

Recurrence time statistics: Versatile tools for genomic DNA sequence analysis

Yinhe Cao^{1,3}, Wen-wen Tung², and J.B. Gao³

¹1026 Springfield Drive, Campbell, CA 95008, USA

²National Center for Atmospheric Research

P.O. BOX 3000, Boulder, CO 80307-3000, USA

³Department of Electrical and Computer Engineering,
University of Florida, Gainesville, FL 32611, USA

Email: {contact@biosieve.com, wwtung@ucar.edu, gao@ece.ufl.edu}

Abstract

With the completion of the human and a few model organisms' genomes, and the genomes of many other organisms waiting to be sequenced, it has become increasingly important to develop faster computational tools which are capable of easily identifying the structures and extracting features from DNA sequences. One of the more important structures in a DNA sequence is repeat-related. Often they have to be masked before protein coding regions along a DNA sequence are to be identified or redundant expressed sequence tags (ESTs) are to be sequenced. Here we report a novel recurrence time based method for sequence analysis. The method can conveniently study all kinds of periodicity and exhaustively find all repeat-related features from a genomic DNA sequence. An efficient codon index is also derived from the recurrence time statistics, which has the salient features of being largely species-independent and working well on very short sequences. Efficient codon indices are key elements of successful gene finding algorithms, and are particularly useful for determining whether a suspected EST belongs to a coding or non-coding region. We illustrate the power of the method by studying the genomes of *E. coli*, the yeast *S. cerevisiae*, the nematode worm *C. elegans*, and the human, *Homo sapiens*. Computationally, our method is very efficient. It allows us to carry out analysis of genomes on the whole genomic scale by a PC.

Introduction

The structure of human genome and genomes of other organisms is very complicated. With the completion of many different types of genomes, especially the human genome, one of the grand challenges for the future genomics research is to comprehensively identify the structural and functional components encoded in a genome [1]. Outstanding structural components include all kinds of repeat-related structures [2, 3], and periodicity and quasi-periodicity, such as period-3, which is considered to reflect codon usage [4], and period 10-11, which may be due to the alternation of hydrophobic and hydrophilic amino acids [5] and DNA bending [6]. Extracting and understanding these structural components will greatly facilitate the identification of functional components encoded in a genome, and the study of the evolutionary variations across species and the mechanisms underlying those variations. Equally or even more important, repeat-related features often have to be masked before protein coding regions along a DNA sequence are to be identified or redundant expressed sequence tags (ESTs) are to be sequenced.

More important than finding repeat-related structures in a genome is the identification of genes and other functional units along a DNA sequence. In order to be successful, a gene finding algorithm has to incorporate good indices for the protein coding regions. A few representative indices are the Codon Bias Index (CBI) [7], the Codon Adaptation Index (CAI) [8, 9], the period-3 feature of nucleotide sequence in the coding regions [10–13] and the recently proposed YZ score [14]. Each index captures certain but not all features of a DNA sequence. The strongest signal can only be obtained when one combines multiple different sources of information [15]. In order to improve the accuracy and simplify the training of existing coding-region or gene identification algorithms (see the recent review articles [16, 17] and the references therein), and to facilitate the development of new gene recognition algorithms, it would be highly desirable to find new codon indices.

Here, we propose a simple recurrence time based method, which has a number of distinct features: i) The method is most convenient for finding out large blocks of insertions or deletions. ii) Computationally it is very efficient: with a computational time proportional to $N \log N$, where N is the size of the sequence, and a memory of $6N$, it can exhaust all repeat-related and periodic or quasi-periodic structures. This feature allows us to carry out genome analysis on the entire genomic scale by a PC. iii) It is model free in the sense that it does not make any assumption about the sequences under study. Instead, it builds a representation of the sequence in such a way that all interesting sequence units are automatically extracted. (iv) The method defines an efficient codon index, which is largely species-independent and works well on very short sequences. This feature makes the method especially appealing for the study of short ESTs. Below, we shall illustrate the power of the method by extracting outstanding structures including insertion sequences (ISs), rRNA clusters, repeat genes, simple sequence repeats (SSRs), transposons, and gene and genome segmental duplications such as inter-chromosomal duplication from genome sequences. We shall also discuss the usefulness of the method for the study of the evolutionary variations across species by carefully studying mutations, insertions and deletions. Finally, we shall evaluate the effectiveness of the recurrence time based codon index by studying all of the 16 yeast chromosomes.

Databases and Methods

Databases

We have studied the DNA sequence data from the following four species: (a) *E. coli* [18], (b) the yeast *S. cerevisiae* [19], (c) the nematode worm *C. elegans* [20], (d) and the human, *Homo sapiens* [2, 21]. These DNA sequence data are available from the following URLs respectively:

E. coli: <http://www.genome.wisc.edu/sequencing/k12.htm>

Yeast: ftp://genome-ftp.stanford.edu/pub/yeast/data_download/

C. elegans: http://www.sanger.ac.uk/Projects/C_elegans/

Human: <http://www.ncbi.nlm.nih.gov/genome/guide/human/>

Except the *E. coli* genome, the other three contain gaps, since they are not yet completely sequenced. Those gaps are deleted before we compute the recurrence times from them. For the yeast *S. cerevisiae*, the sequences of chromosome 1 to chromosome 16 are joined together into a single sequence in that order.

Basic idea of recurrence time statistics

Notations. Let us denote a sequence we want to study by $S = b_1 b_2 b_3 \cdots b_N$, where N is the length of the sequence, b_i , $i = 1, \cdots, N$, are nucleotide bases. For instance, if we take

$$S_1 = \text{ACGAAAAACGATTTTAAA},$$

then $N = 18$, $b_1 = A$, $b_2 = C$, \cdots , $b_{18} = A$. Next we group a consecutive nucleotide bases of window size w together and call that a word of size w . Using maximal overlapping sliding window, we then obtain $n = N - w + 1$ such words. We associate these words with the positions of the original DNA sequence from 1 to n , i.e., $W_i = b_i b_{i+1} \cdots b_{i+w-1}$ is a word of size w associated with the position i along the DNA sequence. Two words are considered equal if all of their corresponding bases match. That is, $W_i = W_j$, if and only if $b_{i+k} = b_{j+k}$, $k = 0, \cdots, w - 1$. $S[u \rightarrow v] = b_u b_{u+1} \cdots b_v$ will denote a subsequence of S from position u to v .

Recurrence time. The recurrence time $T(i)$ of position i for a DNA sequence S is a discretized version of the recurrence times of the second type for dynamical systems introduced recently by Gao [22–24]. It is defined as follows.

Definition: The recurrence time $T(i)$ for a position i along the DNA sequence is the smallest $j - i$ such that $j > i$ and $W_j = W_i$. If no such j exists, then there is no repeat for the word W_i after position i in the sequence S , and we indicate such a situation by $T(i) = -1$.

To analyze the $T(i)$ sequence, we first filter out all those $T(i) = -1$, then denote the remaining positive integer sequence by $R(k)$, $k = 1, \dots, m$, and finally estimate the probability distribution functions for both $R(k)$ and $\log_{10} R(k)$ sequence. These two probability distribution functions are what we mean by the recurrence time statistics. The reason that we also work on $\log_{10} R(k)$ is that the largest $R(k)$ computed from a genomic DNA sequence can be extremely long, hence, it may be difficult to visualize the distribution for $R(k)$ in linear scale.

Let us take S_1 as an example. If $w = 3$, then $n = 16$, and its recurrence time series $T(i)$ is:

$$7, 7, -1, 1, 1, 10, -1, -1, -1, -1, -1, 1, -1, -1, -1, -1$$

Discarding all the -1 terms from the $T(i)$ sequence, we then get the following recurrence time $R(i)$ series:

$$7, 7, 1, 1, 10, 1$$

where $m = 6$. The motivation for introducing the above definition is that the recurrence time sequence $T(i)$, $i = 1, \dots, n$, for a DNA sequence and a completely random sequence will be very different, and that by exploiting this difference, we would be able to exhaustively identify most of the interesting features contained in a DNA sequence.

Recurrence time statistics for completely random (pseudo-DNA) sequences

In order to characterize the difference between a DNA sequence and a completely random sequence in terms of the recurrence times, we study a completely random sequence first. We have the following interesting theorem.

Theorem: Given a sequence of independent words W_i , $i = 1, \dots, n$, where there are a total of K distinct words, each occurs with probability $p = 1/K$, the probability that the recurrence time $T(i)$ being $T \geq 1$ is given by

$$P\{T(i) = T\} \propto [n - T] \cdot p \cdot [1 - p]^{(T-1)} \quad (1 \leq T < n). \quad (1)$$

Proof: It suffices to note that the probability for an arbitrary word W_i , where i is from the positions 1 to $n - T$, to repeat exactly after $T \geq 1$ positions is given by the geometrical distribution, $p \cdot [1 - p]^{(T-1)}$. Since there are a total of $n - T$ such positions or words, while each position along the sequence from 1 to $n - T$ has the same role, the total probability is then proportional to the summation of $n - T$ terms of $p \cdot [1 - p]^{(T-1)}$. This completes the the proof.

If we assume the four chemical bases A, C, T and G to occur completely randomly along a (pseudo) DNA sequence, then there are a total of 4^w words of length w , each occurs with probability $p = 4^{-w}$. Hence, the probability for a word to repeat exactly after $T \geq w$ locations is given by Eq. (1), while the distribution for the log-recurrence time $\log_{10} R(k)$ is given by

$$f(t) = C \cdot T \cdot [n - T] \cdot p \cdot [1 - p]^{(T-1)}, \quad (0 \leq t < \log_{10} n), \quad (2)$$

where $T = 10^t$, and C is a normalization constant. To prove Eq. (2), it suffices to note that $p(T)dT = f(t)dt$. Due to overlapping of adjacent words, we do not have information about the distribution of T when $T \leq w - 1$.

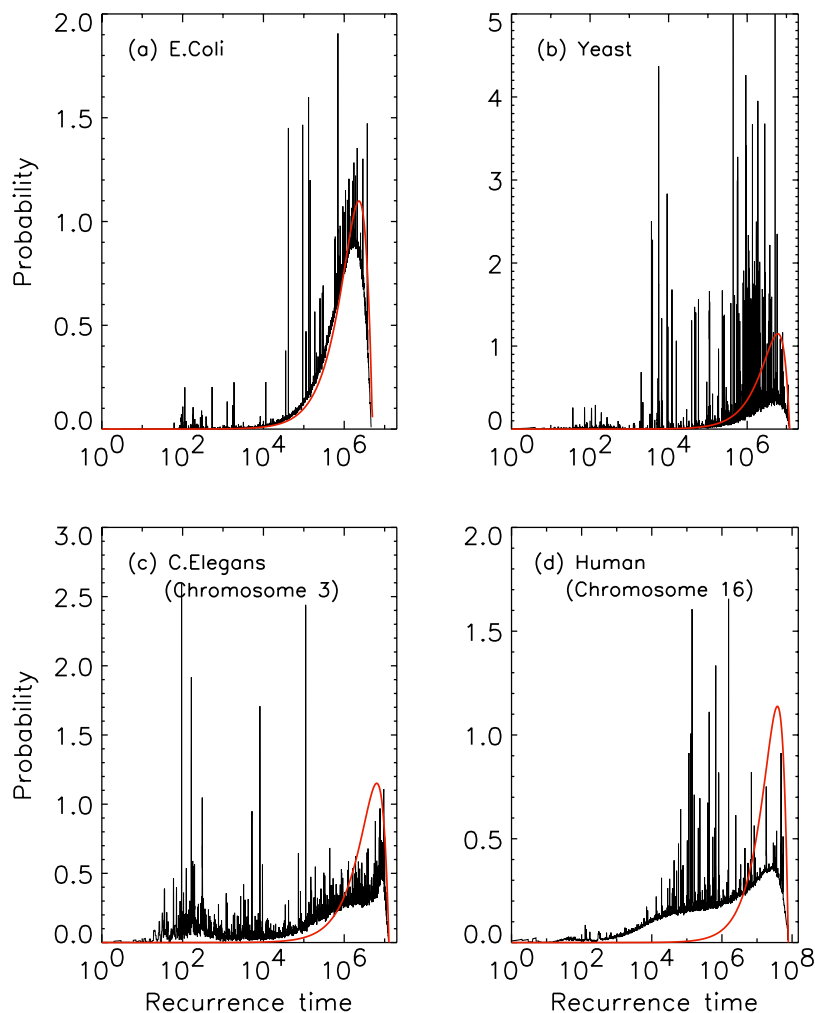


Figure 1: The probability density function (pdf) for the recurrence time $R(i)$ sequence (in log scale) computed from the DNA sequence of (a) *E. coli*, (b) the yeast *S. cerevisiae*, (c) chromosome 3 of the nematode worm *C. elegans*, and (d) chromosome 16 of the human. Red curves are computed from Eq. (2) and represent the situation where the four bases A, C, T, and G occur completely randomly with equal probability.

Recurrence time statistics for DNA sequences

(A) Recurrence time statistics and a novel codon index We have plotted in Fig. 1 the probability density functions (pdfs) of log recurrence time, i.e., $\log_{10}R(i)$, for the DNA sequence data from the four species. The red curves in Fig. 1 are computed according to Eq. (2) and represent those of completely random sequences with their length and the word size chosen to analyze them the same as those of the DNA sequences. The word sizes used are 12, 15, 16, 15 for Fig. 1 (a) to Fig. 1 (d) respectively. We observe two interesting features: (i) the pdfs for the genome sequences are very different from those for the random sequences, as signaled by the many sharp peaks in the curves of the pdfs for the genome sequences; (ii) The degree of this difference varies vastly among the four genomes studied. In fact, Eq. (2) fairly well describes the background distribution for the $\log_{10}R(i)$ sequence for *E. coli*, but very poorly describes that for the chromosome 16 of the human. This suggests that the longer a genome sequence has evolved, the more it deviates from the completely random sequence. Each sharp peak in Fig. 1 may actually represents many sharp peaks

if we plot the pdf for the $R(i)$ sequence instead of that for the $\log_{10} R(i)$ sequence. This is because with logarithmic scale, a whole interval of $R(i)$ will be lumped together. It is important to emphasize that all the sharp peaks indicate distinct features of a genome sequence. To better understand this, let us take an example. A sequence of $(A)_l$, which represents a consecutive sequence of A 's of length l , contributes to a peak at $R = 1$, if l is larger than the word length w . In fact, when $l > w$, $(A)_l$ contributes a total of $l - w$ counts to $R = 1$. Other single base repeats similarly contribute to $R = 1$. As another example, we note that a sequence such as $(AC)_l$ contributes to $R = 2$ a total of $2l - w$ counts.

We are now ready to propose a novel recurrence time based codon index. This index is based on the period-3 feature. To appreciate the idea, we have shown in Fig. 2 the probability distributions for the recurrence times not greater than 40, for the genome sequences of four species, E.Coli, Yeast, C. elegans, and the Human. The black curves are for the coding regions. The red curves in Fig. 2(b-d) are for the non-coding regions. Due to the low percentage of non-coding regions in the E.Coli genome, such a curve is not computed. We observe that the black curves all have very well defined peaks at recurrence times of 3, 6, 9, etc. Also note that the black curves are very similar among the four different species. Such period-3 feature can be conveniently used to define a codon index, which we shall denote by RT_{p3} :

$$RT_{p3} = \sum_{i=1}^m [2p(3i) - p(3i+1) - p(3i+2)] \quad (3)$$

where $p(i)$ is the probability for the recurrence time $T = i$ calculated for a coding or non-coding sequence, n is the number of bases of the coding/non-coding sequence, and m is a cutoff parameter typically chosen not to be larger than 20 so that very short sequences can be studied.

Before we go ahead to evaluate the efficiency of RT_{p3} as a codon index, let us focus on how we can exhaustively find all the repeat-related structures by tracing the peaks in Fig. 1 back to the DNA sequence. This can be easily done.

(B) Computation of exact repeat elements from recurrence times. Let $T(i)$ be the recurrence time for position i of a DNA sequence S , where $i = 1, 2, \dots, n$. For each particular value $T(1 \leq T < n)$ of the recurrence time, we build a list of indices $L(T : S)$ by linearly scanning $T(i)$ from $i = 1$ to $i = n$ and adding i to $L(T : S)$ whenever $T(i) = T$. Denote the index set $L(T : S)$ by

$$L(T : S) = \{i_1, i_2, \dots, i_C\},$$

where $T(i_k) = T$, $k = 1, 2, \dots, C$, $i_k < i_{k+1}$ for $k = 1, 2, \dots, C - 1$, and C is the count of the occurrence of T in the recurrence time series $T(i)$. If we take S_1 as an example, then

$$L(1 : S_1) = \{4, 5, 12\}, L(7 : S_1) = \{1, 2\}, L(10 : S_1) = \{6\}.$$

When the count C is larger than 1, we define the gap between two consecutive indices of $L(T : S)$ as:

$$g_k = i_{k+1} - i_k \quad (k = 1, 2, \dots, C - 1)$$

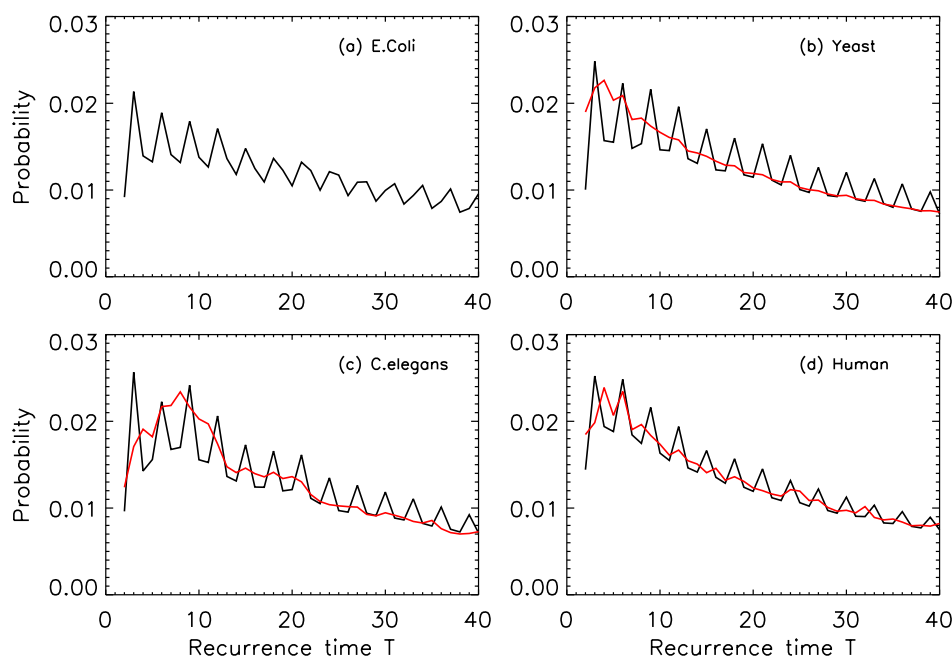


Figure 2: The probability distribution curves computed from the genomes of four organisms studied. The red and black curves are for non-coding and coding sequences, respectively. The window size w is 3 in all the computations.

Let $g^* = w$. What happens when all $g_k \leq g^*$? In this situation, the sequence segment $S_r = S[i_1 \rightarrow i_C - i_1 + w - 1]$ is an exact repeat with period T . To see this, we note that for each term i_k in $L(T : S)$, we have $W_{i_k} = W_{i_k+T}$, or more explicitly, $b_{i_k+j} = b_{i_k+T+j}$ for $j = 0, 1, \dots, w - 1$, and $k = 1, 2, \dots, C$. Hence, we can concatenate b_{i_2} at $b_{i_1+g_1}$ to combine the repeats starting from b_{i_1} and b_{i_2} . More concretely, For $k = 1$, we have $b_{i_1+j} = b_{i_1+T+j}$, $j = 0, 1, \dots, w - 1$; similarly, for $k = 2$, we have $b_{i_2+j} = b_{i_2+T+j}$, $j = 0, 1, \dots, w - 1$. Noting $g_1 \leq w$ means $i_2 - i_1 \leq w$, or $i_2 \leq i_1 + w$, we have $b_{i_1+j} = b_{i_1+T+j}$ for $j = 0, 1, \dots, i_2 - i_1 + w - 1$. Continuing this procedure till $k = C$, we have $b_{i_1+j} = b_{i_1+j+T}$ for $j = 0, 1, \dots, i_C - i_1 + w - 1$. Hence, the sequence $S_r = S[i_1 \rightarrow i_C - i_1 + w - 1]$ is an exact repeat with period T with its length $l = i_C - i_1 + w$. If $T < l$, S_r is then a tandem repeat.

In general, some gaps g_k may be larger than $g^* = w$. When there are P such gaps, we can decompose $L(T : S)$ into $P + 1$ subsets, such that within each subset all of the gaps are not larger than g^* . Then following the procedure detailed in the last paragraph, we see that each subset represents an exact repeat of period T . Taking S_1 as an example again, then from $L(7 : S_1) = \{1, 2\}$ we get an exact repeat ACGA of period 7, and from $L(1 : S_1) = \{4, 5, 12\}$, we get two exact repeats of period 1: the first one is AAAAA which is a simple sequence repeat, the other is TTTT, which is also a simple sequence repeat.

Before we move on, we emphasize that the procedure outlined here makes the recurrence time method largely independent of the word size w : any feature with length longer than w can be re-combined. For simplicity, we shall call this a **re-combination algorithm**. This algorithm, together with the features related to mutations, deletions, and insertions, which are to be discussed shortly, makes the recurrence time based method especially convenient for identifying horizontally transferred genes.

(C) Single Nucleotide Mutation and Single Nucleotide Polymorphism (SNP). Suppose we have two exactly repeating sequence segments, S_{lead} and S_{lag} , where the subscript lead and lag mean S_{lead} appears earlier than S_{lag} in a genome. In the simplest case, each word constructed from segments of S_{lead} has the same period T . In general, however, a few words constructed from segments of S_{lead} may have smaller recurrence times, due to the possibility that those words may find their copies in between S_{lead} and S_{lag} . Now suppose one nucleotide somewhere within S_{lead} is mutated. Since the mutated nucleotide appears in a consecutive w words, each of length w , we see that for those words, the period T will have to be different than T . This means if we plot out the recurrence time vs. the sequence position curve, then we should observe a gap of length w in an otherwise almost constant (T) curve. When this is the case, we can suspect that there may be a single nucleotide mutation at the end of the gap. If the gap corresponds to recurrence times larger than T or equal to -1 (meaning no repeats), then we can conclude that there is a single nucleotide mutation at the end of the gap. This is actually a sufficient condition, since it excludes the possibility that a few words may have copies in between S_{lead} and S_{lag} .

The study of single nucleotide mutation is most relevant to the study of Single Nucleotide Polymorphism (SNP), where DNA sequence variations occur when a single nucleotide (A, T, C, or G) in the genome sequence among different populations is changed, possibly due to evolution. It is clear that if we concatenate two genome sequences for different subjects together, then we can treat SNP as a special type of single nucleotide mutation.

(D) Insertion/Deletion and relations between repeat sequences of different periods. Suppose we have a sequence segment starting from the position i_a . What happens if we insert a sequence of length, say, a few thousand bases, in the middle of that segment, then let the segment with insertion to repeat somewhere in the genome? Equivalently, the original sequence segment can be considered a result of deletion from the longer (i.e., with insertion) sequence segment. This interesting situation is revealed by a jump in the recurrence time vs. sequence position plot, with the height of the jump being the length of the insertion sequence, as we explain below.

Let $T(i)$ denote the recurrence time for the word at the position i of the sequence S . Suppose $i_a < i_b$, $T(i_a) = T(i_{a+1}) = \dots = T(i_b) = T_1 > 0$, and

$$i_b < i_c \leq i_b + w < i_d, \quad T(i_c) = T(i_{c+1}) = \dots = T(i_d) = T_2 > T_1.$$

Then the sequence segment

$$S_a = S[i_a \rightarrow (i_d + w - 1)]$$

can be considered the result of deleting the sequence

$$S_{\text{deletion}} = S[(i_c + T_1 - 1) \rightarrow (i_c + T_2 - 1)]$$

from the sequence segment

$$S_b = S[(i_a + T_1) \rightarrow (i_d + w - 1 + T_2)].$$

Equivalently, S_b is the result of inserting the sequence S_{deletion} into S_a right before the position i_c . Note that the condition of $i_c \leq i_b + w$ comes from the fact that the boundary always affects a

consecutive w words, each of length w . When the first w bases of the deleted sequence segment do not have their copies at the positions starting from i_c , we have $i_c = i_b + w$. Otherwise, we have in-equalities.

Results and Discussion

Extraction of repeat-related structures

We now present examples of structures which can be found by tracing the peaks in Fig. 1 back to the genome sequences. These structures include insertion sequences (ISs), rRNA clusters, repeat genes, simple sequence repeats (SSRs), transposons, and gene and genome segmental duplications such as inter-chromosomal duplication. We shall illustrate most of these structures using the yeast *S. cerevisiae* as an example.

We first study SSRs. SSRs are perfect or slightly imperfect tandem repeats of particular k-mers. They have been extremely important in human genetic studies, because they show a high degree of length polymorphism in human population owing to frequent slippage by DNA polymerase during replication [2]. Any tandem repeat of k-mers, disregarding its exact content, will contribute to the count of occurrence of period $T = k$ in recurrence time statistics, hence can be easily found by following the peak of $T = k$ in Fig. 1. As an example, we note that there are 39 sequence segments contributing to $k = 13$. Three of them are CCACACCCACACA, GGTGTGTGGGTGT, and TACCGACGAGGCT. Note that by Fig. 1(a) we can conclude that *E. coli* has very few SSRs.

One of the more striking features of the yeast *S. cerevisiae* genome is that it contains many copies of transposon yeast (**Ty**) elements. Each **Ty** element is about 6.3 kb long, with the last 330 bp at each end constituting direct repeats, called δ . Those direct end repeats are responsible for the peaks around 5500 in Fig. 1(b), which enable us to find all of those **Ty** elements on both strands of the genome. As two examples, we mention that the transposon Ty3-1 on the Watson strand of chromosome 7 starts at the position 707196 and ends at 712546, and has a period of 5011. The transposon Ty1-1 on the Crick strand of chromosome 1 starts at the position 166162 and ends at 160238, and has a period of 5588.

Gene duplication is an important source of evolutionary novelty. Many duplicate genes have been found in the yeast *S. cerevisiae* genome, and they often seem to be phenotypically redundant [25–27]. Any gene duplication will contribute to one of the sharp peaks in Fig. 1(b). As an example, we note that a gene (standard name MCH2, systematic name YKL221W), which is on chromosome 6 starting from the position 6931, is repeated on chromosome 13, starting from the position 7749.

Genome segmental duplications consist of large blocks that have been copied from one region of the genome to another. They have been found among genomes of many species including the yeast *S. cerevisiae* [25], and the *Homo sapiens* [2, 28]. In fact, they contribute to some of the the sharpest peaks in Fig. 1. An example of such segmental duplications is the inter-chromosomal duplication corresponding to the peak at $T = 5150433$ in Fig. 1(b).

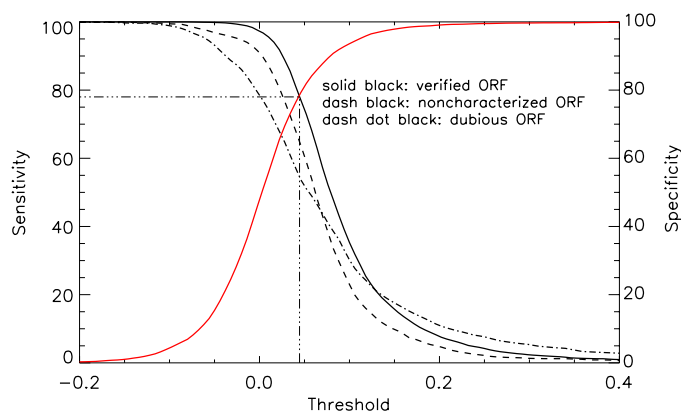


Figure 3: The specificity and sensitivity curves for the RT_{p3} index evaluated on all of the 16 yeast chromosomes.

Evaluation of the novel codon index

In order to evaluate the effectiveness of the RT_{p3} as a codon index, we study all of the 16 yeast chromosomes. Our sample pool is comprised of two sets of DNA segments: the coding set (fully coding regions or exons), which contains 4125 verified ORFs, 1626 uncharacterized ORFs, and 812 dubious ORFs, and the non-coding set, which contains 5993 segments (fully non-coding regions or introns). Some of these coding and non-coding segments are very short. Regardless of their length, each segment is counted as one when calculating the sensitivity and specificity curves. Fig. 3 shows the specificity and sensitivity curves for all of the 16 yeast chromosomes, where the red curve is the cumulative distribution function for RT_{p3} for the non-coding regions, and the black curves are the complementary cumulative distribution function for the coding regions, where for clarity, we have computed such distributions for verified ORFs, uncharacterized ORFs, and dubious ORFs, separately. To understand the meaning of such curves, let us focus on the intersection of the solid black curve and the red curve. When we choose that RT_{p3_0} as a threshold value, then with 78% probability a coding sequence is characterized as coding sequence, while with 78% probability a non-coding sequence is also taken as a non-coding sequence. As expected, this percentage is lower for uncharacterized and dubious ORFs. It is interesting to note that the percentage of accuracy calculated on Human genomes is around 74%, close to 78%. Because of this (see also Fig. 2), we conclude that the method is largely species-independent.

It is interesting to note that the period-3 feature is often quantified by performing the Fourier spectral analysis on fairly long DNA sequences. In order to make such analysis applicable to sequences as short as 162 bases, recently a lengthen-shuffle algorithm is proposed [11]. Fourier spectral analysis together with the lengthen-shuffle algorithm gives about 69% of sensitivity and specificity when evaluated on a prokaryote genome, the *V.cholerae* chromosome I, and about 61% when evaluated on eukaryotic genomes [12]. It is clear that the RT_{p3} index is more accurate. Other features of the recurrence time based method are: (i) DNA sequences as short as 40 bases can be very well studied. Noting that an expressed sequence tag (EST) is usually very short and that little may be known about the genome to which the EST belongs, this feature, together with the species-independent one, makes the method particularly useful for determining whether a suspected EST belongs to a coding or non-coding region. (ii) The method directly works on the DNA sequence.

In contrast, numerical sequences have to be obtained by certain mapping rules in order to use the Fourier spectral analysis based methods.

Discussion

In this paper, we have proposed a simple recurrence time based method for DNA sequence analysis, and shown that the method can conveniently exhaust all repeat-related structures of length greater than an arbitrarily chosen small word of size w in a genome. We have also shown that the method is very convenient for the study of mutations, insertions and deletions, hence, it holds great potential for the study of evolutionary variations across species and the mechanisms underlying it. By characterizing the peaks at multiples of 3, we have defined a very efficient codon index which is largely species independent and works well on very short sequences. We emphasize that one of the more appealing features of RT_{p3} as a codon index is that no priori knowledge about the sequence is used. Hence, the method will be especially convenient for the study of genome sequences that very little is known. This is the case, for example, when a genome sequence is to be sequenced by a few small research groups by studying expressed sequence tags (ESTs).

While the accuracy of 78% for the yeast genome is already satisfactory, we note that it is possible to improve this percentage by designing other indices from the recurrence times. Readers interested in this issue are encouraged to contact the authors for the raw recurrence time data.

References

- [1] Collins, F.S., Green, E.D., Guttmacher, A.E., & Guyer, M.S. (2003) A vision for the future of genomics research, *Nature* **422** (6934): 835-847.
- [2] International Human Genome Sequencing Consortium, Initial sequencing and analysis of the human genome (2001) *Nature* **409** 860-921.
- [3] Jurka, J. (1998) Repeats in genomic DNA: mining and meaning. *Curr Opin Struct Biol* **8** 333-337.
- [4] Guigo, R. (1999) DNA Composition, Codon Usage and Exon Prediction, in *Genetics Databases*, ed. Bishop, M.J. (San Diego, CA : Academic Press, 1999), pp. 53-80.
- [5] Herzel, H. Weiss, D., & Trifonov, E.N. (1999) 10-11 bp periodicities in complete genomes reflect protein structure and DNA folding, *Bioinformatics* **15** (3): 187-193.
- [6] Fukushima, A., Ikemura, T., Kinouchi, M., et al. (2002) Periodicity in prokaryotic and eukaryotic genomes identified by power spectrum analysis, *Gene* **300** (1-2): 203-211.
- [7] Bennetzen, J.L. and Hall, B.D. (1982) Codon selection in yeast, *J. Biol. Chem.*, **257**, 3026-3031.
- [8] Sharp, P.M. and Li, W.-H. (1987) The codon Adaptation Index—a measure of directional synonymous codon usage bias, and its potential applications, *Nucleic Acids Res.*, **15**, 1281-1295.

- [9] Jansen, R., Bussemaker, H.J. and Gerstein, M. (2003) Revisiting the codon adaptation index from a whole-genome perspective: analyzing the relationship between gene expression and codon occurrence in yeast using a variety of models, *Nucleic Acids Res.*, **31**, 2242-2251.
- [10] Tiwari, S., Ramachandran, S., Bhattacharya, A., Bhattacharya, S. and Ramaswamy, R. (1997) Prediction of probable genes by Fourier analysis of genomic sequences, *Comput. Appl. Biosci.*, **13**, 263-270.
- [11] Yan, M., Lin, Z.S. and Zhang, C.T. (1998) A new Fourier transform approach for protein coding measure based on the format of the Z curve, *Bioinformatics*, **14**, 685-690.
- [12] Issac, B., Singh, H. and Kaur, H. (2002) Locating probable genes using Fourier transform approach. *Bioinformatics*, **18**, 196-197.
- [13] Kotlar, D. and Lavner, Y. (2003) Gene prediction by spectral rotation measure: A new method for identifying protein-coding regions, *Genome Res.*, **13**, 1930-1937.
- [14] Zhang, C.T. and Wang, J. (2000) Recognition of protein coding genes in the yeast genome at better than 95% accuracy based on the Z curve, *Nucleic Acids Res.*, **28**, 2804-2814.
- [15] Snyder, M. and Gerstein, M. (2003) Genomics - Defining genes in the genomics era, *Science*, **300**, 258-260.
- [16] Fickett, J.W. and Guig, R., Computational gene identification (1996) In "Internet for the Molecular Biologist", Swindell, S., Miller, R. and Myers, G. (eds.), 73-100, Horizon Scientific Press, Wymondham, UK.
- [17] Zhang, M.Q. (2002) Computational prediction of eukaryotic protein-coding genes, *Nat. Rev. Genet.*, **3**, 698-709.
- [18] Blattner, F.R., et al. (1997) The complete genome sequence of Escherichia coli K-12. *Science* **277** 1453-1474.
- [19] Mewes, H.W., et al. (1997) Overview of the yeast genome. *Nature* **387** 7-8.
- [20] The C. elegans Sequencing Consortium, Genome Sequence of the Nematode Caenorhabditis elegans-A Platform for Investigating Biology (1998) *Science* **282** 2012-2018.
- [21] The Celera Genomics Sequencing Team, The sequence of the human genome (2001) *Science* **291** 1304-1351.
- [22] Gao, J.B. (1999) Recurrence Time Statistics for Chaotic Systems and Their Applications. *Phys. Rev. Lett.* **83**, 3178-3181.
- [23] Gao, J.B. & Cai, H.Q. (2000) On the structures and quantification of recurrence plots. *Phys. Lett. A* **270**, 75-87.
- [24] Gao, J.G. (2001) Detecting nonstationarity and state transitions in a time series. *Phys. Rev. E* **63**, 066202.
- [25] Wolfe, K.H. & Shields, D.C. (1997) Molecular evidence for an ancient duplication of the entire yeast genome. *Nature* **387**, 708-13.
- [26] Seoighe, C. & Wolfe, K.H. (1999) Updated map of duplicated regions in the yeast genome. *Gene* **1** 253-261.
- [27] Glaever, G., et al. (2002) Functional profiling of the *Saccharomyces cerevisiae* genome. *Nature* **418** 387-391.
- [28] Brendan, J., et al. (1999) Genome duplications and other features in 12 Mb of DNA sequence from human chromosome 16p and 16q. *Genomics* **60** 295-308.