

Finding Unusual Medical Time-Series Subsequences: Algorithms and Applications

Eamonn Keogh, Jessica Lin, Ada Waichee Fu, and Helga Van Herle

Abstract—In this work, we introduce the new problem of finding time series discords. Time series discords are subsequences of longer time series that are maximally different to all the rest of the time series subsequences. They thus capture the sense of the most unusual subsequence within a time series. While discords have many uses for data mining, they are particularly attractive as anomaly detectors because they only require one intuitive parameter (the length of the subsequence), unlike most anomaly detection algorithms that typically require many parameters. While the brute force algorithm to discover time series discords is quadratic in the length of the time series, we show a simple algorithm that is three to four orders of magnitude faster than brute force, while guaranteed to produce identical results. We evaluate our work with a comprehensive set of experiments on electrocardiograms and other medical datasets.

Index Terms—Anomaly detection, clustering, time-series data mining.

I. INTRODUCTION

THE previous decade has seen hundreds of papers on time series similarity search, which is the task of finding a time series that is most similar to a particular query sequence [1]. In this work, we pose the new problem of finding the sequence that is least similar to all other sequences. We call such sequences time series discords. Fig. 1 shows a time-series discord found in a human electrocardiogram. The fact that the discord in Fig. 1 coincides with the location annotated by a cardiologist as containing an anomalous heartbeat hints at one possible use of discords. As we shall show, time series discords are superlative anomaly detectors, able to detect subtle anomalies in diverse medical domains.

One reason why discords are particularly suited for the increasingly important problem of anomaly detection is that they only require consideration of a single intuitive parameter, the length of the subsequences. In contrast, many other anomaly detection algorithms require three to seven unintuitive parameters [2]. With so many parameters to set, we need access to huge amounts of training data; even then, avoiding overfitting remains a challenge.

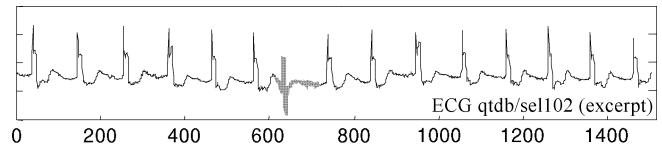


Fig. 1. The time-series discord found in an excerpt of electrocardiogram qtdb/sel102 (marked in bold line). The location of the discord exactly coincides with a premature ventricular contraction.

Time-series discords have other uses. Clustering algorithms can often benefit from removing a handful of “tricky cases,” which can be removed from the dataset before the clustering algorithm is run [3]. We could attempt to define these tricky cases as ones that do not belong to any cluster; however, this opens the possibility of a chicken and egg paradox. This effect has been noted in clustering points in k -dimensional space, but it is also true for time series, and the removal of discords offers a solution.

This paper makes two fundamental contributions in discovering unusual time series subsequences. First, while the idea of the “most unusual subsequence” is intuitive, great care must be taken in creating a workable definition; otherwise, we will be plagued with uninteresting pathological solutions. We introduce such a definition here and validate it in diverse medical domains. Second, the brute-force algorithm to discover the most unusual subsequence requires a quadratic “all to all” comparison, which is untenable for large real-world datasets. We introduce a simple algorithm that can achieve three to four orders of magnitude of speedup on real problems. Our algorithm works by admissibly pruning off some fruitless calculations, and using heuristics to reorder the search such that as many fruitless calculations are pruned as possible.

The rest of the paper is organized as follows. In Section II, we review related work and discuss some background material before introducing our formal definition of time series discords. In Section III, we consider the brute-force algorithm for finding discords, and introduce a general framework for speeding up the search based on admissible pruning and reordering the order in which the search examines the subsequences. Section IV introduces a particular reordering strategy based on examining a symbolic version of the data. We perform an extensive empirical evaluation in Section V to demonstrate both the utility of discords and our ability to find them quickly. Finally, Section VI offers some conclusions and suggestions for future work.

II. RELATED WORK AND BACKGROUND

Our review of related work is exceptionally brief because we are considering a new problem. Most real valued time series

Manuscript received July 20, 2005. This work was supported in part by the National Science Foundation under Grant BS123456.

E. Keogh is with the University of California, Riverside, CA 92521-0144 USA (e-mail: eamonn@cs.ucr.edu).

J. Lin was with University of California, Riverside, CA 92521-0144 USA. She is now with George Mason University, Fairfax, VA 22030 USA (e-mail: jessica@ise.gmu.edu).

A. Fu is with The Chinese University of Hong Kong, Hong Kong (e-mail: adafu@cse.cuhk.edu).

H. Van Herle is with the David Geffen School of Medicine, University of California, Los Angeles, CA 90095 USA (e-mail: hvanherle@mednet.ucla.edu).

Digital Object Identifier 10.1109/TITB.2005.863870

problems such as motif discovery [4], [5], [6], longest common subsequence matching, sequence averaging, segmentation, indexing [1], etc. have approximate or exact analogs in the discrete world, and have been addressed by the text processing or bioinformatics communities. However, time series discords do not appear to have a discrete version. Note that the superficially similar sounding Furthest (Sub)String Problem requires us to build a string, not to find one in the data [7]. As we shall see, one major use of discords is in anomaly detection. This topic has been an area of extensive research in recent years, and we refer the reader to [2], which gives a detailed survey.

A. Notation

For concreteness, we begin with a definition of our data type of interest, the time series.

Definition 1: Time Series: A time series $T = t_1, \dots, t_m$ is an ordered set of m real-valued variables. For data mining purposes, we are often not interested in any of the global properties of a time series [8], [5], [2], [9]; rather, we are interested in local subsections of the time series, which are called subsequences.

Definition 2: Subsequence: Given a time series T of length m , a subsequence C of T is a sampling of length $n \leq m$ of contiguous position from T , that is, $C = t_p, \dots, t_{p+n-1}$ for $1 \leq p \leq m-n+1$.

Since all subsequences may potentially be discords, any algorithm will eventually have to extract all of them; this can be achieved by use of a sliding window.

Definition 3: Sliding Window: Given a time series T of length m , and a user-defined subsequence length of n , all possible subsequences can be extracted by sliding a window of size n across T and considering each subsequence C_p .

Since our task is to find the most distant subsequence under some distance measure $\text{Dist}(C, M)$, we will define distance.

Definition 4: Distance: Dist is a function that has C and M as inputs and returns a nonnegative value R , which is said to be the distance from M to C . For subsequent definitions to work, we require that the function D be symmetric; that is, $\text{Dist}(C, M) = \text{Dist}(M, C)$. We also assume that the two subsequences are of equal length.

While the definition of a distance is obvious and intuitive, we need it to exclude trivial matches. In general, the best matches to a subsequence (apart from itself) tend to be located one or two points to the left or the right of the subsequence in question. Such matches have previously been called trivial matches [4], [5], [6]. As we shall see, it is critical when finding discords to exclude trivial matches; otherwise, almost all real datasets have degenerate and unintuitive solutions. We will therefore take the time to formally define a non-self match.

Definition 5: Non-Self Match: Given a time series T containing a subsequence C of length n beginning at position p and a matching subsequence M beginning at q , we say that M is a non-self match to C at distance of $\text{Dist}(M, C)$ if $|p-q| \geq n$.

We can most easily see the importance of non-self matches for the problem at hand if we consider the analogy of the problem in the discrete world. Consider the following string:

abcabcabcabcXXXabcabcabacabc.

The eye is immediately drawn to the subsequence of “X,” which surely forms the discord here. However, if we assume a sliding window length of three, and that our distance measure is the hamming distance, then the subsequence that is farthest from its nearest neighbor subsequence is “bac.” Below, the string is annotated by subscripts that give the distance to the nearest neighbor for each subsequence of length 3

$$a_0b_0c_0a_0b_0c_0a_0b_0c_0a_0b_1c_1X_1X_1X_1$$

$$a_0b_0c_0a_0b_0c_0a_1b_2a_1c_0a \ b \ c.$$

This unexpected and unintuitive result is caused by allowing trivial matches. While the subsequence XXX may appear unusual, it is only one unit distance from the subsequence XXa, which shares two elements simply shifted by one place. We can see the difference this makes by annotating the string with the non-self match distance to its nearest neighbor subsequence

$$a_0b_0c_0a_0b_0c_0a_0b_0c_0a_0b_1c_2X_3X_2X_1$$

$$a_0b_0c_0a_0b_0c_0a_1b_2a_1c_0a \ b \ c.$$

Here, the results are much more intuitive. While it is a simple and contrived example on discrete data, as we shall see, identical remarks apply to real world, real valued data. Note that the idea that one must exclude “partial self” comparisons in order to create meaningful definitions is well known in the bioinformatics community [10], and increasingly understood in the time series data mining community [4], [5], [6], [9], [11]. We will therefore use the definition of non-self matches to define time series discords:

Definition 6: Time Series Discord: Given a time series T , the subsequence D of length n beginning at position l is said to be the discord of T if D has the largest distance to its nearest non-self match. That is, \forall subsequence C of T , non-self match M_D of D , and non-self match M_C of C , $\min(\text{Dist}(D, M_D)) > \min(\text{Dist}(C, M_C))$.

We will denote the location of the discord as $D.l$ and the distance to the nearest non-self matching neighbor as $D.\text{dist}$. The length of the discord we denote as D_n .

We may be interested in examining the top K discords, which we define as:

Definition 7: K th Time Series Discord: Given a time series T , the subsequence D of length n beginning at position p is the K th-discord of T if D has the K th largest distance to its nearest non-self match, with no overlapping region to the i th discord beginning at position p_i , for all $1 \leq i < K$. That is, $|p-p_i| \geq n$.

We have deliberately omitted naming a distance function up to this point for generality. For concreteness, we will use the ubiquitous Euclidean distance measure throughout the rest of this paper [5], [1].

Definition 8: Euclidean Distance: Given two time series Q and C of length n , the Euclidean distance between them is defined as

$$\text{Dist}(Q, C) \equiv \sqrt{\sum_{i=1}^n (q_i - c_i)^2}.$$

Each time series subsequence is normalized to have mean zero and a standard deviation of one before calling the distance

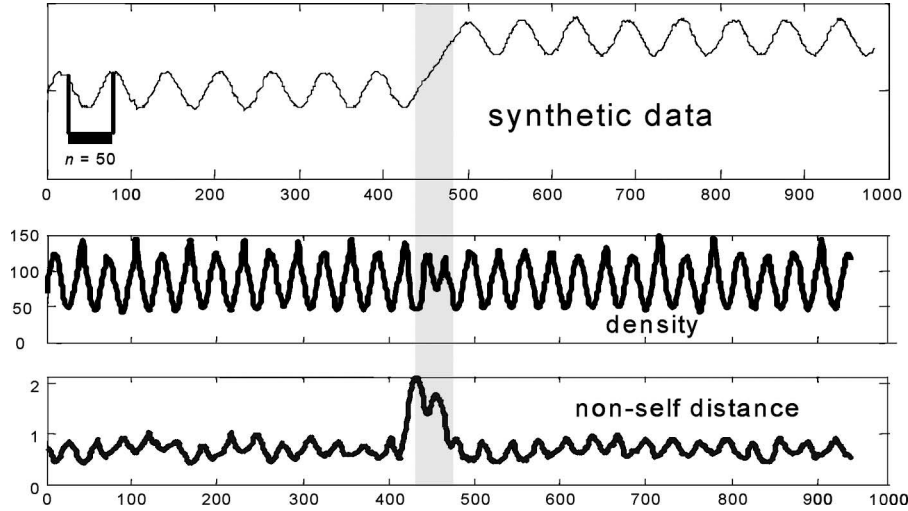


Fig. 2. (Top) A synthetic time series with an obvious anomaly. (Middle). The local density of subsequences of length 50, measured by calculating the number of matching subsequences within a range of 2. (Bottom) The non-self match to the nearest neighbor for all subsequences of length 50.

function, because it is well understood that in virtually all settings, it is meaningless to compare time series with different offsets and amplitudes [1].

B. Some Properties of Time Series Discords

Here, we discuss some properties of time series discords to enhance the readers' understanding of them and to discount some possible research directions for finding algorithms for quickly locating them.

1) *Discords are not Necessarily Found in Sparse Space*: The idea of considering time series subsequences as points in space has long been exploited by dozens of indexing techniques [1], so one might imagine that such a representation would be useful for the task at hand. We could simply project our time series into n -dimensional space and use existing outlier detection methods [12], [3]. The problem with this idea is the unintuitive fact that discords do not necessarily live in sparse areas of n -dimensional space (conversely, repeated patterns do not necessarily live in dense parts of the n -dimensional space [4], [5], [6]). The full explanation has consequences for other problems and is perhaps deserving of a separate paper; however, here, we content ourselves with a visual example and a brief explanation. In Fig. 2, we consider a simple time series consisting of a slightly noisy sine wave. We introduce an “anomaly” of length 50 by shifting the entire second half of the time series.

We can now extract all subsequences of length 50, project them into 50-dimensional space, and measure the local density around each subsequence. Surprisingly, the anomaly is not in the sparsest (or in any other way remarkable) region of space. However, note that the definition of non-self match that is at the heart of time series discords clearly identifies the anomalous region.

The explanation of this unintuitive finding harkens back to the idea of trivial matches. Consider a subsequence C located at t_p that is “simple;” that is to say it has only one or two features such as peaks or valleys. This simple subsequence is very close in n -dimensional space to the subsequences beginning at $t_{p+1}, t_{p-1}, t_{p+2}$, etc. In contrast, consider a subsequence M

located at t_q that is “complex;” that is to say it has many features such as peaks or valleys. This complex subsequence is relatively far from subsequences beginning at $t_{q+1}, t_{q-1}, t_{q+2}$, etc. In other words, simple (and smooth) shapes appear to be in dense neighborhoods because we overcount shifted versions of them. This problem prevents us from using existing density based algorithms to find time series discords. Note that even if current density based algorithms could be adapted to consider non-self distance, most of them degrade to quadratic time complexity for high dimensionality data.

C. Discords Results are Noncombinable

Several generic paradigms for solving problems rely on the ability to decompose a problem into smaller subproblems, which can be solved and admissibly recombined. Depending on the exact definitions, such techniques are variously called dynamic programming, divide and conquer, bottom-up, etc [13]. Unfortunately, as we will show, such ideas are unlikely to help us efficiently find discords.

Imagine that we break a time series T into two sections, A and B , and that we find the discords for both sections, recording their locations as $A.l, B.l$ and values as $A.dist$ and $B.dist$, respectively. Furthermore, imagine that we now concatenate A and B to reproduce the original time series T (for simplicity, let us assume that when the discord for T is discovered; it will not span the end of A and the beginning of B). What can we now say about the discord for T ? Surprisingly, the answer is very little. We cannot assume that it will be either in location $A.l$ or in location $|A| + B.l$, because both of the two previously discovered discords may have good matches in the other section. All we can do is give weak bounds. The value of $T.dist$ is at most $\max(A.dist, B.dist)$. The lower bound of $T.dist$ is a trivial zero (To see this, imagine $A = B$). As to the location of $T.l$, we can say nothing.

If we consider the complementary situation, where we know the discord information $T.l$ and $T.dist$ for T , and we split

TABLE I
BRUTE FORCE DISCORD DISCOVERY

1	Function [dist, loc] = Brute_Force(T, n)
2	best_so_far_dist = 0
3	best_so_far_loc = NaN
4	
5	For $p = 1$ to $ T - n + 1$ // Begin Outer Loop
6	nearest_neighbor_dist = infinity
7	For $q = 1$ to $ T - n + 1$ // Begin Inner Loop
8	IF $ p - q \geq n$ // non-self match?
9	IF $Dist([t_p, \dots, t_{p+n-1}], [t_q, \dots, t_{q+n-1}]) < \text{nearest_neighbor_dist}$
10	nearest_neighbor_dist = $Dist(t_p, \dots, t_{p+n-1}, t_q, \dots, t_{q+n-1})$
11	End
12	End // End non-self match test
13	End // End Inner Loop
14	IF nearest_neighbor_dist > best_so_far_dist
15	best_so_far_dist = nearest_neighbor_dist
16	best_so_far_loc = p
17	End
18	End // End Outer Loop
19	Return[best_so_far_dist, best_so_far_loc]

into two new time series A and B , we are similarly frustrated. Assume that the discord from T happened to fall into A . We can lower bound $A.dist$ as being greater than or equal to $T.dist$, but we cannot provide an upper bound. In addition, we can say nothing about the location of $A.l$. As for $B.dist$ and $B.l$, we can say nothing as well.

A combination of these two results also frustrates any thought of exploiting a sliding window algorithm, since ingesting and egressing a single point can change the location and value of the discords.

The preceding results suggest that existing algorithms/paradigms are of little utility for finding discords. This motivates the introduction of an original algorithm in the next section.

III. FINDING TIME SERIES DISCORDS

The brute force algorithm for finding discords is simple and obvious. We simply take each possible subsequence and find the distance to the nearest non-self match. The subsequence that has the greatest such value is the discord. This is achieved with nested loops, where the outer loop considers each possible candidate subsequence, and the inner loop is a linear scan to identify the candidate's nearest non-self match. For clarity, the pseudocode is shown in Table I.

Note that the algorithm requires exactly one parameter—the length of subsequences to consider. The algorithm is easy to implement and produces exact results. However, it has one fatal flaw for data mining. It has $O(m^2)$ time complexity, which is simply untenable for even moderately large datasets.

The following two observations offer hope to improve the algorithm's running time.

Observation 1: In the inner loop, we don't actually need to find the true nearest neighbor to the current candidate. As soon as we find any subsequence that is closer to the current candidate than the best_so_far_dist, we can abandon that instance of the inner loop, safe in the knowledge that the current candidate could not be the time series discord.

Observation 2: The utility of the preceding optimization depends on the order which the outer loop considers the candidates for the discord, and the order which the inner loop visits the other

TABLE II
HEURISTIC DISCORD DISCOVERY

1	Function [dist, loc] = Heuristic_Search($T, n, \text{Outer}, \text{Inner}$)
2	best_so_far_dist = 0
3	best_so_far_loc = NaN
4	
5	For Each p in T ordered by heuristic <i>Outer</i> // Begin Outer Loop
6	nearest_neighbor_dist = infinity
7	For Each q in T ordered by heuristic <i>Inner</i> // Begin Inner Loop
8	IF $ p - q \geq n$ // non-self match?
9	IF $Dist([t_p, \dots, t_{p+n-1}], [t_q, \dots, t_{q+n-1}]) < \text{best_so_far_dist}$
10	Break // Break out of Inner Loop
11	End
12	IF $Dist([t_p, \dots, t_{p+n-1}], [t_q, \dots, t_{q+n-1}]) < \text{nearest_neighbor_dist}$
13	nearest_neighbor_dist = $Dist(t_p, \dots, t_{p+n-1}, t_q, \dots, t_{q+n-1})$
14	End
15	End // End non-self match test
16	End // End Inner Loop
17	IF nearest_neighbor_dist > best_so_far_dist
18	best_so_far_dist = nearest_neighbor_dist
19	best_so_far_loc = p
20	End
21	End // End Outer Loop
22	Return[best_so_far_dist, best_so_far_loc]

subsequences in its attempt to find a sequence that will allow an early abandon of the inner loop.

While these are simple ideas and only minor modifications of the original algorithm, for concreteness, we will make them clear. The pseudocode is shown in Table II.

Note that the input has been augmented by two heuristics—one to determine the order in which the outer loop visits the subsequences, and one to determine the order in which the inner loop visits the subsequences. Note that the heuristic for the outer loop is used once, but the heuristic for the inner loop takes the current candidate into account, and is thus invoked to produce a new ordering for every iteration of the outer loop.

We have now reduced the discord discovery problem into a generic framework where all one needs to do is to specify the heuristics. Note that we should not attempt to “cheat” the algorithm. We could provide very good heuristic orderings if we are allowed to completely solve the brute force problem each time the heuristic functions are invoked! However, this is simply hiding the time complexity in a different part of the implementation. We must therefore insist that the Outer heuristic (invoked only once) takes at most $O(m)$ to calculate, and the Inner r heuristic (invoked $m - n$ times) takes $O(1)$. Note that this requirement precludes the possibility of using R-trees, K-d trees or other classic indexing algorithms [1], [14].

It is very important to recognize that while we are using heuristics to speed up the search for discords, the results of the algorithm are exact, and completely independent of heuristics used. The heuristics change only the speed of the algorithm.

To gain some intuition into our new algorithm, and to hint at our eventual solution to this problem, let us consider three possible heuristic strategies.

1) Random: We could simply have both the Outer and Inner heuristics randomly order the subsequences to consider. It is difficult to analyze this strategy since its performance is bounded from below by $O(m)$ and from above by $O(m^2)$ (see below for explanation), and depends on the data. However, empirically it works reasonably well. The conditional test on line none of Table II is often true and the inner loop can be abandoned early, considerably speeding up the algorithm.

2) *Magic*: In this hypothetical situation, we imagine that a friendly oracle gives us the best possible orderings. These are as follows: For Outer, the subsequences are sorted by descending order of the non-self distance to their nearest neighbor, so that the true discord is the first object examined. For Inner, the subsequences are sorted in ascending order of distance to the current candidate. For the Magic heuristics, the first invocation of the inner loop will run to completion. Thereafter, all subsequent invocations of the inner loop will be abandoned during the very first iteration. The time complexity is thus one occurrence of $m - n + 1$ steps for the first inner loop, and $m - n$ occurrences of the $O(1)$ step of each subsequent invocation of the inner loop, giving a total time complexity of $O(m) + O(m)$, or just $O(m)$. Note that we have $m \gg n$.

3) *Perverse*: In this hypothetical situation, we imagine that a less than friendly oracle gives us the worst possible orderings. These are identical to the Magic orderings with ascending/descending orderings reversed. In this case, we are back to the original $O(m^2)$ time complexity, and we waste some time in the conditional tests on line nine of Table II.

These results are something of a mixed bag for us. They suggest that a linear time algorithm is possible, but only with the aid of some very wishful thinking. The Magic heuristic requires a perfect ordering of subsequences in the inner loop, and any perfect ordering (i.e., sorting) requires at least $O(m \log m)$. Furthermore, the only known way to produce the perfect ordering of subsequences in the outer loop requires $O(m^2)$ work, but we are only allowed $O(m)$ time. The following two observations, however, offer us some hope for a fast algorithm.

Observation 3: In the outer loop, we do not actually need to achieve a perfect ordering to achieve dramatic speedup. All we really require is that among the first few subsequences being examined, we have at least one that has a large distance to its nearest neighbor. This will give the `best_so_far_dist` variable a large value early on, which will make the conditional test on line 9 of Table II be true more often, thus allowing more early terminations of the inner loop.

Observation 4: In the inner loop, we also do not actually need to achieve a perfect ordering to achieve dramatic speedup. All we really require is that among the first few subsequences being examined we have at least one that has a distance to the candidate sequence being considered that is less than the current value of the `best_so_far_dist` variable. This is a sufficient condition to allow early termination of the inner loop.

We can imagine a full spectrum of algorithms, which only differ by how well they order subsequences relative to the Magic ordering. This spectrum spans {Perverse . . . Random . . . Magic}. Our goal, then, is to find the best possible approximation to the Magic ordering, which is the topic of Section IV.

At the risk of redundancy, we again emphasize that this search problem requires a specialized solution, and we cannot leverage off the huge literature on time series similarity search [8]. Kd-Trees, R-trees, and their many variants require $O(\log(m))$ time per lookup, but we can spare only $O(1)$ time. In any case, these search algorithms support *nearest* neighbor search, whereas all we require here is “*near-enough*” neighbor search, as noted in Observation 4.

TABLE III
A LOOKUP TABLE THAT CONTAINS THE BREAKPOINTS THAT
DIVIDE A GAUSSIAN DISTRIBUTION IN AN ARBITRARY NUMBER
(FROM THREE TO FIVE) OF EQUIPROBABLE REGIONS

β_i a	3	4	5
β_1	-0.43	-0.67	-0.84
β_2	0.43	0	-0.25
β_3		0.67	0.25
β_4			0.84

IV. APPROXIMATIONS TO MAGIC

Before we introduce our techniques for approximating the perfect ordering returned by the hypothetical Magic heuristics, we must briefly review the symbolic aggregate approximation (SAX) representation of time series introduced in [15]. While there are at least 200 different symbolic approximations of time series in the literature, SAX is unique in that it is the only one that allows both dimensionality reduction and lower bounding of L_p norms. Since its relatively recent introduction, SAX has become an important tool in the time series data mining toolbox. It has been used to find time series motifs [5], [11], to mine rules in health data [4], for anomaly detection [2], to extract features from a hepatitis database [6], for visualization [16], [17], and a host of other data mining tasks.

A. A Brief Review of SAX

A time series C of length n can be represented in a w -dimensional space by a vector $\bar{C} = \bar{c}_1, \dots, \bar{c}_w$. The i th element of \bar{C} is calculated by

$$\bar{c}_i = \frac{w}{n} \sum_{j=\frac{n}{w}(i-1)+1}^{\frac{n}{w}i} c_j$$

In other words, to transform the time series from n dimensions to w dimensions, the data is divided into w equal sized “frames.” The mean value of the data falling within a frame is calculated and a vector of these values becomes the dimensionality-reduced representation. This simple representation is known as piecewise aggregate approximation (PAA) [15].

Having transformed a time series into the PAA representation, we can apply a further transformation to obtain a discrete representation. It is desirable to have a discretization technique that will produce symbols with equiprobability [5], [2]. In empirical tests on more than 50 datasets, we noted that normalized subsequences have highly Gaussian distribution [15], so we can simply determine the “breakpoints” that will produce equal-sized areas under Gaussian curve.

Definition 9: Breakpoints: Breakpoints are a sorted list of numbers $B = \beta_1, \dots, \beta_{a-1}$ such that the area under a $N(0, 1)$ Gaussian curve from β_i to $\beta_{i+1} = 1/a$ (β_0 and β_a are defined as $-\infty$ and ∞ , respectively).

These breakpoints may be determined by looking them up in a statistical table. For example, Table III gives the breakpoints for values of a from three to five.

We can now state our Outer heuristic; we scan the rightmost column of the array to find the smallest count mincount (its value is virtually always 1). The indexes of all SAX words that occur mincount times are recorded, and are given to the outer loop to search over first. After the outer loop has exhausted this set of candidates, the rest of the candidates are visited in random order.

The intuition behind our Outer heuristic is simple. Unusual subsequences are very likely to map to unique or rare SAX words. By considering the candidate sequences that map to unique or rare SAX words early in the outer loop, we have an excellent chance of giving a large value to the `best_so_far_dist` variable early on, which (as noted in observation 3) will make the conditional test on line nine of Table II be true more often, thus allowing more early terminations of the inner loop.

C. Approximating the Magic Inner Loop

Our Inner heuristic also leverages off the two data structures shown in Fig. 4. When candidate i is first considered in the outer loop, we look up the SAX word that it maps to, by examining the i th word in the array. We then visit the trie and order the first items in the inner loop in the order of the elements in the linked list index found at the terminal nodes. For example, imagine we are working on the problem shown in Fig. 4. If we were examining the candidate C_{731} in the outer loop, we would visit the array at location 731. Here we would find the SAX word `caa`. We could use the SAX values to traverse the trie to discover that subsequences 1, 3, 731 map here. These three subsequences are visited first in the inner loop (note that line eight of Table I prevents 731 from being compared to itself). After this step, the rest of the subsequences are visited in random order.

The intuition behind our Inner heuristic is also simple. Subsequences that have the same SAX encoding as the candidate subsequence are very likely to be highly similar (this fact is at the heart of more than 20 research efforts [4], [5], [2], [6], [17], [9]). As noted in observation 4, we just need to find one such subsequence that is similar enough (has a distance to the candidate than the current value of the `best_so_far_dist` variable) in order to terminate the inner loop.

1) *Minor Optimizations and Parameter Setting:* There are several minor optimizations we can apply to the heuristic search algorithm. For example, imagine we are considering candidate C_i in the outer loop, and as we traverse through the inner loop, we find that subsequence C_j is close enough to it to allow early abandonment. In addition to saving time with the early termination, we can also delete C_j from the list of candidates in outer loop (if it has not already been visited). The key observation is that since we are assuming a symmetric distance measure, if nearness to C_j disqualifies candidate C_i from being the discord, then the same nearness to C_i would also disqualify candidate C_j from being the discord. Empirically, this simple optimization gives a speed-up factor of approximately two. In addition, there are several well-known optimizations to the Euclidean distance that we can use [1].

As noted previously, we must choose two parameters, the cardinality of the SAX alphabet size a , and the SAX word size w . Recall what it is we want to optimize. We would like the distribution of the SAX words to be highly skewed, so that the discord will map to a SAX word that is unique or rare, and all the other subsequences will map to SAX words that are very frequent. This is the best situation for both our heuristics. If we choose very large values of a and/or w , almost all subsequences will map to unique words; if we choose very small values of a and/or w , all subsequences will map to just a small handful of words. Either of these situations is bad for our heuristics.

The good news is that there is little freedom for the a parameter; extensive experiments carried out by the current authors [5], [2], [16], [15], [17] and dozens of other researchers worldwide [4], [6], [9], [11] suggest that a value of either three or four is best for virtually any task on any dataset. After empirically confirming this on the current problem with experiments on more than 50 datasets, we will simply hardcode $a = 3$ for the rest of this work. Having fixed a , we performed an exhaustive empirical examination of the role of the w parameter. The best value for this parameter depends on the data. In general, relatively smooth and slowly changing datasets favor a smaller value of w , whereas more complex time series favor a larger value of w . The following observations mitigate the problem of parameter setting.

The speedup does not critically depend on w parameter. After empirically finding the best value on a particular data we found we could vary the value of w in the range of 60% to 150% with less than a 12% decrease in speedup.

Once we learn a good setting on a particular data type, say ECGs, that setting will also work well on other datasets of the same type (assuming the sampling rate is the same).

V. EMPIRICAL EVALUATION

We begin by showing the utility of time series discords for a several medical domains, then go on to show that our algorithm is able to find discords very efficiently.

A. The Utility of Time Series Discords

In this paper, we will only demonstrate the utility of discords as anomaly detectors. We have done extensive successful experiments in other tasks, such as improving the quality of clustering and summarization; however, anomaly detection is unique in that it allows immediate and intuitive visual confirmation. The additional experiments for other tasks, together with many extra anomaly detection experiments can be found here [19]. We encourage the interested reader to consult this site for additional examples and larger and more detailed figures of the experiments shown below.

After much reflection, we have decided not to include comparisons to other approaches here. There are two reasons for this. Firstly, it is very difficult to make meaningful comparisons between our method, which requires only one intuitive parameter, and some of the rival methods that require three to seven parameters (see [2] for a detailed discussion of this), including

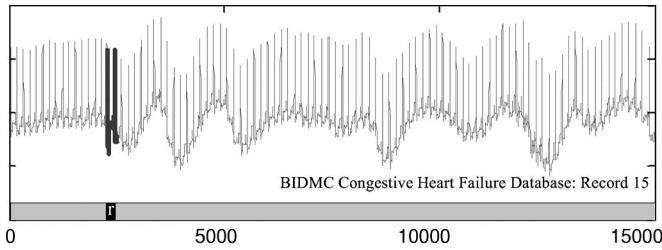


Fig. 5. An ECG that has been annotated by a cardiologist (bottom bar) as containing one premature ventricular contraction. The discord_{256} (bold line) exactly coincides with the heart anomaly.

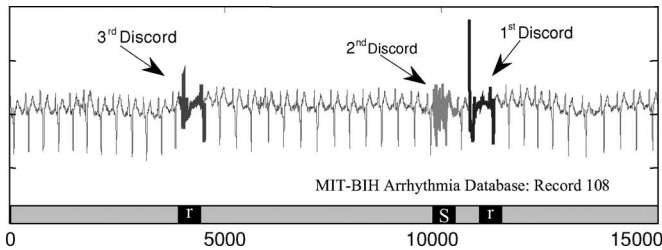


Fig. 6. An excerpt of an ECG that has been annotated by a cardiologist (bottom bar) as containing three various anomalies. The first three discord_{600} (bold lines) exactly coincides with the anomalies.

some parameters for which we may have poor intuition, such as Embedding dimension [20], Kernel function [21], SOM topology, or the number of Parzen windows.

The second reason we do not compare to other anomaly detectors is that most algorithms require a separate training dataset (in order to learn the parameters), whereas our approach finds anomalies while only examining the test dataset. One could easily imagine generalizing the discord discovery algorithm to examine only the test data in the outer loop and only training data in the inner loop. However, we wish to concentrate on first proving our simple intuitive definitions before creating generalizations.

1) *Anomaly Detection in Electrocardiograms:* Electrocardiograms (ECGs) are a time series of the electrical potential between two points on the surface of the body caused by a beating heart. They are arguably the most important time series, and as such, there are many annotated test datasets we can consider. We have already considered the utility of discords in one ECG in Fig. 1. That was a very simple and “clean” example for clarity; however, it is remarkable how varied and complex normal healthy ECGs can be. For example, Fig. 5 shows a very complicated signal with remarkable variability. Surprisingly, this ECG contains only one small anomaly, which is easily discovered by a discord detection algorithm.

In Fig. 6, we consider an ECG that has several different types of anomalies. Here, the first three discords exactly line up with the cardiologist’s annotations. In this figure, we could perhaps spot the anomalies by eye; however, the full time series is much longer, and impossible to scrutinize without a scrollbar and much patience.

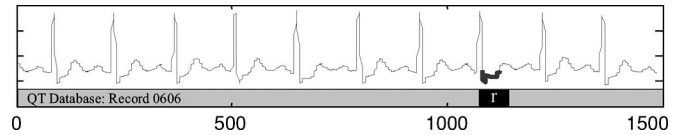


Fig. 7. An ECG that has been annotated by a cardiologist (bottom bar) as containing one premature ventricular contraction. The discord_{40} (bold line) exactly coincides with the heart anomaly.

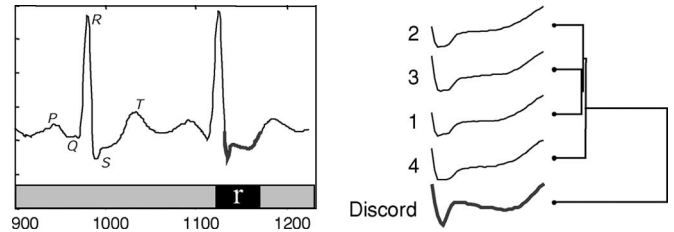


Fig. 8. (Left) A zoom-in of a section of Fig. 7. The first heartbeat has been annotated with the classic PQRS notation. (Right) Five ST waves from Fig. 7 (including the discord) hierarchically clustered.

In the above cases, we simply set the length of the discords to be approximately one full heartbeat (note that the two datasets have different sampling rates). Although we found that we could double or half the parameters without affecting the quality of results, on just a handful of the dozens of ECG datasets we examined, the discords had a harder time finding the anomalous heartbeats. One of the authors of the current work, Helga Van Herle M.D is a cardiologist. She informed us that heart irregularities can sometimes manifest themselves at scales significantly shorter than a single heartbeat. Armed with this knowledge, we searched for discords at approximately $1/4$ the length of a single heartbeat. In Fig. 7, we show the results of a search with the shorter length discords.

While the result is satisfying in that it immediately locates the anomaly, it is not obvious from the figure that the discord is actually different for the other heartbeats. In Fig. 8 (left), we see a zoom-in of the subsequence surrounding the discord, and we can see that the discord falls over the ST wave. In Fig. 8 (right), we manually extracted 4 ST waves from the subsequence in Fig. 7 and clustered them together with the discord. This makes the source of the anomaly apparent. Note that in the four normal ST waves, after the brief descending section, the signal rises monotonically. However, the anomalous ST wave has an additional local peak caused by a premature beat, thus justifying the cardiologist’s diagnosis of premature ventricular contraction.

2) *Change Detection in Patient Monitoring:* The problem of change detection is fundamentally different from anomaly detection. In anomaly detection, the task is to find one or more “different” subsequences that exist in the background of a normal data. In the problem of change detection, we assume that the underlying model that produces the signal changes in some (possibly very subtle) way at various points. The task is to identify the locations of these change points.

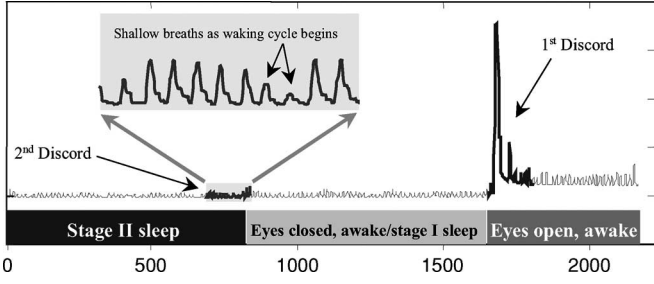


Fig. 9. The first two discords found in a time series of a patient's respiration as they wake up. The annotations show in the boxes at the bottom of the screen are provided by a medical expert.

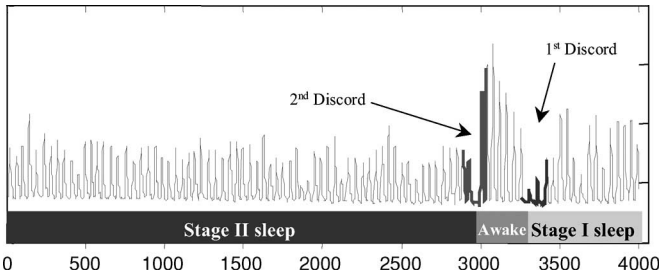


Fig. 10. The first two discords found in a time series of a patient's respiration.

Time series discords do not appear to be likely candidates for change detection, since they look at local patterns, whereas most change detection algorithms consider global (or at least much larger “local”) information. However, we believe that in some situations, the change in underlying global model may produce some unusual local shapes because the local pattern must straddle two different models.

To test this idea, we investigated a time series showing a patient's respiration (measured by thorax extension), as they wake up. A medical expert, Dr. J. Rittweger of the Institute for Physiology, Free University of Berlin, manually segmented the data. We choose a discord length corresponding to 32 s because we want to span several breaths. In Fig. 9, we see the outcome of the first experiment.

The 1st discord is a very obvious deep breath taken as the patient opened their eyes. In contrast, the 2nd discord is much more subtle and difficult to see at this scale. A zoom-in suggests that Dr. Rittweger noticed a few shallow breaths that indicated the transition of sleeping states. In both cases, the discords straddle the change in sleeping cycle. We tested many such datasets with equally positive results. Fig. 10 shows another representative example.

B. The Utility of Heuristic Ordered Search

It is increasingly recognized that comparing algorithm performance by examining wall clock or CPU time invites the possibility of implementation bias [1], which in turn invites the possibility of irreproducible “improvements.” Instead, we measure here the number of times that the distance function is called on line nine in Tables I and II. A simple analysis of the

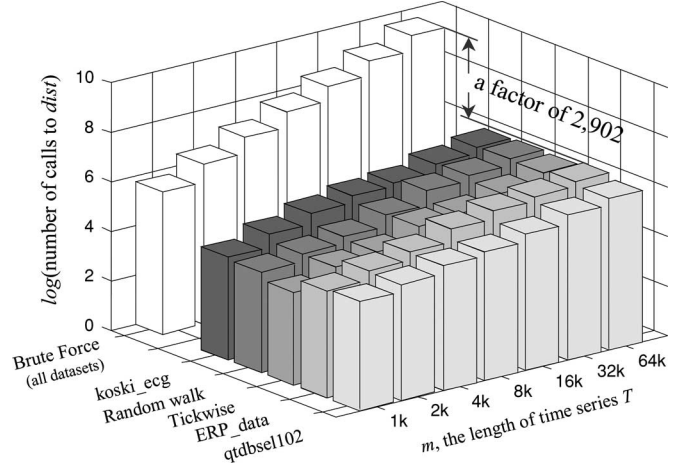


Fig. 11. The number of calls to the distance function required by brute force and heuristic search for discord_{128} over a range of data sizes for five representative datasets.

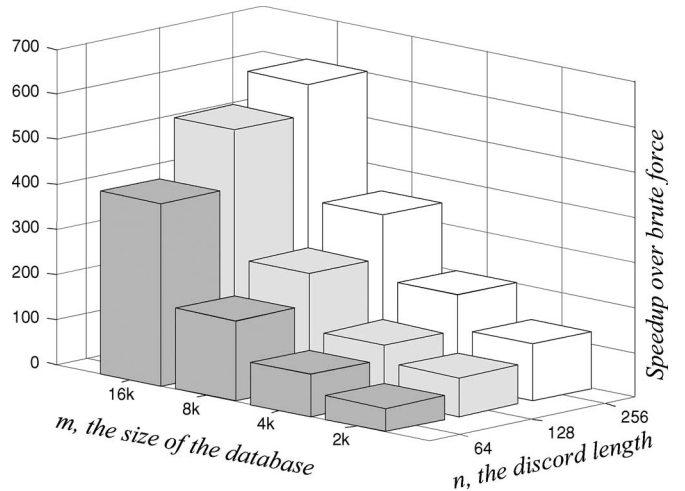


Fig. 12. The speed obtained over brute force search for various discord lengths and database sizes, averaged over 82 diverse datasets.

pseudocode (confirmed with a profiler) tells us that this single line of code accounts for more than 99% of the running time for both algorithms. In addition to fairness and reproducibility, there is another pragmatic reason for this metric. For brute force search, this number depends only on n and m and can simply be computed (recall m is the length of the time series and n is the length of the subsequence). If we had to actually measure the wall clock time for brute force search for all the experiments in this work, it would take several years.

The preceding metric does not include the time it takes to build the data structures discussed in Section 4-2; however, we note that this is a $O(m)$, one time cost. For datasets of a reasonable size (i.e., the datasets shown in Figs. 11 or 12), this overhead takes much less than 0.1% of the total time. Furthermore, as the datasets get larger, it takes an even smaller percentage of time.

In Fig. 11, we compare the brute force algorithm to the heuristic search algorithm in terms of the number of times the Euclidean distance function on line nine is called. For the heuristic search, we averaged the results for each setting of dataset/length over 100 runs on different subsets of the data.

Note that as the data sizes increase, the differences get larger. For a time series of length 64 000, the heuristic algorithm is almost three thousand times faster than brute force for all datasets. This experiment is actually pessimistic in that we made sure that the test data did not have any obvious anomalies or unusual patterns. In general, if there are truly unusual patterns in the time series, the heuristic algorithm is even faster.

In general, these results strongly suggest that we can reasonably expect at least three orders of magnitude of a speedup for most problems. To concretely ground these numbers, consider the following. While our current implementation is in relatively lethargic Matlab, the experiments shown in Fig. 10, Figs. 11, and 12 take a few seconds using heuristic search, but several hours using brute force search.

To make sure that these results were not the result of a happy coincidence of “easy” datasets and the right setting of the single parameter, we repeated the experiment for every dataset in the UCR Time Series Data Mining Archive over a range of values for n . We tested all datasets that have a length of at least 16 000; there are currently 82 such datasets from a diverse set of domains. Fig. 12 shows the results.

This experiment produces pessimistic results in that many of the datasets we averaged over are exceptionally noisy. In addition, the maximum size of the data (16 k) was relatively small to allow us to average over many datasets. Nevertheless, the results support the contention that a minimum speedup of two orders of magnitude can be expected for any combination of dataset and value for n , and even greater speedup can be expected as the datasets get larger.

We also considered the *Perverse* and *Random* heuristics on the preceding problem. As one might expect, perverse has exactly same performance as brute force search. The Random heuristic typically produces an approximately ten fold speedup over brute force, independent of the value of m ; while this is not insignificant, it is completely dwarfed by the *Magic* heuristic.

VI. CONCLUSION

In this work, we have defined time series discords, a new primitive for time series data mining. We introduced a novel algorithm to efficiently find discords and demonstrated their utility on a host of domains.

Many future directions suggest themselves; most obvious among them are extensions to multidimensional time series, to streaming data, and to other distance measures. In addition, for truly massive datasets, even the large speedups obtained may be insufficient for real-time interaction. We therefore plan to investigate an anytime version of our algorithm.

Reproducible results statement: In the interests of competitive scientific inquiry, all datasets used in this work can be found in [7].

ACKNOWLEDGMENT

The authors gratefully acknowledge the other donors of datasets. They also acknowledge insightful comments from C. Ratanamahatana.

REFERENCES

- [1] E. Keogh and S. Kasetty, “On the need for time series data mining benchmarks: A survey and empirical demonstration,” in *Proc. SIGKDD*, 2002, pp. 102–111.
- [2] E. Keogh, S. Lonardi, and C. Ratanamahatana, “Towards parameter-free data mining,” in *Proc. 10th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2004, pp. 206–215.
- [3] Z. Chen, A. Fu, and J. Tang, “On complementarity of cluster and outlier detection schemes,” in *Proc. Data Warehousing and Knowledge Discovery (DaWaK)*, 2004, pp. 234–243.
- [4] J. L. Bentley and R. Sedgwick, “Fast algorithms for sorting and searching strings,” in *Proc. 8th Annu. ACM-SIAM Symp. Discrete Algorithms*, 1997, pp. 360–369.
- [5] B. Chiu, E. Keogh, and S. Lonardi, “Probabilistic discovery of time series motifs,” in *Proc. 9th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2004, pp. 493–498.
- [6] S. Kitaguchi, “Extracting feature based on motif from a chronic hepatitis dataset,” presented at the 18th Annu. Conf. Jpn. Soc. Artificial Intelligence (JSAI), Kanazawa, Japan, Jun. 2–4, 2004.
- [7] J. K. Lanctot, M. Li, B. Ma, S. Wang, and L. Zhang, “Distinguishing string selection problems,” *Inf. Comput.*, vol. 185, pp. 41–55, 2003.
- [8] F. Duchene, C. Garbayl, and V. Rialle, “Mining heterogeneous multivariate time-series for learning meaningful patterns: Application to home health telecare,” Laboratory TIMC-IMAG, Faculté de médecine de Grenoble, Grenoble, France, 2004.
- [9] Y. Tanaka and K. Uehara, “Motif discovery algorithm from motion data,” presented at the 18th Annu. Conf. Jpn. Soc. Artificial Intelligence (JSAI), Kanazawa, Japan, Jun. 2–4, 2004.
- [10] W. L. Ruzzo and M. Tompa, “A linear time algorithm for finding all maximal scoring subsequences,” in *Proc. Int. Conf. Intell. Syst. Mol. Biol.*, 1999, pp. 234–241.
- [11] S. Romblo and G. Terracina, “Discovering representative models in large time series databases,” in *Proc. 6th Int. Conf. Flexible Query Answering Systems*, 2004, pp. 84–97.
- [12] E. Knorr, R. Ng, and V. Tucakov, “Distance-based outliers: Algorithms and applications,” *VLDB J.*, vol. 8, no. 3–4, pp. 237–253, 2000.
- [13] T. H. Corman, C. E. Leiserson, and R. L. Rivest, *Introduction to Algorithms*. New York: McGraw-Hill, 1990.
- [14] C. Ratanamahatana and E. Keogh, “Making time-series classification more accurate using learned constraints,” presented at the 4th SIAM Int. Conf. Data Mining (SDM’04), Lake Buena Vista, FL, Apr. 22–24, 2004, pp. 11–22.
- [15] J. Lin, E. Keogh, S. Lonardi, and B. Chiu, “A symbolic representation of time series, with implications for streaming algorithms,” presented at the 8th ACM SIGMOD Workshop Research Issues in Data Mining and Knowledge Discovery, San Diego, CA, Jun. 13, 2003.
- [16] N. Kumar, N. Lolla, E. Keogh, S. Lonardi, C. A. Ratanamahatana, and L. Wei, “Time-series bitmaps: A practical visualization tool for working with large time series databases,” in *Proc. SIAM Int. Conf. Data Mining (SDM’05)*, Newport Beach, CA, Apr. 21–23, 2005, pp. 531–535.
- [17] J. Lin, E. Keogh, S. Lonardi, J. P. Lankford, and D. M. Nystrom, “Visually mining and monitoring massive time series,” in *Proc. 10th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2004, pp. 460–469.
- [18] K. Sadakane, “Compressed text databases with efficient query algorithms based on the compressed suffix array,” in *Proc. ISAAC, LNCS*, 2000, pp. 410–421.
- [19] E. Keogh (2005). The time series discords page. [Online]. Available: <http://www.cs.ucr.edu/~eamonn/discords/>
- [20] D. Dasgupta and S. Forrest, “Novelty detection in time series data using ideas from immunology,” in *Proc. 5th Int. Conf. Intelligent Systems*, Reno, NV, Jun. 1996, pp. 87–92.
- [21] J. Ma and S. Perkins, “Online novelty detection on temporal sequences,” in *Proc. Int. Conf. Knowledge Discovery and Data Mining*, Aug. 2003, pp. 613–618.
- [22] C. Shahabi, X. Tian, and W. Zhao, “TSA-tree: A wavelet based approach to improve the efficiency of multi-level surprise and trend queries,” in *Proc. Int. Conf. Scientific and Statistical Database Management*, Jul. 2000, pp. 55–68.



Eamonn Keogh received the B.S. degree in computer science from California State University, San Marcos, in 1995, and the Ph.D. degree in computer and information science from the University of California, Irvine, in 2001.

He is an Assistant Professor of Computer Science at the University of California, Riverside. His research interests are in data mining, machine learning, and information retrieval.

Dr. Keogh's several papers have won best paper awards, including papers at SIGKDD and SIGMOD.

He is the recipient of a five year NSF Career Award for "*Efficient Discovery of Previously Unknown Patterns and Relationships in Massive Time Series Databases.*"



Jessica Lin received the B.S. and Ph.D. degrees in computer science from the University of California, Riverside, in 2000 and 2005, respectively.

She is an Assistant Professor of Information and Software Engineering at George Mason University, Fairfax, VA.



Ada Waichee Fu received the B.S. degree from The Chinese University of Hong Kong, Hong Kong, in 1983, and the M.S. and Ph.D. degrees from Simon Fraser University, Burnaby, BC, Canada, in 1986 and 1990, respectively, all in computer science.

In 1993, she joined The Chinese University of Hong Kong. Her research interests include issues in distributed databases, XML data, time series databases, data mining, content-based retrieval in multimedia databases, and parallel and distributed systems.



Helga Van Herle received the B.S. degree in chemical engineering from the University of California, Los Angeles (UCLA), in 1985, the M.S. degree in bioengineering from Columbia University, New York, NY, in 1987, and the M.D. degree from UCLA in 1993. She completed her residency in internal medicine at the New York Hospital, Cornell University, Ithaca, NY, in 1996, and her cardiology Fellowship at UCLA in 2001.

She is an Assistant Clinical Professor of Medicine in the Division of Cardiology, David Geffen School

of Medicine, UCLA.