

FLEXIBLE MUSIC RETRIEVAL IN SUBLINEAR TIME

KIMMO FREDRIKSSON

Department of Computer Science, University of Joensuu, Finland
kfredrik@cs.joensuu.fi

VELI MÄKINEN

Department of Computer Science, University of Helsinki, Finland
vmakinen@cs.helsinki.fi

GONZALO NAVARRO

Department of Computer Science, University of Chile, Chile.
gnavarro@dcc.uchile.cl

Received (received date)

Revised (revised date)

Communicated by Editor's name

ABSTRACT

Music sequences can be treated as texts in order to perform music retrieval tasks on them. However, the text search problems that result from this modeling are unique to music retrieval. Up to date, several approaches derived from classical string matching have been proposed to cope with the new search problems, yet each problem had its own algorithms. In this paper we show that a technique recently developed for multipattern approximate string matching is flexible enough to be successfully extended to solve many different music retrieval problems, as well as combinations thereof not addressed before. We show that the resulting algorithms are average-optimal in many cases and close to average-optimal otherwise. Empirically, they are much better than existing approaches in many practical cases.

Keywords: Music retrieval; Approximate string matching; (δ, γ) -matching; Transposition invariance.

1. Introduction

In this paper we are interested in music retrieval, and in particular, in a recent approach to it where musical scores are regarded as strings and string matching techniques can be used to solve music retrieval problems. In order to map the problem to string matching, the alphabet of the string could simply be the set of notes in the chromatic or diatonic notation, or the set of intervals that appear between notes (for example, pitches may be represented as MIDI numbers and pitch intervals as number of semitones). In both cases, we deal with *numeric* strings. Then, many

music retrieval problems can be converted into string matching problems, that is, find the occurrences of a short string (called the *pattern*) in a longer string (called the *text*). This is usually not enough to fully solve all music retrieval problems, but it provides a useful and efficient filter to leave the most promising candidates for a more profound and costly evaluation. There are also some problems where two long musical pieces are compared, which we do not address in this paper.

Exact string matching cannot be used to find occurrences of a particular melody, because a number of irrelevant distortions could exist between the melody sought and its version stored in the music database. To perform meaningful music retrieval one must resort to diverse forms of *approximate* matching, where a limited amount of *differences* of diverse kinds are permitted between the search pattern and its occurrence in the text. Different versions of the approximate string matching problem arise in different fields [1], yet those of music retrieval are unique of this area [2, 3, 4].

One approximate matching model of use in music retrieval is (δ, γ) -*matching*. In this model, two strings $a_1a_2 \dots a_m$ and $b_1b_2 \dots b_m$ of the same length m match if (i) the absolute differences between corresponding characters do not exceed δ , that is, $|a_i - b_i| \leq \delta$ for all $1 \leq i \leq m$ (or, alternatively, $\max_{1 \leq i \leq m} |a_i - b_i| \leq \delta$), and (ii) the sum of those absolute differences does not exceed γ , that is, $\sum_{1 \leq i \leq m} |a_i - b_i| \leq \gamma$. This model accounts for small differences that may arise between two versions of the same melody, setting a limit for the individual absolute differences, as well as a global limit to the overall differences. Searching for pattern p under (δ, γ) -matching consists of finding all the text positions where a text substring that (δ, γ) -matches p appears. Less popular subproblems are δ -matching and γ -matching, which only enforce one of the two conditions.

A second relevant approximate matching model is the *longest common subsequence (LCS)* and its dual *indel distance*. The former, $LCS(a, b)$, is the maximum length of a string that is subsequence both of a and b , that is, $LCS(a, b) = \max\{|s|, s \sqsubseteq a, s \sqsubseteq b\}$. A string $s = s_1s_2 \dots s_r$ is a *subsequence* of string $a_1a_2 \dots a_m$, $s \sqsubseteq a$, if s can be obtained by removing zero or more characters from a , that is, $s = a_{i_1}a_{i_2} \dots a_{i_r}$ for $1 \leq i_1 < i_2 < \dots < i_r \leq m$. The LCS has been largely used in computational biology to model biological similarity, and it is also relevant to identify musical passages that are similar except for a few extra or missing notes. This is especially relevant because music contains various kind of “decorations”, such as grace notes and ornamentations, that are not essential for matching. The indel distance $id(a, b)$ between strings a and b is the number of characters one has to add or remove to a and b to make them equal, $id(a, b) = |a| + |b| - 2 \cdot LCS(a, b)$. Searching for pattern p under indel distance with tolerance k consists of finding all the text positions where a string p' appears so that $id(p, p') \leq k$. Other variants of indel distance, which are less popular in music retrieval, are Levenshtein or edit distance (where substitutions of characters are also permitted), episode matching (where only insertions in the pattern are permitted), and Hamming distance (where only substitutions are permitted).

Finally, a third similarity concept of relevance in music retrieval is *transposition*

invariance. Two strings $a = a_1a_2 \dots a_m$ and $b = b_1b_2 \dots b_m$ are one the transposed version of the other if there is a constant t such that $a + t = (a_1 + t)(a_2 + t) \dots (a_m + t) = b$. Transposition invariance is very relevant because Western people tend to listen to music analytically, by observing the intervals between consecutive pitch values rather than the actual pitch values themselves. As a result, a melody performed in two distinct pitch levels is perceived as equal regardless of whether it is performed in a lower or higher level of pitches.

As a string matching problem, dealing with transposition invariance alone is trivial because it suffices to represent text and pattern as differences between consecutive notes and then apply exact string matching. However, in most cases of interest the above problems appear in combined form. In particular, transposition invariance is usually combined with longest common subsequence. The longest common transposition invariant subsequence between two strings a and b , $LCTS(a, b)$, permits transposing a or b as necessary to find the longest common subsequence among them, $LCTS(a, b) = \max_{t \in \mathbb{Z}} LCS(a + t, b)$.

In recent years, there has been much activity around developing specific string matching techniques to solve diverse music retrieval problems, mostly consisting of combinations of those outlined above. Several theoretical and practical results of interest have been achieved. We cover these in the next section.

Our contribution in this paper is to show that a particular approach recently developed for multiple approximate string matching [5] is flexible enough to be successfully adapted to solve most of the combinations of problems sketched above. Basically the same search technique, coupled with slightly different pattern preprocessings, yield algorithms that solve each combination. In theoretical terms, we show that many of the resulting algorithms are average-optimal, matching the current lower bound that does not consider transposition invariance. Other algorithms are shown to be close to optimal. That is, including transposition invariance yields no or very little cost on average.

More specifically, we prove that the lower bounds on the average complexity of approximate string matching under several edit distance-like models hold for their transposition-invariant versions too, by deriving new average-optimal transposition-invariant algorithms for them. We also derive lower bounds for (non-transposition-invariant) δ - and γ -matching, as well as average-optimal algorithms for them. We also give almost optimal algorithms for transposition-invariant δ -matching and γ -matching. Finally, we show how to combine δ -matching with edit distance-like models, with a complexity that remains optimal without transposition invariance and almost optimal with it. Obtaining similar complexities for the combination of γ -matching with edit distance-like models remains as an open challenge.

On the practical side, we show experimentally that our technique largely outperforms all the existing ones in most cases of interest. For small to moderate error thresholds our algorithms are substantially faster than previous approaches for all but very short texts. These are the parameter values that are most interesting in most music retrieval applications.

2. Related Work

In which follows, we assume that a long text $T = t_1t_2 \dots t_n$ is searched for a comparatively short pattern $p = p_1p_2 \dots p_m$. Both are sequences over alphabet Σ , a finite contiguous subset of \mathbb{Z} , of size σ .

2.1. (δ, γ) -Matching

Several recent algorithms exist to solve this problem. These can be classified as follows:

Bit-parallel: The idea is to take advantage of the intrinsic parallelism of the bit operations inside a computer word of w bits [6], so as to pack several values in a single word and manage to update them all in one step [7, 8, 9]. The best complexity achieved [9] is $O(n m \log(\gamma)/w)$ in the worst case and $O(n)$ on average.

Occurrence heuristics: Inspired by Boyer-Moore techniques [10], they skip some text characters according to the position of some characters in the pattern [7, 11]. In general, only δ is used to skip characters, while the γ -condition is used to verify candidates. This makes these algorithms weak for large δ and small γ .

Substring heuristics: Based on suffix automata [12], these algorithms skip text characters according to the position of some pattern substrings [11, 9]. In the second article, they use bit-parallelism to filter the text using both δ and γ , unlike previous approaches. This is shown to be the approach examining the least number of text characters.

FFT-related: It is possible to solve the δ -matching and (δ, γ) -matching problems in $O(\delta n \log m)$ time, and γ -matching problem in $O(n\sqrt{m \log m})$ time [13] using Fast Fourier Transform (FFT) based techniques. The $O(n\gamma \log \gamma)$ time algorithm in [14] is faster for small γ . This algorithm is based on bounded divide-and-conquer and non-boolean convolutions. This technique can be also used to solve the δ -matching problem in $O(n \log m \sqrt{\delta})$ time. Other FFT based $o(mn)$ solutions exist for related problems, see e.g. [15], and especially related to δ -matching [16, 17]. Matching under γ -restriction is possible in $O(mn/\log_\sigma n)$ time [18] without using FFT (but using the Four-Russians trick).

In practice, the best current algorithms for (δ, γ) -matching are those in [9], as demonstrated by the experiments in [11, 9]. In [9] they present a plain bit-parallel and a substring heuristic. The first is shown to be the best in most cases, but for short patterns and small δ and γ , the character-skipping technique is better.

The FFT based techniques, although elegant, have considerably large overheads to make them practical. Our preliminary tests show that they only become faster than the naive algorithm on very long patterns. Searching for long patterns is not typical in music retrieval. The solution based on the Four-Russians trick is only practical for small alphabets, much smaller than what is required for music retrieval.

2.2. Transposition Invariant LCS and Indel Distance

Plain (non-transposed) LCS among strings p and T can be computed in $O(mn)$ time using dynamic programming [19]. In general, any LCTS algorithm can be adapted to text searching with indel distance. The LCTS problem was first stated in [20], where $O(\sigma mn)$ time was obtained by trying out all the $2\sigma - 1$ possible transpositions one by one. Further solutions to the problem can be classified as follows.

Brute-force: The idea is to pick any LCS algorithm and try it for all the $2\sigma - 1$ possible transpositions. Apart from the original proposal [20], several others have been attempted considering different practical LCS algorithms based on bit-parallelism [21, 22]. The best complexity achieved is $O(\sigma mn/w)$.

Sparse dynamic programming: An evolution over the above scheme is to notice that the $LCS(a + t, b)$ problem for each transposition t has only a few character matches between a and b , mn in total. Those *sparse* problems are best handled by sparse dynamic programming algorithms. This idea led to several solutions [23, 24, 25]. The best complexity achieved is $O(mn \log \log \min(m, \sigma))$, yet a version with complexity $O(mn \log \sigma / \log w)$ is shown to be better in practice.

Branch and bound: In this case the idea is to search for the best possible transposition t by a backtracking method, recursively dividing the space of $2\sigma - 1$ transpositions into ranges until finding the best one [26]. This yields a best-case complexity of $O((mn + \log \log \sigma) \log \sigma)$, and the method works well in practice. Yet, it cannot be extended to searching with indel distance.

Experiments in [26, 22, 25] demonstrate that the $O(mn \log \sigma / \log w)$ algorithm in [25] is the fastest in practice. This method can be adapted to searching with indel distance. We emphasize that all existing search algorithms for this problem (including transposition invariance) examine all text characters.

3. Optimal Multiple Approximate String Matching

In [5], new algorithms for single and multiple approximate string matching were presented. Those algorithms were not only optimal on average, but also very efficient in practice, even in the more competitive area of single approximate string matching. It was shown that, to search for the occurrences of r patterns of length m in a text of length n , all them uniformly distributed over an alphabet of size σ , the algorithm required $O(n(k + \log_\sigma(rm))/m)$ time on average. Here k is the maximum number of missing, extra, or substituted characters permitted to match a pattern against a text string (searching under edit distance). This average complexity is optimal [27, 28].

We first explain how to search for a single pattern p . We choose a *block length* ℓ , and compute $med(b, p)$ for every possible block $b \in \Sigma^\ell$ (that is, every possible

ℓ -gram). Here, $med(b, p)$ is the minimum edit distance between b and a substring of p ,

$$med(b, p) = \min\{ed(b, p'), \exists x, y, p = xp'y\},$$

being $ed(b, p')$ the edit distance between b and p' .

Now, the text $T = t_1 t_2 \dots t_n$ is scanned as follows. Since the minimum length of an occurrence of $p = p_1 p_2 \dots p_m$ in T with edit distance at most k has length at least $m - k$ (when k deletions occur on p), we slide a window of length $m - k$ along the text. For each window tried, $t_{i+1} t_{i+2} \dots t_{i+m-k}$, we read its ℓ -grams right to left. That is, we read at most $\lfloor (m - k) / \ell \rfloor$ ℓ -grams b_1, b_2 , and so on, so that $b_1 = t_{i+m-k-\ell+1} \dots t_{i+m-k}$ is the rightmost, $b_2 = t_{i+m-k-2\ell+1} \dots t_{i+m-k-\ell}$ precedes b_1 , etc. The invariant is that any occurrence of p starting at positions $\leq i$ has already been reported.

For each such ℓ -gram $b_j = t_{i+m-k-j\ell+1} \dots t_{i+m-k-j\ell}$, we find $med(b_j, p)$ in the precomputed table. If, after reading b_j , we have $med(b_1, p) + med(b_2, p) + \dots + med(b_j, p) > k$, then no possible occurrence of p can contain the text $b_j b_{j-1} \dots b_2 b_1$, thus the window is slid forward to start at the second character of b_j , that is, we set $i \leftarrow i + m - k - j\ell + 1$ (as the new window will start at $i + 1$).

If, on the other hand, all the ℓ -grams of the window are scanned and yet the window cannot be shifted, it must be verified for a real occurrence. At this point, we must check if there is an occurrence p' of p starting at text position $i + 1$. Since the maximum length of an occurrence is $m + k$ (where k insertions occur into p), any potential p' must finish between positions $i + m - k$ and $i + m + k$. So we compute

$$led(p, i) = \min\{ed(p, t_{i+1} \dots t_{i+m-k+d}), 0 \leq d \leq 2k\},$$

which can be done in $O(m^2)$ time by computing $ed(\)$ incrementally in d . If $led(p, i) \leq k$, we report $i + 1$ as the starting position of an occurrence. Finally, we advance the window by one position, $i \leftarrow i + 1$.

We show now that the way we shift the window is safe, that is, no occurrence can start at positions $i + 1$ to $i + m - k - j\ell + 1$. Any such occurrence, of length at least $m - k$, must contain the sequence of ℓ -grams $b_j \dots b_1$. Let $p' = x b_j \dots b_1 y$ be such an occurrence. This is a split of p' into $j + 2$ pieces. The main point is that the edit distance is *decomposable*: For any strings p and p' , given any split $p' = p'_1 \dots p'_{j+2}$, there is a split $p = p_1 \dots p_{j+2}$ such that $ed(p', p) = ed(p'_1, p_1) + \dots + ed(p'_{j+2}, p_{j+2})$. But each such $ed(p'_s, p_s) \geq med(p'_s, p) \geq 0$, by definition of $med(\)$.

Hence, in our particular case, $ed(p', p) \geq med(b_j, p) + \dots + med(b_1, p)$. Thus if the latter exceeds k , there can be no occurrence of p containing $b_j \dots b_1$.

The extension of the algorithm for multiple patterns is trivial. We only have to change the preprocessing so that p is now a set of patterns $p = \{p^1 \dots p^r\}$ and now $med(b, p) = \min_{1 \leq i \leq r} med(b, p^i)$. So $med(b, p)$ is a lower bound to the cost of matching b anywhere inside any pattern of the set.

By appropriately choosing $\ell = \Theta(\log_\sigma(rm))$, we obtain the promised complexity.

3.1. Extensions

Several other improvements are studied in [5]. We briefly review some that are used in our experiments. For more details see [5].

On the windows that have to be verified, we could simply run the verification for every pattern, one by one. A more sophisticated choice is *hierarchical verification* [29]. We form a tree whose nodes have the form $[i, j]$ and represent the group of patterns $p^i \dots p^j$. The root is $[1, r]$, and the leaves have the form $[i, i]$. Every internal node $[i, j]$ has two children $[i, \lfloor (i+j)/2 \rfloor]$ and $[\lfloor (i+j)/2 \rfloor + 1, j]$.

The preprocessing is done first for the leaves, as in the single pattern case, that is, we compute a table for $med(b, p^i)$. The internal nodes contain tables for $\min_{i \leq h \leq j} med(b, p^h)$, computed by minimizing over the two tables of the subtrees. In the filtering phase, we first use the table for the root, corresponding to the full set of patterns, and if the current window has to be verified with respect to a node in the hierarchy, we rescan the window considering the two children of the current node. It is possible that the window can be discarded for both children, for one, or for none. We recursively repeat the process for every child that does not permit discarding the window. If we process a leaf node and still have to verify the window, then we run the verification algorithm for the corresponding single pattern.

The second improvement is to have *bit-parallel counters*. In this case we reserve only $O(\log_2 k)$ bits to accumulate the differences $med(b_j, p)$. This means that if we have a computer word of w bits, we can process $O(w/\log_2 k)$ patterns in parallel. This technique can also be used with the hierarchical verification, to increase the arity of the tree to $O(w/\log_2 k)$.

The third improvement is to use *ordered ℓ -grams*, where each b_j is permitted to match only in the area of p where it could be aligned in an occurrence starting at $i+1$. In an approximate occurrence of $b_j \dots b_1$ inside the pattern, b_i cannot be closer than $(i-1)\ell$ positions to the end of the pattern. Therefore, we compute tables for $med^j(b, p)$, $1 \leq j \leq \lfloor (m-k)/\ell \rfloor$, where $med^j(b, p) = \min\{ed(b, p'), \exists x, y, |y| \geq (j-1)\ell, p = xp'y\}$. This allows us to discard a window whenever $med^1(b_1, p) + med^2(b_2, p) + \dots + med^j(b_j, p) > k$. This reduces verifications but increases preprocessing time and space.

Finally, it is possible to improve the preprocessing time by using a trie of all the possible ℓ -grams to reuse preprocessing work. All the improvements can be combined into a single algorithm.

4. Adapting to Music Retrieval

The method above was designed for multiple string matching under edit distance. Yet its main idea is much more general and can be used to solve many other problems. In this section we demonstrate that the idea solves most of the music retrieval problems we have focused on in this paper. We note that this gives immediately a solution to the multipattern version of the same problems.

4.1. Transposition Invariant Indel Distance and Variants

Let us start with searching with transposition invariant indel distance. For each

ℓ -gram $b \in \Sigma^\ell$, we compute

$$mtid(b, p) = \min\{id(b + t, p'), \exists x, y, p = xp'y, -\sigma < t < \sigma\}. \quad (1)$$

This is the minimum transposition invariant indel distance to match b anywhere inside p . The same algorithm of the previous section is used, and the same argument shows that we cannot discard a window that starts an occurrence of p in T . Indel distance is decomposable just like edit distance, that is, for any split $p' = p'_1 \dots p'_{j+2}$, there is a split $p = p_1 \dots p_{j+2}$ such that $id(p', p) = id(p'_1, p_1) + \dots + id(p'_{j+2}, p_{j+2})$. Assume p matches t the current window $xb_j \dots b_1y$ starting at position $i+1$. That is, there exists a transposition t such that $id(p', p) \leq k$, $p' = (x+t)(b_j+t) \dots (b_1+t)(y+t)$. Now, $id(p', p) \geq id(b_j+t, p_2) + \dots + id(b_1+t, p_{j+1}) \geq mtid(b_j, p) + \dots + mtid(b_1, p)$. Thus if the latter exceeds k we can safely shift the window.

When a window starting at position $i+1$ cannot be shifted, we simply compute $LCTS(p, t_{i+1} \dots t_{i+m-k+d})$ for any $0 \leq d \leq 2k$, and report position $i+1$ if $LCTS(p, t_{i+1} \dots t_{i+m-k+d}) \geq (m+m-k+d-k)/2 = m-k+d/2$ for some d , as this is equivalent to $id(p, t_{i+1} \dots t_{i+m-k+d}) \leq k$ for some transposition t .

Fig. 1 shows simplified pseudocode. The very same algorithm can be used to handle other distances, just by changing the preprocessing. For transposition invariant Levenshtein distance we use edit distance ed instead of indel distance id in Eq. (1). For transposition invariant Hamming distance we use Hamming instead of indel distance in Eq. (1), and let the window length be m . For transposition invariant episode matching we permit only deletions in b in Eq. (1) and use windows of length m . Note that, for Hamming distance, verification of a window only requires to compare it against the pattern.

Search ()	Shift (i, D)
1. $D \leftarrow \mathbf{Preprocess} ()$	1. $M \leftarrow 0$
2. $i \leftarrow 0$	2. $c \leftarrow m - k$
3. While $i \leq n - (m - k)$ Do	3. While $c \geq \ell$ Do
4. $pos \leftarrow \mathbf{Shift} (i, D)$	4. $c \leftarrow c - \ell$
5. If $pos = i$	5. $M \leftarrow M + D[t_{i+c+1} \dots t_{i+c+\ell}]$
6. Verify area $t_{i+1} \dots t_{i+m+k}$	6. If $M > k$ Return $i + c + 1$
7. $pos \leftarrow pos + 1$	7. Return i
8. $i \leftarrow pos$	

Preprocess ()

1. $\ell \leftarrow \Theta(\log_\sigma m)$
2. **For** $b \in \Sigma^\ell$ **Do** $D[b] \leftarrow mtid(b, p)$
3. **Return** D

Fig. 1. Simplified description of the transposition invariant indel algorithm.

4.2. (δ, γ) -Matching

Alternatively, we can search for (δ, γ) -matches of p in T . In this case the window is of length m , as occurrences are all of that length. For each ℓ -gram $b \in \Sigma^\ell$, we compute

$$mdg(b, p) = \min\{\gamma', \exists x, y, p = xp'y, b \text{ } (\delta, \gamma')\text{-matches } p'\}.$$

This is the minimum total number of absolute differences obtained by b inside p , where we restrict those positions to δ -match as well. The same algorithm of the previous section is used with this preprocessing (and the threshold is γ instead of k).

Being γ -matching a cumulative measure, the sum of $mdg(b_j, p)$ values is a lower bound to the γ needed to match the window inside p . Consider window $p' = t_{i+1} \dots t_{i+m} = xb_j \dots b_1$. Assume p' (δ, γ) -matches p . Then, by definition of (δ, γ) -matching, b_1 (δ, γ_1) -matches $p_{m-\ell+1} \dots p_m$, and so on until b_j , which (δ, γ_j) -matches $p_{m-j\ell+1} \dots p_{m-j\ell+\ell}$, so that $\gamma_1 + \dots + \gamma_j \leq \gamma$. As each b_s (δ, γ_s) -matches $p_{m-s\ell+1} \dots p_{m-s\ell+\ell}$, it holds $mdg(b_s, p) \leq \gamma_s$, and $mdg(b_j, p) + \dots + mdg(b_1, p) \leq \gamma$.

When a window $t_{i+1} \dots t_{i+m}$ cannot be shifted, we check whether p (δ, γ) -matches the window in time $O(m)$, and report position $i + 1$ if this is the case.

The pseudocode of Fig. 1 can be easily adapted to this model. One needs only to replace $mtid()$ with $mdg()$, k with γ , and adjust the window size from $m - k$ to m , and verification area from $t_{i+1} \dots t_{i+m+k}$ to $t_{i+1} \dots t_{i+m}$.

4.3. Feasible and Unfeasible Combinations

We can also combine transposition invariant indel distance with δ -matching. In this case we count indels, but two characters match whenever they do not differ by more than δ units. This is easily handled by modifying $mtid(b, p)$ formula so that $id(b + t, p')$ considers matches in the more relaxed way. Transposition invariance can also be combined with (δ, γ) -matching, by using $mtdg(b, p)$ instead of $mdg(b, p)$, so that

$$mtdg(b, p) = \min\{\gamma', \exists x, y, p = xp'y, b + t \text{ } (\delta, \gamma')\text{-matches } p', -\sigma < t < \sigma\}.$$

We cannot directly combine transposition invariant indel distance with (δ, γ) -matching. The reason is that we do not have here a single value to minimize, such as the number of indels or γ , but both of them at the same time. It was possible to combine transposition invariant indel distance with δ -matching because the latter is not a parameter to optimize but a condition for matching. Likewise, it was possible to combine γ -matching with δ -matching to obtain (δ, γ) -matching. Yet, if we want to combine indel distance (even without transposition invariance) with γ -matching, the problem is that each pair (b, p') produces a tradeoff between the number of indels and the sum of differences. It is not a matter of adding up indels or differences over a set of tradeoffs in order to stay below k for the first and below γ for the second. Thus our algorithms work as long as we have a single parameter to optimize.

5. Complexity and Optimality

In this section we analyze the average case behavior of our algorithms and prove the average-optimality of some of them. Those that are not average-optimal are close to it. We assume that text and pattern are sequences of symbols uniformly and independently distributed over σ values.

5.1. Transposition Invariance

As we have described it, our algorithm for transposition invariant indel distance is equivalent to multipattern search with indel distance for the set $\{p^1 = p - (\sigma - 1), p^2 = p - (\sigma - 2), \dots, p^{2\sigma-1} = p + (\sigma - 1)\}$. Since $id(a, b) \geq ed(a, b)$ for any strings a and b , the analysis of [5] on edit distance applies to indel distance and the result is pessimistic (yet tight). According to the analysis in [5], searching for r random patterns in random text yields average complexity $O(n(k + \log_\sigma(rm))/m)$, as long as $k/m \leq 1/2 - O(1/\sqrt{\sigma})$. In our case $r = 2\sigma - 1$, and then the complexity boils down to $O(n(k + \log_\sigma m)/m)$, which is optimal even for one pattern without transposition invariance [30, 27]. Thus our transposition-invariant algorithm is optimal too.

The analysis holds as well for any other distance that upper bounds edit distance, such as episode matching, Hamming distance, and (obviously) the same edit distance. Actually, any distance built over a subset of the edit distance operations (i.e., insertions, deletions, replacements) is covered by the analysis above. In the case of Hamming distance, however, the i -th character of the pattern can only align with the i -th character of the occurrence, and thus the result applies for $k/m \leq 1/2 - O(1/\sigma)$ [31].

All the analysis above assumes that our $2\sigma - 1$ patterns are random. However, this is not the case, as they are all the transpositions of a single random pattern. For example, if $\ell = 1$, then our $2\sigma - 1$ patterns necessarily match any string of length 1, whereas the same number of random patterns do not. We show in Section 5.3 that the average-case analysis is, however, still valid in this case.

Furthermore, we can also search for r patterns permitting transposition invariance in average time $O(n(k + \log_\sigma(rm))/m)$ under any of these models. This is optimal as well [28]. The result is summarized in the following theorem.

Theorem 1 *Our algorithms permit searching for r patterns of length m in a text of length n , both random sequences over an alphabet of size σ , permitting transposition invariance and at most k differences between patterns and their occurrences (the differences being character insertions, deletions, substitutions, or any subset thereof), in average time $O(n(k + \log_\sigma(rm))/m)$ provided $k/m \leq 1/2 - O(1/\sqrt{\sigma})$ (or $k/m \leq 1/2 - O(1/\sigma)$ if only substitutions are permitted). This is average-optimal even when no transpositions are allowed.*

In the worst case the algorithms require, per pattern, $O(n)$ verifications over $O(m)$ characters each, for a total that in no case exceeds $O(rmn \log m)$ [23]. More practical algorithms require $O(rmn\sigma/w)$ in the worst case [21, 22]. Preprocessing time and space is polynomial in rm , as we preprocess all the σ^ℓ different ℓ -grams, for $\ell = \Theta(\log_\sigma(rm))$.

5.2. δ -Matching, γ -Matching, and Combinations

Let us start with δ -matching alone (i.e., no γ restriction nor transposition invariance nor differences). In this case, the probability of a random pattern and text characters matching is $\leq (2\delta + 1)/\sigma$. It is enough to set $\ell \geq 3 \log_{\frac{\sigma}{2\delta+1}} m$ to ensure that the first window ℓ -gram read will δ -match within the pattern with probability $\leq 1/m^2$. Assuming pessimistically that, as soon as the first window ℓ -gram matches the pattern, we traverse the whole window and shift it by one position, those “bad” cases do not contribute more than $O(n/m)$ to the average complexity. “Good” cases (where the first ℓ -gram does not δ -match within the window), make us work $O(\ell)$ and shift $m - \ell$, dominating the overall $O(n \log_{\frac{\sigma}{\delta+1}}(m)/m)$ average time (see [5] if needing more details on this kind of analysis). Note that $2\delta + 1$ must be bounded away from σ for the analysis to hold.

It is easy to see that this complexity is average-optimal: To do plain string matching over a numeric alphabet of size σ , multiply all text and pattern character values by $\delta + 1$ and permit δ -matching over this new alphabet of size $\sigma' = \sigma(\delta + 1)$. If one can do δ -matching in less than $\Omega(n \log_{\frac{\sigma'}{\delta+1}}(m)/m)$ time, then Yao’s [27] bound $\Omega(n \log_{\sigma}(m)/m)$ on plain string matching can be broken over the original alphabet.

We can add transposition invariance to δ -matching by, again, reducing to multi-pattern matching. The resulting complexity is $O(n \log_{\frac{\sigma}{\delta+1}}(\sigma m)/m) = O(n \log_{\frac{\sigma}{\delta+1}}((\delta + 1)m)/m)$. We can also combine δ -matching with any of the distances considered in the previous section, with or without transposition invariance, adding $O(nk/m)$ average time. Finally, we can afford multipattern search for r patterns converting m to rm inside the logarithms. This is easily seen by following the original analysis without δ -matching [28] over an alphabet of size $\sigma/(2\delta + 1)$. We get the following theorem.

Theorem 2 *Our algorithms permit searching for r patterns of length m in a text of length n , both random sequences over an alphabet of size σ , permitting δ -matching in average time $O(n \log_{\frac{\sigma}{\delta+1}}(rm)/m)$. They also combine δ -matching with permitting at most k differences between patterns and their occurrences (the differences being characters insertions, deletions, substitutions, or any subset thereof), in average time $O(n(k + \log_{\frac{\sigma}{\delta+1}}(rm))/m)$, provided $k/m \leq 1/2 - O(1/\sqrt{\sigma})$ (or $k/m \leq 1/2 - O(1/\sigma)$ if only substitutions are permitted). Those complexities are average-optimal for δ -matching. The algorithms can be further combined with transposition invariance at the cost of converting the $O(\log_{\frac{\sigma}{\delta+1}}(rm))$ term into $O(\log_{\frac{\sigma}{\delta+1}}((\delta + 1)rm))$.*

The case of γ -matching is explicitly described and analyzed in [31] (without transposition invariance, and calling the model “accumulated”). It is shown that the average complexity is $O(n(\gamma/\sigma + \log_{\frac{\sigma}{1+\gamma/m}}(rm))/m)$ for $\gamma/m < \sigma/(2e) - O(1)$.

^a It is conjectured in [31] that this is average-optimal; we prove it here for any $\gamma \leq \alpha m\sigma/2$, for any constant $\alpha < 1/2$.

^aActually the analysis is a bit oversimplified in [31], by assuming too quickly the bad case $\sigma = \Theta(\sigma m)$.

Let us start with the $O(\gamma/\sigma)$ additive term. Assume we divide the text into consecutive blocks of length m and just focus on the problem of reporting which of those blocks γ -match the pattern. Let us call X_i the absolute difference between the i -th text and pattern cells. As both cells are independently and uniformly distributed in $[1, \sigma]$, we have $E(X_i) = (\sigma - 1)/3$. Thus, on average we have to add $\Theta(\gamma/E(X_i)) = \Theta(\gamma/\sigma)$ differences so that they add up more than γ [32, page 359]. Hence we need $\Omega(\gamma/\sigma)$ accesses on average per text block in order to discard it, and thus cannot work less than $\Omega(n(\gamma/\sigma)/m)$ on average.

Note that we can “discard” a block in the other way, by noting that it γ -matches the pattern without having completely scanned it. More precisely, if we accumulate γ' differences after examining m' block characters and $\gamma - \gamma' \geq (m - m')\lfloor\sigma/2\rfloor$, then we know that the block will γ -match the pattern no matter what the unseen differences are. Let us call $Y_i = \sigma - 1 - X_i$, thus $E(Y_i) = 2(\sigma - 1)/3$ and $\gamma' = \sum_{i=1}^{m'} X_i$. Now, if $\gamma - \gamma' \geq (m - m')\lfloor\sigma/2\rfloor$ then $\sum_{i=1}^{m'} Y_i \geq (m + m')(\sigma - 1)/2 - \gamma$. Even neglecting m' in the right hand, we have that m' has to be large enough so that $\sum_{i=1}^{m'} Y_i \geq m(\sigma - 1)/2 - \gamma$. By our restriction on γ , the right hand is $\Omega(m\sigma)$. Using again the same result on probability [32, page 359], we need $m' = \Theta(m) = \Omega(\gamma/\sigma)$.

Let us now focus on the logarithmic term. We already know that $O(n \log_{\sigma}(rm)/m)$ is average-optimal for exact multipattern matching [28, 27]. Thus the case $\gamma = O(m)$ is already optimal. Otherwise the base of the logarithm is $\sigma m/\gamma$. We note that any string δ -matching P , for $\delta = \lfloor\gamma/m\rfloor$, will also γ -match it (the limit we have set on γ implies $2\delta + 1 < \sigma$ as required). Therefore, a way to solve δ -matching for that δ value is to run γ -matching and check the δ -condition before reporting any text position as an occurrence. This cannot miss any δ -occurrence, and the cost introduced by the extra δ -check is negligible as all those text characters checked must already have been examined in order to report a γ -match for them^b. Therefore, the lower bound proved on the average complexity of δ -matching holds also for γ -matching with $\gamma = m\delta$. The lower bound is $\Omega(n \log_{\frac{\sigma}{\delta+1}}(m)/m) = \Omega(n \log_{\sigma m/\gamma}(m)/m)$ as promised. The multipattern case follows similarly.

Therefore, the γ -matching algorithm we have described (and that was previously described in [31]) is average-optimal. Let us analyze our transposition-invariant version. Its average cost is $O(n(\gamma/\sigma + \log_{\frac{\sigma}{1+\gamma/m}}(\sigma m))/m) = O(n(\gamma/\sigma + \log_{\frac{\sigma}{1+\gamma/m}}(\gamma + m))/m)$. Thus transposition invariance is included at negligible cost if $\gamma = O(m)$, otherwise an additive term appears which can be as large as $O(\log \sigma)$. Multipattern search can be included as usual, multiplying $(\gamma + m)$ by r inside the logarithm. As explained in Section 4.3, we have not devised a way to combine γ -matching with edit operations.

For (δ, γ) -matching, we will traverse the windows as long as the ℓ -grams we have read both δ -match and γ -match within the patterns ℓ -grams. Thus in each possible window the number of ℓ -grams read (amount of shifting) will be the minimum (maximum) between the corresponding ones for δ - and for γ -matching. Thus a

^bAgain, it is possible to declare a γ -match without having seen all its text characters, if $\gamma - \gamma' \geq (m - m')\lfloor\sigma/2\rfloor$. Yet, this is only significant in terms of complexity if $m' = o(m)$. In this case, even if $\gamma' = 0$ we need $\gamma \geq m\lfloor\sigma/2\rfloor(1 + o(1))$, which is outside our bounds on γ .

conservative complexity for (δ, γ) -matching is the minimum between both. We get the following theorem.

Theorem 3 *The γ -matching algorithm we have described (as well as [31]) permits searching for r patterns of length m in a text of length n , both random sequences over an alphabet of size σ , in average time $O(n(\gamma/\sigma + \log_{\frac{\sigma}{1+\gamma/m}}(rm))/m)$. This complexity is average-optimal if $\gamma \leq \alpha m \sigma$ for any constant $\alpha < 1/2$. Our algorithms permit also transposition-invariant γ -matching in average time $O(n(\gamma/\sigma + \log_{\frac{\sigma}{1+\gamma/m}}(r(\gamma + m)))/m)$. Finally, it is possible to do (δ, γ) -matching, with or without transposition invariance, with the best complexity among those for δ -matching and γ -matching.*

For δ - and/or γ -matching the worst cases are $O(rmn)$ without transpositions and $O(rmn\sigma)$ with transpositions. We note that these can be improved by using the more efficient worst-case algorithms available in the literature. Preprocessing time and space is $O(\sigma^\ell \text{poly}(rm))$. As $\ell = \Theta(\log_{\frac{\sigma}{\delta+1}}(rm))$ for δ -matching and $\ell = \Theta(\log_{\frac{\sigma}{1+\gamma/m}}(rm))$ for γ -matching [31], this is polynomial in rm if we assume that σ is constant or that δ and γ/m are $O(\sigma^\alpha)$ for some constant $0 \leq \alpha < 1$.

5.3. Turning Arbitrary Patterns into Random Patterns

In this section we show that the analysis we have done, reducing transposition invariant search to multipattern search for $O(\sigma)$ patterns, is valid even when those patterns are not random.

We note that the amount of work and amount of shifting in a window depends solely on whether the window ℓ -grams coincide with some pattern ℓ -gram at each possible position in $[1, m - k - \ell + 1]$. More precisely, given our set of patterns $\{p^1, p^2, \dots, p^{2\sigma-1}\}$, let

$$B_j = \{p_j^i \dots p_{j+\ell-1}^i, 1 \leq i \leq 2\sigma - 1\}$$

be the set of pattern ℓ -grams starting at pattern position j , for $1 \leq j \leq m - k - \ell + 1$. The performance of our search algorithm is monotonic with $(B_1, B_2, \dots, B_{m-k-\ell+1})$: If one searches for another pattern set $\{(p')^1, (p')^2, \dots, (p')^r\}$ with B' sets such that $B_j \subseteq B'_j$ for all $1 \leq j \leq m - k - \ell + 1$, then this second search is guaranteed to work more and shift less for every possible text window. The reason is that all the $mtid(\cdot)$ values will be smaller or equal and thus more window ℓ -grams will be examined before surpassing the threshold k .

The set of $O(m\sigma)$ ℓ -grams we produce when converting transposition invariant into multipattern searching is not random. Yet, assume we generate r random patterns with the hope that, for each real ℓ -gram at each position of the real patterns, $b \in B_j$, our set of r random patterns will contain that ℓ -gram at that position for some pattern, $b \in B'_j$.

The search time for the random set will be $O(n(k + \log_\sigma(mr))/m)$, as those are now random patterns. This will be optimal for our problem as long as $r = O(\text{poly}(m\sigma))$.

Take now one specific ℓ -gram from the real set we have to search for, $b \in B_j$. The probability of *not* appearing at the same position j in our r random patterns, $b \notin B'_j$, is that of not appearing in a random choice of ℓ -grams, $(1 - 1/\sigma^\ell)^r$. Since $\ell = O(\log_\sigma m)$, say $\ell \leq c \log_\sigma m$, this probability is $\leq (1 - 1/m^c)^r \leq e^{-r/m^c}$.

Let us call random variable $X_h = 1$ if the h -th real ℓ -gram, $b_h \in B_j$, does not appear in B'_j , and 0 otherwise. Then $E(X_h) = P(X_h = 1) \leq e^{-r/m^c}$.

Now let X be the total number of real ℓ -grams not in their B' set, $X = X_1 + X_2 + \dots + X_{O(m\sigma)}$. Those X_h variables are dependent on each other, but even so, $E(X) = \sum E(X_h) = O(m\sigma e^{-r/m^c})$. Finally, the probability of *some* real ℓ -gram $b_h \in B_j$ not belonging to set B'_j is $P(X \geq 1) \leq E(X) = O(m\sigma e^{-r/m^c})$.

Consider now the following randomized process:

1. Generate r random ℓ -grams.
2. If they happen to contain all the real pattern ℓ -grams at each position, then run our search algorithm over the r random patterns.
3. Otherwise, perform a classical $O(\sigma mn)$ time search for the real patterns one by one.

The process is at least as costly as the real search we do, no matter which of (2) or (3) is chosen, so the average case analysis of this process upper bounds the real one. Case (3) occurs with probability $O(m\sigma e^{-r/m^c})$, so it contributes $O(nm^2\sigma^2 e^{-r/m^c})$ to the average complexity. Case (2) contributes $O(n(k + \log_\sigma(mr))/m)$.

Now, it is sufficient that $r \geq \sigma m^{1+c}$ to make $O(nm^2\sigma^2 e^{-r/m^c}) = O(nm^2\sigma^2 e^{-m\sigma})$, which is $O(n/m)$. The other term of the complexity becomes $O(n(k + \log_\sigma(mr))/m) = O(n(k + (2 + c) \log_\sigma m)/m)$, optimal for any constant c .

This analysis adapts straightforwardly to all the other distances and matching models we have considered.

6. Experimental Results

We have implemented the algorithms in C, compiled using `icc 8.0` with full optimizations. The experiments were run in a 2GHz Pentium 4, with 512MB RAM, running Linux 2.4.18. The computer word length is $w = 32$ bits.

For the text we used a concatenation of 7543 music pieces, whose total length is 1828089 bytes. The file was obtained by extracting the pitch values from MIDI files. The pitch values are in the range $[0 \dots 127]$. A set of 100 patterns were randomly extracted from the text. Each pattern was then searched for separately, and we report the average search times. We measured user times. We have separated the preprocessing and search times, which makes it easier to compare the search performance. Our preprocessing cost is considerably high, but this is amortized by large music collections that arise in practical applications.

6.1. Implementation

Several variants of the optimal multipattern algorithm were considered in [5]. For (δ, γ) -matching without transpositions, we used the basic single pattern algorithm. As the transpositions were implemented as multipattern search, we used bit-parallel counters and hierarchical verification in these cases, which give a considerable speed-up. For indels, we used the IndelMYE algorithm [22] for the final verifications. We ran each experiment with and without ordered ℓ -grams. The former is an order of magnitude faster in many cases, but it has higher preprocessing cost, justified only for large texts.

For all experiments we used $\ell = 2$. Due to the considerably large alphabet size, larger ℓ values were not practical. On the other hand, $\ell = 1$ gives in general poor results, especially combined with transpositions (but note that with bit-parallel counters even 1-grams are not guaranteed to match always, as different transposition ranges are mapped to different counters).

As the alphabet size was large (128), but most of the values occur in the middle of the range, we mapped the alphabet into the range $0 \dots 63$. That is, values $32 \dots 95$ were mapped to $0 \dots 63$, values $0 \dots 31$ to 0, and values $96 \dots 127$ to 95. This mapping allows us to use the original δ values. Verification was done using the original alphabet. This improves the preprocessing times, without worsening the search times.

We note that other alphabet mappings may make sense. In particular, for music applications, it might be acceptable to make the alphabet octave-independent, so that the same notes in different octaves are mapped to the same value.

6.2. Preprocessing Time

Table 1 gives the preprocessing times. For $mtid()$ and $mtdg()$ we have considered hierarchical verification because it gave consistently better results, so the preprocessing timings include all the hierarchy construction. Using ordered ℓ -grams increases the preprocessing cost, but improves the search performance.

Table 1. Preprocessing times in seconds for $\ell = 2$. The second timings are for ordered ℓ -grams.

$mtid(), m = 32$	$mdg(), m = 8$	$mdg(), m = 64$	$mtdg(), m = 32$
0.0699 / 0.2680	0.0048 / 0.0052	0.0067 / 0.0092	0.0936 / 0.5177

6.3. Transposition Invariant Indel Distance

We compared our approach against the LCTS algorithm [25], whose running time is $O(mn \log \sigma / \log w)$. Although the algorithm solves the dual problem, it could be adapted to searching with indel distance as well. We also compared against the bit-parallel dynamic programming algorithm IndelMYE [22], whose running time for a single transposition is $O(mn/w)$. We superimposed [29] all the transpositioned patterns and used hierarchical verification, in the same manner as in [5] with BPM algorithm. This works very well in practice, although the worst case complexity is still $O(\sigma mn/w)$. Fig. 2 shows the results for $m = 8 \dots 64$ and $k = 1 \dots 5$. Our

algorithm is by far the fastest for small k/m . LCTS is competitive only for very large k/m , while IndelMYE is the best choice for moderate k/m . Our algorithm clearly improves with ordered ℓ -grams, at the cost of higher preprocessing effort and memory requirements.

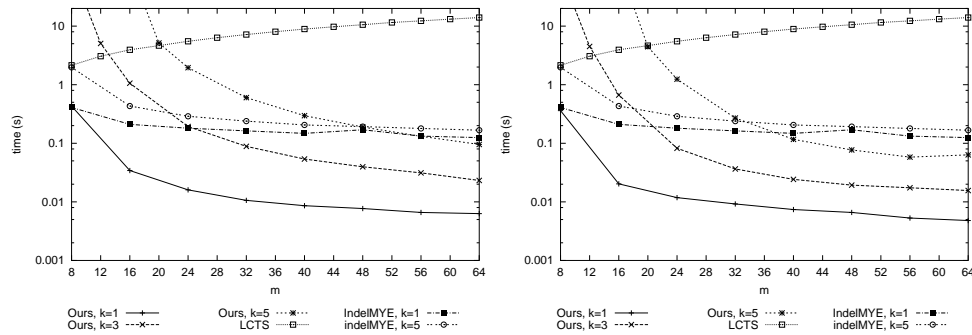


Fig. 2. Left: Search time in seconds for transposition invariant indel/LCS for $m = 8 \dots 64$. Right: The same with ordered ℓ -grams.

Fig. 3 shows the results for $m = 32$, $k = 1 \dots 6$ and $\delta = 0 \dots 2$. The LCTS algorithm cannot be applied for this setting. Being bit-parallel algorithm, IndelMYE can be easily adapted to this case by using classes of characters to implement δ . In this case we are again competitive against IndelMYE for small k/m , but only for very small δ . Ordered ℓ -grams boost the search considerably.

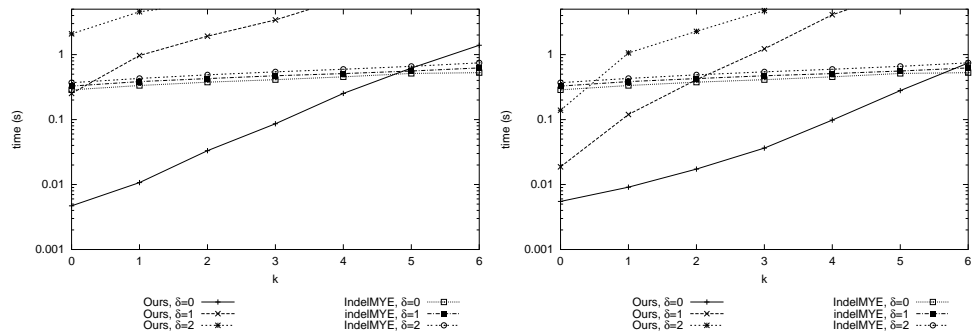


Fig. 3. Left: Search times in seconds for transposition invariant indel for $\delta = 1 \dots 3$, and $m = 32$. Right: The same with ordered ℓ -grams.

6.4. (δ, γ) -Matching

For (δ, γ) -matching we compared against the bit-parallel Forward matching algorithm (Fwd) of [9]. Fig. 4 shows the results for $m = 8 \dots 64$, $\delta = 1 \dots 3$ and $\gamma = m\delta/2$. Our algorithm is much more sensitive to increasing δ than Fwd, but for small δ values we are an order of magnitude faster. Using ordered ℓ -grams makes our algorithm more tolerant for increasing γ (but note that γ/m is constant here).

In [9] they give also bit-parallel backward matching algorithm, that is able to skip some text characters. The implementation restricts the pattern lengths to be at

most $\Theta(w/\log_2(\gamma))$. This means that in this experiment this algorithm is applicable only for the case $m = 8$, $\delta = 1$, and $\gamma = 8 * 1/2 = 4$. The algorithm takes 0.0063s average time, in this case, and marginally beats our algorithm (0.0065s)

Timings for $m = 32$, $\delta = 1 \dots 3$, and $\gamma = 4 \dots 40$ are shown in Fig. 5. (Note that for $\delta = 1$ there is no point for using $\gamma > m$.) Again, Fwd becomes eventually faster for large δ and γ , while our algorithm dominates for small parameter values. Fig. 6 repeats the experiment for transposition invariant (δ, γ) -matching. Note that no competitors exist in this case, although transposition superimposition and hierarchical verification could be applied for some of the existing (δ, γ) matching algorithms. However, observe that our transposition invariant algorithm is faster than Fwd algorithm (without transpositions) for small δ and γ .

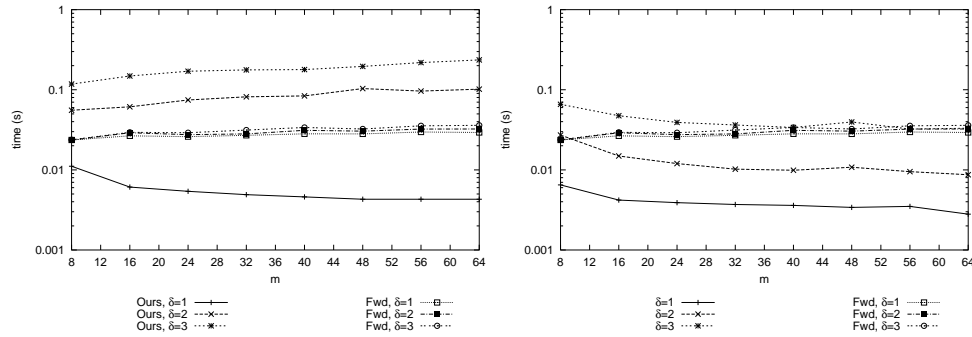


Fig. 4. Left: Search times in seconds for (δ, γ) -matching for $m = 8 \dots 64$ and $\delta = 1 \dots 3$. For each data point $\gamma = m\delta/2$. Right: The same with ordered ℓ -grams.

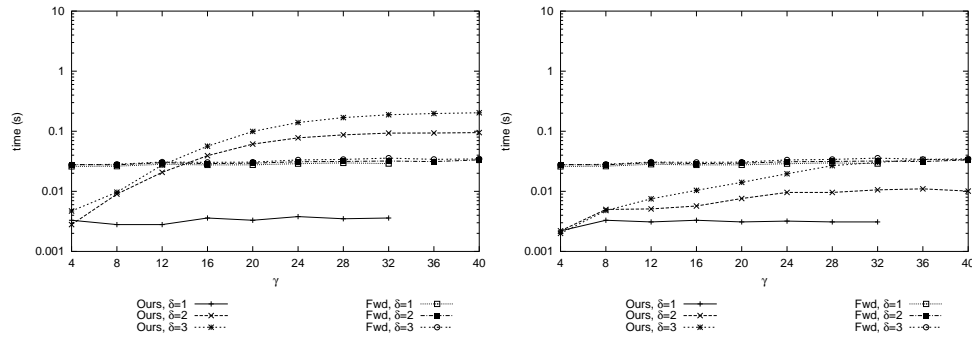


Fig. 5. Left: Search times in seconds for (δ, γ) -matching for $m = 32$, $\delta = 1 \dots 3$, and $\gamma = 4 \dots 40$. Right: The same with ordered ℓ -grams.

6.5. Comparison

We have separated the preprocessing and searching times in presenting the experimental results. This may seem unfair against the competing algorithms, and so it is for short texts. To show that our algorithms *are* competitive, Table 2 gives estimates for the minimum file sizes required to beat the competing approaches for

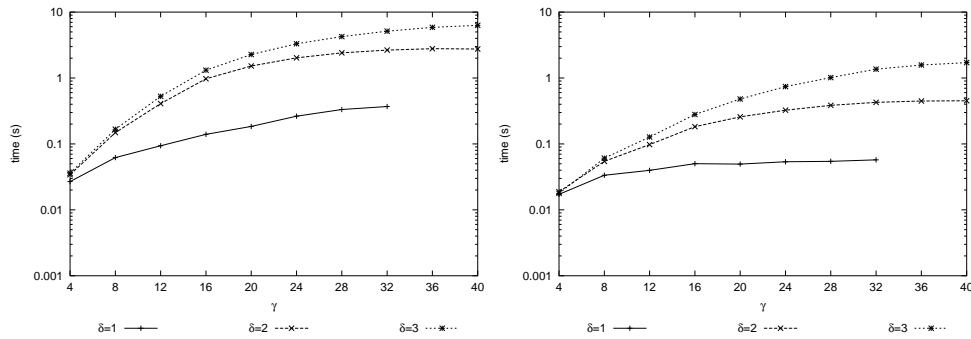


Fig. 6. Left: Search times in seconds for (δ, γ) -matching with transpositions for $m = 32$, $\delta = 1 \dots 3$, and $\gamma = 4 \dots 40$. Right: The same with ordered ℓ -grams.

various problem instances. These limits are quite modest, and for smaller parameter values even shorter files are sufficient.

Table 2. Examples of music file sizes where we begin to win for a few settings. The first row shows the parameter values, and the second row gives an estimate of the minimum file size where our algorithm wins its competitor. For smaller parameters shorter files would suffice. The estimates are for $m = 32$.

Indels		(δ, γ) -matching		
$k = 4, \delta = 0$	$k = 1, \delta = 1$	$(1, \infty)$	$(2, \infty)$	$(3, 24)$
> 0.61 Mb	> 1.77 Mb	> 0.46 Mb	> 0.71 Mb	> 1.52 Mb

7. Conclusions

We have presented new algorithms with applications to music retrieval, where several non-standard string matching problems arise. Many of our new algorithms are average-optimal, and the rest are very close to it. In several cases ours are the first algorithms that do not inspect all the text characters. Our algorithms are also very efficient in practice. The experiments show that for small to moderate error thresholds our algorithms are substantially faster than previous approaches for all but very short texts. These are the parameter values that are most interesting in most music retrieval applications.

In addition, our new algorithms are extremely flexible. We can solve many different problem variants essentially without any modifications to the search algorithms, only preprocessing changes according to the search model. In particular, we are able to solve some variants where no competing algorithms currently exist. These are transposition invariant indel with $\delta > 0$, and transposition invariant (δ, γ) -matching. Moreover, our algorithms can be used for multipattern search as well.

On the other hand, we have shown a basic difficulty of our algorithms to combine γ -matching with edit-like distances. It remains as an interesting open challenge to achieve this combination with optimal average complexity, that is $\Theta(n(k + \gamma/\sigma + \log_{\frac{\sigma}{1+\gamma/m}} m)/m)$. Another open problem is to close the gap between the lower

bounds and the complexities achieved for transposition invariant δ - and γ -matching.

Acknowledgements

This work was partially funded by the Academy of Finland, grant 202281 (Kimmo Fredriksson), and the Millennium Nucleus Center for Web Research, Grant P04-067-F, Mideplan, Chile (Gonzalo Navarro).

References

1. G. Navarro. A guided tour to approximate string matching. *ACM Computing Surveys*, 33(1):31–88, 2001.
2. T. Crawford, C. Iliopoulos, and R. Raman. String matching techniques for musical similarity and melodic recognition. *Computing in Musicology*, 11:73–100, 1998.
3. E. Cambouropoulos, T. Crawford, and C. Iliopoulos. Pattern processing in melodic sequences: Challenges, caveats and prospects. In *Proc. AISB'99*, pages 42–47, 1999.
4. P. Roland and J. Ganascia. Musical pattern extraction and similarity assessment. In E. Miranda, editor, *Readings in Music and Artificial Intelligence*, pages 115–144. Harwood Academic Publishers, 2000.
5. K. Fredriksson and G. Navarro. Average-optimal single and multiple approximate string matching. *ACM J. of Experimental Algorithmics*, 9(1.4), 2004.
6. G. Navarro and M. Raffinot. *Flexible Pattern Matching in Strings*. Cambridge University Press, 2002.
7. E. Cambouropoulos, M. Crochemore, C. Iliopoulos, L. Mouchard, and Y. J. Pinzon. Algorithms for computing approximate repetitions in musical sequences. In *Proc. AWOCA '99*, pages 129–144, 1999.
8. E. Cambouropoulos, M. Crochemore, C. S. Iliopoulos, L. Mouchard, and Y. J. Pinzon. Algorithms for computing approximate repetitions in musical sequences. *J. of Computational Mathematics*, 79(11):1135–1148, 2002.
9. M. Crochemore, C. Iliopoulos, G. Navarro, Y. Pinzon, and A. Salinger. Bit-parallel (δ, γ) -matching suffix automata. *J. of Discrete Algorithms*, 3(2–4):198–214, 2005.
10. R. Boyer and J. Moore. A fast string searching algorithm. *Comm. of the ACM*, 20(10):762–772, 1977.
11. M. Crochemore, C. Iliopoulos, T. Lecroq, Y. J. Pinzon, W. Plandowski, and W. Rytter. Occurrence and substrings heuristics for δ -matching. *Fundamenta Informaticae*, 55:1–15, 2003.
12. M. Crochemore and W. Rytter. *Text algorithms*. Oxford University Press, 1994.
13. P. Clifford, R. Clifford, and C. Iliopoulos. Faster algorithms for (δ, γ) -matching and related problems. In *Proc. CPM'05*, LNCS v. 3537, pages 68–78, 2005.
14. A. Amir, O. Lipsky, E. Porat, and J. Umanski. Approximate matching in the L_1 metric. In *Proc. CPM'05*, LNCS v. 3537, pages 91–103, 2005.
15. R. Cole and R. Hariharan. Verifying candidate matches in sparse and wildcard matching. In *Proc. STOC'02*, pages 592–601, 2002.
16. A. Amir and M. Farach. Efficient 2-dimensional approximate matching of half-rectangular figures. *Information and Computation*, 118(1):1–11, 1995.
17. R. Cole, C. Iliopoulos, T. Lecroq, W. Plandowski, and W. Rytter. On special families of morpishms related to δ -matching and don't care symbols. *Information Processing Letters*, 85(5):227–233, 2003.

18. V. Mäkinen. Sub-quadratic algorithm for weighted k -mismatches problem. Technical Report C-2004-1, Dept. of Computer Science, Univ. of Helsinki, 2004. <http://www.cs.helsinki.fi/u/vmakinen/papers/weightedkmm.ps.gz>.
19. D. Gusfield. *Algorithms on Strings, Trees and Sequences: Computer Science and Computational Biology*. Cambridge University Press, 1997.
20. K. Lemström and E. Ukkonen. Including interval encoding into edit distance based music comparison and retrieval. In *Proc. AISB'00*, pages 53–60, 2000.
21. M. Crochemore, C. Iliopoulos, Y. Pinzon, and J. Reid. A fast and practical bit-vector algorithm for the longest common subsequence problem. *Information Processing Letters*, 80(6):279–285, 2001.
22. H. Hyvrö, Y. Pinzon, and A. Shinohara. New bit-parallel algorithm for approximate string matching under indel distance. In *Proc. WEA '05*, LNCS v. 3503, pages 380–390, 2005.
23. V. Mäkinen, G. Navarro, and E. Ukkonen. Transposition invariant string matching. *J. of Algorithms* 56(2):124–153, 2005.
24. G. Navarro, Sz. Grabowski, V. Mäkinen, and S. Deorowicz. Improved time and space complexities for transposition invariant string matching. Technical Report TR/DCC-2005-4, Dept. of Computer Science, Univ. of Chile, 2005. <ftp://ftp.dcc.uchile.cl/pub/users/gnavarro/mnloglogs.ps.gz>.
25. S. Deorowicz. Speeding up transposition invariant string matching. Technical report, Institute of Computer Science, Silesian University of Technology, Poland, 2005. <http://www-zo.iinf.polsl.gliwice.pl/~sdeor/pub/deo05babs.htm>.
26. K. Lemström, G. Navarro, and Y. Pinzon. Practical algorithms for transposition-invariant string-matching. *J. of Discrete Algorithms*, 3(2–4):267–292, 2005.
27. A. C. Yao. The complexity of pattern matching for a random string. *SIAM J. of Computing*, 8(3):368–387, 1979.
28. G. Navarro and K. Fredriksson. Average complexity of exact and approximate multiple string matching. *Theoretical Computer Science*, 321(2–3):283–290, 2004.
29. R. Baeza-Yates and G. Navarro. New and faster filters for multiple approximate string matching. *Random Structures and Algorithms*, 20:23–49, 2002.
30. W. Chang and T. Marr. Approximate string matching and local similarity. In *Proc. CPM'94*, LNCS v. 807, pages 259–273, 1994.
31. K. Fredriksson, G. Navarro, and E. Ukkonen. Sequential and indexed two-dimensional pattern matching allowing rotations. Technical Report TR/DCC-2003-2, Dept. of Computer Science, Univ. of Chile, May 2003. To appear in *Theoretical Computer Science*.
32. W. Feller. *An Introduction to Probability Theory and Its Applications*. Vol. II, John Wiley, New York, 1966.