

Long-range power-law correlations in condensed matter physics and biophysics

H.E. Stanley^a, S.V. Buldyrev^a, A.L. Goldberger^b, S. Havlin^{a,c},
C.-K. Peng^a and M. Simons^{b,d}

^a*Center for Polymer Studies and Department of Physics, Boston University,
Boston, MA 02215, USA*

^b*Cardiovascular Division, Harvard Medical School, Beth Israel Hospital,
Boston, MA 02215, USA*

^c*Department of Physics, Bar-Ilan University, Ramat-Gan 52900, Israel*

^d*Department of Biology, MIT, Cambridge, MA 02139, USA*

We discuss the appearance of long-range power-law correlations in various systems of interest to condensed matter physicists and biophysicists, with emphasis on the recent discovery of long-range correlations in DNA sequences that contain non-coding regions.

1. Introduction

For what basic physics advances will the twentieth century be remembered? Certainly the first half will be known principally for the discovery of quantum mechanics. The second half witnessed the developed of a myriad of applications of quantum mechanics, without which much of everyday life would not be recognizable. But what are the *basic* advances in fundamental understanding of the workings of nature?

Here, we shall exemplify one such basic advance – the discovery of long-range power-law correlations in a remarkably wide variety of systems. Such long-range power-law correlations are a physical fact that in turn gives rise to the increasingly appreciated “fractal geometry of nature” [1–6]. So if fractals are indeed so widespread, it makes sense to anticipate that long-range power-law correlations may be similarly widespread. Indeed, recognizing the ubiquity of long-range power-law correlations can help us in our efforts to understand nature, since as soon as we find power-law correlations we can quantify these with a critical exponent (called α in this paper). Quantification of different behavior allows us to recognize similarities between different systems, thereby eventually leading to recognizing underlying unifications that might otherwise have gone unnoticed. For example, as soon as phenomena occurring near critical points were quantified with critical exponents, it was recognized that

the entire “zoo” of critical phenomena partitioned itself neatly into a relatively small number of distinct “universality classes”.

Our intuition tells us that correlations should decay exponentially, not as power laws. Consider, e.g., a set of two-state Ising spins in dimension one ($d = 1$) with interactions J between neighbors. If C_1 denotes the correlation function for two spins that are nearest neighbors, then intuition tells us that the correlation function for any two spins separated by a distance r is [7]

$$C_r = (C_1)^r = e^{-r/\xi}, \quad (1a)$$

where the second equality in (1a) serves to define the correlation length $\xi \equiv -1/\log C_1$.

For $d = 1$ site percolation, C_r denotes the pair connectedness, the probability that a site at position r is both occupied and also connected by a string of occupied sites to an occupied site at the origin [8]. Again, (1a) holds, but now with $C_1 = p$, the probability that a site is occupied.

Our simple intuition, that correlations decay exponentially because of the fashion in which order is “propagated”, seems to always work – except at the critical point, where the exponential decay of (1a) gives over to a power law decay

$$C_r \sim (1/r)^{d-2+\eta}. \quad (1b)$$

The difference between (1a) and (1b) is profound: (1a) states that there is a characteristic length ξ fixed by the strength of the nearest-neighbor correlation C_1 , while (1b) states that there is no characteristic length at all.

Can we intuitively understand how it is possible to find a *non-exponential* decay of correlations? At first sight, it might appear that whenever we increment the distance between two spins by one lattice constant, the correlation should decrease by roughly the same factor, but this intuition leads immediately to exponential decay. A possible resolution to this paradox stems from the fact that near a critical point, “information” propagates from a spin at the origin to a spin at position r *not via a single path* (as for $d = 1$), but *rather via an infinite number of paths*; some of these paths are explicitly enumerated in fig. 9.4 of ref. [7]. Ornstein and Zernike [9] recognized this fact, but approximated the fashion in which “order is propagated” and so obtained predictions that today we call “classical” (fig. 7.5 of ref. [7]). Exact enumeration methods, such as high-temperature series expansions, take into account exactly such paths up to a certain length k_{\max} , where k_{\max} is typically 20. To obtain power-law correlations, the exact results for $k < 20$ are extrapolated to obtain an estimate of the behavior for all k . In some sense, although the correlation along each path *decreases* exponentially with the length of the path,

the number of such path *increases* exponentially. Therefore, the net effect is that there arise longer range power-law correlations.

At one time, it was imagined that the “scale-free” case of (1b) was relevant to only a fairly narrow slice of physical phenomena – only to systems that had been “tuned” by exceedingly painstaking experimental work to be exactly at a critical point [7]. Now we appreciate the ubiquity of systems displaying scale-invariant behavior. First of all, any system examined on length scales smaller than the correlation length is likely to display power-law behavior (because all paths between the origin and r are relevant up to the correlation length, and these cancel out the exponential decay for $r < \xi$). Moreover, the number and nature of systems displaying power-law correlations has increased dramatically, including systems that no one might ever have suspected as falling under the umbrella of “critical phenomena”. The latter part of the century has witnessed a veritable expulsion in the study, both experimental and theoretical, of such systems. The 1991 Nobel Prize was awarded to P.-G. de Gennes in part for his recognition that polymer systems behave analogously to systems near their critical points. The 1993 Wolf Prize will be awarded to Benoit Mandelbrot for the recognition of the “fractal geometry of nature”. Another very prestigious Israeli prize, the 1993 Israel Prize, is being awarded this year to Shlomo Alexander, in large part for his discoveries that under appropriate conditions a wide range of systems obey scaling or scale invariance.

Indeed, many systems drive themselves spontaneously toward critical points. One of the simplest systems exhibiting such “self-organized criticality” [10] is invasion percolation, a generic model that has recently found applicability to describing anomalous behavior of rough interfaces [11]. Instead of occupying all sites with random numbers below a pre-set parameter p , in invasion percolation one “grows” the incipient infinite cluster right at the percolation threshold by the trick of occupying always the perimeter site whose random number is smallest. Thus small clusters are certainly not scale-invariant and in fact contain sites with a wide distribution of random numbers. As the mass of the clusters increases, the cluster becomes closer and closer to being scale invariant or “fractal”. One says that such a system drives itself to a “self-organized critical state” [10].

The list of systems in which power law correlations appear has grown rapidly in recent years, including models of turbulence and even earthquakes [12]. What do we anticipate for biological systems? Generally speaking, when “entropy wins over energy” – i.e., randomness dominates the behavior – we find power laws and scale invariance. Biological systems sometimes are described in language that makes one think of a Swiss watch. Mechanistic or

“Rube Goldberg” descriptions must in some sense be incomplete, since it is only some appropriately chosen averages that appear to behave in a regular fashion. The trajectory of each individual biological molecule is of necessity random – albeit correlated. Thus one might hope that recent advances in understanding “correlated randomness” [13–16] could be relevant to biological phenomena. While there have been reports of scale invariant phenomena in isolated biological systems – ranging from the fractal shapes of neurons [17] to long-range correlations in heart beat intervals [18], human writings [19], and the stock market [20] – there has been no systematic study of *biological* system that displays power-law correlations.

Here we will attempt to summarize the key findings of some recent work [21–40] suggesting that under suitable conditions – the sequence of base pairs or “nucleotides” in DNA also displays power-law correlations. The underlying basis of such power-law correlations is not understood at present, but it is least possible that this reason is of as fundamental importance as it is in other systems in nature that have been found to display power-law correlations.

2. Discovery of long-range correlations in DNA

In order to study the scale-invariant long-range correlations of a DNA sequence, we first introduced a graphical representation of DNA sequences, which we term a *fractal landscape* or *DNA walk* [21]. For the conventional one-dimensional random walk model [41], a walker moves either “up” [$u(i) = +1$] or “down” [$u(i) = -1$] one unit length for each step i of the walk [1]. For the case of an uncorrelated walk, the direction of each step is independent of the previous steps. For the case of a correlated random walk, the direction of each step depends on the history (“memory”) of the walker [14–16].

One definition of the DNA walk is that the walker steps “up” [$u(i) = +1$] if a pyrimidine (C or T) occurs at position a linear distance i along the DNA chain, while the walker steps “down” [$u(i) = -1$] if a purine (A or G) occurs at position i . Other definitions are discussed in the caption to fig. 7. The question we asked was whether such a walk displays only short-range correlations (as in an n -step Markov chain) or long-range correlations (as in critical phenomena and other scale-free “fractal” phenomena).

The DNA walk provides a graphical representation for each gene and permits the degree of correlation in the base pair sequence to be directly visualized, as in fig. 1. Fig. 1 naturally motivates a quantification of this correlation by calculating the “net displacement” of the walker after l steps, which is the sum of the unit steps $u(i)$ for each step i . Thus $y(l) \equiv \sum_{i=1}^l u(i)$.

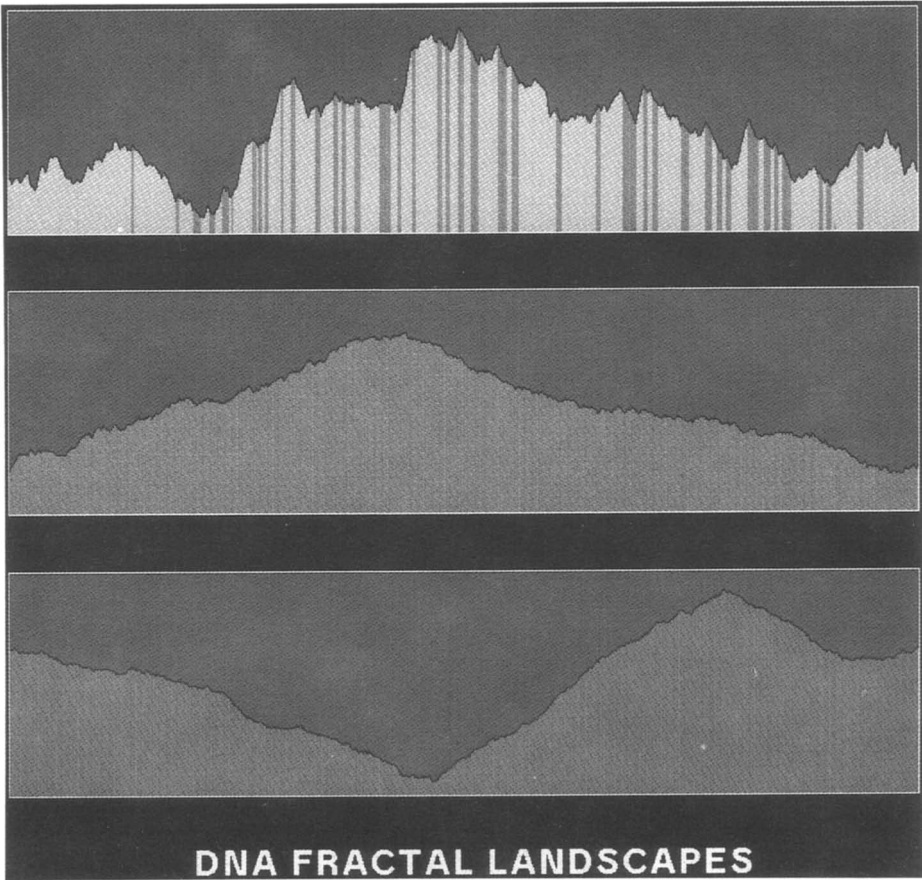


Fig. 1. The DNA walk representations of (top) human β -cardiac myosin heavy chain gene sequence, showing the coding regions as vertical dark bars, (middle) the spliced together coding regions, and (bottom) the bacteriophage lambda DNA which contains only coding regions. Note the more complex fluctuations for (top) compared with the coding sequences (middle) and (bottom). We found that for almost all coding sequences studied that there appear regions with one strand bias, followed by regions of a different strand bias. The fluctuation on either side of the overall strand bias we found to be random, a fact that is plausible by visual inspection of the DNA walk representations. We used different step heights for purine and pyrimide in order to align the end point with the starting point. This procedure is for graphical display purposes only (to allow one to visualize the fluctuations more easily) and is not used in any analytic calculations.

An important statistical quantity characterizing any walk [41] is the root mean square fluctuation $F(l)$ about the average of the displacement; $F(l)$ is defined in terms of the difference between the average of the square and the square of the average,

$$F^2(l) = [\overline{\Delta y(l) - \overline{\Delta y(l)}}]^2 = \overline{[\Delta y(l)]^2} - \overline{\Delta y(l)}^2, \quad (2)$$

of a quantity $\Delta y(l)$ defined by $\Delta y(l) \equiv y(l_0 + l) - y(l_0)$. Here the bars indicate an *average* over all positions l_0 in the gene. Operationally, this is equivalent to (a) taking a set of calipers set for a fixed distance l , (b) moving the beginning point sequentially from $l_0 = 1$ to $l_0 = 2, \dots$ and (c) calculating the quantity $\Delta y(l)$ (and its square) for each value of l_0 , and (d) averaging all of the calculated quantities to obtain $F^2(l)$.

The mean square fluctuation is related to the auto-correlation function $C(l) \equiv \overline{u(l_0) u(l_0 + l)} - \overline{u(l_0)}^2$ through the relation: $F^2(l) = \sum_{i=1}^l \sum_{j=1}^l C(j-i)$. The calculation of $F(l)$ can distinguish three possible types of behavior. (i) If the base pair sequence were random, then $C(l)$ would be zero on average [except $C(0) = 1$], so $F(l) \sim l^{1/2}$ (as expected for a *normal* random walk). (ii) If there were a local correlation extending up to a characteristic range R (such as in Markov chains), then $C(l) \sim \exp(-l/R)$; nonetheless *the asymptotic behavior* $F(l) \sim l^{1/2}$ *would be unchanged from the purely random case*. (iii) If there is no characteristic length (i.e., if the correlation were “infinite-range”), then the scaling property of $C(l)$ would not be exponential, but would most likely to be a power-law function, and the fluctuations will also be described by a power law

$$F(l) \sim l^\alpha \quad (2a)$$

with $\alpha \neq 1/2$. Fig 1(top) shows a typical example of a gene that contains a significant fraction of base pairs that do *not* code for amino acids [42–44]. It is immediately apparent that the DNA walk has an extremely jagged contour, which we shall see corresponds to long-range correlations. Fig. 2 shows double logarithmic plots of the mean square fluctuation function $F(l)$ as a function of the linear distance l along the DNA chain for the three *randomly chosen* sub-sequences (1000 base pairs of each) from fig. 1(top).

The fact that the data are linear on this double logarithmic plot confirms that $F(l) \sim l^\alpha$. A least-squares fit produces a straight line with slope α substantially larger than the prediction for an uncorrelated walk, $\alpha = 1/2$, thus providing direct experimental evidence for the result that there exists long-range correlation.

Peng et al. also addressed the question of whether the long-range correlation properties are different for coding and non-coding regions of a DNA sequence [21], a point that is currently the subject of some continuing debate [24,28]. Fig. 1(middle) shows the DNA walk for a sequence formed by splicing together the coding regions of the DNA sequence of this same gene. Fig. 1(bottom) displays the DNA walk for a typical sequence with only coding regions. In contrast to fig. 1(top), these coding sequences have less jagged contours, suggesting a shorter-range correlation.

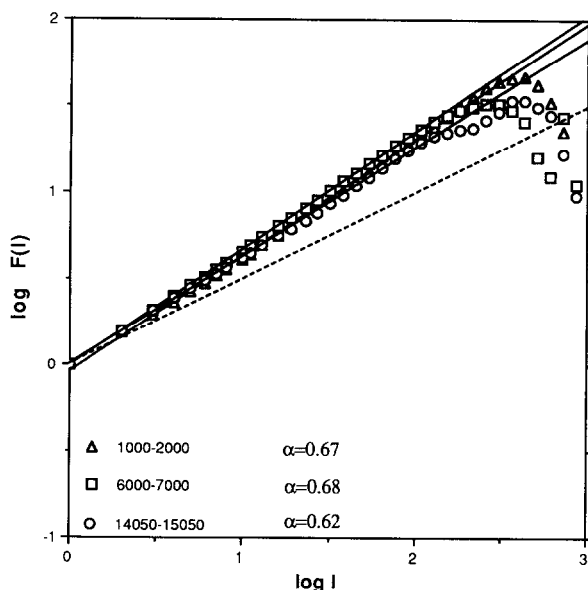


Fig. 2. Double logarithmic plots of the mean square fluctuation function $F(l)$ as a function of the linear distance l along the DNA chain for three *randomly chosen* sub-sequences (1000 base pairs of each) from fig. 1a. The dashed line has slope 0.5, corresponding to the expectation if the correlations were only short-range.

To analyze the middle and bottom parts, we first observe that for almost all sequence that we studied, purine-rich regions (compared to the average concentration over the entire strand) alternate with pyrimidine-rich regions, corresponding to the “up-hill” and “down-hill” portions of the DNA walk. To take into account the fact that the concentrations of purines and pyrimidines are not constant throughout the single strand base pair sequence, each DNA walk representation is partitioned into three segments demarcated by the global maximum (“max”) and minimum (“min”) displacements. Then we analyze the fluctuation within each segment separately. We found that for the middle and bottom parts that $\alpha = 1/2$ to within the level of fluctuation associated with the finite length of chain analyzed.

3. Possible artifacts

Naturally, we worried constantly that there was some possible “artifact” in the analysis that would invalidate our finding that spliced together coding regions as well as sequences containing only coding regions are uncorrelated, while sequences containing non-coding “junk” possess long-range power-law

correlations. Hence we carried out numerous tests, some of which are reported on below. Since the calculation of $F(l)$ for the DNA walk representation thus has the potential of providing a new, *quantitative* method to distinguish coding and non-coding regions, it is particularly important to be certain that there are no artifacts of this method.

3.1. Sampling statistics

In order to see if this scaling behavior is “universal”, we first applied our analysis to more than 100 representative DNA sequences across the phylogenetic spectrum (comprising altogether some 10^7 base pairs analyzed – by contrast, Voss [24] has confirmed our findings using 25 000 DNA sequences). The result of some of this analysis is provided in table 1 of ref. [21]. The results confirm that long-range correlations ($\alpha > 1/2$) are characteristic of DNA containing non-coding base pairs but for coding sequences, $\alpha \cong 0.50 \pm 0.05$.

3.2. Biased but uncorrelated walks

One of the first concerns that we met in presenting our work to others was the confusion that for a biased but uncorrelated random walk, α would also be larger than $1/2$. Much DNA material contains regions that have more purines (or pyrimidines) than 50%; this phenomenon is termed “strand bias”. Therefore, we studied a variety of “artificial” base sequences in which we deliberately introduce a controlled measure of strand bias. These artificial sequences nonetheless all have $\alpha = 1/2$.

To demonstrate this fact explicitly, we constructed fig. 3. Fig. 3a shows an *unbiased* random walk with exactly the same number of steps, 2941, as in the case of the same gene analyzed in fig. 4. The data clearly corroborate the expected result, $\alpha = 1/2$. Fig. 3b shows a 2941-step *biased* random walk, and again the data clearly corroborate the expected result, $\alpha = 1/2$. Fig. 3c shows a 2941-step *correlated* random walk, with correlation parameter 0.61, and now the data corroborate the expected result, $\alpha = 0.61$. Thus, long-range correlations bear no relation whatsoever to strand bias: *The exponent is determined by the correlations, not by the bias*. What can in fact produce artifacts in estimating α is the abrupt change of bias, which is discussed in section 4 below.

3.3. Effect of finite sequence length

To demonstrate how the finite size sample affects the statistical quantity F , we show in fig. 4 plots of F for human metallothionein for three different sizes. First we *randomly* choose a sub-sequence of 300 nucleotides from the entire

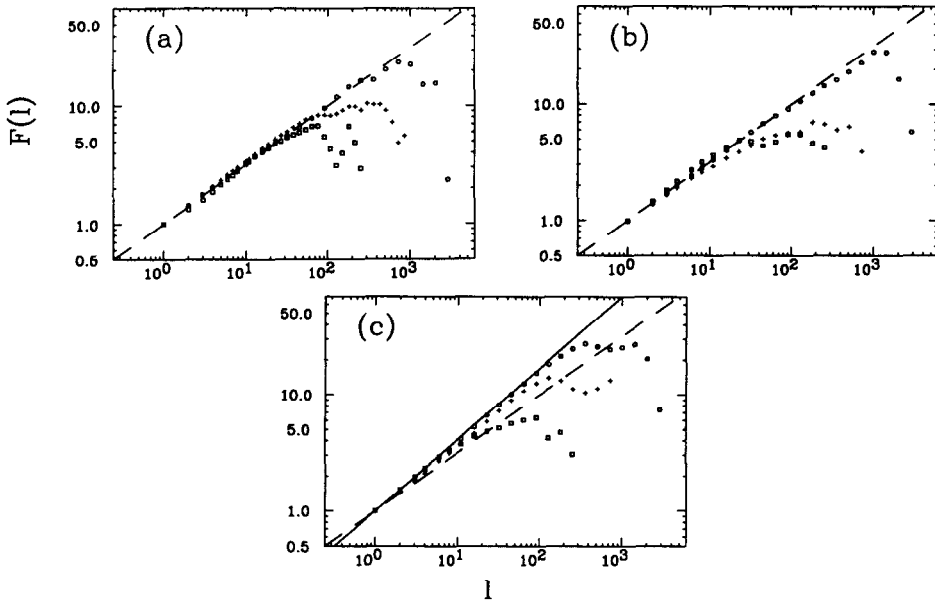


Fig. 3. $F(l)$ versus l for three different type of artificial sequences: (a) *Unbiased* random uncorrelated sequence (i.e., 50% purines); (b) *biased* but still uncorrelated sequence (with 60% purines); and (c) *correlated* sequence with correlation parameter 0.61 (with 50% purines). The dashed line has slope 0.5, while the solid line has slope 0.61. The different symbols represent different size of the sequences: The entire sequence of 2941 nucleotides (circles), the sub-sequences of 1000 nucleotides (crosses) and 300 nucleotides (open squares).

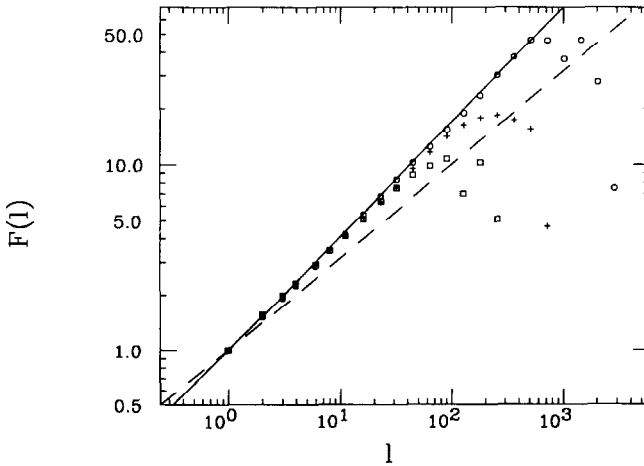


Fig. 4. $F(l)$ versus l for three different sizes of samples: The whole sequence of 2941 nucleotides (circles), the sub-sequences of 1000 nucleotides (crosses) and 300 nucleotides (open squares). The dashed line has slope 0.5, while the solid line has slope 0.61. Note that the linearity in all cases extends up to a fraction of about 1/10 the sample size, a fact familiar to workers involved in statistical analysis of this sort.

gene sequence to calculate the quantity F within this small sample (open square). Second, we choose (again randomly) a piece of the sequence of 1000 nucleotides from the same gene, repeat the analysis and plot the results (crosses). Finally, we analyze the entire sequence (2941 nucleotides) and plot the results (circles). The “fall-off” in the straight line behavior after the distance l reaches approximately one tenth of the size of the sample is typical of all fractal analyses. It is also found in sequences of correlated and random numbers (fig. 3). The trend of obtaining longer regions of straight line behavior for larger size samples is what one expects on statistical grounds.

One can always worry that the long-range features will disappear for longer DNA sequences. As evidence that the long-range feature does not disappear for larger samples, ref. [21] analyzed the entire human beta globin region (73326 bases) and found linearity up to $l \approx 7000$. Recently, Munson et al. [26] analyzed the entire yeast III chromosome (315 000 bases) [45] and found linearity up to $l \approx 31\,500$. In general, the data are linear over a range that is about a factor of ten less than the range of the data. This increase in statistical error when there are less than roughly 10 independent data sets is usually found for analyses of this sort. Thus, e.g., if a gene has 10 000 nucleotides, then there are only 10 *independent* sets of data obtained when the calipers are separated by a distance 1000.

One can worry about the apparent lack of consistency between values of α measured for different genes, or even for different regions within the same gene. Peng et al. have recently carried out a systematic study of the fluctuations in the correlation exponents obtained [35]. They indeed find prominent sample-to-sample variations in the scaling exponent, as well as variations within a single sample. To determine if these fluctuations may result from finite system size, they generate correlated random sequences of comparable length and study the fluctuations in this control system. Peng et al. find that the DNA exponent fluctuations are consistent with those obtained from the control sequences having long-range power law correlations.

3.4. Other methods of measuring long-range correlations

One can also worry that the apparent long-range correlation is some artifact of the DNA walk method itself. To compare the fluctuations of α in our DNA walk method with those found in other methods, Peng et al. used two standard methods to study the correlation property of sequences, namely the correlation function $C(l)$ and the power spectrum $S(f)$. The power spectrum density, $S(f)$, is obtained by (a) Fourier transforming the sequence $\{u(i)\}$ and (b) taking the square of the Fourier component. For a stationary sequence, the power spectrum is the Fourier transform of the correlation function. If the

correlation decays algebraically (not exponentially), i.e., there is no characteristic scale for the decay of the correlation, as we found in the non-coding DNA sequences, then we expect power-law behavior for both the power spectrum and the correlation function,

$$S(f) \sim (1/f)^\beta \quad (2b)$$

and

$$C(l) \sim (1/l)^\gamma. \quad (2c)$$

The correlation exponents α , β and γ are not independent, since [14,15]

$$\alpha = (1 + \beta)/2 = (2 - \gamma)/2. \quad (3)$$

For a typical DNA sequence of finite length, both the correlation function and power spectrum are fairly noisy, but the estimates of β and γ obtained are consistent with those calculated from the DNA walk method (see fig. 5). The reason for the smaller fluctuations of α in the DNA walk method is due to the fact that $F^2(l)$ is a double summation of $C(l)$. Thus it would seem that the original DNA walk method is more useful due to reduced noise.

4. Difference between correlation properties of coding and non-coding regions

Our initial report [21] on long-range (scale-invariant) correlations in DNA sequences has generated contradicting responses. Some [22,24,26] support our initial finding, while some [23,28,32,34] disagree. Furthermore, the conclusions of refs. [24] and [23,28,32] are inconsistent *with one another* in that [23] and [34] doubt the existence of long-range correlations (even in non-coding sequences) while [23] and [28,32] conclude that even coding regions display long-range correlations ($\alpha > 1/2$). Prabhu and Claverie [28] claim that their analysis of the putative *coding* regions of the yeast chromosome III [45] produces a *wide range of exponent values*, some larger than 0.5. The source of these contradicting claims may arise from the fact that, in addition to normal statistical fluctuations expected for analysis of rather short sequences coding regions typically consist of only a few lengthy regions of alternating strand bias. Hence scaling analysis cannot be applied reliably to the entire sequence but only to sub-sequences.

Figure 6a displays our analysis of a typical coding sequence, consisting of two large sub-regions, each with different strand bias; the first (roughly 22 000 nucleotides) is G rich (compared to the average concentration of the entire

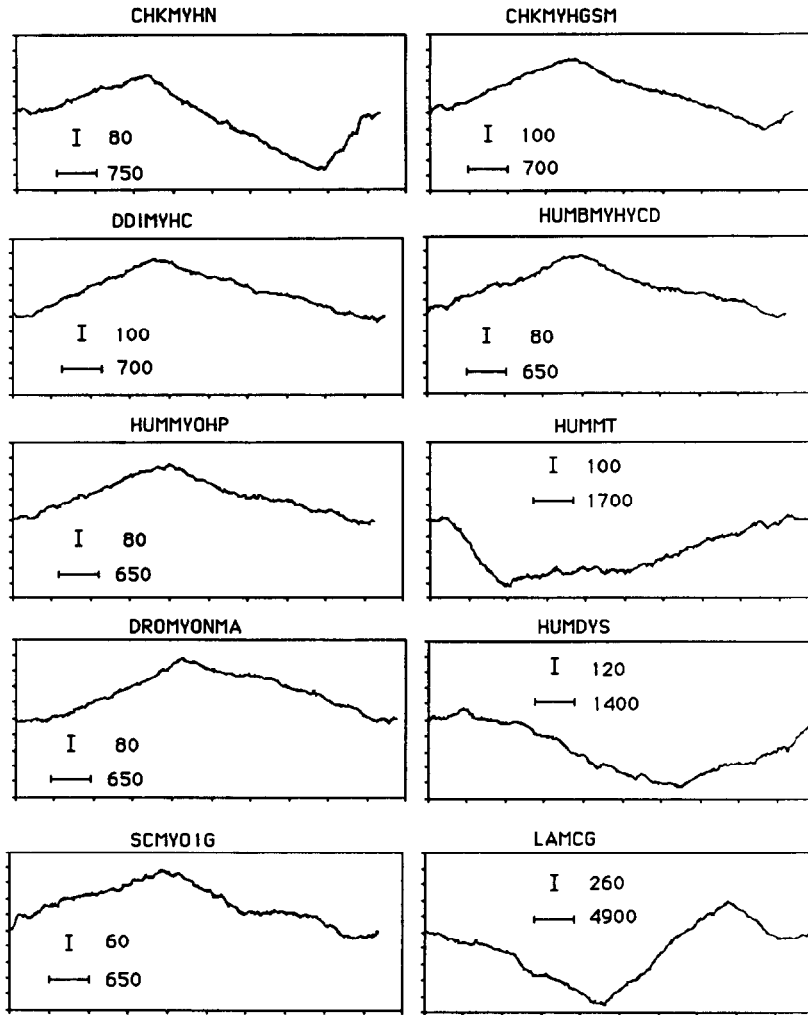


Fig. 5. Ten different coding sequences, plotted in the same form as fig. 1. Note that for all ten cases, there appear regions with one strand bias, followed by regions of a different strand bias. The fluctuation on either side of the overall strand bias is found to be random, a fact that is plausible by looking at these DNA walk representations. The bias introduced by the change in concentration of purine and pyrimidine would not be eliminated by the average term in eq. (2) of the manuscript if pieces of different bias would be analyzed together.

sequence), the second G poor. The scaling analysis of $F(l)$ (fig. 6c) on the *entire* sequence shows a *crossover* behavior, i.e., the log-log plot of the $F(l)$ versus l line has an initial slope 0.5 and curving toward 1 at larger values of l . This crossover behavior is typical of many physical systems having a characteristic scale. In this case, this characteristic length scale is associated with the

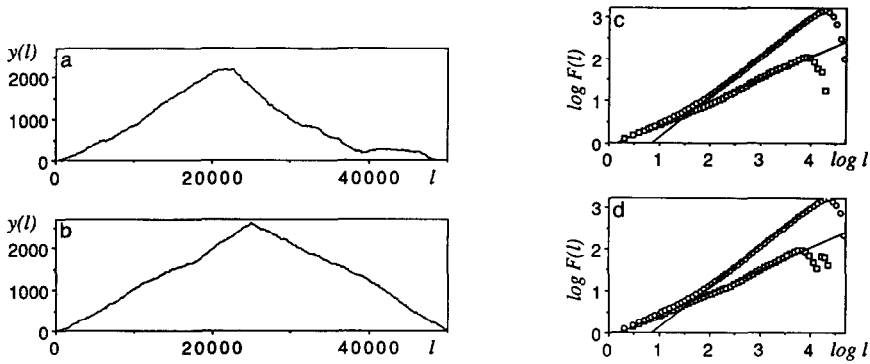


Fig. 6. (a) Landscape representation of the intronless sequence of the complete genome for lambda phage (GenBank: LAMCG, 48502 base pairs). Each “up” step corresponds to a guanine (G) and a “down” step corresponds to any one of the other three nucleotides (A, C or T). For graphical representation we plot the DNA walk such that the end point has the same vertical displacement as the starting point (for the statistical analysis, we use the original definitions). (b) Landscape for a biased random walk, where the bias is similar to that of the DNA sequence: in the first half, $p_{up} = 0.3$ ($>1/4$, the value expected in the absence of strand bias) followed by $p_{up} = 0.2$ ($<1/4$) in the second half. (c) The rms fluctuation in landscape altitude [2], $F(l)$, for the full genome (\circ , slope $\alpha = (\beta + 1)/2 = 0.95$ for $l > 100$) and for the first sub-region (\square , slope 0.54). (d) The rms fluctuation for the biased random walk (slope 0.96) and for the first sub-region (slope 0.54). Similar behavior was observed for the DNA walk with the purine–pyrimidine rule (step “up” for C and T; step “down” for A and G).

length of the two regions of strand bias. However, when the effects of this strand bias are first removed by *separately analyzing the sub-regions*, then we find α close to 0.5, indicating no long-range correlations.

We also calculate $F(l)$ for an artificial “control” sequence consisting of a 50 000-step biased random walk with similar strand biases as in the two regions of the DNA sequence (fig. 6b). We observe the same crossover behavior when the *entire* sequence was analyzed, but obtain the correct exponent $\alpha = 0.5$ when *each sub-region is separately analyzed* (fig. 6d). Figs. 6c, d also show that failure to correct for the crossover due to alternating regions of strand bias gives rise to a larger slope (upper curves) at larger values of l and hence misleadingly large values for the correlation exponents.

The power spectrum $S(f)$ for the *entire* sequence has an initial region of $1/f^\beta$ behavior at low frequency (that could be misinterpreted as indicative of long-range correlation, $\beta \neq 0$) followed by a flat region (indicative of no correlation or “white noise”). However, if the effects of strand bias are first removed by *separately analyzing the sub-regions*, then we find a flat $S(f)$ – and hence the correct correlation exponent $\beta = 0$. We also calculated $S(f)$ for a “control” consisting of a 50 000-step biased random walk with similar strand biases as in the two regions of the gene (fig. 6b), and found a misleading

exponent $\beta \neq 0$ when the *entire* sequence was analyzed, but the correct exponent $\beta = 0$ when *each sub-region was separately analyzed*.

Apart from the reduced noise mentioned above, one additional advantage of the DNA walk method [21] is that to find the exponent characterizing the long-range correlation one need not correct the data by subtracting the white noise, $S(\infty)$ [24]. Since there is no unambiguous method of estimating $S(\infty)$, this need to correct the data introduces an uncontrollable source of uncertainty.

Prabhu and Claverie [28] noted that nonlinear cures were obtained *even some intron-containing sequences*. The DNA walks for intron-containing sequences do not show any apparent length scale of strand bias; as noted above, there seems to be a broad distribution of lengths of strand bias. In principle, for this type of DNA walk, the min-max procedure is not necessary since there is no characteristic length that needs to be taken into account. For most *intron-containing* sequences, our analysis shows that there is a broad range of scaling (with $\alpha > 0.5$) when we study the entire sequence as a whole. For certain intron-containing genes, however, we find that the min-max procedure can extend the scaling region, leading to a larger range of constant- α behavior.

One might also wonder if the selection of “max” and “min” segments gives a bias to the calculations. The “max” and “min” operation is a systematic method to treat all DNA sequences on equal footing without applying any *a priori* knowledge of the sequence itself. Notice that although this procedure eliminates the problem of large scale variation of concentration for the coding sequences, it will not alter the result of a truly self-similar sequence. Fig. 2 shows that we can obtain the same value of α by applying our algorithm for any part of the sequence. Peng et al. [39] have recently applied the “bridge method” to obviate the need to perform the min-max partitioning.

To provide an “unbiased” test of the thesis that non-coding regions possess but coding regions lack long-range correlations, we analyzed [40] several uncorrelated and correlated control sequences of size 10^5 nucleotides using the GRAIL neural net algorithm [46]. The GRAIL algorithm identified about 60 putative exons in the uncorrelated sequences, but only about 5 putative exons in the correlated sequences. We also used the beachcomber method [40], which shows pronounced dips in α in the region where genes are expected (fig. 7).

5. Mosaic structure of DNA: “Patches”

The connection between our finding and the mosaic structure of DNA remains not entirely clear. Some authors [23,28] argue that the observed long-range correlation is a trivial consequence of the heterogeneous “mosaic”

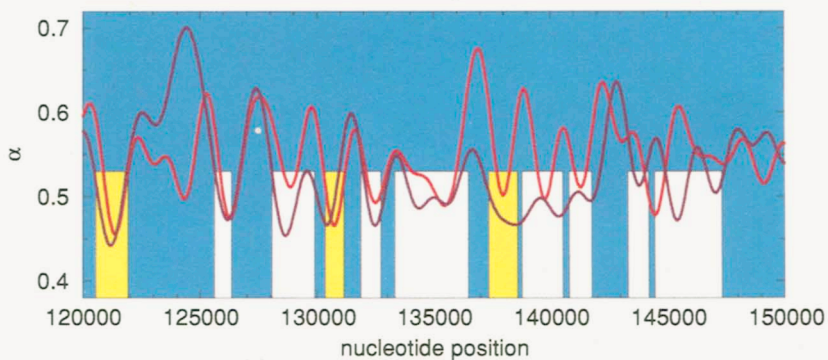
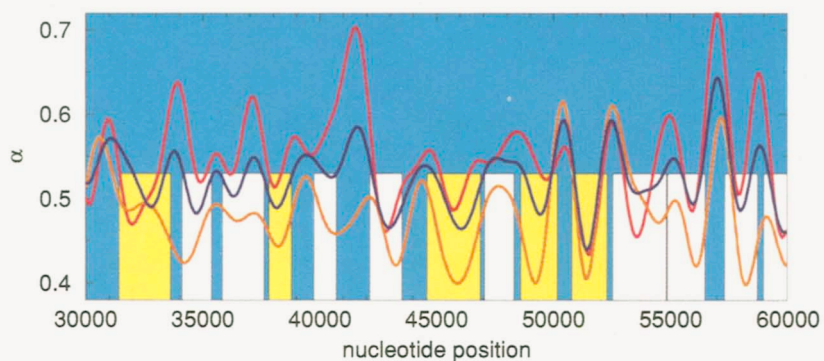


Fig. 7.

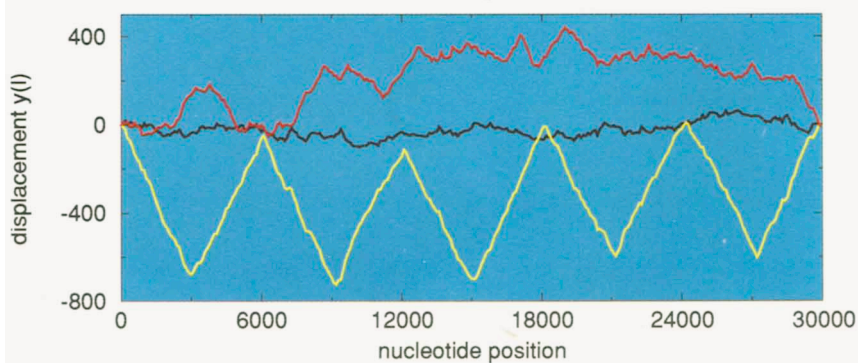
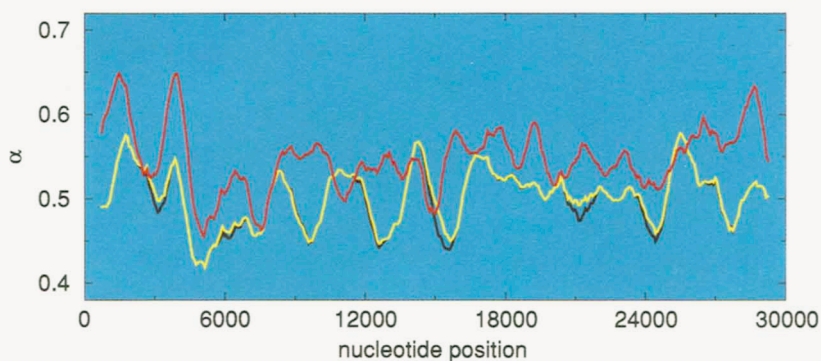


Fig. 9.

structure of DNA. We have argued [36,38,39] that *such mosaic structure can arise directly from long-range correlations*, thus providing a concrete origin for the “patchy” structure of non-coding region of DNA.

The mosaic structure of DNA is visually apparent when one examines the “DNA walk” land-scape (see, e.g., fig. 1). We find that a simulated (control) sequence generated with a well-defined long-range correlation [14–16] leads to a similar mosaic structure [36]. To quantitatively characterize the “patches” comprising the mosaic structure, we studied the histogram giving the number of patches of size s , for both the gene and the control system for the case $\alpha = 0.65$ ($\alpha \neq 0.5$), and found quantitative agreement (fig. 8).

Fig. 7. Top: Beachcomber plot for a typical section containing about 10% of the yeast chromosome III [45] – from base pair #30 000 to #60 000. The vertical yellow bars indicate the set of base pairs forming identified genes (GenBank Release #76), while the white bars indicate less certain “putative genes” determined from analysis of open reading frames [45]. The exponent α is calculated by the beachcomber method [40]: We form an observation box of length 800, place this box at the beginning of the chromosome, and calculate the long-range correlation exponent α for the 800 base pairs lying inside this box. Then we move the box 75 base pairs further along the chromosome, and again calculate α for the 800 base pairs lying inside this box. Iterating this procedure, we obtain $315\,000/75 = 4186$ successive values of α , each giving a “local” measurement of the degree of long-range correlation. The red curve is obtained using rule 1 – a “down” step for A or G (purines) and an “up” step for C or T (pyrimidine). The gold curve is obtained using rule 2 – a “down” step for C or G but an “up” step for A or T. The violet curve is the average of rules 3–6, where rule 3 takes three “up” steps for each A and a “down” step for C, G, or T (with rules 4–6 defined similarly). We see that when the box is covering coding regions, the value of α is generally small, while in between coding regions, there is frequently a peak in α . If α were the same for coding and non-coding regions, we would expect the peaks and dips to occur with no evident correlation in the position of genes. We carried out this analysis for the entire chromosome, and concluded that of the 6 rules, rule 1 was most accurate. Bottom: Beachcomber plot for a different region of yeast chromosome III [45] – from base pair #120 000 to #150 000, showing as the red curve the results of rule 1. Shown for comparison as the violet curve are the results of analyzing by the same rule an artificial sequence which is chosen to be uncorrelated for those base pairs belonging to genes, but to have a power-law correlation (with $\alpha = 0.6$) for base pairs in between genes.

Fig. 9. The *bottom* panel displays three landscapes: (i) in blue an uncorrelated biased random walk (normalized, as in fig. 1, such that the final point has the same altitude as the initial point); (ii) in yellow, a walk which is identical to the first except that after each 3000 base pairs the bias is reversed; and (iii) in red a generalized Levy walk [36] with correlation parameter $\mu = 2.2$ corresponding to $\alpha = 0.9$. The *top* panel shows the beachcomber analysis for these three landscapes. Note that for the uncorrelated landscapes (i) and (ii) the local values of α are almost identical, with the peaks and valleys corresponding solely to the sort of statistical fluctuations one expects from finite-size samples. Since the “Nee patches” have almost the identical beachcomber spectrum as the uncorrelated biased random walk, we conclude that such patches are not sufficient to explain the observed long-range correlations and concomitant large values of α found in DNA sequences with non-coding base pairs. In contrast, for the correlated landscape (iii), there are larger fluctuations that correspond better to the beachcomber plots for real DNA sequences (such as shown in fig. 7).

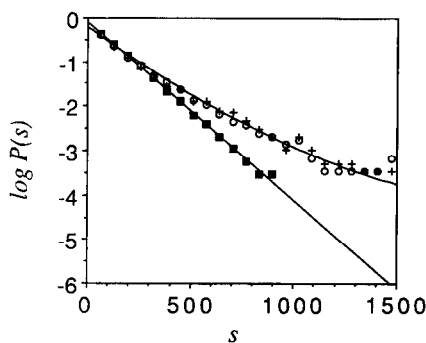


Fig. 8. A semi-log plot of $P(s)$ vs s , where $P(s)$ is the probability to find a “run” of s steps in the same direction; we use the coarse-grained sequence calculated with a window of size 64. The squares display the results for runs of *uncorrelated* variables (solid line). The data from a *correlated* control sequence with $\alpha = 0.75$ (\circ) and from 5 human genomic DNA (+) are in good agreement with each other, and deviate from the random uncorrelated data.

In contrast, simple “patchiness” by itself does not generally lead to a long-range correlation. Suppose, e.g., one randomly divides a sequence of length N into k sub-regions with alternating strand bias (fig. 9). The patch size distribution will display a characteristic length given by N/k . Simulations of this straightforward patch model display a distinct maximum in the local slopes $\alpha(l)$ defined in [31], but the bridge method [39] correctly gives $\alpha = 1/2$. Non-coding DNA sequences rarely display this maximum – $\alpha(l)$ is constant within statistical fluctuation.

Nee [23] raises the possibility that coding regions correspond to uncorrelated *biased* random walks, while non-coding regions correspond to uncorrelated *unbiased* random walks. Since coding and non-coding regions interdigitate, the overall behavior could appear to be a *correlated random walk*. To test this possibility, we studied genes with the coding regions removed. If Nee’s hypothesis were correct, we would have observed $\alpha = 0.5$, as for an uncorrelated unbiased random walk. Instead, we find a comparable value of α (>0.5) as in the full gene (see fig. 1 of [30]).

We have noted [36,38] that Nee’s suggested mechanism could give rise to long-range correlations if the length distributed possesses no characteristic scale. Indeed, long-range correlations can legitimately arise from the mixing of regions (of length l_j) with alternating bias *provided* that the length distribution of these biased regions follows a power-law, $P(l_j) \propto l_j^{-\mu}$, in which case the correlation exponent α is related to μ through $\alpha = 2 - \mu/2$ if $2 < \mu < 3$. A

possible molecular basis for a power-law distribution of biased regions was also discussed [36,38].

6. Discussion

Prabhu and Claverie [28] have also stated that our DNA walk approach cannot be used as a “general method to identify intronless sequences”. We agree with their opinion provided one were to use *only* our algorithm, and provided one were to require “complete certainty of identification”. However, when used in concert with other methods, and when statistical fluctuations are taken into account, our recently developed “beachcomber algorithm” [40] is proving to be quite useful (see, e.g., fig. 7).

What is the biological meaning of our finding? If two nucleotides whose positions differ by 10 000 were uncorrelated, then there might be no meaning. However, if they are correlated there must be a reason! The method we describe points out a new element of DNA structure and suggests a possible fundamental role for the non-coding regions (termed “introns”). Our work may also reveal an interesting feature of the coding regions (termed “exons”).

Next, we return to the discrepancy between our findings and those of Voss [24], who claims that coding regions also display long-range correlations. However, Voss did not take into account the known strand bias (excess of one base pair over another in a large regime) in coding regions. His neglect of crossover effects could be the reason he found exponents larger than ours. This neglect therefore casts some doubt upon Voss’ main conclusion that the correlation exponent behaves in a non-monotonic fashion as evolution proceeds. In fact, we recently tested [38] the Voss conjecture by analyzing 11 distinct myosin heavy chain genes belong to 8 different species – and their spliced together coding regions as well. Our analysis suggests a monotonic increase in fractal complexity for MHC genes with evolution with vertebrate > invertebrate > plant (yeast). We also developed a simple iterative model, based on known properties of polymeric sequences, that generates long-range nucleotide correlations from an initially non-correlated coding region. The DNA walk analysis and this new model support the “intron-late theory” of gene evolution [47–49].

Before concluding, we note that the long-range correlations in DNA sequences are of interest because they may be an indirect clue to its three-dimensional structure [33,37] or a reflection of certain scale-invariant properties of long polymer chains [50]. In any case, the statistically meaningful long-range (scale-invariant) correlations in non-coding regions and their

absence in coding regions will need to be accounted for by future explanations of global properties in gene organization and evolution.

Acknowledgements

We wish to thank F. Sciortino for important contributions in the initial stages of this project, and C.R. Cantor, C. DeLisi, R.D. Rosenberg, J.J. Schwartz, M. Schwartz and N. Shworak for valuable discussions. Partial support was provided to ALG by the G. Harold and Leila Y. Mathers Charitable Foundation, the National Heart, Lung and Blood Institute and the National Aeronautics and Space Administration, to MS by the American Heart Association, and to HES by the National Science Foundation and Office of Naval Research.

Appendix. Scale invariance and power laws

To remove some confusion concerning the relation between the existence of power-law correlation and the concept of scale invariance, we demonstrate that there is a simple mathematical connection between the two. A scale invariant function has the property that if the variable x is increased by an arbitrary factor λ , then the function is changed by a factor λ^p which is independent of the value of x ,

$$f(\lambda x) = \lambda^p f(x) \quad (4a)$$

for all λ . A functional equation, such as (4a), constrains the set of possible functional forms of $f(x)$: any function $f(x)$ satisfying (4a) must be a power law, as may be seen by substituting the choice $\lambda = 1/x$ in (4a),

$$f(x) = Ax^p. \quad (4b)$$

We say that scale invariance (eq. (4a)) implies power-law behavior (eq. (4b)).

Conversely, power-law behavior implies scale invariance, since any function $f(x)$ obeying (4b) also obeys (4a) – one can verify this by substitution. Thus scale invariance is mathematically equivalent to power-law behavior.

References

- [1] B.B. Mandelbrot, *The Fractal Geometry of Nature* (Freeman, San Francisco, 1982).
- [2] A. Bunde and S. Havlin, eds., *Fractals and Disordered systems* (Springer, Berlin, 1991).

- [3] D. Stauffer and H.E. Stanley, *From Newton to Mandelbrot: A Primer in Theoretical Physics* (Springer, Heidelberg, New York, 1990).
- [4] E. Guyon and H.E. Stanley, *Les Formes Fractales* (Palais de la Découverte, Paris, 1991) [English translation: *Fractal Forms* (Elsevier, North-Holland, Amsterdam, 1991)].
- [5] T. Vicsek, *Fractal Growth Phenomena*, 2nd ed. (World Scientific, Singapore, 1992); J. Feder, *Fractals* (Plenum, New York, 1988).
- [6] H.E. Stanley and N. Ostrowsky, eds., *On Growth and Form: Fractal and Non-Fractal Patterns in Physics*, Proc. 1985 Cargèse NATO ASI Series E: Applied Sciences, vol. 100 (Nijhoff, Dordrecht, 1985).
- [7] H.E. Stanley, *Introduction to Phase Transitions and Critical Phenomena* (Oxford Univ. Press, London, 1971).
- [8] H.E. Stanley and N. Ostrowsky, eds., *Correlations and Connectivity: Geometric Aspects of Physics, Chemistry and Biology*, Proc. 1990 Cargèse NATO ASI Series E: Applied Sciences, vol. 188 (Kluwer, Dordrecht, 1990).
- [9] F. Zernike, *Physica* 7 (1940) 565.
- [10] P. Bak, C. Tang, and K. Wiesenfeld, *Phys. Rev. Lett.* 59 (1987) 381.
- [11] S. Havlin, A.-L. Barabási, S.V. Buldyrev, C.-K. Peng, M. Schwartz, H.E. Stanley and T. Vicsek, *Anomalous surface roughening: experiment and models*, in: *Growth Patterns in Physical Sciences and Biology*, Proc. 1991 NATO Advanced Research Workshop, Granada, Spain, October 1991, E. Louis, L. Sander and P. Meakin, eds. (Plenum, New York, 1992); K. Sneppen, *Phys. Rev. Lett.*, submitted.
- [12] P. Bak, in: *Fractals in Science*, A. Bunde and S. Havlin, eds. (Springer, Berlin, 1994).
- [13] H.E. Stanley and N. Ostrowsky, eds., *Random Fluctuations and Pattern Growth: Experiments and Models*, Proc. 1988 NATO ASI, Cargèse (Kluwer, Dordrecht, 1988).
- [14] S. Havlin, R. Selinger, M. Schwartz, H.E. Stanley and A. Bunde, *Phys. Rev. Lett.* 61 (1988) 1438;
S. Havlin, M. Schwartz, R. Blumberg Selinger, A. Bunde and H.E. Stanley, *Phys. Rev. A* 40 (1989) 1717;
R.B. Selinger, S. Havlin, F. Leyvraz, M. Schwartz and H.E. Stanley, *Phys. Rev. A* 40 (1989) 6755.
- [15] C.-K. Peng, S. Havlin, M. Schwartz, H.E. Stanley and G.H. Weiss, *Physica A* 178 (1991) 401;
C.-K. Peng, S. Havlin, M. Schwartz and H.E. Stanley, *Phys. Rev. A* 44 (1991) 2239.
- [16] M. Araujo, S. Havlin, G.H. Weiss and H.E. Stanley, *Phys. Rev. A* 43 (1991) 5207;
S. Havlin, S.V. Buldyrev, H.E. Stanley and G.H. Weiss, *J. Phys. A* 24 (1991) L925;
S. Prakash, S. Havlin, M. Schwartz and H.E. Stanley, *Phys. Rev. A* 46 (1992) R-1724.
- [17] F. Caserta, H.E. Stanley, W. Eldred, G. Daccord, R. Hausmann and J. Nittmann, *Phys. Rev. Lett.* 64 (1990) 95.
- [18] C.-K. Peng, J. Mietus, J. Hausdorff, S. Havlin, H.E. Stanley and A.L. Goldberger, *Phys. Rev. Lett.* 70 (1993) 1343.
- [19] A. Schenkel, J. Zhang and Y.-C. Zhang, *Fractals* 1 (1993) 47.
- [20] R.N. Mantegna, *Physica A* 179 (1991) 232.
- [21] C.-K. Peng, S.V. Buldyrev, A.L. Goldberger, S. Havlin, F. Sciortino, M. Simons and H.E. Stanley, *Nature* 356 (1992) 168.
- [22] W. Li and K. Kaneko, *Europhys. Lett.* 17 (1992) 655.
- [23] S. Nee, *Nature* 357 (1992) 450.
- [24] R. Voss, *Phys. Rev. Lett.* 68 (1992) 3805.
- [25] J. Maddox, *Nature* 358 (1992) 103.
- [26] P.J. Munson, R.C. Taylor and G.S. Michaels, *Nature* 360 (1992) 636.
- [27] I. Amato, *Science* 257 (1992) 747.
- [28] V.V. Prabhu and J.-M. Claverie, *Nature* 357 (1992) 782.
- [29] P. Yam, *Sci. Am.* (Sept. 1992) 23.
- [30] C.-K. Peng, S.V. Buldyrev, A.L. Goldberger, S. Havlin, F. Sciortino, M. Simons and H.E. Stanley, *Physica A* 191 (1992) 25.
- [31] H.E. Stanley, S.V. Buldyrev, A.L. Goldberger, J.M. Hausdorff, S. Havlin, J. Mietus, C.-K. Peng, F. Sciortino and M. Simons, *Physica A* 191 (1992) 1.

- [32] C.A. Chatzidimitriou-Dreismann and D. Larhammar, *Nature* 361 (1993) 212.
- [33] A. Yu. Grosberg, Y. Rabin, S. Havlin and A. Nir, *Biofisika (Russia)* 26 (1993) 1; *Europhys. Lett.*, in press.
- [34] S. Karlin and V. Brendel, *Science* 259 (1993) 677.
- [35] C.-K. Peng, S.V. Buldyrev, A.L. Goldberger, S. Havlin, M. Simmons and H.E. Stanley, *Phys. Rev. E* 47 (1993) 3730.
- [36] S.V. Buldyrev, A.L. Goldberger, S. Havlin, C.-K. Peng, M. Simons and H.E. Stanley, *Phys. Rev. E* 47 (1993) 4514.
- [37] A.S. Borovik, A. Yu. Grosberg and M.D. Frank-Kamenetskii, *Nature*, in press.
- [38] S.V. Buldyrev, A.L. Goldberger, S. Havlin, C.-K. Peng, H.E. Stanley and M. Simons, *Fractal landscapes and molecular evolution: modeling the myosin heavy chain gene family*, *Biophys. J.*, submitted.
- [39] C.-K. Peng, S.V. Buldyrev, S. Havlin, M. Simons, H.E. Stanley and A.L. Goldberger, *The basis of long-range nucleotide correlations in DNA: distinguishing patchiness from fractal organization*, *Science*, submitted.
- [40] S.M. Ossadnik, S.V. Buldyrev, A.L. Goldberger, S. Havlin, C.-K. Peng, M. Simons and H.E. Stanley, *A statistical method for finding coding regions in DNA sequences*, submitted.
- [41] E.W. Montroll and M.F. Shlesinger, *The wonderful world of random walks*, in: *Nonequilibrium Phenomena II. From Stochastics to Hydrodynamics*, J.L. Lebowitz and E.W. Montroll, eds. (North-Holland, Amsterdam, 1984) pp. 1–121.
- [42] S. Tavaré and B.W. Giddings, *Some statistical aspects of the primary structure of nucleotide sequences*, in: *Mathematical Methods for DNA Sequences*, M.S. Waterman, ed. (CRC Press, Boca Raton, 1989) pp. 117–132, and references therein.
- [43] B. Levin, *Genes IV* (Oxford Univ. Press, Oxford, 1990).
- [44] J.D. Watson, M. Gilman, J. Witkowski and M. Zoller, *Recombinant DNA* (Scientific American Books, New York, 1992).
- [45] S.G. Oliver et al., *Nature* 357 (1992) 38.
- [46] E.C. Uberbacher and R.J. Mural, *Proc. Natl. Acad. Sci. USA* 88 (1991) 11261.
- [47] W.F. Doolittle, *Understanding introns: origins and functions*, in: *Intervening Sequences in Evolution and Development*, E. Stone and R. Schwartz, eds. (Oxford Univ. Press, New York, 1990) pp. 42–62.
- [48] W. Gilbert, *Nature* 271 (1978) 501.
- [49] W.-H. Li and D. Graur, *Fundamentals of Molecular Evolution* (Sinauer Associates, Sunderland, MA, 1991).
- [50] P.-G. de Gennes, *Scaling Concepts in Polymer Physics* (Cornell Univ. Press, Ithaca, NY, 1979).