

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/234804842>

# Scaling and Time Warping in Time Series Querying

Article in *The VLDB Journal* · July 2008

DOI: 10.1007/s00778-006-0040-z · Source: DBLP

CITATIONS

168

READS

534

5 authors, including:



[Ada W. Fu](#)

The Chinese University of Hong Kong

179 PUBLICATIONS 11,968 CITATIONS

[SEE PROFILE](#)



[Leo Y. H. Lau](#)

The Chinese University of Hong Kong

2 PUBLICATIONS 236 CITATIONS

[SEE PROFILE](#)



[Chotirat Ann Ratanamahatana](#)

Chulalongkorn University

89 PUBLICATIONS 7,175 CITATIONS

[SEE PROFILE](#)



[Raymond Chi-Wing Wong](#)

The Hong Kong University of Science and Technology

189 PUBLICATIONS 4,651 CITATIONS

[SEE PROFILE](#)

# Scaling and time warping in time series querying

Ada Wai-Chee Fu · Eamonn Keogh ·  
Leo Yung Hang Lau · Chotirat Ann Ratanamahatana ·  
Raymond Chi-Wing Wong

Received: 10 October 2005 / Revised: 27 February 2006 / Accepted: 22 August 2006  
© Springer-Verlag 2007

**Abstract** The last few years have seen an increasing understanding that dynamic time warping (DTW), a technique that allows local flexibility in aligning time series, is superior to the ubiquitous Euclidean distance for time series classification, clustering, and indexing. More recently, it has been shown that for some problems, uniform scaling (US), a technique that allows global scaling of time series, may just be as important for some problems. In this work, we note that for many real world problems, it is necessary to combine both DTW and US to achieve meaningful results. This is particularly true in domains where we must account for the natural variability of human actions, including biometrics, query by humming, motion-capture/animation, and handwriting recognition. We introduce the first technique which can handle both DTW and US simultaneously, our techniques involve search pruning by means of a lower bounding technique and multi-dimensional

indexing to speed up the search. We demonstrate the utility and effectiveness of our method on a wide range of problems in industry, medicine, and entertainment.

**Keywords** Dynamic time warping · Nearest neighbor search · Scaled and warped matching · Subsequence matching · Uniform scaling

## 1 Introduction

We propose to query time series with both the accommodation of a scaling factor and the consideration of time warping effects. In this section we justify our proposal with some examples.

### 1.1 Justifying the need for uniform scaling and DTW

Because time series are near ubiquitous, and are becoming increasingly prevalent as our ability to capture and store them improves, there is increasing interest in the database community in techniques for efficiently indexing large time series collections [9, 27]. It is found that in most domains, it is necessary to match sequences with tolerance of small local misalignments, and dynamic time warping has been shown to be the right tool for this [6, 17, 35, 38, 29]. For example, in speech comparison, small fluctuation of the tempo of the speakers should be allowed in order to identify similar contents. More recently, it has been shown that in many domains it is also necessary to match sequences with the allowance of a global scaling factor [20]. In this work, we argue that for most real world problems, it is necessary to be able to handle both types of distortions simultaneously. In

---

A. W. -C. Fu (✉) · L. Y. H. Lau · R. C. -W. Wong  
The Chinese University of Hong Kong,  
Shatin, Hong Kong  
e-mail: adafu@cse.cuhk.edu.hk

L. Y. H. Lau  
e-mail: yhlau@cse.cuhk.edu.hk

R. C. -W. Wong  
e-mail: cwwong@cse.cuhk.edu.hk

E. Keogh  
University of California, Riverside, USA  
e-mail: eamonn@cs.ucr.edu

C. A. Ratanamahatana  
Computer Engineering Department,  
Chulalongkorn University,  
Bangkok, Thailand  
e-mail: ann@cp.eng.chula.ac.th

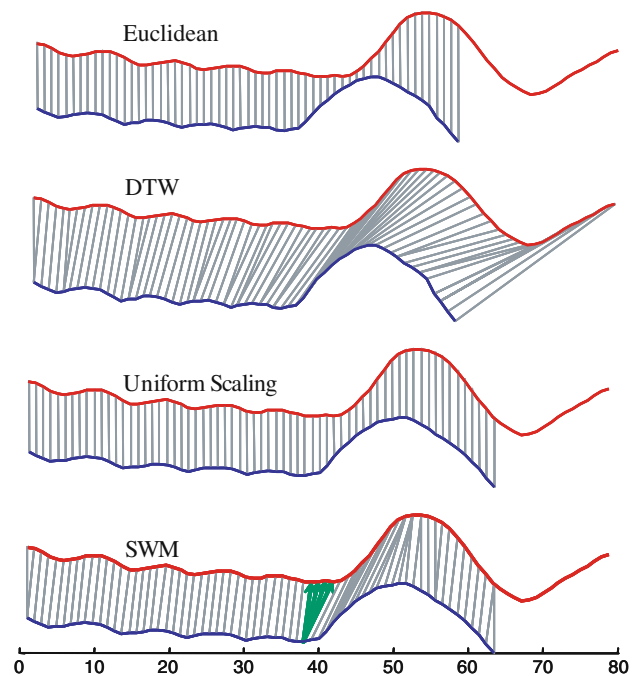
fact, even a casual glance of existing literature confirms this. For example, in query-by-humming systems, it is well-understood that we must allow for uniform scaling in addition to DTW. The current solution is to simply do DTW at many resolutions that span the possible range of tempos. For example, Meek and Birmingham [24] reports “We account for the phenomenon of persons reproducing the same tune at different speeds ... allow(ing) for nine tempo mappings.” However, repeating the query nine times clearly slows the system down. Furthermore, it is possible that the best match occurs somewhere in-between the nine discrete scalings. In [22], the authors also note that in addition to the local problems handled by DTW, “(people can) perform faster or slower than usual.” They again deal with this with multiple scaled queries, achieving reasonable performance only by the use of parallel processing.

Dynamic time warping is frequently used as the basis of gait recognition algorithms, but even in this highly structured domain, it is recognized that uniform scaling is also needed. For example, [16] notes “different gait cycles tend to have unequal lengths.” In fact, even if we discount human variability, it is well-understood that parallax effects from cameras (static or pan-and-tilt) can produce apparent changes in uniform scaling [15].

The need for uniform scaling has been noted in bioinformatics. Moeller-Levet et al. [25] noted that previous work that addressed *only* local scaling (with DTW) is inadequate, and they stressed that “(uniform) scaling factors in the expression level hide similar expressions and have to be eliminated or not considered when assessing the similarity between expression profiles” [16].

Finally, the simultaneous need for both uniform scaling and DTW is well understood in the motion-capture and animation community. For example, Pullen and Bregler [28], explaining their motion-capture editing system, noted “we stretch or compress the real data fragments in time by linearly resampling them to make them the same length as the keyframed fragment ... (then do DTW).” The computational difficulty of dealing with both uniform scaling and DTW at the same time has even led to practitioners abandoning temporal information altogether! Campbell and Bobick [4] used a phase space representation in which the velocity dimensions are removed, thus completely disregarding the time component of the data. This makes the learning and matching of motion patterns simpler and faster, but only at the cost of a massive increase in false positives.

Let us call “DTW with Uniform Scaling” **SWM**, which stands for *Scaled and Warped Matching*. In this paper, we study the combined effects of scaling and time warping in time series querying.



**Fig. 1** Two examples of an athlete’s trajectories aligned with various measures

## 1.2 Motivating examples

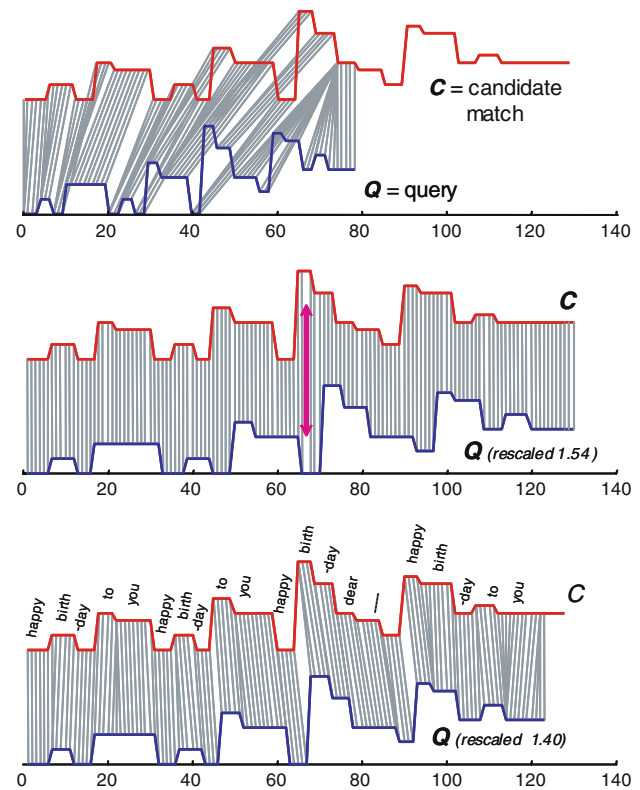
We present two concrete examples that require SWM to produce meaningful and intuitive results.

*Example 1 (Indexing video)* There is increasing interest in indexing sports data, both from sports fans who may wish to find particular types of shots or moves, and from coaches who are interested in analyzing their athletes’ performance over time. As a concrete example, we consider high jump. We can automatically collect the athlete’s center of mass information from video and convert the data into a time series (It is possible to correct for the cameras pan and tilt; see [8]). We found that when we issued queries to a database of high jumps, we got intuitive answers only when doing SWM. It is easy to see why if we look at two particular examples from the same athlete and consider all possible matching options, as shown in Fig. 1. In this figure, we show four different ways to match two time series, the horizontal axis is the time axis. In each case, we have shifted one of the two series upward to show the way the points in the two series are matched. Each vertical line in the diagrams shows the matching of two points. Visually, we can say that the two time series are similar, and hence the distance between them should be small. We want to see which of the four measurements can generate a result that gives a small distance as expected. From top to bottom:

- If we attempt simple Euclidean matching (after truncating the longer sequence), we get a large distance (which we can consider as error) because we are mapping part of the *flight* of one sequence to the *takeoff drive* in the other.
- If we simply use DTW to match the entire sequences, we get a large error because we are trying to explain part of the sequence in one attempt (the bounce from the mat) that simply does not exist in the other sequence. This problem can be corrected by constraints such as the Sakoe-Chiba Band, but without scaling, the matching will be poor.
- If we attempt just uniform scaling, we get the best match when we stretch the shorter sequence by 112%. However, the local alignment, particularly the *takeoff drive* and *up-flight*, is quite poor.
- Finally, when we match the two sequences with SWM, we get an intuitive alignment between the two sequences. The global stretching (once again at 112%) allows DTW to align the small local differences. In this case, the fact that DTW needed to map a single point in a time series onto 4 points in the other time series suggests an important local difference in one of these sequences. Inspection of the original videos by a professional coach suggests that the athlete misjudged his approach and attempted a clumsy correction just before his *takeoff drive*.

**Example 2 (Query by Humming)** The need for both local and global alignment when working with music has been extensively demonstrated [5,23,24,38]. For completeness, we will briefly review it here. Finding similar sequences of music has applications in copyright infringement detection, analyzing the evolution of music styles [5], automatic annotation, etc. (It is interesting to note that these studies are not confined to human endeavors; similar techniques have been used in animal “music” such as humpback whales and songbirds [23]). However, the vast majority of research in this area is used to support query by humming.

The basic idea of query by humming is to allow users to search large music collections by providing an example of the desired content, by humming (or singing, or tapping) a snippet. Clearly, humans cannot be expected to reproduce an exact fragment of a song, so the system must be invariant to certain distortions. Some of these are trivial to deal with. For example, the query can be made invariant to key by normalizing both the query and the database to a standard key. However, two types of errors are more difficult to deal with; users may perform the query at the wrong tempo, and users may insert or delete notes. The former corresponds with uniform



**Fig. 2** Two performances of *Happy Birthday to You* aligned with different metrics. Both performances were performed in the same key, but are shifted in the Y-axis for visual clarity

scaling, the latter with DTW. The music retrieval community has traditionally dealt with these two problems in two ways. The first is to do DTW multiple times, at different scalings [24]. However, this clearly produces scalability problems. The other common approach is to only do DTW with relatively short song snippets as queries believing that short sequences are less sensitive to uniform scaling problems than long sequences. While this is undoubtedly true, short snippets also have less discriminating power.

In Fig. 2, we demonstrate the problems with the universally familiar piece of music, *Happy Birthday to You*. For clarity of illustration, the music was produced by the fourth author on a keyboard and converted into a pitch contour; however, similar remarks apply to other music representations. From top to bottom:

- Because the query sequence was performed at a much faster tempo, direct application of DTW fails to produce an intuitive alignment.
- Rescaling the shorter performance by a scaling factor of 1.54 seems to improve the alignment, but note for example that the higher pitched note produced on the third “birth...” of the candidate is forced to

align with the lower note of the third “happy...” in the query.

- Only the application of *both* uniform scaling and DTW produces the correct alignment.

### 1.3 Contributions

This subsection summarizes the main contributions of our research.

1. We reviewed existing time series similarity measures, including Euclidean distance, dynamic time warping (DTW), and uniform scaling (US). We showed that these measures are inappropriate or insufficient for many applications.
2. We proposed a new time series similarity measure, scaled and warped matching (SWM), that combines the power of DTW and US to solve these real world problems.
3. We derived a lower bounding function for SWM that speed up time series matching, based on previous work in lower bounding DTW and US.
4. We experimentally showed the lower bounding function was effective in pruning most of the data in processing SWM query.
5. We proposed an optimization on the lower bounding function that further reduced query time. We showed the reduction in query time after applying the proposed optimization by extensive experiments.
6. We proposed the use of index to further improve the performance of time series matching under the proposed SWM distance.
7. We evaluated the usefulness and performance of the index in handling both one-nearest neighbor queries and  $k$ -nearest neighbor queries.

### 1.4 Organization

The rest of this article is organized as follow. Having developed the intuition for DTW and US, and having demonstrated the need to handle both types of distortions simultaneously, we will next define the problem of similarity measurement under SWM more formally in the following section. Section 3 explains dynamic time warping (DTW), uniform scaling (US), as well as the lower bounding distances (LB), on which this work was built. Section 4 introduces our work on scaled and warped matching (SWM). Section 5 suggests an optimization to speed up the computation of SWM. Section 6 explains the use of an index to further shorten query time. Section 7 reports on the empirical study on the

proposed mechanisms. Section 8 concludes this work and suggests possible future work.

## 2 Problem definition

Assume that we are given a database  $D$  which contains  $N$  time series (note that this formulation does not preclude the subsequence matching case, since it may be trivially transformed into this formulation). Further, assume that we are given a query  $Q$ , and a scaling factor  $l$ , where  $l \geq 1$ , which represents the maximum allowable stretching of the time series. The maximum allowable shrinking is implicitly set to  $1/l$ .<sup>1</sup> Hence when we compare a query sequence with a data sequence, either the query or the data can be shortened by a factor of up to  $1/l$  or lengthened by a factor of up to  $l$ .<sup>2</sup> Note that while  $1/l$  is bounded below by zero, and  $l$  is bounded from above by infinity, such loose bounds would allow pathological solutions to certain problems, and in any case are surely impossible to support efficiently. We therefore restrict our attention to scaling factors in the range  $0.5 \leq 1/l \leq 1 \leq l \leq 2$ . Note that this range encompasses the necessary flexibility documented in virtually every domain we are aware of. For instance, in [24], the authors reported excellent results with a query-by-humming system that allows “a (maximum) tempo scaling of 1.25.” Ref. [24] notes that in their experience, amateur singers can speed up their rendition of a song by as much as 200% or slow down to as little as 50%.

Recently, it has been shown that for nearly all types of time series data, using appropriate global constraints always improves the classification or clustering accuracy and the precision and recall of indexing [29]. Therefore a global constraint is typically enforced to limit the warping path to a roughly *diagonal portion* of the warping matrix.

Given  $N$  variable-length data sequences and a query sequence  $Q$ , we would like to find all data sequences that are “similar” to  $Q$ . Suppose the query sequence is  $Q = Q_1, Q_2, \dots, Q_m$ , where  $Q_i$  is a numerical value. We are interested in tackling the following problem.

**Problem** Assume the data sequences can be longer than the query sequence  $Q$ . Find the best match to  $Q$  in database, for any rescaling in a given range, under the dynamic time warping distance with a global constraint.

<sup>1</sup> Such formulation assumes the maximum allowable stretching and shrinking is symmetric. If this is not the case for a specific application, it is trivial to add as an extra parameter: the maximum allowable shrinking  $s$ .

<sup>2</sup> In our remaining discussions we sometimes talk about scaling the data, while at other times we talk about scaling the query.



By best match we mean the data sequence with the smallest distance from  $Q$ .

This problem has never been considered in the literature before. This problem is realistic in applications such as query by humming.

Before proceeding to review the existing distance measures, we note that, in some literature, uniform scaling may also refer to scaling of the values of a time series (scaling of the amplitude axis) [1,34]. However, this is not the focus of this research.

### 3 Preliminaries

In this section, we review separately time series querying with time warping distance and also querying with the scaling effect. For each case, we can apply a lower bounding technique for pruning the search space.

#### 3.1 Time warping distance

Intuitively, *dynamic time warping* is a distance measure that allows time series to be *locally* stretched or shrunk before applying the base distance measure. Definition 1 formally defines time warping distance.

**Definition 1** [Time warping distance (DTW)] Given two sequences  $C = C_1, C_2, \dots, C_n$  and  $Q = Q_1, Q_2, \dots, Q_m$ , the time warping distance DTW is defined recursively as follows:

$$\begin{aligned} \text{DTW}(\phi, \phi) &= 0 \\ \text{DTW}(C, \phi) &= \text{DTW}(\phi, Q) = \infty \\ \text{DTW}(C, Q) &= D_{\text{base}}(\text{First}(C), \text{First}(Q)) \\ &\quad + \min \begin{cases} \text{DTW}(C, \text{Rest}(Q)) \\ \text{DTW}(\text{Rest}(C), Q) \\ \text{DTW}(\text{Rest}(C), \text{Rest}(Q)) \end{cases} \end{aligned}$$

where  $\phi$  is the empty sequence,  $\text{First}(C) = C_1$ ,  $\text{Rest}(C) = C_2, C_3, \dots, C_n$ , and  $D_{\text{base}}$  denotes the distance between two entries.

Several metrics were used as the  $D_{\text{base}}$  distance in previous literature, such as Manhattan distance [37] and squared Euclidean distance [17,32]. We will use squared Euclidean distance as the  $D_{\text{base}}$  measure. That is,

$$D_{\text{base}}(C_i, Q_j) = (C_i - Q_j)^2$$

Note, we deliberately omit the final square root function in our distance definitions. Such optimization speeds up computations without altering the relative ranking given by these distances, which is more important than the actual value in most applications. The same optimization has been used before in [20]. However, if such

**Table 1** An example warping matrix aligning the time series  $\{1, 2, 2, 4, 5\}$  and  $\{1, 1, 2, 3, 5, 6\}$

5	27	27	13	5	<b>1</b>	<b>2</b>
4	11	11	4	<b>1</b>	2	6
2	2	2	<b>0</b>	1	10	26
2	1	1	<b>0</b>	1	10	26
1	<b>0</b>	<b>0</b>	1	5	21	46
	1	1	2	3	5	6

The warping path is highlighted in bold

**Table 2** An illustration of the relationship between each element and its adjacent elements in a warping matrix

$\text{DTW}(C, \text{Rest}(Q))$	$\text{DTW}(C, Q)$
$\text{DTW}(\text{Rest}(C), \text{Rest}(Q))$	$\text{DTW}(\text{Rest}(C), Q)$

The top right element  $\text{DTW}(C, Q)$  can be computed by looking up the values of the top left, bottom left and bottom right elements, which would have been computed already before the top right element

optimization is not desired, we can also consistently insert the final square root function without altering the essence of this work.

It is well-known that dynamic time warping distance can be computed by filling a *warping matrix* using a *dynamic programming algorithm* directly derived from the definition of time warping distance.

Table 1 shows an example warping matrix aligning the time series  $\{1, 2, 2, 4, 5\}$  and  $\{1, 1, 2, 3, 5, 6\}$ . Table 2 illustrates the relationship between each element and its adjacent elements. A *warping path* can be identified by tracing the elements in the warping matrix that were used to compute the time warping distance. Formally, a warping path  $W$  for two sequences  $Q$  and  $C$  is a sequence of elements  $w_1, w_2, \dots, w_p$  so that  $w_k = (i_k, j_k)$  is an entry in the warping matrix, where  $i_k \geq i_{k-1}$  and  $j_k \geq j_{k-1}$ ,  $\max(|Q|, |C|) \leq |W| \leq |Q| + |C| - 1$ .<sup>3</sup>

#### 3.2 Constraints and lower bounding

In the previous section we have explained with examples the importance of having global constraints on time warping in order to give meaningful results. Ref. [17] suggested a lower bounding measure based on such *global constraints* on time warping. Two commonly used global constraints exist. The *Sakoe–Chiba Band* limits the warping path to a band enclosed by two straight lines that are parallel to the diagonal of the warping matrix [32]. The *Itakura Parallelogram* [14] limits the warping path to be within a parallelogram whose major diagonal is the diagonal of the warping matrix.

<sup>3</sup>  $|X|$  denotes the length of a sequence  $X$ .

Ref. [17] viewed a global constraint as a constraint on the warping path entry  $w_k = (i, j)_k$  and gave a general form of global constraints in terms of inequalities on the indices to the elements in the warping matrix

$$j - r \leq i \leq j + r$$

where  $r$  is a constant for the Sakoe–Chiba Band and  $r$  is a function of  $i$  for the Itakura Parallelogram.

Incorporating the global constraint into the definition of dynamic time warping distance, Definition 1 can be modified as follows.

**Definition 2** [Constrained DTW (cDTW)] Given two sequences  $C = C_1, C_2, \dots, C_n$  and  $Q = Q_1, Q_2, \dots, Q_m$ , and the time warping constraint  $r$ , the constrained time warping distance cDTW is defined recursively as follows:

$$\begin{aligned} \text{Dist}_r(C_i, Q_j) &= \begin{cases} D_{\text{base}}(C_i, Q_j) & \text{if } |i - j| \leq r \\ \infty & \text{otherwise} \end{cases} \\ \text{cDTW}(\phi, \phi, r) &= 0 \\ \text{cDTW}(C, \phi, r) &= \text{cDTW}(\phi, Q, r) = \infty \\ \text{cDTW}(C, Q, r) &= \text{Dist}_r(\text{First}(C), \text{First}(Q)) \\ &\quad + \min \begin{cases} \text{cDTW}(C, \text{Rest}(Q), r) \\ \text{cDTW}(\text{Rest}(C), Q, r) \\ \text{cDTW}(\text{Rest}(C), \text{Rest}(Q), r) \end{cases} \end{aligned}$$

where  $\phi$  is the empty sequence,  $\text{First}(C) = C_1$ ,  $\text{Rest}(C) = C_2, C_3, \dots, C_n$ , and  $D_{\text{base}}$  denotes the distance between two entries.

The *upper bounding sequence*  $UW$  and the *lower bounding sequence*  $LW$  of a sequence  $C$  of size  $m$  are defined using the time warping constraint  $r$  as follows.

**Definition 3** [Enveloping sequences by Keogh [17]] Let  $UW = UW_1, UW_2, \dots, UW_m$  and  $LW = LW_1, LW_2, \dots, LW_m$ ,

$$UW_i = \max(C_{i-r}, \dots, C_{i+r}) \text{ and}$$

$$LW_i = \min(C_{i-r}, \dots, C_{i+r})$$

Considering the boundary cases, the above can be rewritten as

$$UW_i = \max(C_{\max(1, i-r)}, \dots, C_{\min(i+r, n)}) \text{ and}$$

$$LW_i = \min(C_{\max(1, i-r)}, \dots, C_{\min(i+r, n)})$$

These two sequences form an *envelope* which encloses the sequence  $C$ , as shown in Fig. 3.

The lower bounding measure by Keogh [17] bounds the time warping distance between two sequences  $Q$  and  $C$  by the Euclidean distance between  $Q$  and the envelope of  $C$  ( $D_{\text{base}}$  being the Euclidean distance).

Equation (1) below formally defines the lower bounding distance:

$$\text{LB}_W(Q, C) = \sum_{i=1}^m \begin{cases} (Q_i - UW_i)^2 & \text{if } Q_i > UW_i \\ (Q_i - LW_i)^2 & \text{if } Q_i < LW_i \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

### 3.3 Uniform scaling

Consider a query sequence  $Q = Q_1, \dots, Q_m$  and a candidate sequence  $C = C_1, \dots, C_n$ .

We assume that  $m$  is not greater than  $n$  ( $m \leq n$ ); hence, the query is typically shorter than the candidate sequence. We assume that the data can scale up or down by a factor of at most  $l$ , where  $l \geq 1$ . The entry  $Q_m$  may be matched to  $C_{lm}$  when the data is expanded by a factor of  $l$ . To simplify our discussion we shall assume that  $lm \leq n$ .

In order to scale time series  $C = C_1, \dots, C_q$  to produce a new time series  $C' = C'_1, \dots, C'_m$  of length  $m$ , we use the formula

$$C'_j = C_{\lceil j \cdot q/m \rceil} \quad \text{where } 1 \leq j \leq m$$

This is similar to the formula used in [20]. We target to find a scaled prefix in  $C$  to compare with  $Q$ . With a scaling factor of  $l$ ,  $q$  can range from  $\lceil m/l \rceil$  to  $lm$ .

**Definition 4** [Uniform scaling (US)] Given two sequences  $Q = Q_1, \dots, Q_m$  and  $C = C_1, \dots, C_n$  and a scaling factor bound  $l$ , where  $l \geq 1$ . Let  $C(q)$  be the prefix of  $C$  of length  $q$ , where  $\lceil m/l \rceil \leq q \leq lm$  and  $C(m, q)$  be a rescaled version of  $C(q)$  of length  $m$ ,

$$C(m, q)_i = C(q)_{\lceil i \cdot q/m \rceil} \quad \text{where } 1 \leq i \leq m$$

$$\text{US}(C, Q, l) = \min_{q=\lceil m/l \rceil}^{\min(lm, n)} D(C(m, q), Q)$$

where  $D(X, Y)$  denotes the Euclidean distance between two sequences  $X$  and  $Y$ .

Note that the ceiling function in the definition of  $\mathcal{C}(p, q)$  may be replaced by the floor function. The whole definition of  $\mathcal{C}(p, q)$  may also be replaced by some interpolation on the values of  $\mathcal{C}(q)_{\lfloor \cdot \rfloor}$  and  $\mathcal{C}(q)_{\lceil \cdot \rceil}$ .

### 3.4 Lower bounding uniform scaling

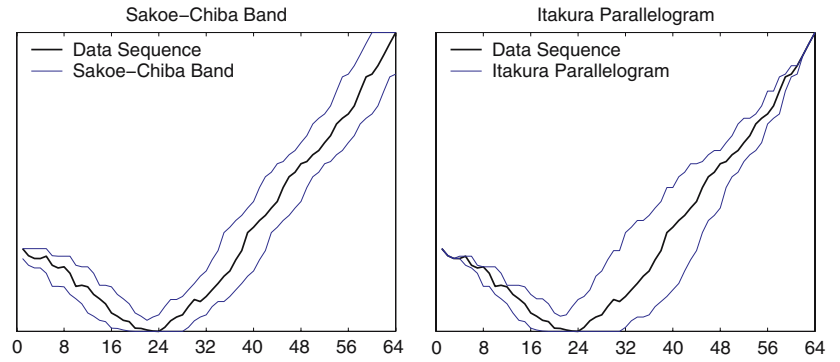
We create two sequences  $UC = UC_1, \dots, UC_m$  and  $LC = LC_1, \dots, LC_m$ , such that

$$UC_i = \max(C_{\lceil i/l \rceil}, \dots, C_{\lceil i \cdot l \rceil})$$

$$LC_i = \min(C_{\lceil i/l \rceil}, \dots, C_{\lceil i \cdot l \rceil})$$

These sequences bound the points of the time series  $C$  that can be matched with  $Q$ .

**Fig. 3** Enveloping sequences derived from two different constraints



The lower bounding function, which lower bounds the distance between  $Q$  and  $C$  for any scaling  $\rho$ , where  $1 \leq \rho \leq l$ , can now be defined as

$$LB_S(Q, C) = \sum_{i=1}^m \begin{cases} (Q_i - UC_i)^2 & \text{if } Q_i > UC_i \\ (Q_i - LC_i)^2 & \text{if } Q_i < LC_i \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

**Lemma 1** For any two sequences  $Q$  and  $C$  of length  $m$  and  $n$ , respectively, for any scaling constraint on the warping path  $w_k = (i, j)_k$  of the form  $j/l \leq i \leq lj$ , the value of  $LB_S(Q, C)$  lower bounds the distance between  $C$  and  $Q$  under a scaling of  $C$  between  $1/l$  and  $l$ , where  $l \geq 1$ .

*Proof* We can assume a matching path  $w_k = (i, j)_k$  which defines a mapping between the indices  $i$  and  $j$ , so that each such mapping constitutes a term  $(Q_i - C_j)^2$  to the required distance. We will show that each term  $t_{lb}$  in the square root of our lower bounding distance  $LB_S(Q, C)$  can be matched with a term  $t$  resulted from the one-to-one mapping, with  $t_{lb} \leq t$ .

Based on the constraints on the scaling factor, we have the constraint  $j/l \leq i \leq lj$  between  $i$  and  $j$  in  $w_k = (i, j)_k$ . From this, we have  $i/l \leq j \leq il$  and by definition

$$UC_i = \max(C_{\lceil i/l \rceil}, \dots, C_{\lceil il \rceil})$$

$$LC_i = \min(C_{\lceil i/l \rceil}, \dots, C_{\lceil il \rceil})$$

thus  $UC_i = \max(C_{\lceil i/l \rceil}, \dots, C_j, \dots, C_{\lceil il \rceil}) \geq C_j$ , or

$$Q_i - UC_i \leq Q_i - C_j$$

If  $Q_i > UC_i$  then  $Q_i - UC_i > 0$ , hence

$$(Q_i - UC_i)^2 \leq (Q_i - C_j)^2$$

Similarly, we can show that if  $Q_i < LC_i$  then

$$(Q_i - LC_i)^2 \leq (Q_i - C_j)^2$$

#### 4 Scaling and time warping

Having reviewed time warping, uniform scaling, and lower bounding, this section introduces scaling and time warping (SWM)[11].

**Definition 5** [Scaling and time warping (SWM)] Given two sequences  $Q = Q_1, \dots, Q_m$  and  $C = C_1, \dots, C_n$ , a bound on the scaling factor  $l$ , where  $l \geq 1$ , and the Sakoe-Chiba Band time warping constraint  $r$  which applies to a sequence of length  $m$ . Let  $C(q)$  be the prefix of  $C$  of length  $q$ , where  $\lceil m/l \rceil \leq q \leq \min(lm, n)$  and  $C(m, q)$  be a rescaled version of  $C(q)$  of length  $m$ ,

$$C(m, q)_i = C(q)_{\lceil i \cdot q/m \rceil} \quad \text{where } 1 \leq i \leq m$$

$$SWM(C, Q, l, r) = \min_{q=\lceil m/l \rceil}^{\min(lm, n)} cDTW(C(m, q), Q, r)$$

To simplify our discussion we shall assume that  $lm \leq n$ . We are interested in being able to scale the sequence and also to find nearest neighbor or evaluate range query by means of time warping distance. As noted in [20], a naive search for the uniform scaling problem alone requires  $O(|D| \cdot (a - b))$  time, where  $[b, a]$  is the range of lengths resulting from scaling. Time warping computation alone requires  $O(n^2)$  time for time series of length  $n$ . Hence we need to find a more efficient technique for the SWM problem.

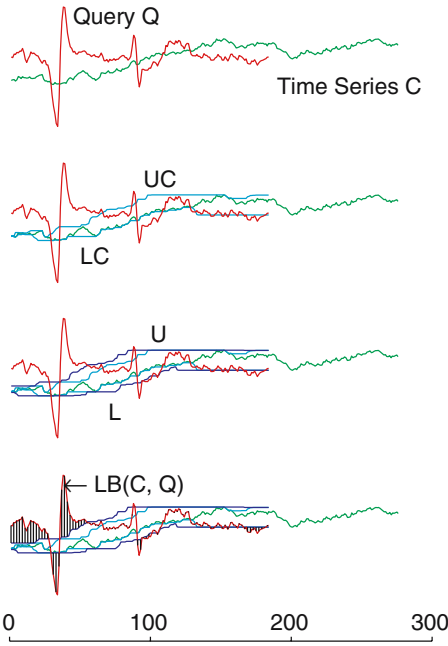
In previous sections, we reviewed the lower bounding technique for each sub-problem. Here, we propose to combine these lower bounds to form overall lower bounds for the querying problem. Figure 4 illustrates this graphically.<sup>4</sup>

We apply time warping on top of scaling, i.e., we scale the sequence first, and then measure the time warping distance of the scaled sequence with the query. Typically, time warping with Sakoe-Chiba Band constrains the warping path by a fraction of the data length, which

<sup>4</sup> In this example, the scaling factor is  $l = 1.5$ , the time warping constraint is  $r' = 10\%$  of the length of  $C$ .

□





**Fig. 4** An illustration of the SWM envelopes. From *top to bottom*: A time series  $C$  and a query  $Q$ ; the series  $C$  bounded from above and below, respectively, by  $UC$  and  $LC$ , the envelope for scaling; the series  $UC$  bounded above by  $U$  and  $LC$  bounded below by  $L$ , forming the overall envelope for scaling and time warping; and the lower bounding distance  $LB$  derived from the overall envelope

is translated into a constant  $r$ . Hence, if the fraction is 10%, then  $r = 0.1|C|$ . If the length of  $C$  is changed according to the scaling fraction  $\rho$ , that is,  $C$  is changed to  $\rho C$ , then the Sakoe–Chiba Band time warping constraint is  $r = 0.1|\rho C|$ . Hence, we have  $r = r'\rho$ , where  $r'$  is the Sakoe–Chiba Band time warping constraint on the unscaled sequence, and  $\rho$  is the scaling factor.

The lower envelope  $L_i$  and upper envelope  $U_i$  on  $C$  can be deduced as follows. Recall that the upper and lower bounds for uniform scaling between  $1/l$  and  $l$  are given by the following:

$$UC_i = \max(C_{\lceil i/l \rceil}, \dots, C_{\lceil il \rceil})$$

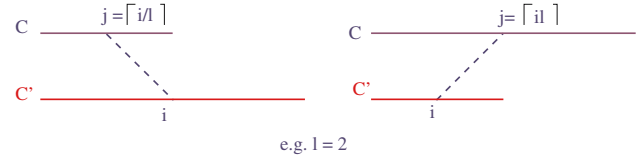
$$LC_i = \min(C_{\lceil i/l \rceil}, \dots, C_{\lceil il \rceil})$$

and the upper and lower bounds for a Sakoe–Chiba Band time warping constraint factor of  $r$  for a point  $C_i$  are given by

$$UW_i = \max(C_{\max(1, i-r)}, \dots, C_{\min(i+r, n)})$$

$$LW_i = \min(C_{\max(1, i-r)}, \dots, C_{\min(i+r, n)})$$

Therefore, when we apply time warping on top of scaling the upper and lower bounds will be



**Fig. 5** An illustration of the scaling effect, given a sequence  $C$ ,  $C'$  is the result after scaling. Note that the Sakoe–Chiba Band time warping constraint  $r'$  applies to  $C$ . Hence the range of  $j$  is given by  $\lceil i/l \rceil - r', \lceil il \rceil + r'$

$$U_i = \max(UW_{\lceil i/l \rceil}, \dots, UW_{\lceil il \rceil})$$

$$= \max(C_{\max(1, \lceil i/l \rceil - r')}, \dots, C_{\min(\lceil i/l \rceil + r', n)}, \dots,$$

$$C_{\max(1, \lceil il \rceil - r')}, \dots, C_{\min(\lceil il \rceil + r', n)})$$

$$= \max(C_{\max(1, \lceil i/l \rceil - r')}, \dots, C_{\min(\lceil il \rceil + r', n)}) \quad (3)$$

$$L_i = \min(LW_{\lceil i/l \rceil}, \dots, LW_{\lceil il \rceil})$$

$$= \min(C_{\max(1, \lceil i/l \rceil - r')}, \dots, C_{\min(\lceil i/l \rceil + r', n)}, \dots,$$

$$C_{\max(1, \lceil il \rceil - r')}, \dots, C_{\min(\lceil il \rceil + r', n)})$$

$$= \min(C_{\max(1, \lceil i/l \rceil - r')}, \dots, C_{\min(\lceil il \rceil + r', n)}) \quad (4)$$

The lower bound function which lower bounds the distance between  $Q$  and  $C$  for any scaling in the range of  $\{1/l, l\}$  and time warping with the Sakoe–Chiba Band constraint factor of  $r'$  on  $C$  is given by

$$LB(Q, C) = \sum_{i=1}^m \begin{cases} (Q_i - U_i)^2 & \text{if } Q_i > U_i \\ (Q_i - L_i)^2 & \text{if } Q_i < L_i \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

**Lemma 2** For any two sequences  $Q$  and  $C$  of length  $m$  and  $n$ , respectively, given a scaling constraint of  $\{1/l, l\}$  (see Sect. 2 on problem definition), where  $l \geq 1$ , and a Sakoe–Chiba Band time warping constraint of  $r'$  on the original (unscaled) sequence  $C$ , the value of  $LB(Q, C)$  lower bounds the distance of  $SWM(C, Q, l, r')$ .

*Proof* The matching warping path  $w_k = (i, j)_k$  defines a mapping between the indices  $i$  and  $j$ . Each such mapping constitutes a term  $t = (Q_i - C_j)^2$  to the required distance. We will show that the  $i$ -th term  $t_{lb}$  in our lower bounding distance  $LB(Q, C)$  can be matched with the term  $t$  resulting in a one-to-one mapping, with  $t_{lb} \leq t$ . For the  $i$ -th term  $t_{lb}$ , if  $Q_i > U_i$ , then  $t_{lb} = (Q_i - U_i)^2$ ; if  $Q_i < L_i$ , then  $t_{lb} = (Q_i - L_i)^2$ , otherwise  $t_{lb} = 0$ , which is always  $\leq t$ .

For scaling plus time warping, as illustrated in Fig. 5, the effective constraint on the range of  $j$  is given by:  $\lceil i/l \rceil - r' \leq j \leq \lceil il \rceil + r'$

By Eq. (3) and (4),

$$U_i = \max(C_{\max(1, \lceil i/l \rceil - r'), \dots, C_{\min(\lceil i/l \rceil + r', n)})$$

$$L_i = \min(C_{\max(1, \lceil i/l \rceil - r'), \dots, C_{\min(\lceil i/l \rceil + r', n)})$$

thus  $U_i \geq C_j$ , or  $Q_i - U_i \leq Q_i - C_j$

Hence if  $(Q_i > U_i)$  then  $Q_i - U_i > 0$  and we have

$$(Q_i - U_i)^2 \leq (Q_i - C_j)^2$$

Similarly, we can show that when  $(Q_i < L_i)$

$$(Q_i - L_i)^2 \leq (Q_i - C_j)^2$$

□

#### 4.1 Tightness of lower bounds

In this section, we show that the lower bounds we have described are tight. In general, to show that a lower bound is *tight*, we need to only find a case where the exact value is equal to the lower bound value. However, this is not exactly applicable in our scenario. For the three lower bounds  $LB_W(Q, C)$ ,  $LB_S(Q, C)$ , and  $LB(Q, C)$  we have discussed so far, the formulae have a similar pattern:

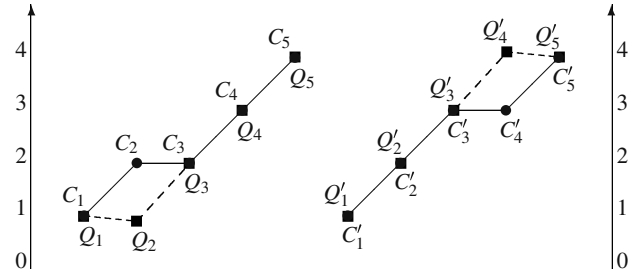
$$LB(Q, C) = \sum_{i=1}^m \begin{cases} (Q_i - U_i)^2 & \text{if } Q_i > U_i \\ (Q_i - L_i)^2 & \text{if } Q_i < L_i \\ 0 & \text{otherwise} \end{cases}$$

For each point in the query sequence, we have a lower envelope value, e.g.,  $L_i$  and an upper envelope value, e.g.,  $U_i$ , so that the sequence  $Q$  can compare in order to calculate the lower bounds. The values of  $L_i$  and  $U_i$  determine the lower bound value. We want to show that both  $L_i$  and  $U_i$  are “tight”. It can happen that for certain pairs of  $Q, C$ , the exact distance is equal to  $LB(Q, C)$  but in the computation of  $LB(Q, C)$  not both of  $L_i$  and  $U_i$  are used, and hence we cannot be sure that both  $L_i$  and  $U_i$  are set as tight as possible.<sup>5</sup> Hence we have the following definition for *tightness*.

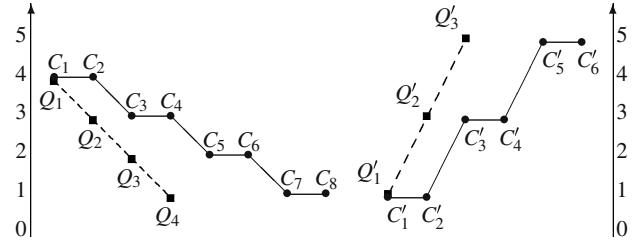
**Definition 6** Consider a lower bound  $LB(Q, C)$  for a distance  $D(Q, C)$  of the form

$$LB(Q, C) = \sum_{i=1}^m \begin{cases} (Q_i - U_i)^2 & \text{if } Q_i > U_i \\ (Q_i - L_i)^2 & \text{if } Q_i < L_i \\ 0 & \text{otherwise} \end{cases}$$

<sup>5</sup> As an example, suppose  $L_i$  is set to be the smallest possible value, so that  $Q_i < L_i$  is always false, the corresponding  $LB(Q, C)$  is still an lower bound, but obviously  $L_i$  is not set tightly. Yet we can have a pair of  $Q, C$ , where  $Q = C$  so that  $LB(Q, C) = D(Q, C)$ . Similarly, we can set  $U_i$  to be the greatest possible value so that  $Q_i > U_i$  is always false and similarly we can find  $Q = C$  so that  $LB(Q, C) = D(Q, C)$ . We can even set both  $L_i$  to be small and  $U_i$  to be big so that neither is tight but  $Q = C$  so that  $LB(Q, C) = D(Q, C)$ .



**Fig. 6** Example sequence pairs  $(Q, C)$  in Lemma 3



**Fig. 7** Example sequence pairs  $(Q, C)$  in Lemma 4

We say that the lower bound is *tight*, if there exists a set of sequence pairs so that for each pair  $\{Q, C\}$  in the set,  $D(Q, C) = LB(Q, C)$ , and the  $U_i$  and  $L_j$  values for some  $i, j$  are used [in the  $(Q_i - U_i)^2$  or  $(Q_j - L_j)^2$  term] at least once in computing the lower bounds in the set.

**Lemma 3** The lower bound  $LB_W(Q, C)$  for the DTW distance with the Sakoe–Chiba Band constraint is tight.

*Proof* Consider DTW with a Sakoe–Chiba Band constraint of  $r = 1$ . Hence in the warping path entry  $(i, j)$ ,  $j - 1 \leq i \leq j + 1$ .

Select two pairs of  $\{Q, C\}$  as follows (illustrated in Fig. 6):

$$Q = \{1, 0.9, 2, 3, 4\}, C = \{1, 2, 2, 3, 4\}$$

$$Q' = \{1, 2, 3, 4.1, 4\}, C' = \{1, 2, 3, 3, 4\}$$

It is easy to see that  $D(Q, C) = LB_W(Q, C)$ , and  $D(Q', C') = LB_W(Q', C')$ .

For  $Q, C$ ,  $Q_2 < LW_2$  and hence  $LW_2$  is used in the computation of  $LB_W(Q, C)$ .

For  $Q', C'$ ,  $Q'_4 > UW'_4$ , hence  $UW'_4$  is used in the computation of  $LB_W(Q', C')$ . □

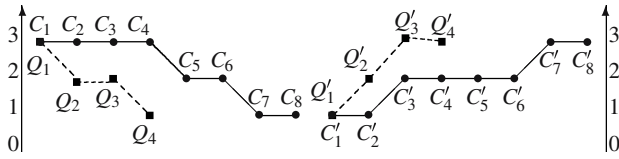
**Lemma 4** The lower bound  $LB_S(Q, C)$  for the distance between  $Q, C$  with a scaling factor between  $1/l$  and  $l$  is tight.

*Proof* Consider scaling between 0.5 and 2. Hence  $l = 2$ .

Select two pairs of  $\{Q, C\}$  as follows (illustrated in Fig. 7):

$$Q = \{4, 3, 2, 1\}, C = \{4.1, 4.1, 3.1, 3.1, 2.1, 2.1, 1.1, 1.1\}$$

$$Q' = \{1.1, 3.1, 5.1\}, C' = \{1, 1, 3, 3, 5, 5\}$$



**Fig. 8** Example sequence pairs  $(Q, C)$  in Lemma 5

It is easy to see that  $D(Q, C) = LB_S(Q, C)$ , and  $D(Q', C') = LB_S(Q', C')$ .

For  $Q, C$ ,  $LC_i > Q_i$  and all  $LC_i$  are used in the computation of  $LB_S(Q, C)$ .

For  $Q', C'$ ,  $UC'_i < Q'_i$  and all  $UC'_i$  are used in the computation of  $LB_S(Q', C')$ .  $\square$

**Lemma 5** The lower bound  $LB(Q, C)$  for the distance between  $Q, C$  with a scaling factor bound  $l$  and time warping with the Sakoe–Chiba Band constraint  $r'$  is tight.

*Proof* Consider a Sakoe–Chiba Band constraint of  $r' = 1$  and a scaling factor between 0.5 and 2. Hence  $l = 2$ .

Select two pairs of  $\{Q, C\}$  as follows (illustrated in Fig. 8):

$$Q = \{3, 1.9, 2, 1\}, C = \{3, 3, 3, 3, 2, 2, 1, 1\}$$

$$Q' = \{1, 2, 3, 1, 3\}, C' = \{1, 1, 2, 2, 2, 2, 3, 3\}$$

It is easy to see that  $SWM(Q, C, l, r') = LB(Q, C)$ , and  $SWM(Q', C', l, r') = LB(Q', C')$ .

For  $Q, C$ ,  $Q_2 < L_2$  and  $L_2$  is used in the computation of  $LB(Q, C)$ .

For  $Q', C'$ ,  $Q'_3 > U'_3$  and  $U'_3$  is used in the computation of  $LB(Q', C')$ .  $\square$

## 5 A faster and more flexible approach

This section suggests an optimization that both speeds up the lower bounding distance computation and increases the flexibility of the index built on top of this optimization. As reported in Sect. 7, this optimization can achieve considerable speed up compared to the original method.

### 5.1 The enveloping sequences revisited

The use of enveloping sequences to define lower bounding distances has been an increasingly popular technique in speeding up time series query processing. However, previous work that used this technique [17, 20, 38] were inconsistent when choosing whether to envelope the query sequence or the data sequence. To the best of our knowledge, previous literature did not compare the relative merits and demerits of enveloping the query

sequence or the data sequence. This may be partly because authors tend to assume that the query sequence and the data sequence take similar roles in the consideration of the distances.

Ref. [17] and [38] chose to envelope the query sequence. This has the advantage that, since the construction of the enveloping sequences is deferred until the actual processing of a query, the index construction algorithm and the index constructed are independent of the parameters controlling the envelopes, such as the time warping constraint in DTW and the scaling factor in uniform scaling. A single index can be built to support any queries that specify possibly different time warping constraint and scaling factor. The drawback of this approach is that the enveloping sequences are specific to individual query and must be rebuilt each time a query is processed.

Ref. [20] chose to envelope the data sequence, instead. One obvious reason is that they considered the length of the data sequence to be longer than the length of the query sequence, and the enveloping sequences are made of the same length as the query sequence. The immediate advantage of enveloping the data is that the enveloping sequences can be pre-computed. Instead of indexing the data sequences, these envelopes can also be indexed in any spatial index as MBRs. However, using this approach, the index is bounded to a particular envelope for each data sequence. Multiple index structures are required to answer queries specifying different time warping constraint or scaling factor.

### 5.2 Speeding up lower bounding distance computation

Following [20], the original approach to solving the SWM problem involves enveloping the data sequences. From the above discussion, it is also possible to envelope the query sequence, and as shown in [17] and [38], there can be considerable benefits. However, in the definition of SWM, the roles of  $Q$  and  $C$  are interchanged, instead of  $SWM(C, Q, l, r)$ , we compute  $SWM(Q, C, l, r)$ . The definition of the lower bounding distance based on this envelope is also slightly different.

Given the enveloping sequences,  $U, L$  of  $Q$ ,

$$U_i = \max(Q_{\max(1, \lceil i/l \rceil - r')}, \dots, Q_{\min(\lceil i/l \rceil + r', m)}) \quad (6)$$

$$L_i = \min(Q_{\max(1, \lceil i/l \rceil - r')}, \dots, Q_{\min(\lceil i/l \rceil + r', m)}) \quad (7)$$

The lower bound function which lower bounds the distance between  $Q$  and  $C$  for any scaling factor in  $SF = 1/l, l$  and time warping with the Sakoe–Chiba Band

constraint factor of  $r'$  on  $C$  is given by

$$LB(Q, C) = \sum_{i=1}^{\lceil m/l \rceil} \begin{cases} (C_i - U_i)^2 & \text{if } C_i > U_i \\ (C_i - L_i)^2 & \text{if } C_i < L_i \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

Note that in the above equation, instead of computing the sum of squared differences from  $i = 1$  up to  $i = m$ , the sum is computed only up to  $i = \lceil m/l \rceil$ . This modification is necessary because, in the extreme case, the query sequence can best match with the data subsequence  $C_1, \dots, C_{\lceil m/l \rceil}$  (after stretching of the query by a factor of  $l$ ).

We next show that this lower bound is tight by the two sets of sequences  $(Q, C), (Q', C')$ :

$Q = \{3, 2, 2, 2, 2\}$ ,  $C = \{3, 1.9, 2, 2, 2\}$  where  $m = 6, l = 2, r = 1$ . Obviously,  $LB(Q, C) = D(Q, C)$  and  $C_2 < L_2$ , so that  $L_2$  is used in the term of  $(C_i - L_i)^2$  in the computation.

$Q' = \{1, 1, 3, 3, 3\}$ ,  $C' = \{1, 1, 3.1, 3, 3\}$ , where  $m = 6, l = 2, r = 1$ .  $LB(Q', C') = D(Q', C')$  and  $C'_3 > U'_3$ , so that  $U'_3$  appears in the term of  $(C'_i - U'_i)^2$  in the computation.

Since the number of terms in the summation in  $LB(Q, C)$  have been reduced from  $m$  to  $\lceil m/l \rceil$ , we may expect that this lower bound is not as selective as the lower bound for enveloping the data sequences. In general, the best match between  $Q$  and  $C$  may make use of more than  $m/l$  entries of  $C$  and the entries that are not included in the summation introduce more inaccuracy. From our experiments this modification affects the pruning power of the lower bounding distance, but it also reduces the computation time significantly. This can be justified because  $m - \lceil m/l \rceil$  additions, subtractions, and multiplications have been saved from each lower bound function computation. Since the computation is performed for each and every sequence, the aggregate saving is more than enough to compensate for the time required to process the increased false alarms. Furthermore, as the nearest neighbor search proceeds, the distance to the current nearest neighbor converges quickly as soon as a few actual SWM distances, which are much better estimates to the actual nearest neighbor distance, are computed. In fact, similar results on DTW have been reported in the literature [29, 30].

## 6 Indexing for SWM

Next, we investigate the use of an index in speeding up the computation of SWM queries. Both scaling and time warping are to be handled at the same time. Furthermore, we would like our index to be able to support subsequence matching, in addition to whole sequence

matching. That is, we would like to find the best matching time series subsequences in the database for a given query.

### 6.1 Related work

The following subsections review existing subsequence matching algorithms. These algorithms worked under the Euclidean distance metric.

#### 6.1.1 Fast subsequence matching

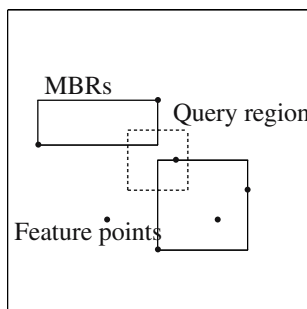
Faloutsos et al. [9] proposed the now ubiquitous method of indexing time-series subsequences under the Euclidean distance metric. The idea is to place a sliding window on every possible position of every data sequence. For each such placement, a subsequence is extracted. Dimensionality reduction technique such as Discrete Fourier Transform is applied to reduce the subsequence to a feature point in the  $f$ -dimensional space. Feature points from nearby windows are grouped together to span a minimum bounding rectangle (MBR). These MBRs are stored in an ordinary spatial index such as the R-Tree [12] and the R\*-Tree [2]. MBRs are stored instead of individual feature points because the number of feature points can be as high as  $O(nl)$  where  $n$  is the number of original data sequence and  $l$  is the length of each data sequence.

To process a range query with range being  $\epsilon$ , the query sequence is divided into consecutive disjoint subsequences. Each subsequence is again transformed to a feature point in the  $f$ -dimensional space. A range query is performed for each feature point. The union of the results from these range queries form the candidate set. It is proven that if the original tolerance (range) is  $\epsilon$  and the query is divided into  $p$  subsequences, no false dismissal is generated even if each range query has a tolerance as low as  $\epsilon/\sqrt{p}$ .

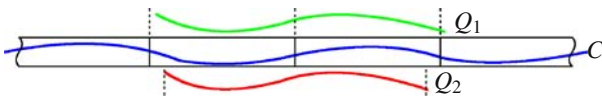
#### 6.1.2 Duality-based subsequence matching

More recently, Moon et al. [27] proposed another approach for subsequence indexing that is a dual approach to the one described above. This is a dual approach because the roles of the data sequence and the query sequence are exchanged. Each data sequence is divided into consecutive disjoint subsequences. Each subsequence can be indexed using a spatial index.

In order to guarantee no false dismissal, a sliding window must be placed on a query sequence. Naively, a range query can be performed on each subsequence obtained. However, this requires repeated accessing of the index structure, incurring extra page accesses.



**Fig. 9** It is less likely that a query region will overlap with the feature points (compared to overlapping with an MBR)



**Fig. 10** Query must be sufficiently long to include at least one disjoint subsequence at every alignment with the data sequence  $C - Q_1$  is long enough but  $Q_2$  is too short

Consequently, instead of performing range query individually, several subsequences can be grouped together to span an MBR. These MBRs are appropriately enlarged by the user specified tolerance at each dimension. A query is performed based on each resulting enlarged MBR. The union of the results form the candidate set. The actual sequences corresponded to the entries in the candidate set are retrieved to compute the actual distance from the query.

The advantage of Moon's approach over Faloutsos' is that individual feature points, instead of MBRs, are stored. It is anticipated that fewer false alarms will be generated because it is less likely that a query region will overlap with the feature points (compared to overlapping with an MBR), as illustrated in Fig. 9. Our enveloping technique on the query translates naturally to query regions (or query MBRs). It is intuitive that enclosing the query regions in an MBR would produce fewer false alarms than enclosing the data with another MBR.

However, since a query may match a data subsequence at any offset but only disjoint subsequences are stored in the database, a query matching a data subsequence may *not* necessarily contain a subsequence corresponding to the disjoint subsequence. To guarantee the correctness of Moon's approach, a query must be sufficiently long to include at least one disjoint subsequence at every alignment with the data sequences as shown in Fig. 10.

Formally, Moon et al. [26,27] proved the following lemma about the number of disjoint subsequences included by a query of a given length.

**Lemma 6** *If a sequence  $C$  is divided into disjoint subsequences of length  $\omega$ , the minimum number of disjoint subsequences  $p$  included by a query of length  $l$  is given by  $p = \lfloor (l + 1)/\omega \rfloor - 1$*

As a consequence of Lemma 6, Moon et al. [26,27] also proved the following lemma about the relationship between the maximum length of a disjoint subsequence and minimum length of a query holds.

**Lemma 7** *If the minimum length of a query sequence is given by  $\min(Q)$ , then the maximum length of a disjoint subsequence is given by  $\lfloor (\min(Q) + 1)/2 \rfloor$*

### 6.1.3 Dimension reduction

It is well-known in the database community that most tree indexing structures suffer from “the curse of dimensionality”, that is, the performance of these tree indices start to degrades increasingly rapidly as the number of dimensions of data increases. And, eventually, it becomes so slow that even a linear scan of data is faster, rendering the index useless. Fortunately, it is still possible to make use of indices to speed up query by indexing a dimension-reduced version of the data. Previous work [17,18,21,36] showed promising results that only a very small dimension was necessary for an index to effectively speed up query computation, making the index very compact.

Several dimension reduction techniques exist. Among them, the Piecewise aggregate approximation, otherwise known as piecewise constant approximation or segmented means, or PAA in short, was shown to be an intuitive and extremely fast dimension reduction method that is especially suitable for time series [18,36,38]. Zhu and Shasha [38] proved that the lower bounding property for time warping distances is preserved when defined over the PAA version of the enveloping sequences. With similar arguments the lower bounding property is also preserved for SWM with PAA.

### 6.2 Proposed indexing for SWM

This subsection describes an index that supports scaled and warped matching of subsequences.

In essence, the enveloping technique enables the construction of indices supporting scaled and warped matching (and in general, any measures that do not follow the triangular inequality) by introducing a lower bounding distance that satisfies the triangular inequality. Consequently, almost any spatial access methods such as the R-Tree [12], the R\*-Tree [2], and the X-Tree [3] can be used as the underlying indexing structure. Applying different envelopes on the query sequence allows the



index to support arbitrarily-defined measures for which the enveloping technique is applicable, such as scaled and warped matching.

### 6.2.1 Index construction algorithm

Algorithm 1 shows the index construction algorithm. Each data sequence is divided into  $\lfloor |C|/\omega \rfloor$  disjoint subsequences each of length  $\omega$ , where  $|C|$  denotes the length of a sequence  $C$ . Piecewise aggregate approximation is used to reduce the dimension of these subsequences by constructing the PAA representation for each subsequence. The PAA representation of these subsequences is inserted into the index structure, together with its ID. Note that only disjoint subsequences are stored, so the index can be kept to a reasonable size. However, these subsequences are stored directly into the index; no extra MBRs are created, leaving the responsibility of grouping feature points back to the underlying spatial access method.

We assume that the index is a tree structure, and that non-leaf nodes keep pointers and bounding rectangles for child nodes, while a leaf node contains the raw data, which is a set of subsequences. Typically, a leaf node is not designed to hold a single subsequence since it is preferred to utilize a page for each node.

---

#### Algorithm 1 Index construction

---

```

1: Initialize the index
2: for all data sequences  $C$  do
3:   Divide  $C$  into  $\lfloor |C|/\omega \rfloor$  disjoint subsequences, each of length
      $\omega$ 
4:   for all disjoint subsequences do
5:     Construct the PAA representation of the subsequence
6:     Insert the PAA representation of the subsequence into
       the index, together with its ID
7:   end for
8: end for

```

---

### 6.2.2 Utilizing the index for range queries

The key to the index for scaled and warped matching lies mostly on the query processing algorithm, as shown in Algorithm 2.

Given a range query sequence  $Q$  with tolerance  $\epsilon$ , the minimum number of disjoint subsequences  $p$  included by  $Q$  is computed. The lower and upper bounding sequences are constructed from  $Q$ . Lower and upper bounding subsequence pairs are then extracted by placing a sliding window on every possible position of the lower and upper bounding sequences. The PAA representation for each subsequence is constructed.

Each pair of subsequences spans an MBR. Enlarge the MBRs by the tolerance  $\epsilon/\sqrt{p}$  (see Sect. 7.3, p 18). Range queries<sup>6</sup> are performed using the enlarged MBRs and the results are stored in the candidate set. For each resulting candidate, the actual sequence  $C$  is retrieved to perform post-processing step. For each data suffix of  $C$ , find the best match by trying all possible scalings. The sequence together with the offset and the best scaling factor is returned if the actual distance lies within the given tolerance  $\epsilon$ .

---

#### Algorithm 2 Query utilizing the index for range query

---

**Require:**  $|Q| \geq 2\omega - 1$

```

1: Compute the minimum number of disjoint subsequences  $p$ 
   included by  $Q$ 
2: Construct the lower and upper enveloping sequences from the
   query sequence  $Q$ 
3: Extract lower and upper bounding subsequences by placing
   a sliding window on every possible position of the sequences.
   Each pair of lower and upper bounding subsequences spans
   an MBR
4: Construct the PAA representation for each subsequence
5: Enlarge the MBRs by the tolerance  $\epsilon/\sqrt{p}$ 
6: for all enlarged MBRs do
7:   Perform range query on the index using the enlarged MBR
8:   Store the query result in the candidate set
9: end for
10: for all candidate in the candidate set do
11:   Retrieve the actual sequence  $C$ 
12:   for offset = 0 to  $(|C| - |Q|)/\text{maximum scaling factor}$  do
13:     Find the best match by trying all possible scalings
14:     Return the sequence together with the offset and best
       scaling factor lying within the given tolerance  $\epsilon$ 
15:   end for
16: end for

```

---

Algorithm 2 returns all the sequences that match the query sequence at a certain offset at a certain scaling factor. However, a large number of range queries are required at line 6. This can be reduced by grouping several query MBRs to form a larger query MBR.

### 6.2.3 Nearest neighbor search

It has been noted that nearest neighbor search is more useful in many occasions [7,13,31,33]. Therefore, we propose to enhance the query algorithm to also handle this type of query.

Algorithm 3 shows our proposed  $k$ -nearest neighbor query algorithm under the SWM distance. Line 1 constructs the enveloping sequences from the query sequence. The query regions can then be constructed from the enveloping sequences at line 2 by placing a

<sup>6</sup> The bounding box defines a high-dimensional range in search space for locating data.

sliding window on every possible position of the enveloping sequences. For simplicity, the query regions for the sliding windows are merged into one rectangle, *Query Rect*. The number of nearest neighbors reported is then initialized to zero at line 3.

Next, two priority queues are created at line 4–5. The first queue *Candidate* stores subsequences that are potential  $k$ -nearest neighbors while the second queue, named *Queue*, stores the non-leaf nodes, the leaf nodes, and the subsequences stored at the leaf nodes during the tree traversal. *Queue* is the same priority queue used in [13]. The queues are sorted by the stored *SWM* or lower bounding distances of the elements.

The root node is first enqueued into the queue *Queue* at line 6 before entering the main loop between line 7 and line 29. At line 8–19, the best-first search strategy is used to traverse the tree while enqueueing non-leaf nodes, leaf nodes, and subsequences into the queue *Queue*. However, when a subsequence is encountered at line 9, it is enqueued into the *Candidate* queue for processing at line 20–28. The value of LB (*QueryRect*, *Subsequence*) is computed for non-leaf nodes. This LB is the minimum distance between the two rectangles, which is defined as follows. Let  $A$  and  $B$  be two  $n$ -dimensional rectangles. Let the range of  $A$  at dimension  $i$  be  $[LA_i, uA_i]$ , and the range for  $B$  be  $[LB_i, uB_i]$ , then

$$LB(A, B) = \sum_i^n |r_i|^2 \text{ where}$$

$$\text{if } uA_i < LB_i \text{ then } r_i = LB_i - uA_i$$

$$\text{if } uB_i < LA_i \text{ then } r_i = LA_i - uB_i$$

otherwise,  $r_i = 0$ .

At line 23, the *SWM* distance between the *Query* and the *Sequence* is computed to replace the lower bounding estimate and the *Sequence* is re-enqueued into the *Candidate* queue with the *SWM* distance, which is also called the key in the algorithm. Therefore, when an element is inserted in the *Candidate* queue, we mark whether the distance (key) is *Estimated* or *NotEstimated* (meaning the actual *SWM* distance).

After this re-enqueuing, if the weighted distance (key) comparison between the first elements of the two queues says that the *Candidate* queue has a smaller distance, we are certain that there are no other subsequences that are nearer to the *Query* sequence than the first element of the *Candidate* queue. And the element is output at line 25 and the corresponding nearest neighbor count is incremented at line 26.

In most other  $k$  nearest neighbor search algorithms such as [10, 31, 33], the actual distance  $d_{\max}$  of the current  $k$ -th nearest neighbor is used to prune other candidates.

### Algorithm 3 $k$ -nearest neighbor query under *SWM*

```

1: Envelope  $\leftarrow$  CONSTRUCTENVELOPINGSEQUENCES(Query)
2: QueryRect  $\leftarrow$  MAKEQUERYRECT(Envelope)
3: NNcount  $\leftarrow$  0
4: Candidate  $\leftarrow$  NEWPRIORITYQUEUE()
5: Queue  $\leftarrow$  NEWPRIORITYQUEUE()
6: ENQUEUE(Queue, R-tree.RootNode, 0)
7: while not ISEMPTY(Queue) and NNcount <  $k$  do
8:   Element  $\leftarrow$  DEQUEUE(Queue)
9:   if Element is a Subsequence of Sequence then
10:    ENQUEUE(Candidate, Subsequence, LB(QueryRect,
11:      Subsequence), Estimated)
12:   else if Element is a leaf node then
13:     for all entry (Subsequence) in leaf node Element do
14:       ENQUEUE(Queue, [Subsequence], LB(QueryRect,
15:         Subsequence))
16:   end for
17: else {Element is a non-leaf node}
18:   for all entry (Node, Rect) in node Element do
19:     ENQUEUE(Queue, Node, LB(QueryRect, Rect))
20:   end for
21: end if
22: while not(ISEMPTY(Candidate)) and NNcount <  $k$  and
23:   FIRST(Candidate).Key  $\leq$  FIRST(Queue).Key  $\cdot |Query|/\omega$  do
24:   Element  $\leftarrow$  DEQUEUE(Candidate)
25:   if ESTIMATED(Element) then
26:     ENQUEUE(Candidate, Element.subsequence,
27:       SWM(Query, Element.subsequence), NotEstimated)
28:   else
29:     Report Element as the next nearest object
30:     NNcount  $\leftarrow$  NNcount + 1
31:   end if
32: end while
33: end while

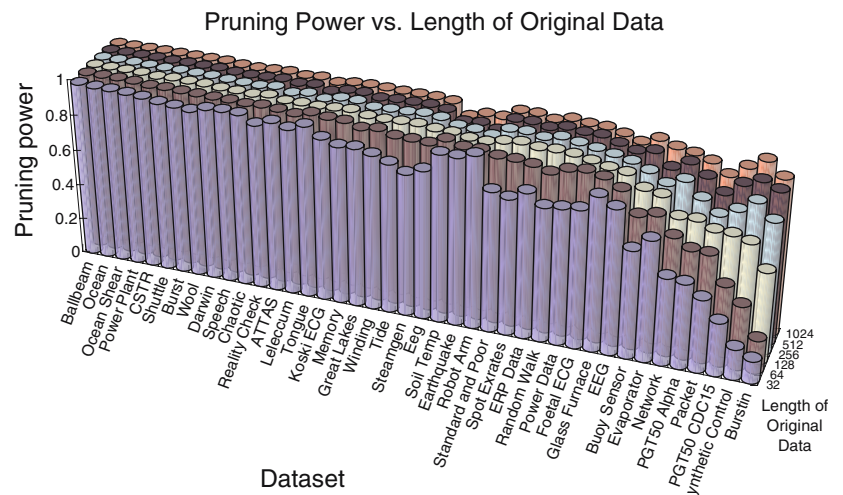
```

$d_{\max}$  is dynamically updated as nearer neighbors are uncovered and whenever a candidate is dequeued, the actual distance (*SWM* in our case) is computed. In our algorithm, when an element is dequeued, the *SWM* is not computed unless it is the smallest in terms of LB. We do not keep track of  $d_{\max}$  and use it for pruning. Instead, the nearest neighbors are returned from the head of the *Candidate* queue one at a time when it satisfies the condition at line 20. In this way we may reduce the number of *SWM* computations.

## 7 Experimental evaluation

This section describes the experiments carried out to verify the effectiveness of the proposed lower bounding distance. The experiments were executed on an Intel Xeon 2.2 GHz Linux PC with 1 GB RAM. The source code for the experiments was written in C Language. MATLAB was also used for pre-processing the raw data.

To evaluate the effectiveness of the proposed lower bounding distance, and thus the proposed solution, an objective measure of the quality of a lower bounding

**Fig. 11** Pruning power versus length of original data

distance is required. The *Pruning power*  $P$  is defined in [17] as follows.

$$P = \frac{\text{Number of objects that do not require full DTW}}{\text{Number of objects in database}}$$

The Pruning power is an objective measure because it is free of implementation bias and choice of underlying spatial index. This measure has become a common metric for evaluating the efficiency of lower bounding distances. Therefore, it was adopted in evaluating the proposed lower bounding distance, by replacing DTW by SWM in the above formula.

Extensive experiments were conducted on as many as 41 different datasets. These datasets, which represent time series from different domains, were obtained from “The UCR Time Series Data Mining Archive” [19].

As the datasets came from a wide variety of different domains, they differed significantly in size and in the length of individual data sequences. In order to produce meaningful results, both parameters must be controlled. Thus, from each original dataset, we derived six sets of data, each containing 1,024 data sequences, with variable lengths of 32, 64, 128, 256, 512 and 1,024, respectively. Short sequences were produced by using only prefixes of the original datasets while long sequences were produced by concatenating original sequences. After that we also normalized the derived datasets using MATLAB by subtracting the mean value from the entries of each sequence and dividing them by their standard deviation, so that each data sequence has a mean zero and a standard deviation of one. All experiments were conducted on these derived datasets.

### 7.1 Enveloping data

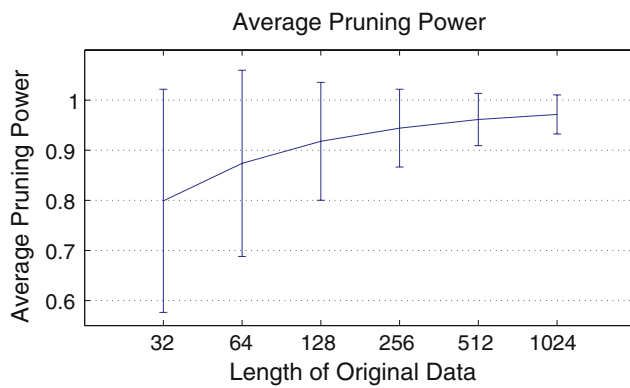
To compute the pruning power of the proposed lower bounding distance, the one-nearest neighbor search was

performed using the linear-scan algorithm. A random subsequence was chosen from the dataset to act as the query, and the remaining 1,023 sequences acted as the data. The search was repeated for 50 trials using a different subsequence as query. The actual dynamic time warping distance did not need to be calculated if the lower bounding measure gave a value larger than the time warping distance of the current nearest neighbor. The fraction of sequences that did not require calculation of actual time warping distance became the pruning power of the lower bounding measure in that query. The average of the 50 queries was reported as the pruning power of that particular dataset.

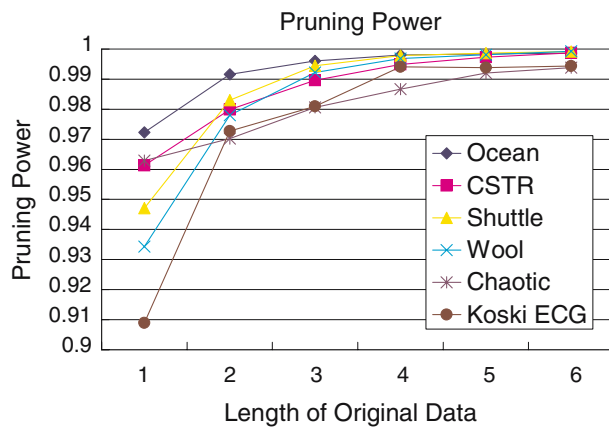
Unless stated otherwise, in all experiments, the length of data was 1,024 data points; the scaling factor was between 1.5 and its reciprocal; the length of query was set so that the longest rescaled query is at most as long as the data, or the length of the query is not longer than the shortest rescaled data; and the width of the Sakoe–Chiba Band was set to 10% of the length of the query. In fact, recent evidence suggests that this is a pessimistic setting, and real world problems benefit from even tighter constraints [29].

Figure 11 shows how the pruning power of the proposed lower bounding measure varies as the lengths of data change on different datasets.<sup>7</sup> For a majority of datasets, the pruning power increased with the length of data, suggesting that the proposed algorithm is likely to perform well in real-life environment, in which long sequences of data are collected for a long period of

<sup>7</sup> The comprehensiveness of our experiments means that some of our figures are very dense and somewhat hard to read. To help the interested reader, we have created webpage which contains large versions on all the figures in addition to tables of raw numbers. This additional information will be maintained on the websites of the first two authors.



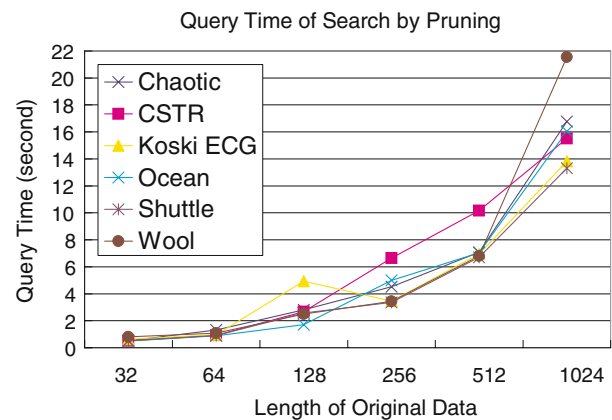
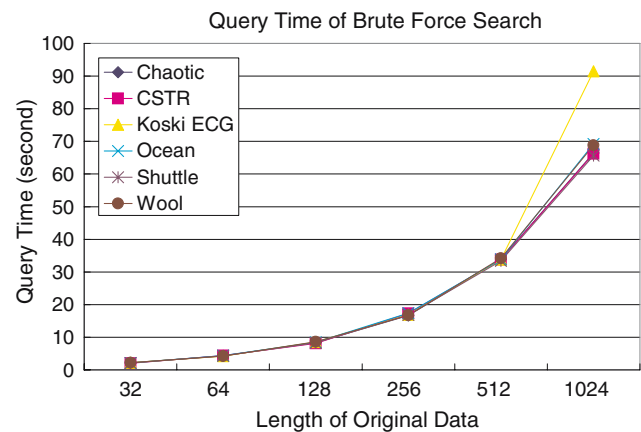
**Fig. 12** Average pruning power versus length of original data



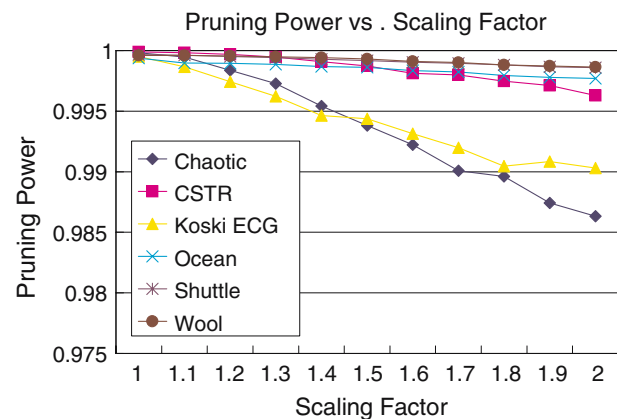
**Fig. 13** Some significant applications

time. More than 78% (32 out of 41) of the datasets obtained a pruning power above 90%. All but three of the datasets exhibited a pruning power of over 90% at length 1,024. Even at length 32, over 75% pruning power was achieved in 80% (33 out of 41) of the datasets. Figure 12 shows the pruning power averaged over all datasets; 97% of data sequences of length 1024 and 80% of data sequences of length 32 did not require computation of the actual time warping distances. Figure 11 may contain too much information so we pick six of the more significant applications to show the pruning power for them more clearly in Fig. 13. The applications include CSTR (speech), ECG, Ocean, Shuttle, Wool and Chaotic.

The promising pruning power will greatly reduce the querying time. We conducted experiments to measure the time required for query evaluation in all the 41 datasets. We compare the brute force approach to the pruning approach. In the timing we included both the time spent on the pruning and the post-processing where the SWM distances for remaining sequences are actually computed. Figure 14 shows the results. The time is consistently reduced, down to about 13% of the time



**Fig. 14** Query time comparison

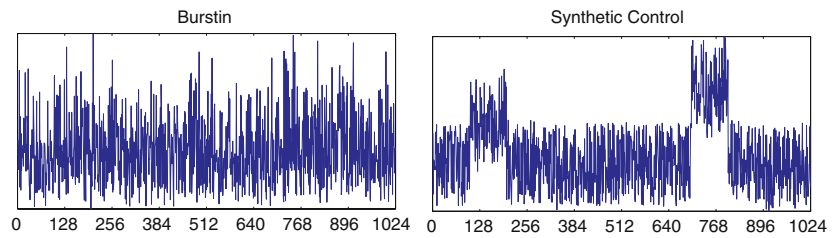
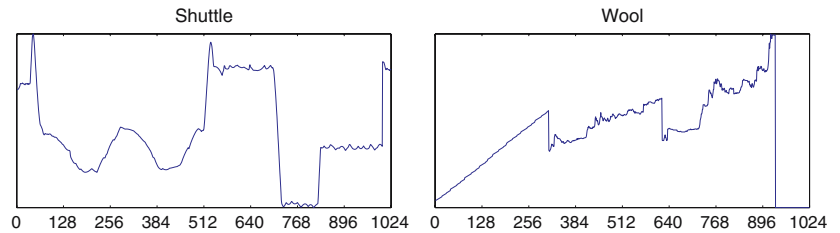


**Fig. 15** Varying the scaling factor

required by brute force search. We have repeated this with some other parameters and the results are similar.

Figure 15 shows the effect of varying the range of allowed scaling factors on pruning power. Note the  $x$ -axis indicates the upper bound range of allowed scaling factor. The lower bound range of allowed scaling factor is the reciprocal of the upper bound. For instance, the label 2 indicates that the range of allowed scaling factor is between  $1/2 = 0.5$  and 2. In particular, the label 1

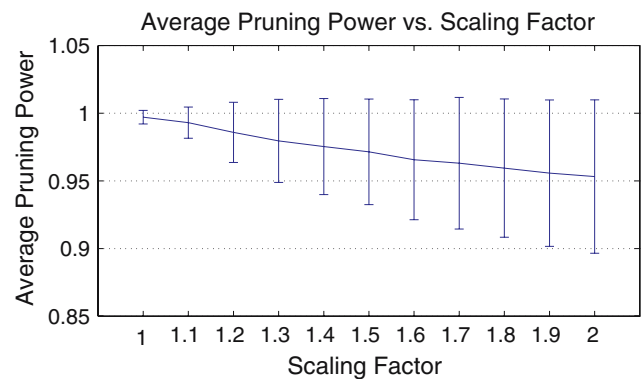


**Fig. 16** Data giving the lowest pruning power**Fig. 17** Data giving the highest pruning power

indicates that the time warping distance was calculated without scaling. It also implies that the size of the range was not increasing linearly. Although we show only six of the significant applications, we have experimented on all 41 sets of real data, the important observation is that for all sizes, a pruning power of over 90% was achieved in nearly 83% (34 out of 41) of the datasets. For all datasets (of length 1,024), the pruning powers never dropped below 80%. The downward trend of the pruning power as the scaling range (factor) increases can be explained by the increasing(decreasing) values of  $U_i(L_i)$ , since the value of  $U_i(L_i)$  is the maximum(minimum) of a greater range of  $C_j$ 's. This will decrease the value of  $LB(Q, C)$  and hence the pruning power.

A more detailed look into the actual data provided some insights as to why most datasets give very high pruning power, and why the few other datasets result in less pruning power. Figure 16 shows sample sequences from the two datasets that give the lowest pruning power, and Fig. 17 shows the sample sequences from the two datasets that give the highest pruning power. The difference between them is rather obvious visually. The sequences giving the lowest pruning power are those that fluctuate vigorously. The sequences giving the highest pruning power are those that are rather smooth. This is because with vigorous data fluctuation, the lower and upper bound envelope will be loose, and the pruning power will be weakened.

Nevertheless, we note that vigorously fluctuating datasets are far less common than smooth datasets. Figure 18 illustrates this claim by showing the pruning power averaged over all the datasets, as the range of allowed scaling factor changes. For all scaling factors, the average pruning powers are always above 95%. Even if we allow for one standard deviation margin below the average, the pruning power is still above 90% in general.

**Fig. 18** Average pruning power versus. scaling factor

In conclusion, the result shows that the proposed lower bounding measure effectively speeds up the query evaluation process. It also confirms the applicability of the lower bounding technique, even when a tight lower bound may not be readily obtainable.

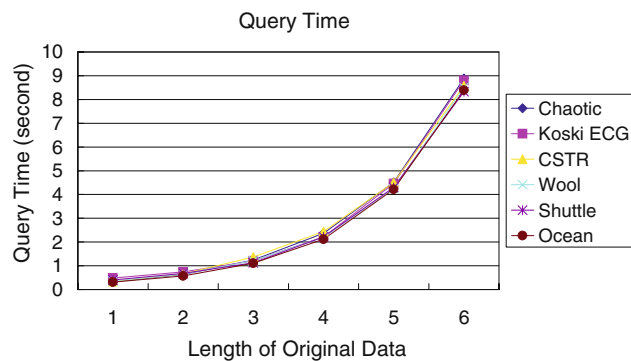
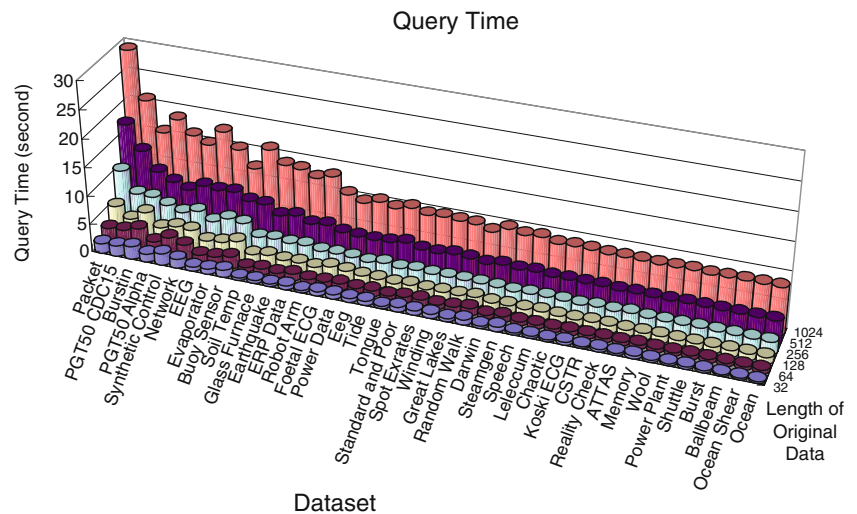
## 7.2 Enhancement by enveloping query

Next, we experiment with the enhancement described in Sect. 5 where the query subsequences are being enveloped instead of the data subsequences. We measure the query time and the pruning power.

### 7.2.1 Query time comparison

The query evaluation time required was measured. Figure 19 shows the results. Compared to the original approach, a general speed up can be seen across all datasets of any lengths. The proposed approach to envelope the query is *six times* faster in the extreme case and is about *twice* faster in the average case. Figure 20 shows six of the more significant applications more clearly.

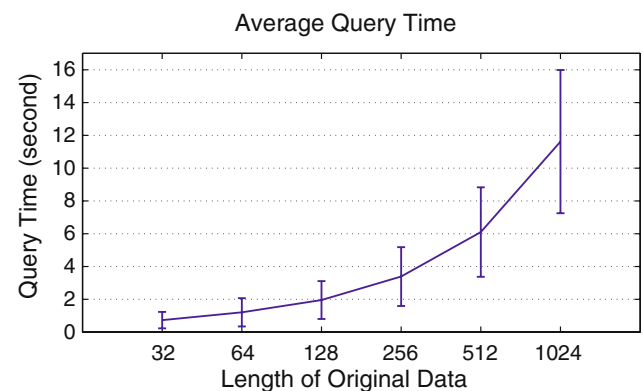


**Fig. 19** Query time**Fig. 20** Query time of selected applications

Only less than half a second was required for these six applications at length 32. Even at length 1,024, no more than 9 seconds was required. Figure 21 shows the average query evaluation time. For the longest data at length 1,024, about 20s was required by the original approach, but only about 12s was required by the proposed approach, confirming that the proposed approach speeded up query evaluation by about twice.

### 7.2.2 Effect on pruning power

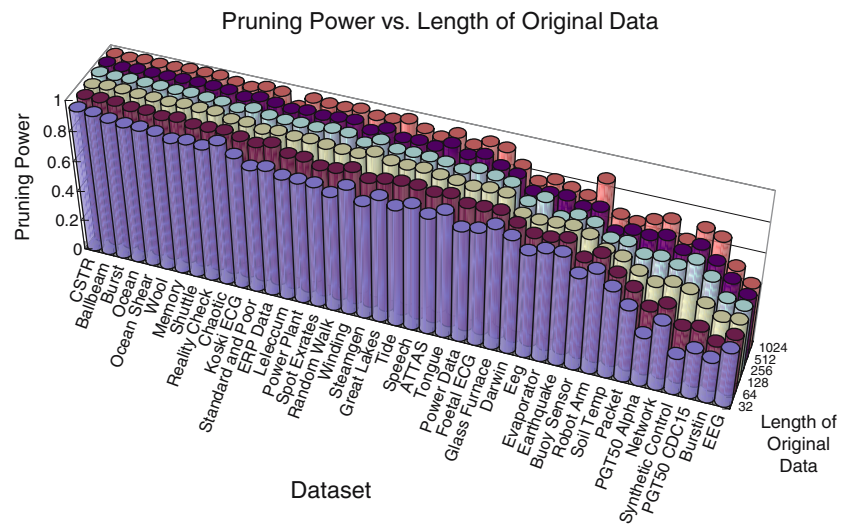
Figure 22 shows the pruning power of the lower bounding distance by enveloping query. Despite the anticipated reduction in pruning power, 26 of the 41 datasets keep their pruning power high above 90% at length 1,024. All but two of the datasets exhibited a pruning power of over 60%. Even at length 32, over 70% pruning power was achieved in three-fourths (30 out of 41) of the datasets. Figure 23 shows the pruning power for six of the more significant applications. The pruning power was above 86% at all lengths. And the pruning power at length 1,024 is even kept high above 96%. Figure 24

**Fig. 21** Average query time

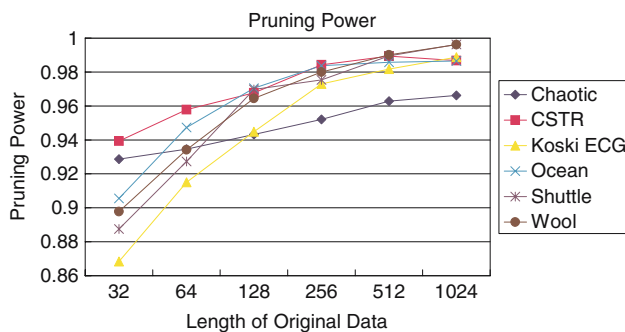
shows the pruning power averaged over all datasets. About 87% of data sequences of length 1,024 and 73% of data sequences of length 32 did not require computation of the actual time warping distances.

### 7.3 Performance of indexing

To investigate the effect of using our proposed index, indices were built on the datasets and various queries were issued. In order to handle the high dimensionality, an X-tree [3], a variant of the R-tree family, was used. Table 3 summarizes the parameters used in the experiments. Experiments were conducted to study the effect of varying the dimension of the feature space. Increasing the dimension would increase computation cost while decreasing it would decrease the pruning power. The dimension reported here represented a balance between the two factors. Figure 25 shows the size of the resulting index against different lengths of original data. It shows that the indices are compact at all lengths and fits into

**Fig. 22** Pruning power versus length of original data**Table 3** A list of X-tree parameters used in the experiments

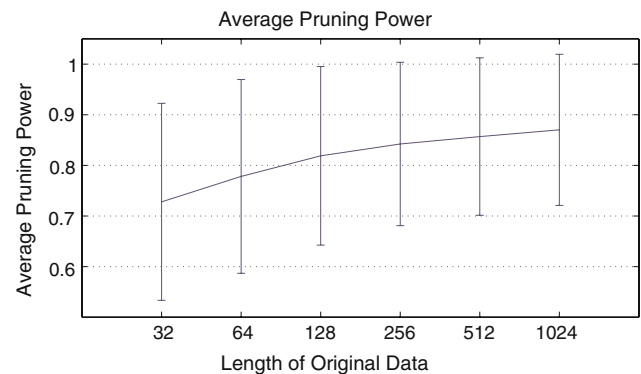
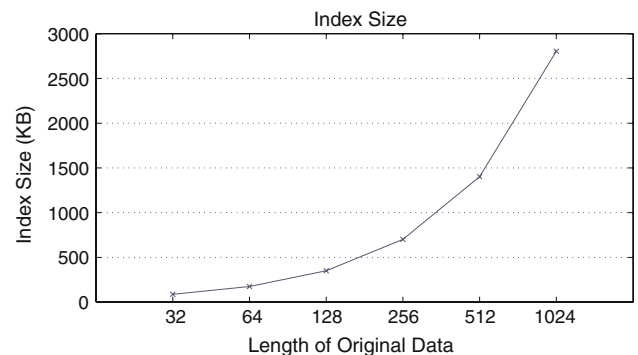
Name	Value
Dimension of feature space	4
Minimum number of entries per node	29
Maximum number of entries per node	59
Minimum query length	32

**Fig. 23** Some significant applications

memory nicely. Note that the index size merely increases linearly with the length of original data.

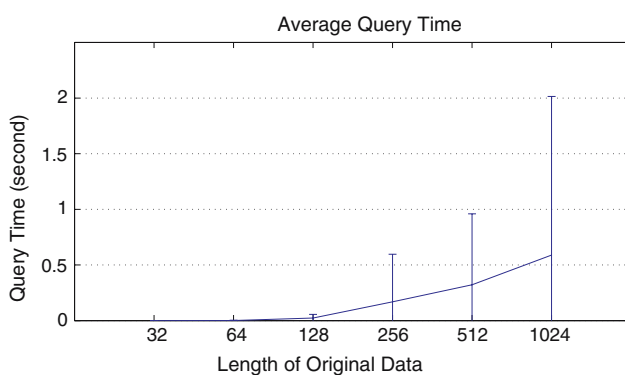
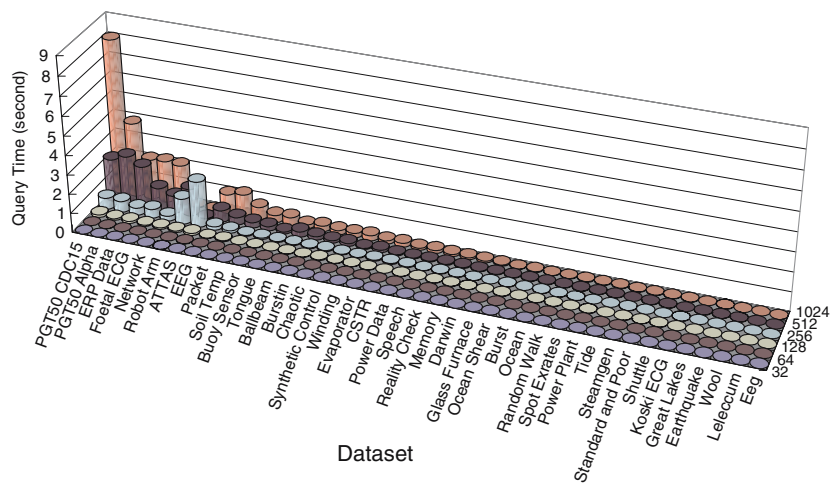
### 7.3.1 One nearest neighbor search

Figure 26 shows the query time using the proposed index. It shows that the index is very efficient for all kinds of datasets and all lengths of data we tested. In particular, the query time required by all datasets of length 128 or below was always less than a second. For length 256, the query time was 2.4 s in the worst case. It was only 2.5 s for length 512 and just above 8 s for length 1,024. Figure 27 shows the average query time, which

**Fig. 24** Average pruning power versus length of original data**Fig. 25** Index size

shows that the query time for most datasets is shorter than the worst case. At length 1,024, the average query time was just 0.59 s. That was only less than 0.4 s for all other lengths.

As a sanity check, we demonstrate that the index structure, rather than the dimension reduction, is the underlying cause of the speed up by repeating our near-

**Fig. 26** Query time using index**Fig. 27** Average query time using index

est neighbor search, but without employing any index. That is, we repeated the nearest neighbor search using linear scan and we pre-computed dimension reduction for the sequences when computing the lower bounding distance function.

Figure 28 compares the query time using linear search with pre-computed PAA with the query time using index. It shows that the index is effective in reducing the query time in almost all cases. Depending on the nature and dimension of the dataset, the query time using index can be faster by tens to thousands of times, representing a very huge speed up. On average, an order of magnitude speed up can be expected for all lengths of data, as illustrated in Table 4. The speed up was more notable at length 32, with an average speed up by 497 times. But the average speed up of 20 times at length 1,024 was also a big improvement. There was no particular trend in the level of speed up across different datasets. This may be because the order of the data sequences presented to the algorithm is different among the datasets (and is controlled by the index). If a match is found early, it is possible to prune away most of the data. And

**Table 4** Average speed up

Length of original data	Mean	Standard deviation
32	496.8026	483.6499
64	454.6806	538.4734
128	11.7928	14.8641
256	10.5620	15.0454
512	9.4354	13.9982
1,024	19.6865	47.6764

this effect is amplified when an index is used because a whole sub-tree of data in the index can be pruned.

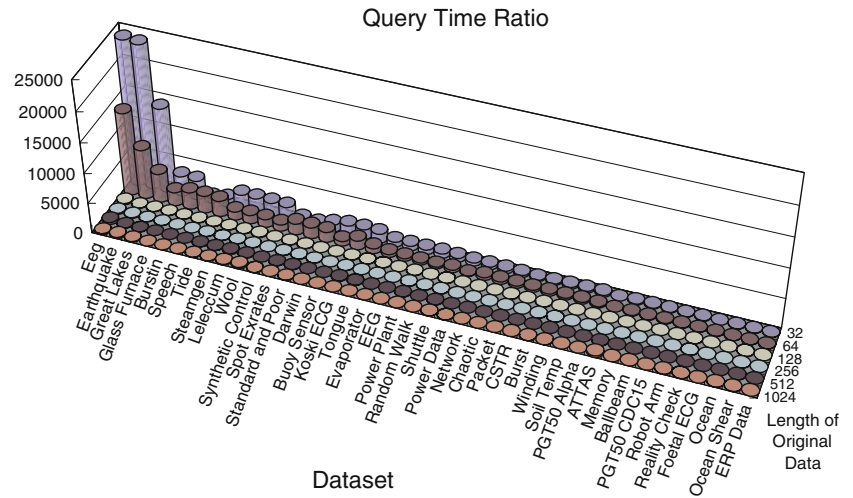
### 7.3.2 *k*-nearest neighbor search

In practice, it is more often that people are interested in more than one nearest neighbor. Thus, experiments have been conducted to evaluate how well the proposed method performs as the number of nearest neighbor  $k$  increases.

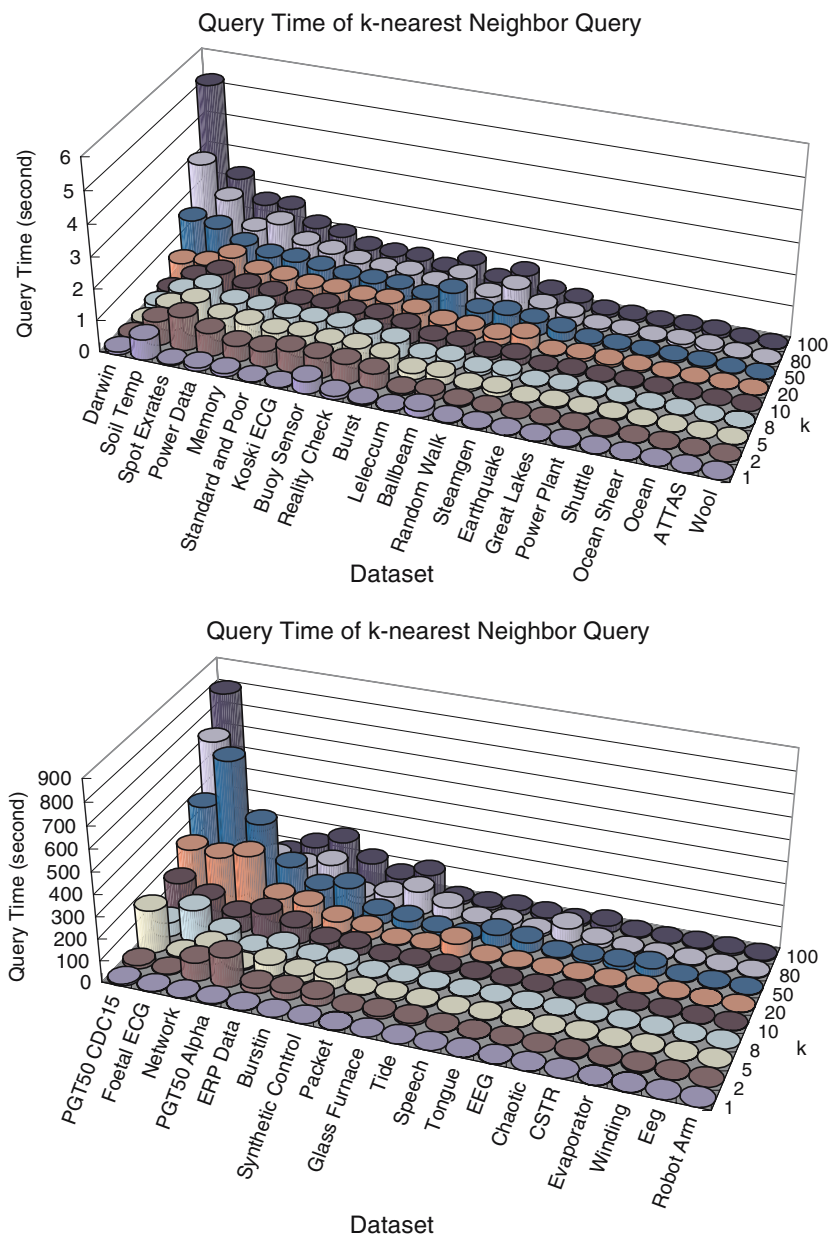
Figure 29 shows how the query time of  $k$ -nearest neighbor search using index varies as the number of nearest neighbor  $k$  increases. It shows that the query time varies substantially across different datasets and different number of nearest neighbor  $k$  required. However, the query time for more than half (27 out of 41) of the datasets was within 10 s when  $k = 100$ , while the query time for most (33 out of 41) of the datasets was within 10 s when  $k = 10$ , and within 6 s when  $k = 5$ .

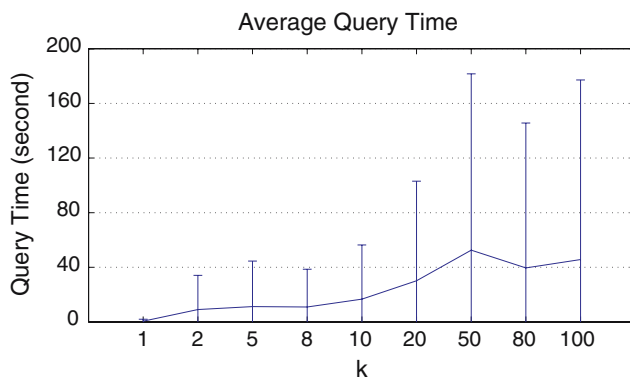
The numbers quoted above clearly demonstrated that the average query time, as plotted on Fig. 30, was a pessimistic estimate of the expected query time. Even so, it is worth noting that the average query time was about 40 s for queries up to  $k = 20$  and it was just above 40 s for queries up to  $k = 100$ . Also, the large standard deviation, as shown by the long error bars, also supported

**Fig. 28** Query time using linear search versus using index

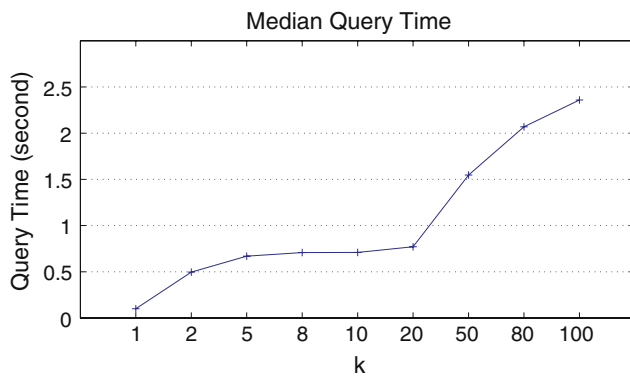


**Fig. 29** Query time of  $k$ -nearest neighbor search





**Fig. 30** Average query time of  $k$ -nearest neighbor search



**Fig. 31** Median query time of  $k$ -nearest neighbor search

the claim that the average query time was biased toward a few unusually long running queries.

The median query time may better reflect the trend of the query time as the number of nearest neighbor  $k$  increases and thus how the index scales on increasing  $k$ . Figure 31 shows that the median query time was just a few seconds, even when  $k = 100$ . And the median query time was below 1 s when  $k = 10$ . Hence, when the number of nearest neighbor  $k$  increased by ten times, the median query time only increased mildly by a few times.

## 8 Conclusion

In this article, we reviewed Euclidean distance, dynamic time warping (DTW), and uniform scaling (US). We motivated our work by showing that these time series similarity measures found in previous literature are inappropriate or insufficient for many applications.

In view of the weakness of previous distance measures, we proposed to combine DTW and US to solve these real world problems. Through this complementary merger, we amplified the power of both DTW and US to

give a better time series similarity measure, scaled and warped matching (SWM).

Although SWM is inherently expensive to compute, we showed that the lower bounding technique is applicable to SWM in drastically improving the query performance. In particular, experiments showed that as many as 97% of the SWM computation could be saved by using the proposed lower bounding function, reducing the query time of one nearest neighbor search to 13% of time required when the lower bounding function was not used.

Based on the first lower bounding function, we proposed an optimization to squeeze query time further. We showed that applying the proposed optimization consistently reduces the query time by extensive experiments.

Moreover, we employed an index to improve the performance of time series matching under the proposed SWM distance. This approach is especially useful when the size of data is large. We evaluated the usefulness and performance of the index in handling both one-nearest neighbor queries and  $k$ -nearest neighbor queries.

For future work, it is observed that the performance of time series similarity matching algorithms can be highly data-dependent. It is of practical values to further analyze the characteristics of data from different domains, and to investigate how these characteristics of data affect the performance of time series similarity matching algorithms.

**Acknowledgments** This research was supported by the RGC Research Direct Grant 03/04, and the RGC Earmarked Research Grant of HKSAR CUHK 4120/05E.

## References

1. Agrawal, R., Lin, K.I., Sawhney, H.S., Shim, K.: Fast similarity search in the presence of noise, scaling, and translation in time-series databases. In: VLDB'95, Proceedings of 21st International Conference on Very Large Data Bases, pp. 490–501. Morgan Kaufmann, Zurich (1995)
2. Beckmann, N., Kriegel, H.P., Schneider, R., Seeger, B.: The R\*-tree: an efficient and robust access method for points and rectangles. In: SIGMOD '90: Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data, pp. 322–331. ACM Press, Atlantic City (1990)
3. Berchtold, S., Keim, D.A., Kriegel, H.P.: The X-tree: An index structure for high-dimensional data. In: VLDB'96, Proceedings of 22nd International Conference on Very Large Data Bases, pp. 28–39. Morgan Kaufmann, Mumbai (Bombay) (1996)
4. Campbell, L., Bobick, A.: Recognition of human body motion using phase space constraints. In: Proceedings of International Conference on Computer Vision, pp. 624–630. IEEE Computer Society, Washington (1995)
5. Chai, W., Vercoe, B.: Folk music classification using hidden markov models. In: Proceedings of International Conference on Artificial Intelligence (2001)



6. Chan, F., Fu, A.: Efficient time series matching by wavelets. In: *Proceedings of the 15th International Conference on Data Engineering*. IEEE Computer Society, Sydney (1999)
7. Cheung, K.L., Fu, A.W.C.: Enhanced nearest neighbour search on the R-tree. *ACM SIGMOD Reco.* **27**(3), 16–21 (1998)
8. Dalal, N., Horaud, R.: Indexing key positions between multiple videos. In: *Proceedings of IEEE Workshop on Motion and Video Computing*, pp. 65–71. IEEE Computer Society, Orlando (2002)
9. Faloutsos, C., Ranganathan, M., Manolopoulos, Y.: Fast subsequence matching in time-series databases. In: *SIGMOD '94: Proceedings of the 1994 ACM SIGMOD International Conference on Management of Data*, pp. 419–429. ACM Press, Minneapolis (1994)
10. Fu, A., Chan, P., Cheung, Y.L., Moon, Y.: Dynamic vp-tree indexing for n-nearest neighbor search given pair-wise distances. *VLDB J.* **9**(2), 154–173 (2000)
11. Fu, A.W.C., Keogh, E., Lau, Y.H., Ratanamahatana, C.A.: Scaling and time warping in time series querying. In: *VLDB 2005, Proceedings of 31st International Conference on Very Large Data Bases*. Morgan Kaufmann, Trondheim (2005) (in press)
12. Guttman, A.: R-trees: a dynamic index structure for spatial searching. In: *SIGMOD '84: Proceedings of the 1984 ACM SIGMOD International Conference on Management of Data*, pp. 47–57. ACM Press, Boston (1984)
13. Hjaltason, G.R., Samet, H.: Distance browsing in spatial databases. *ACM Trans. Database Syst.* **24**(2), 265–318 (1999)
14. Itakura, F.: Minimum prediction residual principle applied to speech recognition. *IEEE Trans. Acoust. Speech Signal Process.* **23**(1), 67–72 (1975)
15. Kale, A., Chowdhury, R., Chellappa, R.: Towards a view invariant gait recognition algorithm. In: *Proceedings of the IEEE International Conference on Advanced Video and Signal based Surveillance (AVSS)*. IEEE Computer Society, Miami (2003)
16. Kale, A., Cuntoor, N., Yegnanarayana, B., Rajagopalan, A., Chellappa, R.: Gait analysis for human identification. In: *Proceedings of the 3rd International conference on Audio and Video Based Person Authentication* (2003)
17. Keogh, E.: Exact indexing of dynamic time warping. In: *VLDB 2002, Proceedings of 28th International Conference on Very Large Data Bases*, pp. 406–417. Morgan Kaufmann, Hong Kong (2002)
18. Keogh, E., Chakrabarti, K., Pazzani, M., Mehrotra, S.: Dimensionality reduction for fast similarity search in large time series databases. *Knowl. Inf. Syst.* **3**(3), 263–286 (2001)
19. Keogh, E., Foliass, T.: The UCR Time Series Data Mining Archive. University of California, Computer Science & Engineering Department, Riverside. (Available at <http://www.cs.ucr.edu/~eamonn/TSDMA/index.html> (2002))
20. Keogh, E., Palpanas, T., Zordan, V.B., Gunopulos, D., Cardle, M.: Indexing large human-motion databases. In: *VLDB 2004, Proceedings of 30th International Conference on Very Large Data Bases*, pp. 780–791. Morgan Kaufmann, Toronto (2004)
21. Keogh, E., Ratanamahatana, C.A.: Exact indexing of dynamic time warping. *Knowl. Inf. Syst.* **7**(3), 358–386 (2004)
22. Kosugi, N., Sakurai, Y., Morimoto, M.: SoundCompass: a practical query-by-humming system; normalization of scalable and shiftable time-series data and effective subsequence generation. In: *SIGMOD '04: Proceedings of the 2004 ACM SIGMOD International Conference on Management of Data*, pp. 881–886. ACM Press, New York (2004)
23. Matessi, G., Pilastro, A., Marin, G.: Variation in quantitative properties of song among European populations of the reed bunting (*Emberiza schoeniclus*) with respect to bill morphology. *Can. J. Zool.* **78**, 428–437 (2000)
24. Meek, C., Birmingham, W.: The dangers of parsimony in query-by-humming applications. In: *Proceedings of International Symposium on Music Information Retrieval* (2003)
25. Moeller-Levet, C., Klawonn, F., Cho, K.H., Wolkenhauer, O.: Fuzzy clustering of short time series and unevenly distributed sampling points. In: *Proceedings of IDA* (2003)
26. Moon, Y.S., Whang, K.Y., Loh, W.K.: Efficient time-series subsequence matching using duality in constructing windows. Technical Report 00-11-001. Advanced Information Technology Research Center (AITrc), KAIST, Taejeon, Korea (2000)
27. Moon, Y.S., Whang, K.Y., Loh, W.K.: Duality-based subsequence matching in time-series databases. In: *Proceedings of the 17th International Conference on Data Engineering*, pp. 263–272. IEEE, IEEE Computer Society, Heidelberg, Germany (2001)
28. Pullen K., Bregler C.: (2002) Motion capture assisted animation: Texturing and synthesis. *ACM Trans. Graph. (Proc SIGGRAPH 2002)* **22**(3)
29. Ratanamahatana, C.A., Keogh, E.: Everything you know about dynamic time warping is wrong. In: *Third Workshop on Mining Temporal and Sequential Data*, in conjunction with the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-2004), Seattle, WA (2004)
30. Ratanamahatana, C.A., Keogh, E.: Three myths about dynamic time warping data mining. In: *Proceedings of the 5th SIAM International Conference on Data Mining*, Newport Beach, CA (2005)
31. Roussopoulos, N., Kelley, S., Vincent, F.: Nearest neighbor queries. In: *SIGMOD '95: Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data*, pp. 71–79. ACM Press, New York (1995)
32. Sakoe, H., Chiba, S.: Dynamic programming algorithm optimization for spoken word recognition. *IEEE Trans. Acoust. Speech Signal Process.* **26**(1), 43–49 (1978)
33. Seidl, T., Kriegel, H.P.: Optimal multi-step k-nearest neighbor search. In: *SIGMOD '98: Proceedings of the 1998 ACM SIGMOD International Conference on Management of Data*, pp. 154–165. ACM Press, Seattle (1998)
34. Vlachos, M., Kollios, G., Gunopoulos, D.: An discovering similar multidimensional trajectories. In: *Proceedings of the 18th International Conference on Data Engineering*, pp. 673–684. IEEE, IEEE Computer Society, San Jose, CA (2002)
35. Wong, T.S.F., Wong, M.H.: Efficient subsequence matching for sequences databases under time warping. In: *Proceedings of the Seventh International Database Engineering and Applications Symposium*. IEEE, IEEE Computer Society, Hong Kong, SAR (2003)
36. Yi, B.K., Faloutsos, C.: Fast time sequence indexing for arbitrary  $L_p$  norms. In: *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases*, pp. 385–394. Morgan Kaufmann, Cairo (2000)
37. Yi, B.K., Jagadish, H.V., Faloutsos, C.: Efficient retrieval of similar time sequences under time warping. In: *Proceedings of the 14th International Conference on Data Engineering*. IEEE, IEEE Computer Society, Orlando, Florida (1998)
38. Zhu, Y., Shasha, D.: Warping indexes with envelope transforms for query by humming. In: *SIGMOD '03: Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data*, pp. 181–192. ACM Press, San Diego (2003)