

Accelerating Dynamic Time Warping Subsequence Search with GPUs and FPGAs

Doruk Sart, Abdullah Mueen, Walid Najjar,
Eamonn Keogh
University of California, Riverside
{sartd, mueen, najjar, eamonn}@cs.ucr.edu

Vit Niennattrakul
Department of Computer Engineering
Chulalongkorn University, Thailand
g49vnn@cp.eng.chula.ac.th

Abstract—Many time series data mining problems require subsequence similarity search as a subroutine. While this can be performed with any distance measure, and dozens of distance measures have been proposed in the last decade, there is increasing evidence that Dynamic Time Warping (DTW) is the best measure across a wide range of domains. Given DTW’s usefulness and ubiquity, there has been a large community-wide effort to mitigate its relative lethargy. Proposed speedup techniques include early abandoning strategies, lower-bound based pruning, indexing and embedding. In this work we argue that we are now close to exhausting all possible speedup from software, and that we must turn to hardware-based solutions if we are to tackle the many problems that are currently untenable even with state-of-the-art algorithms running on high-end desktops. With this motivation, we investigate both GPU (Graphics Processing Unit) and FPGA (Field Programmable Gate Array) based acceleration of subsequence similarity search under the DTW measure. As we shall show, our novel algorithms allow GPUs, which are typically bundled with standard desktops, to achieve two orders of magnitude speedup. For problem domains which require even greater scale up, we show that FPGAs costing just a few thousand dollars can be used to produce four orders of magnitude speedup. We conduct detailed case studies on the classification of astronomical observations and similarity search in commercial agriculture, and demonstrate that our ideas allow us to tackle problems that would be simply untenable otherwise.

Keywords- *time series; similarity search; dynamic time warping; FPGA; GPU;*

I. INTRODUCTION

Subsequence similarity search, the task of finding a region of much longer time series that matches a specified query time series within a given threshold, is a fundamental subroutine in many higher level data mining tasks such as motif discovery [19], anomaly detection [4], association discovery, and classification [20][1][33].

More than one hundred different distance measures for time series have been proposed in the last decade [12]; however there is increasing empirical evidence that Dynamic Time Warping (DTW) (which includes Euclidean Distance as a special case) is the best measure across a wide range of domains [7]. Given DTW’s usefulness and ubiquity, there has been a large community-wide effort to mitigate its relative lethargy in the last decade. Speedup techniques for general search under DTW include various indexing methods [1][11][7]. Speedup techniques for the special case of a subsequence similarity search that we are considering here

include early abandoning strategies, embedding and various “computation reuse” strategies [27]. A recent paper has shown that much of the *apparent* progress made in recent years is fatally flawed [4]. In particular, the speedup comes at the cost of allowing arbitrary false dismissals (we will expand on this surprising result in Section III.B).

Even if the apparent recent results had been correct, there still exist problems for which no current algorithms running on standard hardware can hope to solve in a reasonable amount of time. As a concrete example (which we expand to a case study in Section VII.A), entomologists need to examine telemetry gathered from insects for the occurrence of certain patterns known to be indicative of destructive (to host plants) behaviors. Entomologists at the University of California have created an archive of four hundred million data points of this data in the last four years, as part of an effort to understand and ultimately control just one insect, the Glassy-winged Sharpshooter (*Homalodisca vitripennis*). This insect causes tens of millions of dollars of damage to the grape industry. Moreover, these entomologists are adding one to two million data points *per day* to this archive. Searching this archive under the DTW distance for a single (relatively short) query pattern of length 360 takes nine days on a high-end desktop, using state-of-the-art algorithms. Similar stories can be told for astronomy (cf. Section VII.B), computational finance, motion capture processing data [5], industrial and medical domains.

After surveying and testing the current software solutions, and talking to several domain experts and practitioners, we have come to the conclusion that we are now close to exhausting all possible speedup from software approaches, and that we must turn to hardware-based solutions if we are to tackle the problems faced by real world practitioners. With this motivation, we investigate both GPU and FPGA based acceleration of subsequence similarity search under the DTW measure. The use of specialized hardware to allow subsequence similarity search requires a detailed understanding of both the hardware strengths and limitations, and of the DTW computation itself.

As we shall show, our novel algorithms allow GPUs, which are typically bundled with standard desktops and are thus essentially free, to achieve two orders of magnitude speedup. We show that if a domain practitioner is motivated enough to purchase an FPGA, which can cost as little as a few thousand dollars, our algorithm can achieve a speedup of four orders of magnitude.

It is important to note that we see our work as going beyond the claim that “*we have made an important algorithm faster*”. A factor of say, two, speedup for an

important algorithm is useful, but unlikely to make a significant difference to the community. However, a speedup factor of a thousand or more *really* has the potential to make a significant difference, because it allows problems to be tackled that are otherwise unimaginable. To consider one concrete example, time series motif discovery is a useful tool with applications in dozens of domains. A recent paper introduced a technique to find motifs in datasets containing millions of objects in just hours, a significant speed-up [19]. This method explicitly assumes the Euclidean Distance; however, for the related problem of classification, it is well-known that DTW is significantly more accurate [7][25][33]. Could we do motif discovery under DTW instead? Even if we assume optimistically that the DTW could somehow be made to obey the triangular inequality¹, using the MK algorithm in [19] with DTW instead of Euclidean Distance would push the execution time to years instead of hours.

It is also important to note that we see our work as a service to the time series data mining community. Many research projects use DTW similarity search as a subroutine, and could greatly benefit from significantly improved performance. As such, all of the code developed for both GPU and FPGA will be freely available at [39] in perpetuity.

The rest of this paper is organized as follows. In Section II we consider some related work, followed by the definition of and background to the problem at hand in Section III. Our GPU and FPGA acceleration methods are described in Section IV and V, respectively. Section VI shows an experimental comparisons of different methods. Finally, in Section VII, we conduct three detailed case studies in classification of star light curves, commercial entomology and data condensing for nearest neighbor classification, and demonstrate that our ideas allow us to tackle problems that would be simply untenable otherwise.

II. RELATED WORK

The problem of subsequence similarity search using the DTW measure has been extensively studied in recent years. One of the most cited methods is SPRING [27], where a query time series is searched in a larger streaming time series. The authors achieve significant speed-up by reusing computations. Unfortunately, this reuse means that the method allows *false negatives*, a problem we will elaborate on in Section II.B.

Our paper joins the growing literature on hardware acceleration techniques for important database/data-mining algorithms, such as FFT [9], relational join [10], etc. However to the best of our knowledge, this work is the first to present hardware acceleration techniques for similarity search in streams under the DTW measure.

¹ To be clear, the DTW does not obey the triangular inequality and is thus a distance *measure*, not a distance *metric*, in spite of more than a dozen papers that claim the latter. However, some of the lower bounds to DTW are metrics, and thus some triangular-inequality exploiting algorithms can be adapted to work with DTW.

III. DEFINITION AND BACKGROUND

For concreteness we begin with a formal definition of the problem and a discussion of why the current solutions are inadequate. We begin by defining the time series:

A *time series* T is a sequence of real numbers t_1, t_2, \dots, t_n representing n uniform samples of a measurement. A *subsequence* $C_{s,m}$ of a time series T is any contiguous set of m samples starting at s . (i.e., $t_s, t_{s+1}, \dots, t_{s+m-1}$, where $1 \leq s \leq n-m+1$).

Before we compare two time series under any distance measure, it is critical that we normalize them to have the same mean and variance². As noted in [12], “*without normalization time series similarity has essentially no meaning*”.

The *z-normalization* of a time series T is defined as $\hat{T} = \hat{t}_1, \hat{t}_2, \dots, \hat{t}_n$, where $\hat{t}_i = (t_i - \mu) / \sigma$. Here, μ and σ are the sample mean and the sample standard deviation of T, respectively.

It is critical at this point to clarify a naive misunderstanding which is replete in the literature. If we are doing subsequence similarity search with our z-normalized query Q for the best matching subsequence in a much longer time series T, we *cannot* simply z-normalize T once and proceed. Instead, we must z-normalize each and every subsequence we extract from T. Note that in the case that T is not a batch dataset residing in its entirety in memory (or disk), but in a *data stream*, it would not even be logically possible to z-normalize it all, even if doing so gave meaningful results.

While DTW is defined to allow for the comparison of two time series of possibly different lengths, without losing the generality (see [25], Section 2), we will define it assuming time series of equal lengths.

Suppose we have two time series, $C = c_1, c_2, \dots, c_i, \dots, c_m$ and $Q = q_1, q_2, \dots, q_j, \dots, q_m$. The *Dynamic Time Warping* (DTW) distance between Q and C is denoted by $D(C, Q)$ and defined as below.

$$D(C, Q) = d(m, m)$$

$$d(i, j) = |c_i - q_j| + \min \begin{cases} d(i-1, j) \\ d(i, j-1) \\ d(i-1, j-1) \end{cases}$$

$$d(0, 0) = 0; d(i, 0) = d(0, j) = \infty; i = 1, 2, \dots, m; j = 1, 2, \dots, m$$

In the definition, we use absolute difference for distances between individual samples. Our methods also work for the squared difference with no difference in speed.

The m -by- m matrix, d , is called the *warping matrix*. In a warping matrix, each cell uses a value from any of the three previously computed neighbors. Thus, if we trace back the values used to compute the DTW (i.e. $d(m, m)$), we get the *warping path* that describes the optimal alignment of T and Q (Figure 1).

² Some papers have suggested doing [0,1] or [-1,1] normalization instead. However, the authors do not seem to appreciate how brutally sensitive this method is to even small amounts of noise or a single outlier.

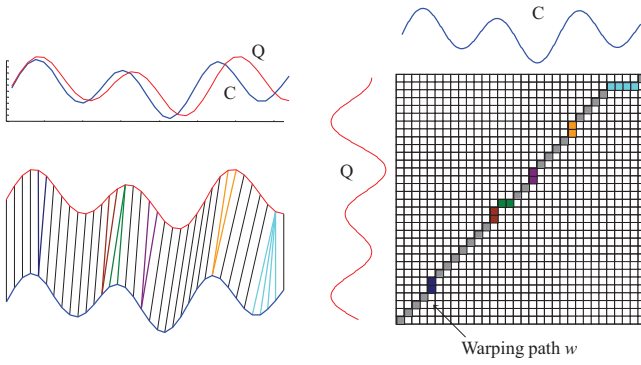


Figure 1 (left-top) Two time series sequences which are similar but out of phase. (right) To compute the dynamic time warping distance and to align the sequences, we construct a warping matrix and find a warping path, represented by solid squares.

The time complexity to compute the $D(C,Q)$ is $O(m^2)$, and the space complexity is also $O(m^2)$. If we only need the value of the distance (i.e. $d(m,m)$) we can delete the trace of the warping path, and thus, the space complexity can be reduced to $O(m)$ by storing only two columns of the matrix.

A. Definition of the Problem

We are now in a position to define the subsequence search problem:

Given a time series $T = t_1, t_2, \dots, t_n$ and a query $Q = q_1, q_2, \dots, q_m$, find the subsequence $C_{s,m}$ of T such that $D(C_{s,m}, Q)$, $1 \leq s \leq n-m+1$, is minimum.

Given the above definition, we could devise a brute force algorithm shown in Table 1, which takes $O(nm^2)$ time and $O(nm)$ space. For completeness, we also show the pseudocode for computing the DTW distance in Table 2.

TABLE 1: SUBSEQUENCE SEARCH ALGORITHM

Procedure SubsequenceSearch(T, Q)	
	T : A time series of n points Q : Query time series of m points
1	z-Normalize(Q)
2	for $s = 1$ to $n-m+1$
3	z-Normalize($C_{s,m}$)
4	Compute $D(C_{s,m}, Q)$
5	Update minimum if necessary

TABLE 2: DTW ALGORITHM

Procedure $D(C, Q)$	
	C : A time series of n points, $C(0) = \infty$ Q : A time series of m points, $Q(0) = \infty$
1	$s = 0$
2	for $i = 0$ to m
3	$d(i, s) = C(1) - Q(i) $
4	$s = s \oplus 1$ // xor operation
5	for $j = 2$ to n
6	for $i = 0$ to m
7	$d(i, s) = C(j) - Q(i) +$ $\min(d(i-1, s), d(i, s \oplus 1), d(i-1, s \oplus 1))$
8	$s = s \oplus 1$
9	return $d(n, s \oplus 1)$

We have chosen the simplest possible problem definition with one query, one time series and the same subsequence length (m). There are more general subsequence search problems where many queries [29] and time series are involved, or where rotation/phase invariance is required under DTW [34][26]. However, all such problems can benefit directly from a speedup of the simple definition.

B. Why Current Software Solutions Are Not the Answer

As we hinted at above, the several *apparent* software solutions to the task at hand contain a serious error. We can best demonstrate this with a simple experiment.

Suppose we task a DTW subsequence search with the simple task of detecting the heartbeats of an individual, using one of that same individual's heartbeats. It is difficult to imagine a simpler problem.

We begin by downloading a long ECG sequence from a 61-year-old female and manually extracting a typical beat as our query [39]. We also manually extract some additional adjacent beats and compare them to our query, finding them to be an average distance of about 20.0, so we set our beat detector at a conservative threshold of 30.0. Figure 2 shows the beats detected in the first 1,800 datapoints, as we can see, the majority of the beats are missed. How could this be?

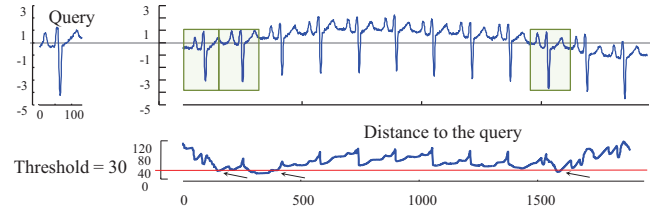


Figure 2: A query heartbeat (left) is scanned across an ECG trace. (top-right) Only three of the twelve beats are detected. Plotting the distance from the query to the relevant subsequence (bottom) reveals that slight differences in a subsequence's mean value (offset) completely dominate the DTW distance calculation, dwarfing any contribution from the similarity of the shape.

Note that while the local mean of the ECG trace starts at about zero, which is also the approximate mean of the query, the trace slowly rises to have a local mean of about 1.0, then descends below zero (allowing the detection of a single beat at about 1,500 as the mean crosses zero).

The problem is that the SPRING algorithm [27] does not (and, more critically *cannot*) normalize the offset or amplitude of the subsequences of the longer time series. It is therefore implicitly assuming that the query will happen to have the same offset and amplitude as the matching subsequence. However, in virtually every domain that assumption is unwarranted. For example, virtually all ECGs wander up and down as in our example, the effect is known as a *wandering baseline* [18]. Similar problems are observed in motion capture [16], astronomy, entomology, industrial process telemetry, EEGs, etc.

It is important to recognize that *there is no simple fix for this problem*. The SPRING [27] algorithm achieves its speedup by exploiting the redundancy of calculations in a sliding DTW matrix, but if each subsequence is z-normalized, as it must be to obtain meaningful results, then there will be no redundant calculations to exploit. For brevity

we will conclude the discussion of the reasons why current software solutions do not work here. The interested reader can view [39], where we have many additional examples of the problem in real domains and a detailed discussion of why the obvious possible fixes for the problem we have just pointed out will not solve it. In our view, this problem is very unlikely to yield to a software solution that improves the time complexity of the brute force algorithm in Table 1. We note in passing there are at least one dozen other works that have the same problem, we pointed to SPRING only because it is the most cited work on the topic.

IV. ACCELERATION BY GPU

The GPU is a computing device that serves as a coprocessor for the CPU. It has its own device memory on the card and can execute many threads in parallel [22]. In this work we use the ubiquitous NVIDIA CUDA architecture, where multiple threads running on multiple processing cores execute the same program on separate data. This Single Instruction Multiple Data (SIMD) architecture allows us to map each normalization and DTW distance computation (Lines 3-4 in Table 1) to be executed in parallel on different segments of the time series.

Each CUDA function (i.e. kernel) is executed by an array of threads. Each of these threads is assigned an ID that it uses to determine memory addresses (i.e. the segment of the time series) it should operate on. The hardware is free to determine the mapping and scheduling of these threads on the available processing cores. A thread block is defined as a batch of threads that are guaranteed to run simultaneously and cooperate with each other through shared resources. The size of a thread block can be specified at runtime. The NVIDIA CUDA thread architecture is shown in Figure 3.

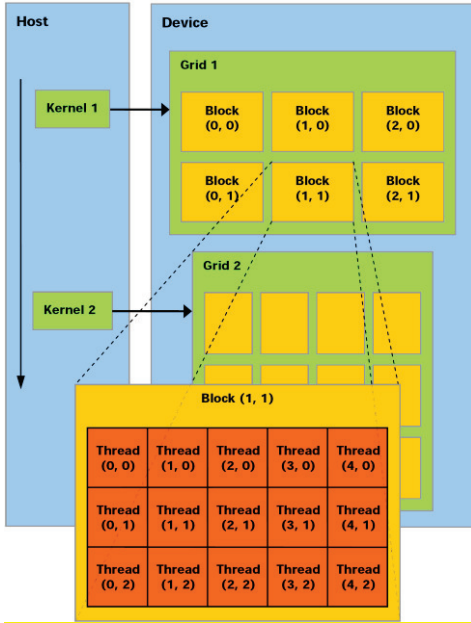


Figure 3: NVIDIA CUDA Architecture: Each GPU kernel is executed as an array of threads. Blocks are batches of threads that run concurrently.

The selection of the number of threads per block is an important parameter to maximize the utilization of the processor cores. Unfortunately, it depends on the architecture of the GPU. Thus, obtaining the optimum value is not trivial. The important variables to consider are the total number of processing cores and their internal pipeline architectures. For a more detailed discussion, see [22].

The GPU implementation of Dynamic Time Warping consists of three main stages:

- The CPU copies the values to the GPU memory
- The CPU calls the GPU kernel
- The CPU copies the output from the GPU

In the first step, the CPU copies the whole time series (T in Table 1) to the global memory of the GPU. If the time series is larger than the available device memory, the CPU splits it into small batches and processes one batch at a time. This process introduces latency in the output but does not hamper the real-time processing, as the time to copy the data is in the range of milliseconds. Therefore, copying batches one at a time can tolerate a data arrival rate of hundreds of hertz without overflowing any buffer. Since the query is fixed and good for all of the batches, we copy it to the global memory in the beginning and keep it there throughout the execution.

In the second step, the CPU calls the kernel in the GPU. Every kernel thread operates on a specific sliding window in two steps: first, accessing the sliding window to compute the mean and variance, and second, computing the normalized DTW distance to the query. For both the steps, each kernel thread accesses a contiguous segment of m numbers from the time series T in the global memory. If we batch the threads responsible for successive sliding windows in a thread block, the memory accesses by these threads will result into coalesced accesses. For example in Figure 4, a block of four threads is shown where the first memory accesses by these four threads require one read from the memory instead of four, because of the threads operating on contiguous locations in the memory.

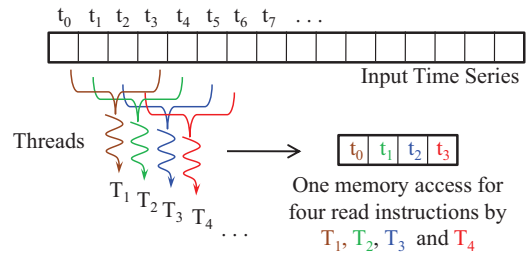


Figure 4: Division of work among threads. Memory accesses are coalesced by overlapping threads (for $m=4$).

In the second step, when the mean and variance are ready, each thread computes the DTW distance between the query and the normalized subsequence by Table 2. Note that in Table 2, the query Q is accessed $O(m^2)$ times whereas the subsequence C is accessed $O(m)$ times. It is important to maintain this distinction between Q and C although swapping Q and C in Table 2 would produce correct results.

The reason is that Q is a fixed time series, and by the problem definition it is not changed during execution. Furthermore Q is much smaller than T so it can fit in the shared memory of the GPU device. Shared memory is a special memory locally available to the processor cores in a multi-processor of the GPU device. Shared memory is 150x faster than the global memory which is available to all of the cores in all of the multi-processors. Typically, shared memory is much smaller than the global memory and thus is the ideal place for the query time series.

In the beginning of the DTW computation, the kernel threads copy the query (Q) into shared memory from the global memory. The internal data structure for computing the DTW distance is two column vectors (d in Table 2) of size m . These vectors are best stored in the shared memory if they all fit together (e.g., in the case of very small m). Otherwise, they are stored in the global memory. With all of the variables in place, the DTW computation is performed. Each thread stores the computed distance in a global array indexed by the thread ID.

Finally in the third step, when all of the threads terminate, the CPU copies the array back to system memory. Although the algorithm looks for the *minimum* distance, it is much simpler and more efficient to copy back *all* of the distances to the CPU instead of computing the minimum in the GPU. Maintaining a variable for the minimum distance computed by all of the threads requires an atomic instruction in all of the threads that update the variable. Such an instruction eventually makes all of these threads update the variable serially and reduces the parallelism. Therefore, atomic updates are not recommended.

V. ACCELERATION BY FPGA

A. Data Representations for FPGAs

While GPUs can handle essentially any data that a normal CPU can, FPGAs require reducing the cardinality (precision) of the data. This is because in modern FPGAs, floating point arithmetic does not scale well with larger applications [15] due to the additional complexity for handling the mantissa and exponent separately. Therefore, we select 8-bit integers for representing values in the time series. Reducing the number of bits not only increases the performance but also increases the area utilization. Furthermore, for a given size of the FPGA, we *may* need to impose an upper limit on the dimensionality (length) of the query.

Given that we may need to reduce the cardinality and/or dimensionality of the data, it is natural to ask if this will reduce the accuracy of FPGA-based DTW similarity search. We can strongly affirm that these changes to raw data will *not* affect the accuracy of DTW in any meaningful (or indeed *perceptible*) way. To see this we can perform some simple experiments. We took all twenty datasets from the UCR Time Series Classification Archive [13] and tested to see what effect reductions in the dimensionality or cardinality would have. We began by running DTW one-nearest-neighbor classification on the raw data, confirming the published results [13]. We then took the data, reduced the

precision from eight *bytes* to just eight *bits*, and reran the DTW one-nearest-neighbor classification. We also took the raw data, and reduced the *dimensionality* of all datasets from their original lengths to exactly 128 (the original lengths for the twenty problems had a maximum of 637 and a mean of 310.8). While there are sophisticated interpolation algorithms for resizing time series, we simply used the trivial MATLAB command:

$$x = x(1:\text{length}(x)/128:\text{end})$$

to resize the data. The results of this experiment can be visualized in Figure 5.

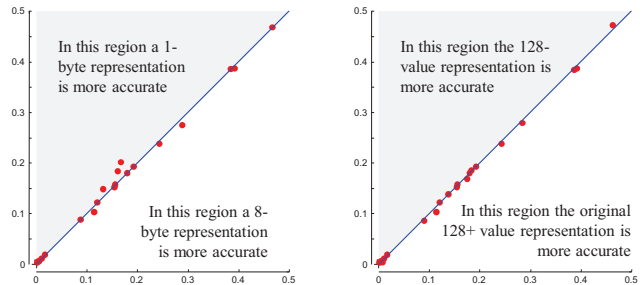


Figure 5: The error rate of twenty time series classification problems conducted with two reduced representations (*left*) The original data vs. a reduced precision data. (*right*) The original data vs. a reduced dimensionality data.

A paired T-test³ of the null hypothesis that there is no difference between the two mean accuracies produces a value of 49.2%, suggesting what is visually obvious from the figures: there is simply no significant difference made by the change of representation.

We briefly note that if there was a reason to insist that the FPGA returns the *exact* same answer as the original data, we believe that this would be possible to achieve by producing a *bound* on the difference between the distances measured in the raw space, and the distances measured in the reduced space, and then returning all of the sequences that match within the bound for additional inspection. Note that the downsampling could be achieved by Piecewise Aggregate Approximation (PAA), and a bound (LB_PAA) between DTW on the raw data and DTW on the PAA representation exists [36]. Likewise, the reduction in cardinality could be achieved by Symbolic Aggregate Approximation (SAX); a bound between DTW on the raw data and DTW on the SAX representation also exists [6]. We do not further explore such considerations, given that Figure 5 strongly suggests it is not a fruitful problem.

B. FPGA Implementation

The design of an FPGA configuration requires programming using a hardware description language (i.e., VHDL/Verilog). In this work, we use an open source C to VHDL compiler system, *ROCCC* [37], which allows us to describe the hardware in C language and generate the VHDL code automatically. *ROCCC* also optimizes the design

³ More appropriate and powerful measures of significance (see [28]), which conforms the null hypothesis.

mainly in three ways. First, it maximizes the throughput by exploiting loop and instruction level parallelism. Second, it reuses the data, and third, it generates a pipelined datapath to minimize the number of clock cycles [29].

Our FPGA design consists of two major blocks: *Normalizer* and *Warper*, to normalize the input data and run the actual DTW matrix calculations, respectively (Figure 6). Input data streamed into the system are first given to a First-In-First-Out (FIFO) buffer. The size and input ratio of this FIFO can be adjusted according to the FPGA interconnection mechanism. However, the output of the FIFO generates one sample (8 bits) every clock cycle. Next, the output of the FIFO is fed into the Normalizer module. Initially, Normalizer waits until the first window is received. Every following normalization operation reuses $m-1$ operands of the previous operation, where m is the query length. After the first output is produced, a new output is generated every clock cycle. This output is given to another FIFO, which acts as the intermediate memory component between the Normalizer and the Warper.

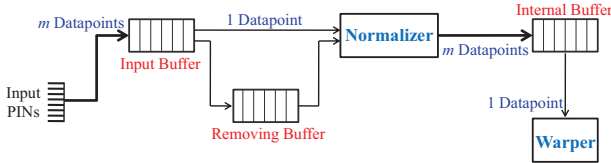


Figure 6: FPGA Block Diagram. Thick lines are for m -point wide connections. Thin lines are for one-point wide connections. Buffers are simple FIFOs.

Internally, a trivial Normalizer module stores m -partial sliding windows. In every clock cycle, it updates statistics for all of the partial windows and outputs the window for which the normalization is complete. Thus, it needs quadratic $O(m^2)$ space in the FPGA and does not scale with larger query lengths. In order to support larger query lengths, we implemented an *online Normalizer*, which does not remember intermediate states. It computes the mean (μ) and standard deviation (σ) online and normalizes exactly one window in every clock cycle. Thus, it needs linear $O(m)$ space in the FPGA. Although the trivial Normalizer has shown better performance in speed due to less overhead, it does not make any difference in the overall system performance. The reason for this is that the Warper module is the real performance bottle-neck as described later.

The Online Normalizer consists of three sub-units, as shown in Figure 7. The first unit calculates the sum and sum of squares of all the inputs in a sliding window fashion, by adding the new value while subtracting the oldest value to be removed from the sum. Initially, “to remove” value is zero until the first window is completely received through the “Datapoint” input. When switching to the next window, the very first value of the first window is given to “to remove” and the sum for the second window is obtained at the output. This output is also given to the Normalize Divider sub-unit, where the mean and the standard deviation of the latest window are obtained. The input stream is provided to the third unit through a buffer. The size of this buffer depends on the delay of the first two modules. The third module must

wait until the corresponding mean and standard deviation values are available for a given window. This delay is provided by the Datapoint Buffer, which is automatically added by *ROCCC*. The unit then runs the actual normalization function. The generated normalized data is provided to the systolic array (warper) through a buffer, as shown in Figure 6.

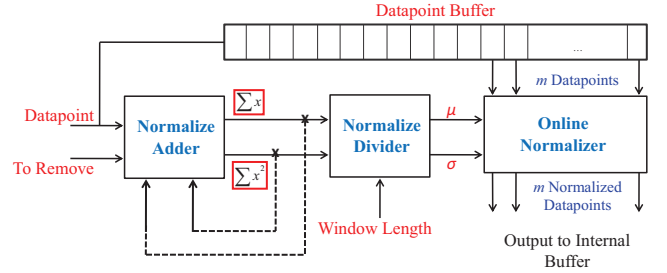


Figure 7: Online Normalization Unit. The sum and sum of squares are obtained in a sliding window approach, by adding the new input and subtracting the oldest value. The input datapoints are delayed through the Datapoint Buffer, to make sure that the correct mean and standard deviation are used.

The Warper module is implemented as a systolic array [3]. A systolic array consists of data processing units connected in a matrix fashion. These data processing units (i.e. cells) share the information with their neighbors immediately after processing. Using *ROCCC*’s built-in systolic array generator, we simply obtain the hardware description of the Warper module. Structurally, the Warper module is the same for any window size, except for the size of the systolic array. This size can be adjusted in the *ROCCC* code by tuning a parameter. A Warper module generates one DTW distance between the normalized sliding window and the query time series in every m clock cycle where m is the window size/query length. Since the normalization unit is m times faster than the Warper unit, we place multiple Warper units to operate on separate normalized windows generated by the normalization unit. Ideally, if we had unlimited FPGA area, we could place m Warper modules to get the maximum processing speed of one DTW distance in every clock cycle. When multiple Warper modules are in place, the Internal Buffer output is fed into them in a round robin fashion.

VI. EVALUATION

In this section, we show the performances for the DTW subsequence search problem in different hardware settings. We use the following platforms:

- Software:** Intel Xeon E5540 CPU at 2.53 GHz
- SSE:** Intel i7- 920 CPU at 2.66 GHz
- GPU:** NVIDIA Tesla C1060 with 240 cores
- FPGA:** Xilinx Virtex 5 LX-330

The SSE (Streaming Single Instruction Multiple Data (SIMD) Extensions) is an instruction set extension to Intel’s x86-architecture. It makes use of 128-bit SSE registers and can merge four 32-bit data to operate concurrently. The software implementation proposed in Table 1 can be parallelized by making use of data independencies. However,

it is practically impossible to realize this parallelization on the DTW matrix calculations. We can only execute SSE instructions on the normalization operations. The performance improvement is therefore not significant compared to the software-based solution.

In Figure 8, we show the time required to answer a query of length 128 by different hardware settings. We achieve the highest speedup over the software through FPGA acceleration, which is 4000 times faster in the best case scenario. GPU acceleration is 36.3 times faster, on average. All of the results reported here use 8-bit integers to represent the values in the time series.

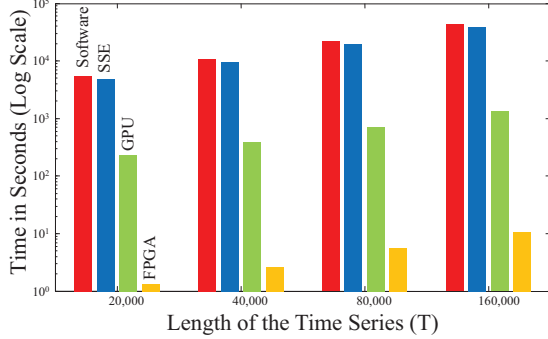


Figure 8: Comparison of execution times with different lengths of the input time series, in LOG scale. The query length is 128. For GPU, block size is 512.

The FPGA performance results are obtained after placement and routing operations. We use one trivial Normalizer unit and eight Warper units. The Normalizer for window size 128 requires 13% of the target device and runs at 180MHz. Each Warper unit requires 7% of the area and run at 240MHz. The Normalization unit provides an output for each clock cycle. The Warper requires 128 clock cycles. Therefore, the Warper unit runs more slowly than the Normalizer. At 240 MHz, one window of length 128 can be processed through the Warper in 128/240M seconds. However, using 8 Warpers reduces this to 16/240M seconds. In other words, 15 million samples (windows) can be processed per second.

In Figure 9, we show the execution times for a fixed window size of 1024, the FPGA provides a maximum speedup of 4500x and the GPU achieves a speedup of 29x over software. For the window size 1024 on the FPGA, we use the online Normalizer, as shown in Figure 7. This module runs at a clock frequency of 180MHz and requires 83% of the FPGA logic. The Warper module runs at 250MHz and requires 9% of the logic. In this case, we can only place one Warper module safely. Although the area is dominated by the Normalizer, the throughput of the system is still determined by the Warper module. The Warper module requires 1024 clock cycles per cycle. One sample (window) can be processed in 1024/250M seconds. This results in a throughput of 244 thousand samples per second. Note that in the streaming case there are probably few data sources that could produce data from a *single* stream at anything like this rate. However, because the overhead for switching between streams is so low, we could use a single

FPGA to handle say two hundred streams at 1,000Hz, by simply multiplexing between (slightly buffered) streams.

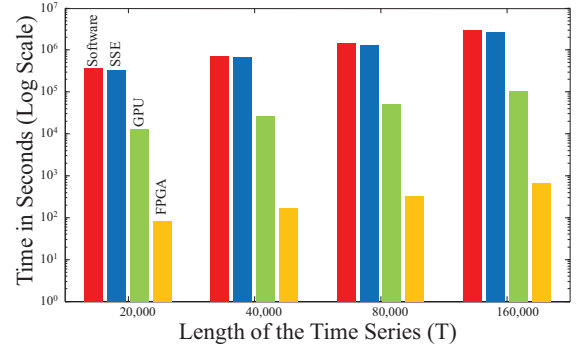


Figure 9: Comparison of execution times with different lengths of the input time series, in LOG scale. The query length is 1024.

In Figure 10 we show the responses of different methods while varying the size of the query. Recall the methods have the same time complexity of $O(nm^2)$. The responses show a clear quadratic growth for software and SSE methods. Our hardware acceleration techniques are much slower in growth because of the parallelism our techniques achieve. The trends in the figure clearly show that our techniques will remain tenable for larger window size while the software methods are already intractable.

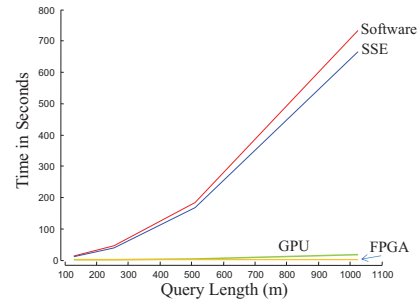


Figure 10: Execution times for a fixed time series of length 40,000 with varying query sizes (m).

VII. EXPERIMENTAL CASE STUDIES

In this section we conduct case studies that forcefully confirm all of the claims and assumptions made in this work. Recall that we have claimed:

- Time series subsequence search is an important problem with scientifically and commercially important applications.
- For at least some problems, Euclidean distance, while tractable, is simply not accurate/robust enough.
- There exist problems of a scale which makes them untenable with conventional hardware. We propose special hardware-based solutions to make them yield.

A. Case Study in Entomology

Many species of insects feed by inserting their stylet (a long needle-like mouthpart) into a plant and sucking out sap. While this behavior in itself is generally not harmful to the plants, if one plant has a disease, the insects will transmit it from plant to plant.

It is impossible to overstate the commercial impact of this insect behavior. For example, if we consider just *one* species of insect, the Aster leafhopper (*macrosteles quadrilineatus*), in *one* American state, Wisconsin, and *one* crop, the carrot, it is estimated that losses exceed two million dollars a year [8]. However, the Aster leafhopper is just one of perhaps two thousand insects that exhibit this harmful piercing/sucking behavior. The worldwide commercial losses to agriculture are hard to estimate, but for the wine industry alone the cost is conservatively estimated at several billion US dollars.

In order to control these insects, their behavior must be understood. However, direct visual observation of a tiny and fast-moving insect is very difficult, to say the least. In a 1964 paper in *Nature*, McLean and Kinsey [17] proposed attaching a small wire to the insect with conductive glue, and creating a circuit through the ground in which the plant is potted. As soon as the insect's stylet penetrates the plant, the circuit is completed and a fluctuating voltage, called the Electrical Penetration Graph (EPG) signal, occurs and can then be amplified and recorded. Figure 11 shows an example EPG.

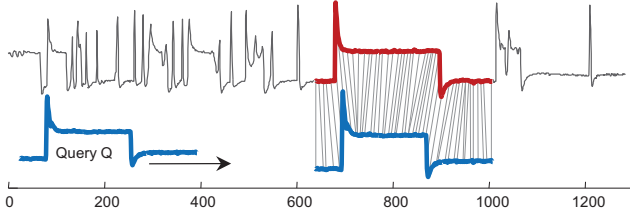


Figure 11: A short extract of an EPG graph of a beet leafhopper (*Circulifer tenellus*). A query Q can be searched against the longer time series to find examples of similar behavior.

Since the invention of the McLean-Kinsey apparatus, EPGs have become the main tool for researchers working on these insects, and some groups produce several million data points of time series each day. One critical task researchers perform is to search for patterns in these long traces [1]. For example, an expert on the beet leafhopper (*Circulifer tenellus*) who has identified the pattern that indicates the insects' probing behavior (the query Q in Figure 11) may wish to know if the less common Black-faced Leafhopper (*Graminella nigrifrons*) from a sister Genera *Graminella*⁴ also exhibits this behavior.

A survey of the literature tells us that such questions are frequently asked, and surprisingly, mostly answered by a *manual search* of vast archives [1]. Some attempts have been made to computerize the search process, but it has been noted that "...it still remains difficult to distinguish some waveform features and especially, in establishing the accurate time of transitions between some subsequent waveforms..." [14] (our emphasis). This variability in time of transition is of course simply *warping* of the time between (sub) behaviors. In Figure 12 we show a zoom-in of Figure 11 which clearly shows that pattern matching in this domain requires time warping.

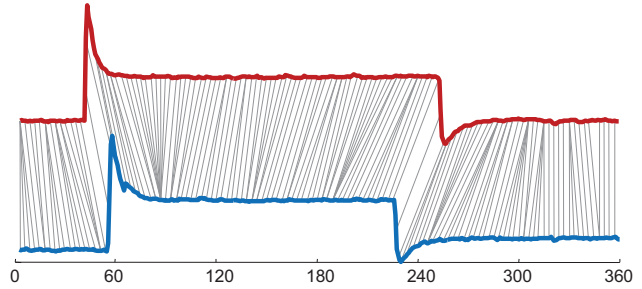


Figure 12: A zoom-in from Figure 11 showing that a query behavior (*bottom*) corresponding to stylet-insertion matches to a subsequence of a different insect's behavior only after warping.

To test the scalability of our ideas in this domain, we searched for the query shown in Figure 11 and Figure 12 in an EPG trace of length 1,499,000. The GPU method took 80.39 seconds, tenable given that collecting this trace took about one day of effort. The FPGA took only 2.24 seconds, allowing the entomologist the luxury of interactive querying. The classic software solution took 48.8 minutes, which is simply so slow that it is never likely to be used by an entomologist. Recall that these times are for a single query-by-content search; however, many higher level data mining algorithms, such as motif discovery [19], anomaly detection[4], density estimation, etc. require multiple searches as a subroutine. If we needed to do 1,000 such searches, the FPGA would require 37 minutes, but the classic software solution would take 34 days!

B. Case Study in Astronomy

A star light curve is a graph which shows the brightness of a stellar object over a period of time. Figure 13 shows two examples. The reasons why the stars change their (actual or apparent) brightness include planetary transits, self-occluding binary systems, cataclysmic or explosive events (nova or supernova) and unknown reasons. The study of light curves has led to the discovery of pulsars, extra-solar planets, supernovae, the rate of expansion of the universe, etc. [24][26].

It is difficult to overstate how many star light curves exist. Looking backwards, there are over 500,000 glass photographic plates at Harvard University that were exposed between 1885 and 1993 [38]. These are currently being digitized to yield millions of light curves. Looking ahead, this year sees work starting on the Large Synoptic Survey Telescope (LSST), a wide-field "survey" reflecting telescope that will photograph the available sky every three nights. It is estimated that LSST will produce billions of light curves in the next decade.

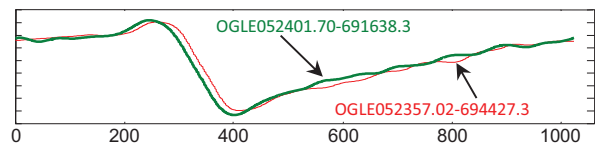


Figure 13: Two star light curves that happen to be very similar. If we happen to know the class label of one, we may predict that the other is in the same class.

⁴ Both insects are in the *Deltocephalinae* subfamily

As both old and new light curves come online, an obvious thing to do is to classify them [24]. Astronomers do have a large number of classified light curves; in some cases they can obtain what is effectively *ground truth* by obtaining extra features for relatively close stars.

While it is possible to extract a single light curve cycle, there is no well-defined starting point. Astronomers have an algorithm called *universal phasing* to produce a canonical alignment for light curves, but bemoan the fact that this is “...an operation that scales poorly to massive data sets”. However, as we shall see, in addition to poorly scaling, the universal phasing algorithm does not work as well as astronomers believe.

We obtained a three-class star light curve dataset which had been universally phased by astronomers at Time Series Center at Harvard University. Because we wanted some experiments with classic CPUs to finish, we created a test set with just 128 objects, and a training set of 1024 objects. Each light curve was normalized to have a length of 1024 (as is the practice in astronomy).

We measured the accuracy of Euclidean distance and DTW, obtaining accuracies of 80.47% and 86.72%, respectively. This tells us that “warping” is useful in this domain, something that had been at least suspected before [34]. However, rather than stopping here, we decided to test the universal phasing assumption. Suppose we ignored it and tested DTW for all possible alignments/shifts. To our knowledge this has never been attempted before, presumably because the rotation invariance version of DTW (rDTW) is $O(n^3)$, which is untenable for a normal CPU. After testing the rotation-invariant versions of both Euclidean distance and DTW, we found that the accuracies jumped to 81.2% and 91.4%, respectively. Clearly, universal phasing does not produce perfect alignments.

Why has this not been noted before? Our entire rDTW experiment took about 22.7 hours using the GPU, but on conventional hardware it would take 1.4 months.

TABLE 3: ACCURACIES AND TIMINGS FOR CLASSIFYING 128 LIGHT CURVES AGAINST A TRAINING SET OF 1,024, WITH ALL SEQUENCES OF LENGTH 1,024. VALUES IN PARENTHESES ARE ESTIMATED BASED ON SHORTER EXPERIMENTS.

	Accuracy	Time FPGA	Time GPU	Time CPU
ED	80.47%	<1.0 seconds	<1.0 seconds	2.5 seconds
rED	81.25%	<1.0 seconds	55.3 seconds	43.6 minutes
DTW	86.72%	<1.0 seconds	43.6 seconds	35.4 minutes
rDTW	91.41%	9.54 minutes	22.7 hours	(42 days)

Finally, we note that there are some algorithmic “tricks” that can speed up both rED and rDTW [34]. However, since their speedup is not dependent of the hardware on which they are implemented, we ignored them here.

C. Case Study in Data Condensing

Our results in the previous section suggest what is becoming increasingly evident in the literature [7][33], that the DTW nearest neighbor classifier (DTWnnc) is

exceptionally difficult to beat. This makes our results useful and broadly applicable. However, there are some domains in which our results *appear* to have no impact. It is often useful to do classification directly on low-powered devices or on sensors. Such hardware is by definition resource limited. Surprisingly, our work can have an impact even here.

The simplest and most obvious way to mitigate the computational and memory demands of DTWnnc is to discard a large fraction of the training data. The general version of this idea is known as data editing/condensing/prototype selection/instance pruning/numerosity reduction etc. [23][31][33].

It is well known that for general numerosity reduction algorithms, if we are very careful in choosing which objects we discard, we can significantly reduce the classification time while maintaining high accuracy, in some cases actually *improving* the accuracy [31]. Accuracy can improve because a judicious algorithm will start by removing any training instances that are mislabeled or are simply outliers that tend to misclassify new instances.

While there are many generic algorithms for numerosity reduction, Xi et al. explicitly considers a version for DTWnnc, called AWARD [33]. On the Two-Pattern dataset, which consists of 1,000 training objects of length 128, they report that AWARD takes 2.31 days to finish. This is a long time, but remarkably they show that their algorithm can throw away about 90% of the data while maintaining the same accuracy (had they thrown away a *random* 90% of the data the accuracy would have decreased significantly [33]).

These results seem very useful, but have had little impact in the literature. We suspect that part of the reason is the relative lethargy of AWARD. Can the results in this work improve the situation? The AWARD algorithm essentially has to do leaving-one-out K times with (slightly different subsets of) the K objects in the training set. This means it must do DTW $O(K^3)$ times.

We re-implemented that algorithm and found that it took 1.49 days. This is faster than the 2.31 days reported in [33], but this slight difference presumably just reflects four years of progress in CPU performance. However, for a GPU the algorithm only takes 50 minutes, and for the FPGA it takes a mere 4.3 minutes.

It is important to recognize that these results promise more than just the speedup of some existing sluggish algorithms. The lead author of [31] bemoaned to us that it was difficult to design the search algorithm that eventually became AWARD [32]. Every attempt at adjusting the search operators required several days to evaluate. Perhaps, if he had the ability to test the algorithm in minutes instead, he could have found an even better numerosity reduction strategy. We leave such considerations to future work.

VIII. CONCLUSION

We have shown that while subsequence similarity search is an important problem that has attracted great interest, the current software based solutions cannot provide adequate speed to tackle important problems in many diverse domains. We have shown that hardware based solutions offers the only realistic chance to bridge the gap between current academic

state-of-the-art and pragmatic needs of practitioners. We have designed (*correctly normalizing*) DTW similarity search algorithms for both GPUs and FPGAs and placed all code in the public domain [39]. We have demonstrated with three detailed case studies that our algorithms can expand the purview of DTW to problem previously thought untenable.

Future work includes revisiting current algorithms that use DTW as a subroutine to see if they can be improved in the light of a thousand-fold decrease in cost.

ACKNOWLEDGMENT

Dr. Najjar's work was funded by NSF CCF0905509 and CCF0811416. Dr. Keogh's work was funded by NSF 0803410 and 808770.

REFERENCES

- [1] Athitsos, V., Papapetrou, P., Potamias, M., Kollios, G. and Gunopulos, D. *Approximate Embedding-based Subsequence Matching of Time Series*. SIGMOD Conference 2008: 365-378.
- [2] Backus, E. A., and Bennett, W. H. *The AC-DC Correlation Monitor: New EPG Design with Flexible Input Resistors to Detect Both R and emf Components for any Piercing-sucking Hemipteran*. Journal of Insect Physiology. 55: 869-884.
- [3] Buyukkurt B. and Najjar W. *Compiler Generated Systolic Arrays For Wavefront Algorithm Acceleration on FPGAs*. International Conference on Field Programmable Logic and Applications, 2008: 655-658.
- [4] Chandola, V., Cheboli, D. and Kumar, V. *Detecting Anomalies in a Time series Database*. CS Technical Report 09-004, January 2009, Computer Science Department, University of Minnesota.
- [5] Chen, Y., Chen, G., Chen, K. and Ooi, B. C. *Efficient Processing of Warping Time Series Join of Motion Capture Data*. ICDE 2009.
- [6] Chu, S., Narayanan, S. and Jay Kuo, C.-C. *Efficient Rotation Invariant Retrieval of Shapes using Dynamic Time Warping with Applications in Medical Databases*. In IEEE International Symposium on Computer-Based Medical Systems (CBMS), Special Track on Data Mining, 2006.
- [7] Ding, H., Trajcevski, G., Scheuermann, P., Wang, X. and Keogh, E. *Querying and Mining of Time Series Data: Experimental Comparison of Representations and Distance Measures*. VLDB 2008.
- [8] Frost, K. E. and Groves, R. L. *Seasonal Ineffectivity of Aster Leafhoppers in Carrot*. Technical Report Department of Entomology, University of Wisconsin-Madison. www.entomology.wisc.edu/vegento/xtras/proc/2009_asterLeafhoppers.pdf
- [9] Govindaraju, N. K., Lloyd, B., Dotsenko, Y., Smith, B. and Manferdelli, J. *High Performance Discrete Fourier Transforms on Graphics Processors*. Proc. of ACM/IEEE SuperComputing 2008.
- [10] He, B., Yang, K., Fang, R., Lu, M., Govindaraju, N. K., Luo, Q. and Sanderv P. *Relational Joins on Graphics Processors* Proc. of ACM SIGMOD 2008.
- [11] Keogh, E. (2002): Exact Indexing of Dynamic Time Warping. VLDB 2002: 406-417
- [12] Keogh, E. J. and Kasetty, S. *On the Need for Time Series Data Mining Benchmarks: A Survey and Empirical Demonstration*. Data Min. Knowl. Discov. 7(4): 349-371 (2003).
- [13] Keogh, E., Xi, X., Wei, L. and Ratanamahatana, C. A. (2006). *The UCR Time Series Classification/Clustering Homepage*: www.cs.ucr.edu/~eamonn/time_series_data/
- [14] Kindt, F., Joosten, N. N. and Tjallingii, W. F. *Electrical Penetration Graphs of Thrips Revised: Combining DC- and AC-EPG Signals*. Journal of Insect Physiology 52: 1-10.
- [15] Lienhart, G., Kugel A., Manner R. *Using floating-point arithmetic on FPGAs to accelerate scientific N-body simulations*. FCCM 2002.
- [16] Matsunaga, M. and Zordan, V. B. *A Dynamics-based Comparison Metric for Motion Graphs*, Computer Graphics International (CGI) 2007.
- [17] McLean, D. L. and Kinsey, M. D. *A Technique for Electronically Recording Aphid Feeding and Salivation*, Nature 202 (1964), pp. 1358-1359.
- [18] Mneimneh, M. A., Yaz, E. E., Johnson, M. T. and Povinelli, R. J. *An Adaptive Kalman Filter for Removing Baseline Wandering in ECG Signal*. Computers in Cardiology, vol. 33, pp.253-256, 2006.
- [19] Mueen, A., Keogh, E. J. and Bigdely-Shamlo, N. *Finding Time Series Motifs in Disk-Resident Data*, ICDM 2009: 367-376.
- [20] Nam, H., Lee, K. and Lee, D. *Identification of Temporal Association Rules from Time-Series Microarray Data Sets*. BMC Bioinformatics, vol. 10 (Suppl 3):S6, March 2009.
- [21] Niennattrakul, V. and Ratanamahatana, C. A. *Meaningful Subsequence Matching under Time Warping Distance for Data Stream*. PAKDD 2009: 1013-1020.
- [22] NVIDIA CUDA Programming Guide. Version 2.3. http://developer.download.nvidia.com/compute/cuda/2_3/toolkit/docs/NVIDIA_CUDA_Programming_Guide_2.3.pdf
- [23] Pekalska, E., Duin, R. P. W., and Paclik, P. (2006). *Prototype Selection for Dissimilarity-Based Classifiers*. Pattern Recognition, 39:2, pp. 189-208.
- [24] Protopapas, P., Giammarco, J. M., Faccioli, L., Struble, M. F., Dave, R. and Alcock, C. *Finding Outlier Light Curves in Catalogues of Periodic Variable Stars*. Mon. Not. R. Astron. Soc. 369(2), 677-696 (2006).
- [25] Ratanamahatana, C. and Keogh, E. J. *Three Myths about Dynamic Time Warping Data Mining*. SDM 2005.
- [26] Rebbapragada, U., Protopapas, P., Brodley, C. E. and Alcock, C. *Finding Anomalous Periodic Time Series: An Application to Catalogs of Periodic Variable Stars*. Machine Learning, Vol. 74, Issue 3, p. 281, 2009.
- [27] Sakurai, Y., Faloutsos, C. and Yamamuro, M. *Stream Monitoring under the Time Warping Distance*. ICDE 2007: 1046-1055.
- [28] Salzberg, S. L. *On Comparing Classifiers: Pitfalls to Avoid and a Recommended Approach*. Data Mining and Knowledge Discovery, 1(3), 1997.
- [29] Villarreal, J, Park, A., Najjar, W. and Halstead, R. *Designing Modular Hardware Accelerators in C With ROCCC 2.0*, in The 18th An. Int. IEEE Symp. On Field-Programmable Custom Computing Machines (FCCM), Charlotte, NC, May 2010.
- [30] Wei, L., Keogh, E. J., Van Herle, H. and Mafrá-Neto, A. *Atomic Wedgie: Efficient Query Filtering for Streaming Times Series*. ICDM 2005: 490-497.
- [31] Wilson, D. R. and Martinez, T. R. (1997). *Instance Pruning Techniques*. ICML'97, Morgan Kaufmann, pp. 403-411.
- [32] Xi, X, (2010) Personal Communication. (note to reviewers, this necessary citation of a personal communication does not violate our anonymity, Dr. Xi has agreed to preserve our anonymity during the review process)
- [33] Xi, X., Keogh, E. J., Shelton, C. R., Wei, L. and Ratanamahatana, C. A. *Fast Time Series Classification Using Numerosity Reduction*. ICML 2006: 1033-1040.
- [34] Yankov, D., Keogh, E. J., Wei, L., Xi, X. and Hodges, W. L. *Fast Best-Match Shape Searching in Rotation-Invariant Metric Spaces*. IEEE Transactions on Multimedia 10(2): 230-239 (2008).
- [35] Ye, L. and Keogh, E. J. *Time Series Shapelets: A New Primitive for Data Mining*. KDD 2009: 947-956.
- [36] Zhu, Y. and Shasha, D. *Warping Indexes with Envelope Transforms for Query by Humming*. SIGMOD Conference 2003: 181-192
- [37] <http://roccc.cs.ucr.edu>
- [38] <http://hea-www.harvard.edu/DASCH/index.php>
- [39] Supporting webpage: For code and data. http://www.cs.ucr.edu/~mueen/GPU_DTW/index.html