 condition holds.

116 The works commented on Section 2 suggest that when it is necessary to find a  
 117 quality approximation to the median string, applying modifications one by one would

118 appear to be a better strategy. The theoretical results in Jiang and Bunke (2002) and  
119 Martínez-Hinarejos (2003) show that the approximation computed by the algorithm  
120 proposed in Martínez-Hinarejos et al. (2003) is very close to the lower bound obtained  
121 for the value of  $SOD(\hat{S})$  for the true median.

### 122 3.1. Computing the approximate median string

123 The algorithm in Martínez-Hinarejos et al. (2003) tests every possible operation in  
124 each position of the partial solution, and it might therefore be very useful to study how  
125 to reduce the size of the search space without spoiling the quality of results, which is  
126 one of the principal motivations of this work. The proposed algorithm is based on two  
127 main ideas:

- 128 - selecting the appropriate modification by paying attention to certain statistics  
129 from the computation of the edit distance from the partial solution to every string  
130 in the set.
- 131 - applying modifications one by one.

132 Heuristic information could help to avoid to test a number of useless solutions,  
133 which would reduce the amount of times that  $SOD(\hat{S})$  is evaluated. Another distinctive  
134 feature is that if the best operation according to the goodness index does not lead to an  
135 improvement, other low ranked operations can be tested.

136 The *AppMedianString* procedure outlines how to compute the approximate median  
137 string.

### 138 3.2. Selecting the best edit operation

139 In our case, the most suitable edit operation in step  $t$  will be selected by examining  
140 two approaches. The first simply implies ranking operations by their *frequency* while  
141 computing the edit distance from the partial solution to strings in the set, as in Fischer  
142 and Zell (2000). Note that the selected operation is that with the best overall ranking,

---

```

Function AppMedianString(S,R) : $\hat{S}$ 
/* S: instance set to compute the approximate median. */
/* R: initialization string. */
R' = R;
repeat
   $\hat{S} = R'$ ;
  foreach instance  $s_i \in S$  do
    compute  $D(R', s_i)$ ;
    obtain that  $Q_{s_i}^{R'}$  is the minimum cost edit sequence needed to transform
      R' into  $s_i$ ;
    update statistics for the operation in each position  $j$  of  $R'$ ;
  end foreach
  let  $O_p$  be an operation queue sorted by its goodness index;
  /* Generate new candidates  $R'$  while none of them improve  $\hat{S}$  */
  while  $\sum_{s_i \in S} D(\hat{S}, s_i) \leq \sum_{s_i \in S} D(R', s_i)$  and  $O_p \neq \emptyset$  do
     $op_i = O_p.dequeue$ ;
    obtain a new candidate  $R'$  applying  $op_i$  to  $\hat{S}$ ;
  end while
until no operation  $op_i$  applied to  $\hat{S}$  improve the result;
return  $\hat{S}$ ;

```

---

143 not the most frequent in a specific position. However, under a more general weighting  
144 scheme for edit operations, the frequency might not be the best assessment of how  
145 promising a transformation is. We therefore propose the use of  $Frequency \times Cost$  as  
146 a goodness index. For example, let  $\hat{S}^t$  be the candidate solution and  $S = \{S_1, S_2, S_3\}$ .  
147 Without loss of generality, let us suppose that the best ranked edit operation ( $op_1$ ) is  
148 a substitution with a frequency of 2, and cost of 1. Let us also suppose that there  
149 is another substitution ( $op_2$ ) with a frequency of 1 but with a cost of 3. From the  
150 results in Bunke et al. (2002) we obtain that an  $\hat{S}^{t+1}$  built by applying  $op_1$  will satisfy  
151  $D(\hat{S}^{t+1}, S_1) = D(\hat{S}^t, S_1) - 1$  and  $D(\hat{S}^{t+1}, S_2) = D(\hat{S}^t, S_2) - 1$ . Regardless of the value  
152 of  $D(\hat{S}^{t+1}, S_3)$  it can be expected that  $SOD(\hat{S})$  will decrease by 2. A similar analysis  
153 shows that the application of  $op_2$  leads to a reduction of 3.

154 *3.3. An illustrative example*

155 The following example illustrates the algorithm's behavior. Let  $\hat{S}^t = \{5, 5, 0\}$ ,  $S_1 =$   
156  $\{3, 1, 1, 2\}$  and  $S_2 = \{0, 6, 1, 6\}$ . The substitution of a symbol  $a$  for  $b$  obtain the cost  
157  $\min\{|a-b|, 8-|a-b|\}$ , while insertions and deletions obtain the cost of 2. Table 1 shows  
158 the computation of the edit distance from  $\hat{S}^t$  to  $S_1$  and  $S_2$ . In the first case, this results  
159 in one of the optimal edit sequences  $\{s(5, 3), s(5, 1), s(0, 1), i(2)\}$ .  $D(\hat{S}^t, S_2)$  results in  
160  $\{s(5, 0), s(5, 6), s(0, 1), i(6)\}$ . Table 2 shows an edit operation ranked by its frequency.  
161 Note how a different goodness index leads to a different ranking. Applying the best  
162 operation  $s(0, 1)$  in position 3 results in  $\hat{S}^{t+1} = \{5, 5, 1\}$ , which improves  $SOD(\hat{S})$   
163 since  $D(\hat{S}^{t+1}, S_1) = 8$  and  $D(\hat{S}^{t+1}, S_2) = 6$ . If the best operation does not lead to  
164 an improvement, then the second best option must be tested, and so on. Note that in  
165 the list of perturbations there may be different operations related to the same position.  
166 This option does not occur in Fischer and Zell (2000) and Mollineda (2004). The  
167 process is repeated by starting from the new solution while some operations lead to  
168 a better approximation. The example above also shows how ranking by *Frequency*  $\times$   
169 *Cost* can lead to better results. As was explained previously, by applying  $s(0, 1)$  we  
170 obtain  $SOD(\hat{S}^{t+1}) = 14$ . The last column in the Table 2 shows that the operations may  
171 be ranked differently. In this case,  $s(5, 1)$  in position 2 is the operation with the best  
172 goodness index. If it were to be applied, then  $\hat{S}^{t+1} = \{5, 1, 0\}$  and thus  $D(\hat{S}^{t+1}, S_1) = 5$   
173 and  $D(\hat{S}^{t+1}, S_2) = 5$ , which is  $SOD(\hat{S}^{t+1}) = 10$ .

174 *3.4. Computational cost analysis*

175 The procedure used to compute the approximate median string needs to compute  
176 the distance from the partial solution to every string in the set. Under the Levenshtein  
177 edit distance this can be carried out in time  $O(l^2)$  by using the dynamic programming  
178 algorithm presented in Wagner and Fischer (1974), where  $l$  is the length of the longest  
179 string. The **foreach** statement loops  $N$  times, and the first stage of the algorithm thus



Table 1: Computation of the edit distance cost from  $\hat{S}^t = \{5, 5, 0\}$  to  $S_1 = \{3, 1, 1, 2\}$  and  $S_2 = \{0, 6, 1, 6\}$ . Substitutions of a symbol  $a$  by a symbol  $b$  have cost  $\min\{|a - b|, 8 - |a - b|\}$  while deletions and insertions have cost of 2. An optimal path is shaded in order to follow the best cost operations easily and visually.

(a)					(b)						
		3	1	1	2			0	6	1	6
	0	2	4	6	8		0	2	4	6	8
5	2	2	4	6	8	5	2	3	3	5	7
5	4	4	6	8	9	5	4	5	4	6	6
0	6	6	5	7	9	0	6	4	6	5	7

Table 2: Ranking of edit operations

Operation	Position	Frequency	Frequency $\times$ Cost
s(0,1)	3	2	2
s(5,0)	1	1	3
s(5,1)	2	1	4
s(5,6)	2	1	1
s(5,3)	1	1	2
i(2)	3	1	2
i(6)	3	1	2

180 requires a time that is proportional to  $O(N \times l^2)$ . Assuming that no perturbations im-  
 181 prove the solution, the inner **while** loop needs to examine the whole queue  $O_p$ .

182 Let  $|\Sigma|$  be the size of the alphabet;  $\min\{N, |\Sigma|\}$  substitutions are possible for each  
 183 of the  $l$  symbols in  $\hat{S}$ , this is the maximum number of substitutions, and there are  
 184 thus  $O(l \times \min\{N, |\Sigma|\})$  potential substitutions. The same result holds for insertions.  
 185 Only  $l$  deletions are possible. A pessimistic upper bound to  $|O_p|$  is therefore  $O(2 \times$   
 186  $l \times \min\{N, |\Sigma|\} + l)$ . In the worst case, each operation in  $O_p$  involves computing the  
 187 distance from  $R'$  to all the strings, which requires  $O(N \times l^2)$ . Under these assumptions,  
 188 inner **while** takes a time proportional to  $O(N \times l^3 \times \min\{N, |\Sigma|\})$ . Let  $k$  be the number of  
 189 times that the outer **repeat** loops, thus the algorithm requires  $O(k \times N \times l^3 \times \min\{N, |\Sigma|\})$ ,  
 190 which is the same time required by the algorithm described by Martínez-Hinarejos et al.  
 191 (2001). However, in practice the proposed approach behaves much better as suggest  
 192 results which will be discussed in Section 4.

#### 193 4. Experimental results

194 Experiments were carried out to evaluate the performance of the proposed approach  
195 when computing an approximate median string. To ensure independent results with  
196 regard to the alphabet, the strings over two sets of symbols were tested. In the first case,  
197  $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, \lambda\}$ , corresponding to the directions of Freeman chain codes  
198 (Freeman, 1974) where  $\lambda$  denotes the empty symbol used for deletions and insertions.  
199 Edit operation costs were fixed in a manner similar to that of Rico-Juan and Micó  
200 (2003), that is, a cost of 2 for deletions and insertions and  $\min\{|a - b|, 8 - |a - b|\}$   
201 for substitutions. The strings in each set are not randomly generated but are a chain  
202 code representation of the contours from two widely known 2D shape databases, the  
203 *NIST-3 Uppercase Letters* and the *USPS Digits*, (Jain and Zongker, 1997; García-Díez  
204 et al., 2011; Rico-Juan and Iñesta, 2012), with 26 and 10 classes, respectively. Four  
205 independent samples of 20 instances per class were drawn for a total of 144 different  
206 sets. Our approach was used to compute an approximate median for each of them. The  
207 proposed algorithm, referred to as *JR-S*, was compared to the methods proposed by  
208 Fischer and Zell (2000) and Mollineda (2004) which performs several modifications  
209 at the same time, and that of Martínez-Hinarejos (2003) which modifies the partial  
210 solution in a one by one manner.

211 In a second test, strings were drawn from the chromosomes dataset used by Martínez-  
212 Hinarejos et al. (2003). This time  $\Sigma = \{a, b, c, d, e, =, A, B, C, D, E, \lambda\}$ , and the cost of  
213 each operation was computed as in Martínez-Hinarejos et al. (2003). Four samples of  
214 20 instances were again selected for each of the 22 classes.

215 Tables 3 and 4 show the results for each set in the respective databases. In order  
216 to facilitate the comparison of the results of different algorithms and datasets, in each  
217 case we compute the ratio  $\frac{SOD(\hat{S})}{SOD(S^M)}$ , where  $S^M$  is the set median. The lower it is, the  
218 better the approximation to the true median found by the algorithm is. In each case “ $\varepsilon$ ”,  
219 “ $S^M$ ” or “ $S^G$ ” refer to the initial string, that is, the empty string, the set median and

220 the greedy initialization proposed by Casacuberta and Antonio (1997). Since all the  
221 algorithms in the test work in an iterative manner, the number of distances computed  
222 by each approach that evaluates  $SOD(\hat{S})$  was also studied. The graphics in Figure 1 and  
223 2 show the average value for  $\frac{SOD(\hat{S})}{SOD(S^M)}$  and the average number of distances computed  
224 by each approach in all the experiments.

225 Besides, a third experiment was carried out to compare the results with respect to  
226 the true median. In this case we collect four sets of 20 random generated strings over  
227 the alphabet  $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, \lambda\}$  with length varying from 3 to 8. Operation  
228 costs were fixed as explained before. Table 5 shows results on this simple database.

229 As mentioned previously, the results confirm that applying perturbations to the par-  
230 tial solution one by one leads to a much better quality approximation to the true median  
231 in terms of  $SOD(\hat{S})$ . In every set, either the proposed approach or Martínez-Hinarejos  
232 (2003) provides the most precise approximation. In general, the solutions computed  
233 with *JR-S* are equivalent to or even better than those attained with Martínez-Hinarejos  
234 (2003) but, as Tables 3 and 4 show, the proposed approach is, on average, about 10  
235 times faster than Martínez-Hinarejos (2003) in terms of the computed distances. In  
236 some cases ranking the operations by  $Frequency \times Cost$  instead  $Frequency$  can lead to  
237 slightly better approximations, but in general, it also requires the computation of addi-  
238 tional distances. On the other hand, although its results are not so good in terms of the  
239 approximate median quality in the methods of Fischer and Zell (2000) and Mollineda  
240 (2004), only a few distances are needed to notably improve the set median. In both  
241 cases it would appear that the algorithm gets stuck in a local minimum after a small  
242 number of iterations.

243 A comparison in terms of running time was also included, as Figure 3 shows. The  
244 experiments were performed in a computer with an Intel X5355-2.66 GHz CPU (4  
245 cores) and 8 Gb RAM. It can be observed that algorithms introduced by Fischer and  
246 Zell (2000) and Mollineda (2004) are in average about 30 times faster than ours. On the

247 other hand, the proposed approach runs near 8 times faster than the methods described  
 248 by Martínez-Hinarejos et al. (2003).

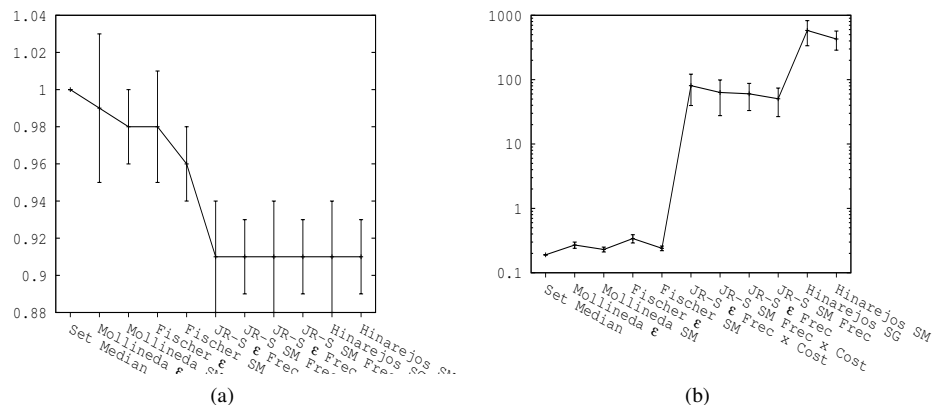


Figure 1: 1a shows the average for  $\frac{SOD(\hat{S})}{SOD(S^M)}$  in all experiments. This measure represents the quality of the results. The chart in 1b shows the average number of distances (in thousands) (Freeman chain codes set). In both cases, less value is better.

## 249 5. Conclusions and Future work

250 A new approach with which to compute a quality approximation to the median  
 251 string has been presented. The algorithm builds an approximate median through the  
 252 successive refinements of a partial solution. Modifications are applied one by one in a  
 253 manner similar to that of Martínez-Hinarejos et al. (2003), and empirical results show  
 254 that this approach leads to better approximations than those methods which apply sev-  
 255 eral perturbations simultaneously, although the latter runs much faster. Comparisons  
 256 with Martínez-Hinarejos (2003) show that the proposed algorithm is able to compute  
 257 high-quality approximations to the true median string, but requires significantly less  
 258 computation, and is about 10 times faster, which makes it highly suitable for applica-  
 259 tions that require a precise approximation. As was pointed in Section 2, an operation  
 260  $op_i$  determines two subsets  $S^{YES}$  and  $S^{NO}$  from  $S$ . Applying  $op_i$  to  $\hat{S}$  results in new  
 261 string  $\hat{S}'$  such as the distance from strings in  $S^{YES}$  to  $\hat{S}'$  will decrease. Further research

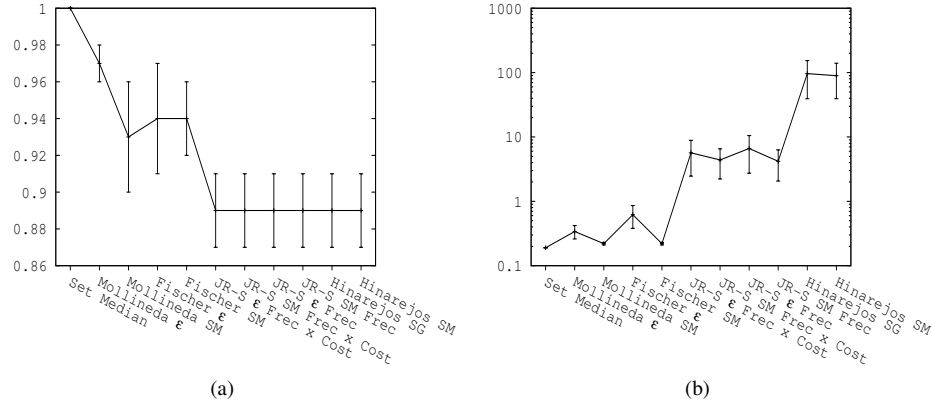


Figure 2: 2a shows the average for  $\frac{SOD(\hat{S})}{SOD(\hat{S}^M)}$  in all experiments. This measure represents the quality of the results. The chart in 2b shows the average number of distances (in thousands) (Copenhagen chromosomes set). In both cases, less value is better.

262 may address to better characterize how behaves the distance from  $\hat{S}'$  to strings in  $S^{NO}$   
 263 without computing those distances, but using information gathered when computing  
 264 the distances to  $\hat{S}$ . This can help to select the better operation to reduce the number of  
 265 distances computed without spoiling the approximation quality. Another subject of in-  
 266 terest is that of analysing how the choice of a different optimal path will affect results,  
 267 since a different ranking might be obtained.

268 **Acknowledgements**

269 This work is partially supported by the Spanish CICYT under project DPI2006-  
 270 15542-C04-01, the Spanish MICINN through project TIN2009-14205-CO4-01 and by  
 271 the Spanish research program Consolider Ingenio 2010: MIPRCV (CSD2007-00018).

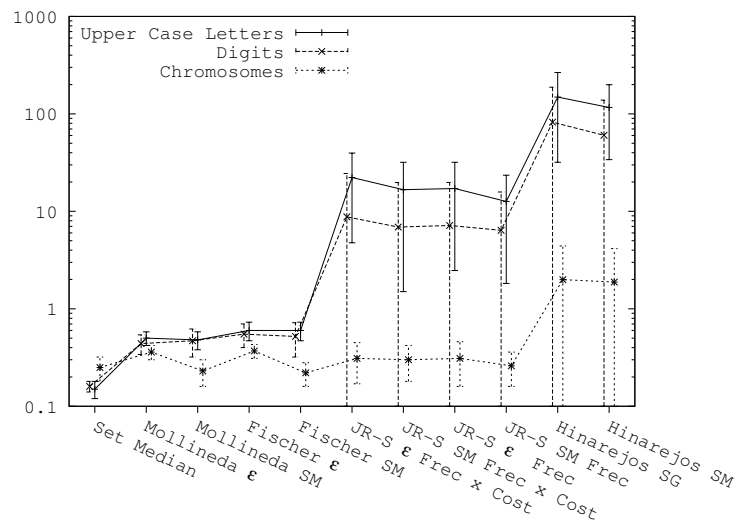


Figure 3: Average running time, in seconds, for each algorithm in each database experiments.

Table 3: Average distance from the approximated median to each string in the set. (Freeman chaincodes)

Class	Set Median	Mollineda $\epsilon$	Mollineda $S^M$	Fischer $\epsilon$	Fischer $S^M$	JR-S $\epsilon$	JR-S $S^M$	JR-S $\epsilon$	JR-S $S^M$	Hinarejos $S^G$	Hinarejos $S^M$
A	102.0±5.0	98.7±6.0	98.0±6.0	99.3±5.0	96.4±5.0	92.2±5.0	92.4±5.0	92.1±5.0	92.2±5.0	92.3±5.0	92.1±5.0
B	120.5±4.0	122.3±7.0	119.7±4.0	120.8±5.0	119.0±4.0	110.3±5.0	110.1±5.0	110.2±5.0	110.2±5.0	110.2±5.0	110.1±5.0
C	116.2±13.0	110.7±10.0	111.6±12.0	109.9±11.0	109.9±11.0	103.6±10.0	103.5±10.0	103.3±10.0	103.3±10.0	103.3±10.0	103.3±10.0
D	126.3±29.0	131.5±29.0	126.3±29.0	128.8±30.0	126.3±29.0	117.1±28.0	117.2±28.0	117.1±28.0	117.0±28.0	117.0±28.0	117.2±28.0
E	181.8±12.0	191.1±19.0	176.0±8.0	174.7±7.0	172.9±7.0	168.5±12.0	165.9±8.0	168.3±12.0	165.5±8.0	168.6±12.0	165.5±8.0
F	154.2±8.0	149.0±10.0	148.1±8.0	146.2±9.0	145.9±7.0	137.0±7.0	137.1±7.0	136.9±7.0	137.0±7.0	137.0±7.0	137.2±7.0
G	190.1±17.0	184.7±17.0	184.8±16.0	180.9±17.0	180.5±16.0	171.0±16.0	170.2±16.0	170.8±16.0	170.6±16.0	170.6±16.0	170.4±16.0
H	193.3±19.0	192.5±20.0	193.1±19.0	189.2±19.0	189.0±19.0	176.7±18.0	176.5±18.0	176.6±18.0	176.6±19.0	176.6±19.0	176.6±18.0
I	141.6±17.0	155.1±20.0	141.6±17.0	150.0±19.0	141.6±17.0	137.3±18.0	137.3±18.0	137.7±18.0	137.6±18.0	137.5±18.0	137.3±18.0
J	184.9±14.0	180.1±13.0	182.0±17.0	182.6±8.0	178.6±13.0	167.0±12.0	167.2±12.0	166.9±11.0	166.9±11.0	167.0±11.0	166.9±11.0
K	197.1±17.0	195.4±21.0	190.7±17.0	186.8±17.0	186.6±15.0	175.1±14.0	175.0±15.0	175.1±15.0	175.5±14.0	179.3±20.0	175.2±15.0
L	89.7±5.0	89.4±4.0	86.9±6.0	87.8±5.0	85.9±6.0	82.4±5.0	82.7±5.0	82.4±5.0	82.4±5.0	82.5±5.0	82.5±5.0
M	225.6±16.0	220.2±11.0	218.4±11.0	214.6±14.0	214.1±12.0	201.4±11.0	201.2±11.0	201.2±11.0	201.5±11.0	201.6±12.0	201.4±11.0
N	191.8±11.0	192.1±10.0	190.7±10.0	188.8±10.0	187.9±11.0	178.8±9.0	178.8±9.0	178.9±9.0	178.6±9.0	179.0±9.0	178.7±9.0
O	75.4±11.0	78.5±16.0	74.8±12.0	76.8±14.0	74.6±12.0	70.7±12.0	70.7±12.0	70.8±12.0	70.6±12.0	70.6±12.0	70.4±12.0
P	100.5±10.0	100.9±13.0	100.4±10.0	99.1±13.0	99.5±11.0	92.2±10.0	92.1±10.0	92.0±10.0	92.2±10.0	92.3±10.0	92.3±10.0
Q	109.7±9.0	115.0±10.0	109.4±8.0	107.0±8.0	107.3±8.0	100.4±7.0	100.5±7.0	100.2±7.0	100.3±7.0	100.2±7.0	100.2±7.0
R	150.5±10.0	149.0±11.0	150.0±10.0	145.9±12.0	147.5±9.0	135.6±9.0	135.4±10.0	135.5±9.0	135.7±10.0	135.4±9.0	135.4±9.0
S	143.3±7.0	141.7±6.0	138.8±6.0	137.8±7.0	136.8±6.0	130.3±6.0	130.2±7.0	130.3±7.0	130.2±6.0	130.3±7.0	130.2±7.0
T	144.7±11.0	149.9±11.0	144.4±12.0	146.0±14.0	144.6±11.0	136.3±11.0	136.4±11.0	136.4±11.0	136.5±11.0	136.2±11.0	136.2±11.0
U	171.0±6.0	179.0±4.0	170.7±6.0	177.7±4.0	170.3±5.0	172.7±21.0	163.0±4.0	172.9±21.0	163.0±4.0	169.1±13.0	162.9±4.0
V	146.6±14.0	142.8±14.0	142.6±15.0	141.3±14.0	140.7±12.0	134.5±13.0	134.6±13.0	134.9±13.0	134.7±13.0	135.0±13.0	134.4±13.0
W	226.8±19.0	221.5±14.0	221.5±13.0	220.4±14.0	217.4±16.0	207.1±14.0	206.0±13.0	206.8±14.0	206.2±13.0	208.3±17.0	206.0±13.0
X	149.7±11.0	152.4±12.0	149.2±12.0	149.2±12.0	147.8±10.0	139.8±12.0	140.0±12.0	140.0±12.0	139.9±12.0	139.6±12.0	139.6±12.0
Y	147.3±4.0	150.3±6.0	147.3±4.0	146.2±5.0	144.9±6.0	136.7±5.0	137.0±4.0	136.6±4.0	136.9±4.0	136.7±5.0	136.6±5.0
Z	172.6±11.0	172.8±11.0	170.8±9.0	171.6±8.0	167.1±9.0	157.7±9.0	157.8±9.0	157.8±9.0	157.9±9.0	157.8±9.0	157.9±9.0
0	54.1±2.0	51.6±3.0	51.6±3.0	53.7±2.0	51.6±3.0	48.1±3.0	48.2±3.0	48.4±3.0	48.0±3.0	48.0±3.0	47.9±3.0
1	46.6±2.0	45.6±2.0	45.6±2.0	49.4±2.0	45.8±2.0	42.8±2.0	42.9±1.0	42.7±1.0	42.8±1.0	43.0±1.0	42.6±2.0
2	122.1±8.0	114.1±7.0	115.1±7.0	116.7±7.0	113.6±8.0	107.9±6.0	107.7±7.0	107.8±6.0	107.9±7.0	107.8±6.0	107.6±6.0
3	122.7±3.0	112.2±4.0	110.8±4.0	114.1±2.0	111.4±4.0	105.5±3.0	105.8±3.0	105.6±3.0	105.5±3.0	105.7±3.0	105.5±3.0
4	147.0±12.0	145.1±11.0	145.1±11.0	144.0±11.0	143.6±13.0	134.9±11.0	135.0±11.0	135.1±11.0	135.0±11.0	134.8±11.0	135.4±11.0
5	155.4±6.0	144.4±6.0	144.9±5.0	145.8±4.0	143.5±5.0	134.9±4.0	135.3±4.0	135.4±4.0	135.1±4.0	135.1±4.0	135.1±4.0
6	87.8±3.0	83.0±2.0	82.5±2.0	87.5±2.0	82.4±2.0	76.7±2.0	77.2±2.0	77.1±2.0	77.1±2.0	76.9±2.0	77.0±2.0
7	102.4±5.0	98.7±5.0	96.7±6.0	101.4±4.0	97.5±5.0	91.3±4.0	91.5±4.0	91.7±4.0	91.3±5.0	91.4±4.0	91.3±5.0
8	103.5±7.0	102.8±8.0	99.5±9.0	101.8±7.0	100.0±8.0	92.8±7.0	93.1±8.0	92.9±8.0	92.9±8.0	93.1±7.0	93.0±7.0
9	85.5±5.0	83.8±4.0	83.2±5.0	84.4±4.0	81.5±5.0	77.0±4.0	77.1±5.0	76.9±4.0	77.0±4.0	76.8±5.0	76.8±5.0
Average	138.3	137.4	135.1	135.5	133.4	126.2	125.9	126.3	125.9	126.3	125.8
$\sigma$	44.0	44.1	43.5	42.2	42.4	40.6	40.1	40.6	40.2	40.7	40.2

Table 4: Average distance from the approximated median to each string in the set. (Copenhagen Chromosomes set)

Class	Set Median	Mollineda $\epsilon$	Mollineda $S^M$	Fischer $\epsilon$	Fischer $S^M$	JR-S $\epsilon$ Freq $\times$ Cost	JR-S $S^M$ Freq $\times$ Cost	JR-S $\epsilon$ Freq	JR-S $S^M$ Freq	Hinarejos $S^G$	Hinarejos $S^M$
chromo1	47.8±3.0	49.8±3.0	44.8±2.0	57.6±4.0	44.6±2.0	42.4±3.0	42.0±3.0	42.2±3.0	42.1±3.0	42.1±3.0	42.0±3.0
chromo2	42.8±1.0	47.5±3.0	40.5±2.0	65.3±5.0	40.8±2.0	37.7±2.0	37.9±2.0	37.8±2.0	37.7±2.0	38.0±2.0	37.8±2.0
chromo3	38.7±1.0	45.4±6.0	36.4±0.5	50.4±5.0	36.2±1.0	34.2±1.0	34.2±1.0	34.3±1.0	34.3±1.0	34.5±1.0	34.5±1.0
chromo4	36.2±1.0	37.0±2.0	33.1±1.0	52.6±1.0	33.6±1.0	31.8±1.0	31.9±1.0	31.9±1.0	31.8±1.0	31.8±1.0	31.8±1.0
chromo5	33.1±1.0	37.1±0.3	30.7±1.0	51.8±3.0	31.1±2.0	28.9±1.0	28.8±1.0	28.8±1.0	28.8±1.0	28.8±1.0	28.7±1.0
chromo6	33.8±2.0	36.7±1.0	32.3±2.0	46.5±7.0	30.5±1.0	29.8±1.0	29.8±1.0	29.9±1.0	29.8±1.0	29.8±1.0	29.8±1.0
chromo7	29.5±1.0	37.0±3.0	27.1±0.2	42.2±5.0	27.2±1.0	25.9±1.0	25.9±1.0	25.9±1.0	25.9±1.0	25.9±1.0	25.9±1.0
chromo8	27.6±1.0	32.1±1.0	24.8±0.1	44.7±1.0	25.1±1.0	24.0±1.0	23.9±1.0	23.9±1.0	23.8±1.0	23.9±1.0	23.9±1.0
chromo9	28.2±2.0	28.2±1.0	28.2±2.0	28.2±3.0	28.2±2.0	28.2±1.0	28.2±1.0	28.2±1.0	28.2±1.0	28.2±1.0	28.2±1.0
chromo10	25.9±1.0	30.2±1.0	25.5±1.0	39.5±6.0	24.9±1.0	23.2±1.0	23.2±1.0	23.2±1.0	23.2±1.0	23.3±1.0	23.2±1.0
chromo11	24.6±2.0	26.4±5.0	21.0±2.0	33.3±5.0	22.3±2.0	21.3±2.0	21.3±2.0	21.3±2.0	21.3±2.0	21.2±2.0	21.3±2.0
chromo12	25.4±1.0	27.0±1.0	24.5±1.0	33.8±4.0	23.9±1.0	22.6±1.0	22.7±1.0	22.6±0.5	22.7±0.5	22.6±0.5	22.8±0.5
chromo13	20.4±1.0	24.9±3.0	19.5±1.0	35.2±3.0	18.7±1.0	18.2±1.0	18.1±1.0	18.1±1.0	18.1±1.0	18.0±1.0	18.1±0.5
chromo14	23.2±1.0	22.0±0.5	21.8±0.5	31.1±5.0	21.4±0.4	20.5±1.0	20.5±1.0	20.6±1.0	20.6±1.0	20.5±1.0	20.6±1.0
chromo15	21.0±1.0	23.4±2.0	19.4±1.0	30.7±2.0	19.8±1.0	18.6±1.0	18.6±1.0	18.6±1.0	18.6±1.0	18.6±1.0	18.6±1.0
chromo16	17.9±0.5	18.0±3.0	16.3±0.3	29.4±2.0	16.6±0.5	15.9±1.0	15.9±1.0	15.9±1.0	15.9±1.0	15.9±1.0	15.9±1.0
chromo17	21.3±1.0	22.5±1.0	19.7±2.0	30.2±3.0	19.8±1.0	18.6±1.0	18.6±1.0	18.6±1.0	18.6±1.0	18.6±1.0	18.6±1.0
chromo18	18.5±1.0	21.5±2.0	17.8±1.0	26.1±2.0	16.9±1.0	16.2±1.0	16.3±1.0	16.1±1.0	16.1±1.0	16.1±1.0	16.1±1.0
chromo19	12.4±1.0	16.2±3.0	11.6±1.0	19.5±2.0	12.0±1.0	11.7±1.0	11.7±1.0	11.7±1.0	11.7±1.0	11.7±1.0	11.7±1.0
chromo20	15.2±2.0	21.3±1.0	14.9±2.0	24.9±1.0	14.6±2.0	14.1±1.0	14.1±1.0	14.1±1.0	14.1±1.0	14.1±1.0	14.1±1.0
chromo21	10.6±1.0	12.7±2.0	10.5±1.0	14.7±3.0	10.2±1.0	10.0±1.0	10.0±1.0	10.0±1.0	10.0±1.0	10.0±1.0	10.0±1.0
chromo22	13.4±0.3	15.9±2.0	13.1±0.4	23.1±1.0	12.7±0.2	12.4±0.3	12.4±0.2	12.4±0.3	12.4±0.3	12.3±0.3	12.4±0.3
Average	26.6	30.0	24.9	38.4	24.9	23.7	23.7	23.7	23.7	23.7	23.7
$\sigma$	9.9	10.5	9.2	13.5	9.2	8.6	8.6	8.6	8.6	8.6	8.6



Table 5: Comparison of the average distance from the approximated median to each string in the set respect the true median. (Synthetic data)

Set	Exact Median	Set Median	Mollineda $\varepsilon$	Mollineda $S^M$	Fischer $\varepsilon$	Fischer $S^M$	JR-S $\varepsilon$ Freq $\times$ Cost	JR-S $S^M$ Freq $\times$ Cost	JR-S $\varepsilon$ Freq	JR-S $S^M$ Freq	Hinarejos $S^G$	Hinarejos $S^M$
Synthetic 1	6.5	6.9	7.7	6.8	6.9	6.9	6.5	6.5	6.8	6.5	6.5	6.5
Synthetic 2	7.9	8.4	8.6	8.3	8.4	8.4	8.1	8.1	8.1	8.1	8.1	8.1
Synthetic 3	8.0	8.5	8.3	8.4	8.5	8.5	8.2	8.0	8.0	8.0	8.0	8.0
Synthetic 4	7.3	7.6	7.6	7.6	7.6	7.6	7.3	7.3	7.4	7.3	7.3	7.3
Average	7.4	7.8	8.0	7.8	7.9	7.8	7.5	7.4	7.6	7.4	7.5	7.5
$\sigma$	0.6	0.7	0.4	0.6	0.7	0.7	0.7	0.7	0.5	0.7	0.7	0.7

272 **References**

- 273 Bunke, H., Jiang, X., Abegglen, K., Kandel, A., 2002. On the weighted mean of a pair  
274 of strings. *Pattern Analysis & Applications* 5, pp. 23-30.
- 275 Casacuberta, F., Antonio, M., 1997. A greedy algorithm for computing approximate  
276 median strings. In: VII Simposium Nacional de Reconocimiento de Formas y  
277 Análisis de Imágenes (AERFAI), pp. 193-198.
- 278 Fischer, I., Zell, A., 2000. String averages and self-organizing map for strings. In:  
279 Proceedings of the Neural Computation, Canada / Switzerland, (ICSC). Academic  
280 Press, pp. 208-215.
- 281 Freeman, H., 1974. Computer processing of line-drawing data. *Computer Surveys* 6,  
282 pp. 57-96.
- 283 García-Díez, S., Fouss, F., Shimbo, M., Saerens, M., 2011. A sum-over-paths extension  
284 of edit distances accounting for all sequence alignments. *Pattern Recognition*, Vol  
285 44(6), pp. 1172-1182.
- 286 González-Rubio, J., Casacuberta, F. 2010. On the use of median string for multi-source  
287 translation. In: 20th International Conference on Pattern Recognition (ICPR), pp.  
288 4328-4331.
- 289 Jain, A., Zongker, D., 1997. Representation and recognition of handwritten digits using  
290 deformable templates. *IEEE Trans. on Pattern Analysis and Machine Intelligence*,  
291 Vol. 19, pp. 1386-1390.
- 292 Jiang, X., Bunke, H., 2002. Optimal lower bound for generalized median problems in  
293 metric spaces. *Structural, Syntactic, and Statistical Pattern Recognition*, LNCS, Vol.  
294 2396, pp. 143-151.
- 295 Jiang, X., Schiffmann, L., Bunke, H., 2000. Computation of median shapes. In: 4th  
296 Asian Conf. on Computer Vision, pp. 300-305.

297 Jiang, X., Bunke, H., Csirik, J., 2004. Median Strings: A Review Data Mining in Time  
298 Series Databases. World Scientific, Vol. 57, pp. 173-192

299 Jiang, X., Wentker, J., Ferrer, M., 2012. Generalized median string computation by  
300 means of string embedding in vector spaces. Pattern Recognition Letters, Vol. 33,  
301 pp. 842-852.

302 Kohonen, T., 1985. Median strings. Pattern Recognition Letters 3, pp. 309–313.

303 Kohonen, T., 1998. Self-organizing maps of symbols strings. Neurocomputing, Vol.  
304 21, pp. 19-30.

305 Kruskal, J., 1983. An overview of sequence comparison. time warps, string edits and  
306 macromolecules. SIAM Reviews, Vol. 2, pp. 201-2037.

307 Kruzslisz, F., 1999. Improved greedy algorithm for computing approximate median  
308 strings. Acta Cybernetica, Vol. 14, pp. 331-339.

309 Levenshtein, V. I., 1966. Binary codes capable of correcting deletions, insertions, and  
310 reversals. Tech. Rep. Vol. 8.

311 Lourenço, A., Fred, A., 2005. Ensemble methods in the clustering of string patterns.  
312 In: 7th IEEE Workshops on Application of Computer Vision (WACV/MOTIONS),  
313 Vol. 1, pp. 143-148.

314 Martínez-Hinarejos, C., Juan, A., Casacuberta, F., 2003. Median strings for k-nearest  
315 neighbour classification. Pattern Recognition Letters, Vol. 24(1-3), pp. 173-181.

316 Martínez-Hinarejos, C. D., 2003. La cadena media y su aplicación en reconocimiento  
317 de formas. Ph.D. thesis.

318 Martínez-Hinarejos, C., Juan, A., Casacuberta, F., Mollineda, Ramón 2002. Reducing  
319 the computational cost of computing approximated median strings. Lecture Notes in  
320 Artificial Intelligence, SSPR&SPR 2396, pp. 47-55.

- 321 Martínez-Hinarejos, C. D., Alfons, J., Casacuberta, F., 2001. Improving classification  
322 using median string and NN rules. In: IX Spanish Symposium on Pattern Recogni-  
323 tion and Image Analysis. Benicàssim. Spain, (AERFAI), Vol. 2, pp. 307-314.
- 324 Mollineda, R. A., 2004. A learning model for multiple-prototype classification of  
325 strings. In: 17th Int. Conf. on Pattern Recognition (ICPR). Vol. 4, pp. 420-423.
- 326 Nicolas, F., Rivals, E., 2005. Hardness results for the center and median string problems  
327 under the weighted and unweighted edit distances. *Journal of Discrete Algorithms*  
328 3 (2-4), 390–415, combinatorial Pattern Matching (CPM) Special Issue. The 14th  
329 annual Symposium on combinatorial Pattern Matching.
- 330 Olivares, C., Oncina, J., 2008. A stochastic approach to median string computation.  
331 *Lecture Notes in Computer Science, SSPR&SPR 5342*, 431–440.
- 332 Rico-Juan, J., Micó, L., 2003. Some results about the use of tree/string edit distances  
333 in a nearest neighbour classification task. *Pattern Recognition and Image Analysis,*  
334 *LNCS, Vol. 2652*, pp. 821-828.
- 335 Rico-Juan, J. R., Iñesta, J. M. I., 2012. New rank methods for reducing the size of the  
336 training set using the nearest neighbor rule. *Pattern Recognition Letters*, Vol. 33(5),  
337 pp. 654 - 660.
- 338 Ristad, E., Yianilos, P., 1998. Learning string-edit distance. *IEEE Trans. on Pattern*  
339 *Analysis and Machine Intelligence*, Vol. 20, pp. 522-532.
- 340 Wagner, R., Fischer, M., 1974. The string-to-string correction problem. *Journal of the*  
341 *ACM*, Vol. 21, pp. 168-173.