

Scale-invariant correlations in the biological and social sciences

By H. E. STANLEY†, L. A. N. AMARAL†, J. S. ANDRADE, JR‡, S. V. BULDYREV†, S. HAVLIN§, H. A. MAKSE¶, C.-K. PENG||, B. SUKI¶¶ and G. VISWANATHAN†

† Center for Polymer Studies and Department of Physics, Boston University, Boston, Massachusetts 02215, USA

‡ Departamento de Física, Universidade Federal do Ceará, 60451-970 Fortaleza, Ceará, Brazil

§ Minerva Center for the Physics of Mesoscopics, Fractals and Neural Networks and Department of Physics, Bar-Ilan University, Ramat Gan, Israel

|| Cardiovascular Division, Harvard Medical School, Beth Israel Hospital, Boston, Massachusetts 02215, USA

¶¶ Department of Biomedical Engineering, Boston University, Boston, Massachusetts 02215, USA

ABSTRACT

In this opening introductory paper, we discuss the possibility that scale-invariant correlations may be a feature of biological and possibly even social systems. We illustrate this possibility by reviewing recent work at Boston University. Specifically, we focus first on the apparent scale-invariant correlations in non-coding deoxyribonucleic acid (DNA) and show that this feature can be used to distinguish coding and non-coding DNA. We argue that the inflating a degassed lung is characterized by a cascade of avalanches, as the airways successively open, and that distribution functions characterizing this cascade are scale invariant. Moving from the lung to the heart, we find that the sequence of interbeat intervals is characterized by scale-invariant correlations in health, but not in disease. Moving from individual organs to entire organisms, we discuss recent experimental evidence that the foraging behaviour of the wandering albatross is governed by a scale-invariant Levy distribution. Finally, we enquire whether scale invariance describes not only animal behaviour but also human behaviour. To this end, we analyse data on urban growth patterns, on finance and on economics. For all cases, we find empirical evidence of scaling behaviour. We conclude by asking why such complex systems might display scale invariance.

§ 1. INTRODUCTION

Here we are going to look at some examples of scale-invariant correlations that are of interest to biological scientists and, possibly, to social scientists.

At one time, it was imagined that the ‘scale-free’ phenomena are relevant to only a fairly narrow slice of physical phenomena (Stanley 1971). However, the range of systems that apparently display power-law and hence scale-invariant correlations has increased dramatically in recent years, ranging from base pair correlations in deoxyribonucleic acid (DNA) (Peng *et al.* 1992), lung inflation (Suki *et al.* 1994, Barabási *et al.* 1996) and interbeat intervals of the human heart (Peng *et al.* 1993, 1995) to complex systems involving large numbers of interacting subunits that display ‘free will’, such as govern city growth (Makse *et al.* 1995) and even economics (Mantegna and Stanley 1995, Stanley *et al.* 1996).

§ 2. SCALE INVARIANCE IN NON-CODING DNA

2.1. *What is the puzzle? Why should we care?*

In the human cell, 97% of the DNA is not coding for protein. What is non-coding DNA doing? In the past, we were taught simply that DNA codes for protein. Now we know that actually only 3% of DNA codes for protein. The rest of the DNA seems to be doing nothing; sometimes it is even referred to as ‘junk’ DNA. However, we might say that, if this DNA is ‘junk’, it is not like the junk that we throw away, but more like the junk we do not know what to do with at present; so we store it in the attic.

There seems to be little agreement among biologists about why the non-coding DNA is present. So why should we care about this puzzle?

The practical reason to care is that worldwide the scientific community is spending the equivalent of three gigadollars to sequence the three gigabases, one dollar per base, of the human genome, that is to uncover the sequence of DNA bases in the entire 46 chromosomes of a human being. If only 3% of these are coding for protein, we could reduce our time and budget by a factor of 30 if we examined only the 90 megabases that are actually coding for protein.

The scientific reason is familiar in biology: if it is there, it usually has a purpose. We have two eyes and know why, and two ears and know why. Now 97% of our DNA is in a non-coding form, and we do not understand why, but it is easy to be tempted to hypothesize that it is there for a reason—that it has some function.

2.2. *What have we done?*

We have modest results to report on both the practical and the scientific sides. First, we have taken advantage of the fact that there are long-range correlations in the non-coding DNA by devising an algorithm that enables us to identify the non-coding and, thus, by implication, the coding parts of DNA, and with a statistical accuracy comparable to other methods. Second, on the biological side we have some work that cautiously suggests that non-coding DNA has features in common with the statistical properties of a structured language.

If we unwind a tiny piece of the double helix that constitutes DNA, we see that it consists of two strands, and that these strands have bases (sometimes called nucleotides) that have four letters of the alphabet: C, A, T and G. These four have the remarkable property that the base on one strand always ‘marries’ a given base on the other strand. In every instance, C is paired with G, and A is paired with T. Because of this pairing, DNA is able to replicate itself: by breaking the hydrogen bonds, followed by each strand making a complementary copy of itself. In the case of multiplying viruses, this takes place approximately every 20 min. Overnight, 30 doublings take place; so each virus leads to 10^9 children.

If we want to describe all the hereditary information contained in, for example, a hair cell, we end up with a real sequence of letters that anyone can access from the GenBank using the internet. (Everyone in the Genome Project has agreed to make every sequence available to the entire world.) In about 41 years, the entire human genome will be available—about a million solid pages of data. What do we do with all that information?

If we push a system to very near its critical point, we shall see scale-invariant fluctuations of all length scales up to and including the wavelength of light. If we shine a laser at a binary mixture near its critical point, it will scatter light, a

phenomenon discovered by Andrews (1869) in England about 100 years ago that we call 'critical opalescence'. Andrews interpreted the phenomenon correctly, deducing that, near the system's critical point, fluctuations of all length scales were present and also that, since among those present were *huge* fluctuations, there must be correlations between the constituents making up the binary mixture. If there were no correlations at all, there would be no way to explain the presence of the huge (near the wavelength of light) fluctuations.

The study of what is termed critical phenomena has become a very active field of scientific inquiry. The simple random walk, the motion of a small particle undergoing Brownian motion, was first accurately measured experimentally by Jean Perrin (1915), French Nobel Prize Winner and pre-founder of the concept of fractals. He recognized that this trajectory had something in common with this critical point; there occur fluctuations on all scales. In Perrin's words, 'Fresh irregularities appear every time I increase the magnification.' Over the last 20 years or so, Sir Sam Edwards has carefully studied the simple random walk using field theory and other methods. This area has also received much attention by numerous experimentalists in Europe and the USA, and the whole field of polymer physics, a field that seems to have nothing to do with the concept of the critical point, was found to have parallel properties.

So it should not be too surprising that we find in DNA something in common with critical point phenomena. What do we do with the one million pages of data describing bases? We start by attempting some visual representation of the data. We can construct a 'visual mountain range' from the DNA sequence; we take an 'up step' each time that we have a C or a T and a 'down step' each time that we have an A or a G. When we do this, we get a landscape. This landscape for the muscle protein myosin represents 30 000 bases and differs in appearance from a landscape made up of an uncorrelated sequence of bases (also called an uncorrelated random walk). Although we can clearly distinguish by eye between the correlated and the uncorrelated sequence, it is still necessary to analyse the data (Li and Kaneko 1992, Peng *et al.* 1992). In the case of critical phenomena, Buckingham, Fairbanks and Kellers were able to analyse data in such a way that the scale invariance was demonstrated quantitatively. In their now-classic graph, the specific heat as a function of temperature is plotted directly on three scales: degrees, millidegrees and microdegrees. The three graphs are fairly similar. A log-log plot gives a straight line over three or four orders of magnitude.

To have an analogous form of quantitative measurement of this landscape, we must find some measure of the fluctuation of this landscape and observe how it depends on the length scale over which the fluctuations are measured. We usually find that the variation is power law in nature, and that the exponent α characterizing that power-law behaviour has a value of $1/2$ for the uncorrelated sequence and for sequences with short-range correlations, and a value greater than or less than $1/2$ when there are long-range correlations.

While everyone can easily obtain DNA sequences, it is difficult to analyse their statistical properties. The sequence of base pairs displays a huge amount of non-stationarity; there are patches of the DNA molecule where there is an excess of one kind of base and other patches where there is an excess of another kind of base. We cannot assume that we have 'independent identically distributed random variables', which textbook methods handle well. So it seems that we have two possibilities at this point: either to give up or to make up new methods. We have devised a method

that is sensitive to non-stationarity: 'detrended fluctuation analysis'. It consists of making 'window-boxes' of a fixed size L , finding the straight-line regression that fits the data and measuring the fluctuations around that regression line (Peng *et al.* 1994). This is the key to avoiding the bad effect of non-stationarity. The analysis is repeated for successively larger window boxes, and a plot is made of the rms fluctuation around that trend line as a function of the window-box size. If the original sequence had no correlations, we would expect to find a straight line on log-log paper with a slope $1/2$. What we actually find is a straight line with a slope that is closer to $2/3$. This straight line extends over approximately three orders of magnitude, similar to that in the work by Buckingham, by Fairbanks and by Kellers on critical phenomena.

In biology, unlike physics, one plot does not a discovery make. We need to study more than one gene, and we need to use more than one method. When we directly measure the correlations or the power spectrum we find comparable exponents consistent with the exponent $2/3$.

Our work was independently verified by Richard Voss (1992), who repeated our work, but for the entire GenBank (about 25 000 sequences at that time), and not just for the 82 genes that we had first worked with, and more recently by Buldyrev *et al.* (1995), who studied even more sequences 3 years later. A group at the National Institutes of Health acquired the first sequence of an entire chromosome and made the analogous analysis for that chromosome; the graph and the fluctuations exhibit the same exponent of $2/3$, but the linearity extends not over three decades, but four.

The upshot is that very careful tests by a variety of research groups have confirmed that there is long-range correlation in DNA. This still leaves the question about coding against non-coding parts. If we look at the actual chromosome, we see that a chromosome consists of some parts that are coding interrupted by other parts that are non-coding. Even an individual gene inside has parts that are coding interrupted by parts that are non-coding. What kind of analysis would enable us to distinguish between the coding and non-coding?

If we return to our original landscape, make a heavy red line for the subset of each gene that is actually coding and analyse the correlation properties of just the coding parts; stitching together all the little coding regions, we see a landscape that differs greatly from the original. This new and different landscape shows no long-range correlations. (Note that these two landscapes are of a sample of non-human DNA and exhibit a 25%: 75% coding: non-coding ratio, unlike human DNA which has a 3%: 97% coding: non-coding ratio.) When we do this procedure repeatedly for these same 82 genes, we find that the full gene has a slope of $2/3$, but that the coding regions have a slope of only $1/2$.

What principle allows us to 'stitch together' coding regions in this analysis? Can we obtain data from regions that are 100% coding? Fortunately, in some simpler forms of life, one finds that the DNA is almost 100% coding. When we test that DNA, we find that the slope is indeed $1/2$.

Our discovery of a difference between coding and non-coding DNA has been confirmed recently by Arneodo *et al.* (1995). They analysed the same sequences as we did but compensated for the non-stationarity using wavelet methods. They found again that the non-coding DNA displays long-range power-law correlations and the coding DNA does not.

If there is a difference between the statistical properties of coding and non-coding DNA, we ought to be able to build an algorithm to distinguish between these two in

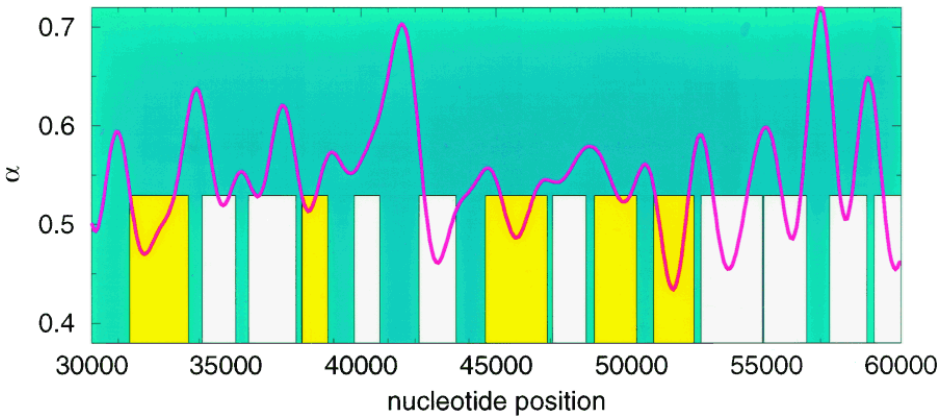


Figure 1. Beachcomber plot for a typical section containing about 10% of the yeast chromosome III, from base pair 30 000 to 60 000. The vertical yellow bars indicate the set of base pairs forming identified genes, while the white bars indicate less certain 'putative genes' determined from analysis of open reading frames. The exponent α is calculated by the beachcomber method (Ossadnik *et al.* 1994). 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 farther 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, namely a 'down' step for A or G (purines) and an 'up' step for C or T (pyrimidine). 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.

an unknown sequence of DNA. Why would we want to do that? Because the human genome project sequences DNA by machine; the machine spits out the letters and does not distinguish between coding and non-coding parts. Our algorithm slides down a DNA sequence with a window-box (covering, say, 1000 bases) that pauses after each successive base pair, allowing the computer to calculate the long-range exponent α for that window box. If the window-box overlaps a coding region, we would anticipate that $\alpha = 1/2$. If it does not overlap a coding region, we would anticipate that $\alpha > 1/2$. That indeed is what one finds; there is a signal that oscillates up and down that tends to 'dip' when there are coding regions, in this case a chromosome in which each of the genes is almost 100% coding (figure 1). This statistical device is not perfectly accurate (although its accuracy is comparable with that of other statistical methods) but it is less sensitive to errors. Occasional mess-ups do not change the outcome much.

2.3. What is the junk's purpose?

Three years ago at the Bar-Ilan Conference, the long-range correlation properties of natural languages were reported (Schenkel *et al.* 1993). The fact that long-range correlation properties are present in both non-coding DNA and in natural languages of course does not mean that non-coding DNA has a language, but it does suggest

that DNA has at least one property in common with natural languages. Because of Mantegna's work with our group (Mantegna *et al.* 1995), we have systematically applied methods of statistical linguistics to the coding and non-coding DNA. One conclusion is that non-coding DNA does indeed display some language-like features: a larger redundancy than coding DNA, and a power-law frequency-rank plot (in contrast, coding DNA has a logarithmic dependence of n-tuple frequency on rank).

§ 3. SCALE INVARIANCE OF THE HEALTHY HEARTBEAT

The interesting music that you heard before this paper was presented was created by a high-school student, Z. Davids, who worked in our research group one recent summer. In the composition process he allowed the time interval between each beat of a particular human heart (using an actual sequence of beats from an electrocardiogram) to determine the next pitch of the piece; if there was a long period between two beats, the algorithm would select a large change in pitch and, if the period was short, a small change in pitch. Interestingly, in the average human lifetime the heart seems to have the capacity to beat roughly three gigabeats without replacement or even repair. The time intervals between each beat of these three gigabeats also show a long-range power-law correlation, actually an anticorrelation (in the sense that, if there is, for whatever reason, an increase in beat rate at some point in time, there will be a corresponding decrease in rate in the future).

Usually, when we consider our heartbeat rate, we only pay attention to the average number of beats during some given time interval. For example, a nurse takes our pulse and tells us our heartbeat rate is 60 beats per minute.

However, just as in critical phenomena, in which there is information not only in the magnetization (the net number of spins up) but also in the fluctuations in the magnetization (which are directly proportional to the susceptibility of the magnet), in heartbeat rates there is information in the fluctuations of the time interval between each successive beat. If the heart beats 60 times per minute, then the interval is roughly 1 s per beat, but some intervals will be 0.95 s and others 1.05 s, and so on.

Using the same analysis methods that we used for DNA, we find that, in the healthy heartbeat sequence of a healthy heart, long-range power-law correlations are present (Peng *et al.* 1992, 1995; Ivanov *et al.* 1996, 1997). These correlations show considerable scale invariance and extend out for as long a period of time as the data record, typically one day or about 10^5 heartbeats.

§ 4. SCALE INVARIANCE IN LUNG INFLATION

In contrast with compact objects, scale-invariant or 'fractal' objects have a very large *surface* area. In fact, they are composed almost entirely of 'surface'. This observation explains why fractals are ubiquitous in biology, where surface phenomena (Bunde and Havlin 1994, 1996, Barabási and Stanley 1995) are of crucial importance.

Lungs exemplify this feature (figure 2). The surface area of a human lung is almost as large as a tennis court. The mammalian lung is made up of self-similar branches with many length scales, which is the defining attribute of a fractal surface. The efficiency of the lung is enhanced by this fractal property, since with each breath

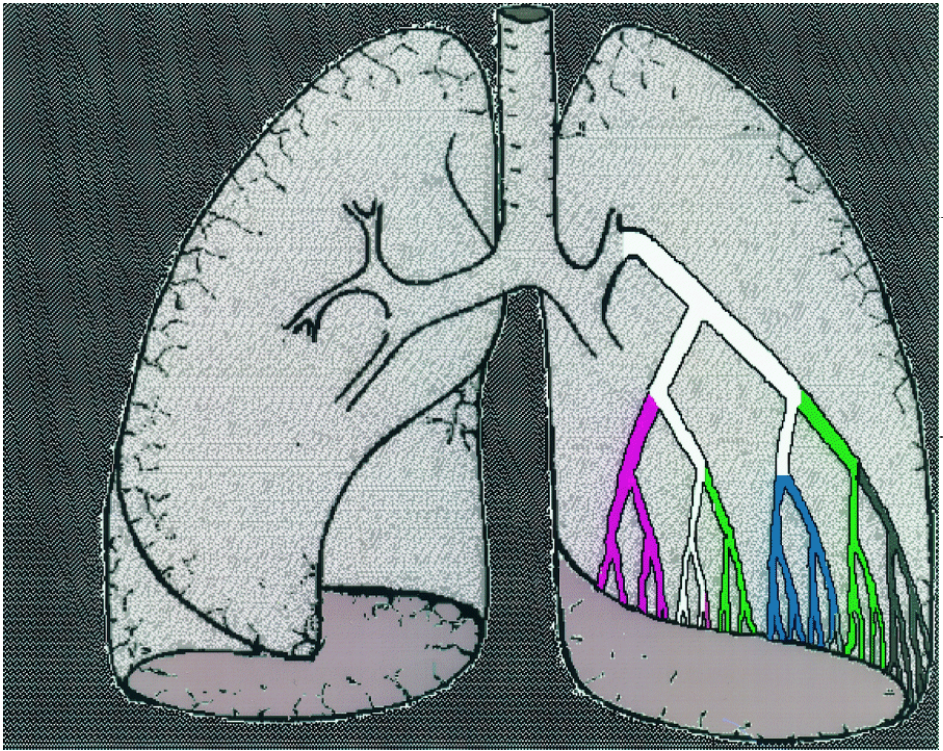


Figure 2. The dynamic mechanism responsible for filling the lung involves ‘avalanches’ or ‘bursts’ of air that occur in all sizes, instead of an exponential distribution, one finds a power-law distribution (Suki *et al.* 1994). The underlying cause of this scale-free distribution of avalanches is the fact that every airway in the lung has its own threshold below which it is not inflated. Shown here is a diagram of the development of avalanches in the airways during airway opening. The left side shows the classic view of the lung, but the right side is colour coded to show the successive avalanches. Note that the last avalanche (blue) opens up about 20% of the total lung volume, thereby significantly increasing the total surface area available for gas exchange.

oxygen and carbon dioxide have to be exchanged at the lung surface. The structure of the bronchial tree has been quantitatively analysed using fractal concepts (Shlesinger and West 1991). In particular, fractal geometry could explain the power-law decay of the average diameter of the airways with the generation number, in contrast with the classical model which predicts an exponential decay (Weibel and Gomez 1962).

Not only is the geometry of the respiratory tree described by fractal geometry, but also so are the time-dependent features of inspiration. Specifically, Suki and co-workers (Suki *et al.* 1994, 1997, Barabási *et al.* 1996, Sujeer *et al.* 1997) studied airway opening in isolated collapsed dog lungs. During constant-flow inflations, they found that the terminal airway resistances decreases in discrete jumps, and that the probability distribution function $\Pi(x)$ of the relative size x of the jumps and the probability distribution $\Pi(t)$ of the time intervals t between these jumps follow a power law over nearly two decades of x and t with exponents of 1.8 and 2.5 respectively. To interpret these findings, they developed a branching airway model in which airways, labelled ij , are closed with a uniform distribution of opening

threshold pressures P . When the ‘airway-opening’ pressure P_{ao} exceeds P_{ij} of an airway, that airway opens along with one or both of its daughter branches if $P_{ij} < P_{ao}$ for the daughters. Thus, the model predicts ‘avalanches’ of airway openings with a wide distribution of sizes, and the statistics of the jumps agree with those of $\Pi(x)$ and $\Pi(t)$ measured experimentally. They concluded that power-law distributions, arising from avalanches triggered by threshold phenomena, govern the recruitment of terminal airspaces.

Recently, it has become possible to solve numerically the full Navier–Stokes equations for an arbitrary geometry, using the FLUENT software package. Andrade *et al.* (1997a,b) have used this package to solve for a range of fluid flow problems. In particular, they have found a potential explanation to the open question concerning the morphogenesis of the lung structure. A commonly held belief is that the asymmetric structure of the lung arises solely from geometrical constraints, but Andrade *et al.* suggested a possible different origin for this structure, since the asymmetry of the bronchial tree can be a result of the fluid flow asymmetry combined with the requirement of homogeneous ventilation (figure 3).

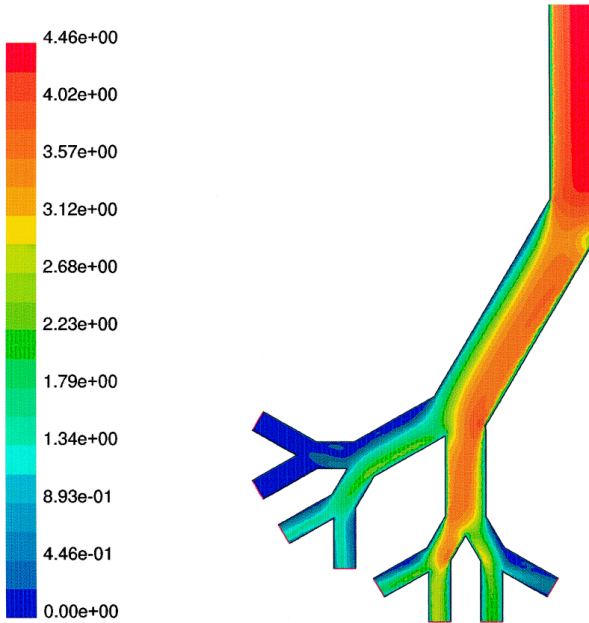


Figure 3. Contour plot of the stream function in a five-generation tree for a value of Reynolds number $Re = 4800$ in the range of normal breathing. The distribution of streamlines at the outlet branches is uniform at low Re , but highly non-uniform at high Re . The trachea is at the top of this figure, and only half of the lung is shown for reasons of structural symmetry. The colours correspond to the value of the air velocity. One might at first expect that the flow would bifurcate equally at each structural bifurcation but, as the simulations demonstrate, the flow ‘remembers’ the flow in its grandparents in the sense that it is larger for the branch of the bifurcation that is parallel to its grandparent, owing to the effects of inertia which become larger at high Re . The text discusses possible implications of this finding; in order to achieve homogeneous ventilation, it is possible that the lung evolved with a structural asymmetry such that the branches that receive small flows subtend proportionally small sub-regions.



Figure 4. Two wandering albatross, whose wing span is up to 4m, and which can circle the entire globe at low latitudes, travelling 8000m before returning to its 'family'. The flight patterns of these giant seabirds are found to have well defined statistical properties corresponding to those of a Levy flight (Peterson 1996, 1997, Viswanathan *et al.* 1996, 1997).

§ 5. SCALE INVARIANCE OF ANIMAL BEHAVIOUR

The wandering albatross, a giant seabird, was recently the subject of a popular work written by Peterson (1996, 1997) (figure 4). The occasion was an analysis done in collaboration with three workers at the British Antarctic Service, who have been leg-banding these birds with tracking devices (Viswanathan *et al.* 1996). On analysing the data, we found that the migratory paths of these birds obey Levy flight statistics, and recently we found that other foraging animals obey well defined statistical rules (Viswanathan *et al.* 1997).

Recently, Keitt and Stanley (1997) have applied to a 30 year data set on bird populations the same sort of techniques used to describe long-term data sets on economics and finance. They find statistical properties that are remarkably similar and consistent with the idea that 'every bird species interacts with every other bird species', just as the economic analysis supports the notion that 'every firm interacts with every other firm'. This empirical result is not without interest, since it serves to cast doubt on models of bird population (and of economic systems) in which one partitions the entire data set into strongly interacting and weakly interacting subsets and then ignores or oversimplifies the interactions in the weakly interacting subset.

§ 6. SCALE INVARIANCE IN HUMAN BEHAVIOUR: URBAN GROWTH PATTERNS

Predicting urban growth is important for the challenge that it presents to theoretical frameworks for cluster dynamics (Benguigui and Daoud 1991, Batty and Longley 1994, Benguigui 1995). Recently, the model of diffusion-limited aggregation

(DLA) has been applied to describe urban growth (Batty and Longley 1994) and results in tree-like dendritic structures which have a core or 'central business district' (CBD). The DLA model predicts that there exists only one large fractal cluster that is almost perfectly screened from incoming 'development units' (people, capital, resources, etc), so that almost all the cluster growth occurs in the extreme peripheral tips. In recent work (Makse *et al.* 1995) an alternative model to DLA that better describes the morphology and the area distribution of systems of cities, as well as the scaling of the urban perimeter of individual cities, has been developed. The results agree both qualitatively and quantitatively with actual urban data. The resulting growth morphology can be understood in terms of the effects of interactions among the constituent units forming a urban region and can be modelled using the correlated percolation model in the presence of a gradient.

In the model, one takes into account the following points:

- (i) Urban data on the population density $\rho(r)$ of actual urban systems are known to conform to the relation (Clark 1951) $\rho(r) = \rho_0 \exp(-\lambda r)$, where r is the radial distance from the CBD situated at the core, and λ is the density gradient. Therefore, in our model the development units are positioned with an occupancy probability $p(r) \equiv \rho(r)/\rho_0$ that behaves in the same fashion as is known experimentally.
- (ii) In actual urban systems, the development units are not positioned at *random*. Rather, there exist *correlations* arising from the fact that, when a development unit is located in a given place, the probability of adjacent development units increases naturally, that is each site is not independently occupied by a development unit but is occupied with a probability that depends on the occupancy of the neighbourhood.

In order to quantify these ideas, we consider the *correlated* percolation model (Coniglio *et al.* 1977, Prakash *et al.* 1992). In the limit where correlations are so small as to be negligible, a site at position \mathbf{r} is occupied if the occupancy variable $u(\mathbf{r})$ is smaller than the occupation probability $p(\mathbf{r})$; the variables $u(\mathbf{r})$ are uncorrelated random numbers. To introduce correlation among the variables, we convolute the uncorrelated variables $u(\mathbf{r})$ with a suitable power-law kernel (Prakash *et al.* 1992) and define a new set of random variables $\eta(\mathbf{r})$ with long-range power-law correlations that decay as $r^{-\alpha}$, where $r \equiv |\mathbf{r}|$. The assumption of power-law interactions is motivated by the fact that the 'decision' for a development unit to be placed in a given location decays gradually with the distance from an occupied neighbourhood. The correlation exponent α is the only parameter to be determined by empirical observations.

To discuss the morphology of a system of cities generated in the present model, we performed simulations of correlated urban systems for a fixed value of the density gradient λ , and for different degree of correlations. The correlations have the effect of agglomerating the units around a urban area. In the simulated systems the largest city is situated in the core, which is regarded as the attractive centre of the city, and is surrounded by small clusters or 'towns'. The correlated clusters are nearly compact near their centres and become less compact near their boundaries, in qualitative agreement with empirical data on actual large cities such as Berlin, Paris and London (Batty and Longley 1994, Frankhauser 1994).

So far, we have argued how correlations between occupancy probabilities can account for the irregular morphology of towns in a urban system. As can be seen in

figure 5, the towns surrounding a large city such as Berlin are characterized by a wide range of sizes. We are interested in the laws that quantify the town size distribution $N(A)$, where A is the area occupied by a given town or ‘mass’ of the agglomeration; so we calculate the actual distribution of the areas of the urban settlements around Berlin and London and find that, for both cities, $N(A)$ follows a power law.

This new result of a power-law area distribution $N(A)$ of towns can be understood in the context of our model. Insight into this distribution can be developed by first noting that the small clusters surrounding the largest cluster are all situated at distances r from the CBD such that $p(r) < p_c$ or $r > r_f$. Therefore, we find $N(A)$, the cumulative area distribution of clusters of area A , to be

$$N(A) \equiv \int_0^{p_c} n(A, p) dp \sim A^{-(\tau+1/d_f v)}.$$

Here, $n(A, p) \sim A^{-\tau} g(A/A_0)$ is defined to be the average number of clusters containing A sites for a given p at a fixed distance r , and $\tau = 1 + 2/d_f$. Here,

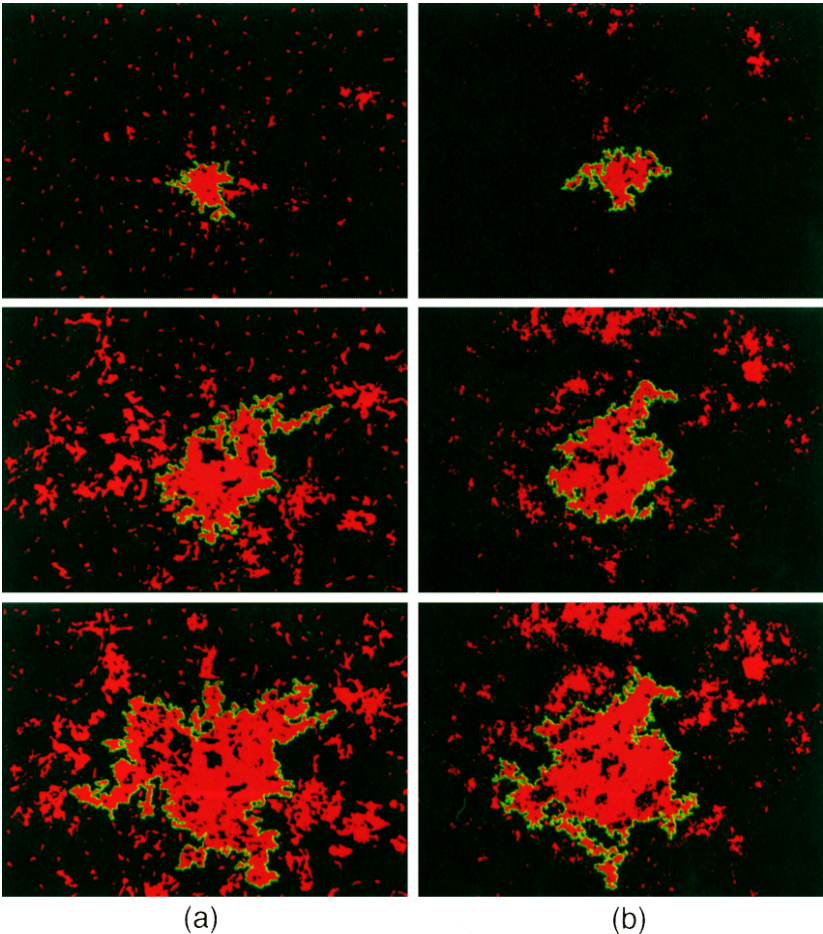


Figure 5. Qualitative comparison between actual urban data and the predictions of the correlated gradient percolation model. (a) Three steps of the growth with time of Berlin and surrounding towns. Data are shown for the years 1875, 1920 and 1945 (from top to bottom). (b) Dynamical urban simulations of the model.

$A_0(r) \sim |p(r) - p_c|^{-d\nu}$ corresponds to the maximum typical area occupied by a cluster situated at a distance r from the CBD, while $g(A/A_0)$ is a scaling function that decays rapidly (exponentially) for $A > A_0$. The exponent $\nu = \nu(\alpha)$ is defined by $\xi(r) \sim |p(r) - p_c|^{-\nu}$, where $\xi(r)$ is the connectedness length that represents the mean linear extension of a cluster at a distance $r > r_f$ from the CBD.

§ 7. SCALE INVARIANCE OF HUMAN BEHAVIOUR: FINANCE AND ECONOMICS

About 35 years ago, Benoit Mandelbrot (1963) wrote an article about fluctuations in cotton prices. This proved to be a seminal work and has been described in many popular books about fractals. In it, he points out the possibility of scaling in financial indices. We have extended his analysis to data sets available now (Mantegna and Stanley 1995) and confirmed the presence of scale invariance (figure 6). Furthermore, it appears that the distribution function conforms to a truncated Levy flight distribution (a Levy distribution with an exponential truncation in the wings) (Mantegna and Stanley 1994). Recently, the general approach of Mantegna and Stanley has been extended to study the scale invariance of one measure of the volatility of a financial index (Cizeau *et al.* 1997, Liu *et al.* 1997).

Economics is different from finance, and we have also looked at economic data. Specifically, in collaboration with a card-carrying economist, Michael Salinger, we

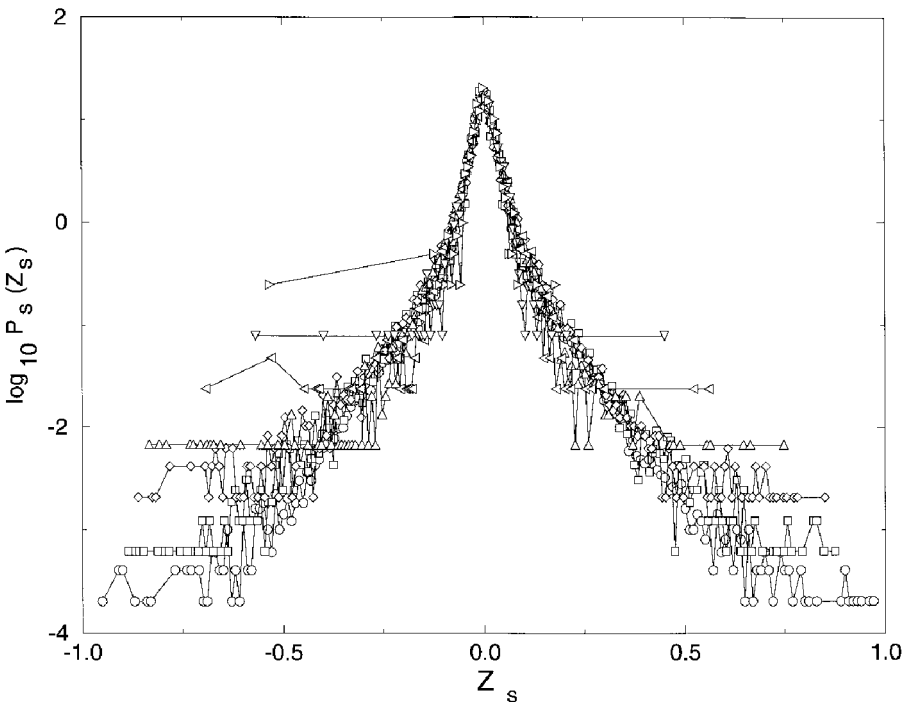


Figure 6. Demonstration of scale invariance in financial time series. Mantegna and Stanley (1995) analysed the probability distribution $P(Z)$ of the S&P Index variations $Z(t)$ observed at time intervals Δt , which ranges from 1 to 1000 min. By increasing Δt , a spreading of the probability distribution characteristic of a random walk is observed. Shown is a scaled plot of the probability distributions shown. All the data collapse onto the $\Delta t = 1$ min distribution by using the scaling transformations appropriate to those of a Levy distribution, with $\alpha = 1.40$. The points outside the average behaviour define the noise level of that specific distribution.

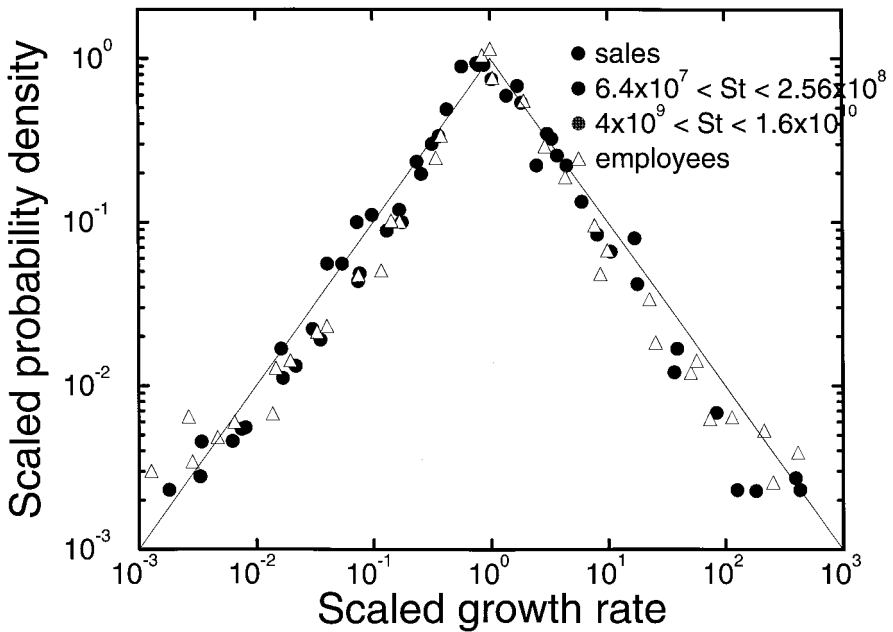


Figure 7. Stanley *et al.* (1996) analysed the fluctuations in the 1 year growth rates of the sales and of the 1 year growth rates of the number of employees as functions of the initial values. They calculated histograms for the probability that a firm grows at a rate r from one year to the next. They found distinct histograms depending on the size of the firm, with larger firms having histograms with a smaller standard deviation. Shown is the data collapse obtained when the histograms are scaled by this standard deviation.

studied the possibility that all the companies in a given economy might interact, more or less, like an Edwards–Anderson spin glass. As in an Edwards–Anderson spin glass, each spin interacts with another spin, but not with the same coupling and not even with the same sign.

If the sales in a given company x decreases by, for example, 10%, it will have repercussions in the economy. Some of the repercussions will be favourable; company y , who competes with x , may experience an increase in market share. Others will be negative; service industries that provide personal services for company x employees may experience a drop-off in sales as employee salaries decline. There are positive and negative correlations for almost any economic change. The notion that we can view the economy as a complicated Ising system is possibly as old as Mandelbrot's first work in this area.

To approach this interesting bit of statistical 'poetry' and make sense of it, we first located and secured a database that listed the actual size of every firm in the USA. With this information, we did an analysis to determine how the distribution of firm size changes from one year to the next. We then made a histogram for each of three characteristic firm sizes. The largest firms have a very narrow distribution, plausible because the percentage of size change from year to year for the largest firms cannot be that great. On the other hand, a tiny company or a garage-based start-up can radically increase (or decrease) in size from year to year. The histograms have a width determined by the size of the firm. When this width is plotted on the y axis of log–log paper as a function of the size of the firm on the x axis, the data are

approximately linear *over eight orders of magnitude*, from the tiniest firms in the database to the largest. The width scales as the firm size to an exponent β , with $\beta \approx 1/6$. We can therefore normalize the growth rate and show that all the data collapse onto a single curve, demonstrating the scaling of this measure of firm size (figure 7).

Why does this occur? We are working on that. We model this firm structure as an approximate Cayley tree, in which each subunit of a firm reacts to its directives from above with a certain probability distribution. This model, developed primarily by Sergey Buldryev, seems to be consistent with the critical exponent $-1/6$ (Buldyrev *et al.* 1997). More recently, Amaral *et al.* (1997a,b) have proposed a microscopic model.

§ 8. DISCUSSION

What is the point of this paper? Just to show that many different systems appear to develop scale-invariant correlations? If so, how do we understand this empirical fact?

Bak's idea that systems self-organize themselves such that they are in effect near a critical point is an appealing unifying principle. Near a critical point, there is a delicate balance between the exponentially growing number of different one-dimensional paths connecting any two faraway subunits and the exponentially decaying correlations along each one-dimensional path (this concept is illustrated, for example, in figure 9.4 in the book by Stanley (1971)). If the control parameter (say, the coupling constant) is too small, the correlations die out so fast along each one-dimensional path that subunits far from one another are not well correlated. However, at a critical point, the exponentially large number of paths connecting each pair of subunits is sufficient to balance out the exponential decay along each path and the 'correction factor' wins out; this correction factor is the power law that governs the total number of one-dimensional paths connecting two distant subunits. The exponent in this correction factor depends primarily on the system dimension, and not at all on the actual arrangement of the subunits (lattice or no lattice).

Could it be that somehow biological and social systems push themselves 'up to the limit', just as a sandpile is pushed to the limit before an avalanche starts, an image that has attracted recent attention in the debate between 'self-organized criticality' and 'plain old criticality (for example Vespignani and Zapperi (1997) and references therein)? For example, in economics every subunit plays according to rules and pushes itself up against the limits imposed by these rules, but social systems display a variety of rich forms of 'order', far richer than we anticipate from studies of ferromagnets and antiferromagnets (see, for example, some of the papers appearing in the work by Knobler *et al.* (1997)). Could such orderings arise from the complex nature of the interactions or from the range of different 'sizes' of the constituent subunits in the same way as, for example, one finds ordering in sandpiles when sand particles of two different grain sizes are dropped onto a heap (for example Makse *et al.* (1997a, 1997b))? These are questions that occupy us now, and questions that I would be delighted to discuss with any readers.

ACKNOWLEDGEMENTS

We are grateful to many individuals, and most especially to A. L. Goldberger, H. Leschhorn, P. Maass, R. N. Mantegna, M. E. Matsa, S. M. Ossadnik, M. A. Salinger, F. Sciortino, M. Simons and M. H. R. Stanley for major contributions to

those results reviewed here that represent collaborative research efforts. We also wish to thank C. Cantor, C. DeLisi, M. Frank-Kamenetskii, A. Yu. Grosberg, G. Huber, I. Labat, L. Liebovitch, G. S. Michaels, P. Munson, R. Nossal, R. Nussinov, R. D. Rosenberg, J. J. Schwartz, M. Schwartz, E. I. Shakhnovich, M. F. Shlesinger, N. Shworak and E. N. Trifonov for valuable discussions. Partial support was provided by the National Science Foundation, by Conselho Nacional de Desenvolvimento Científico e Tecnológico, by National Institutes of Health (NIH) (Human Genome Project), by the G. Harold and Leila Y. Mathers Charitable Foundation, by the National Heart Lung and Blood Institute, by the National Aeronautics and Space Administration, by the Israel–USA Binational Science Foundation, by the Israel Academy of Sciences, and (to C.-K.P.) by an NIH First Award.

REFERENCES

- AMARAL, L. A. N., BULDYREV, S. V., HAVLIN, S., LESCHHORN, H., MAASS, P., SALINGER, M. A., STANLEY, H. E., and STANLEY, M. H. R., 1997a, *J. Phys., Paris*, **1**, 7, 621.
- AMARAL, L. A. N., BULDYREV, S. V., HAVLIN, S., SALINGER, M. A., and STANLEY, H. E., 1997b, *Phys. Rev. Lett.*, **80**, 1385.
- ANDRADE, J. S., ALENCAR, A. M., ALMEIDA, M. P., MENDES FILHO, J., BULDYREV, S. V., ZAPPERI, S., STANLEY, H. E., and SUKI, B., 1997a, *Phys. Rev. Lett.* (submitted).
- ANDRADE, J. S., JR, ALMEIDA, M. P., MENDES FILHO, J., HAVLIN, S., SUKI, B., and STANLEY, H. E., 1997b, *Phys. Rev. Lett.*, **79**, 3901.
- ANDREWS, T., 1869, *Phil. Trans.*, **159**, 575.
- ARNEODO, A., BACRY, E., GRAVES, P. V., and MUZY, J. F., 1995, *Phys. Rev. Lett.*, **74**, 3293.
- BARABÁSI, A.-L., BULDYREV, S. V., STANLEY, H. E., and SUKI, B., 1996, *Phys. Rev. Lett.*, **76**, 2192.
- BARABÁSI, A.-L., and STANLEY, H. E., 1995, *Fractal Concepts in Surface Growth* (Cambridge University Press).
- BATTY, M., and LONGLEY, P., 1994, *Fractal Cities* (San Diego, California: Academic Press).
- BENGUIGUI, L., 1995, *Physica A*, **219**, 13.
- BENGUIGUI, L., and DAUD, M., 1991, *Geogr. Anal.*, **23**, 362.
- BULDYREV, S. V., AMARAL, L. A. N., HAVLIN, S., LESCHHORN, H., HAASS, P., SALINGER, M. A., STANLEY, H. E., and STANLEY, M. H. R., 1997, *J. Phys., Paris*, **1**, 7, 635.
- BULDYREV, S. V., GOLDBERGER, A. L., HAVLIN, S., MANTEGNA, R. N., MATSA, M. E., PENG, C.-K., SIMONS, M., and STANLEY, H. E., 1995, *Phys. Rev. E*, **51**, 5084.
- BUNDE, A., and HAVLIN, S. (editors), 1994, *Fractals in Science* (Berlin: Springer); 1996, *Fractals and Disordered Systems*, second edition (Berlin: Springer).
- CIZEAU, P., LIU, Y., MEYER, M., PENG, C.-K., and STANLEY, H. E., 1997, *Physica*, **245**, 441.
- CLARK, C., 1951, *J. R. statist. Soc. A*, **114**, 490.
- CONIGLIO, A., NAPPI, C., RUSSO, L., and PERUGGI, F., 1977, *J. Phys. A*, **10**, 205.
- FRANKHAUSER, P., 1994, *La Fractalité des Structures Urbaines* (Paris: Collection Villes, Anthropos).
- IVANOV, P. CH., AMARAL, L. A. N., GOLDBERGER, A. L., and STANLEY, H. E., 1997, *Phys. Rev. Lett.* (submitted).
- IVANOV, P. CH., ROSENBLUM, M. G., PENG, C.-K., MIETUS, J., HAVLIN, S., STANLEY, H. E., and GOLDBERGER, A. L., 1996, *Nature*, **383**, 323.
- KNOBLER, C. M., ROBLEDO, A., and STANLEY, H. E., 1997, *Statistical Mechanics in the Physical, Biological, and Social Sciences: Festschrift in Honor of Benjamin Widom on the Occasion of his 70th Birthday* (Amsterdam: Elsevier).
- KEITT, T., and STANLEY, H. E., 1997, *Nature* (submitted).
- LI, W., and KANEKO, K., 1992, *Europhys. Lett.*, **17**, 655.
- LIU, Y., CIZEAU, P., MEYER, M., PENG, C.-K., and STANLEY, H. E., 1997, *Physica*, **245**, 437.
- MAKSE, H. A., CIZEAU, P., and STANLEY, H. E., 1997a, *Phys. Rev. Lett.*, **78**, 3298.
- MAKSE, H. A., HAVLIN, S., KING, P. R., and STANLEY, H. E., 1997b, *Nature*, **386**, 379.
- MAKSE, H. A., HAVLIN, S., and STANLEY, H. E., 1995, *Nature*, **377**, 608.
- MANDELBROT, B. B., 1963, *J. Business*, **36**, 394.

- MANTEGNA, R. N., BULDYREV, S. V., GOLDBERGER, A. L., HAVLIN, S., PENG, C.-K., SIMONS, M., and STANLEY, H. E., 1995, *Phys. Rev. E*, **52**, 2939.
- MANTEGNA, R. N., and STANLEY, H. E., 1994, *Phys. Rev. Lett.*, **73**, 2946; 1995, *Nature*, **376**, 46.
- OSSADNIK, S. M., BULDYREV, S. V., GOLDBERGER, A. L., HAVLIN, S., MANTEGNA, R. N., PENG, C.-K., SIMONS, M., and STANLEY, H. E., 1994, *Biophys. J.*, **67**, 64.
- PENG, C.-K., BULDYREV, S. V., GOLDBERGER, A. L., HAVLIN, S., SCIORTINO, F., SIMONS, M., and STANLEY, H. E., 1992, *Nature*, **356**, 168.
- PENG, C.-K., BULDYREV, S. V., HAVLIN, S., SIMONS, M., STANLEY, H. E., and GOLDBERGER, A. L., 1994, *Phys. Rev. E*, **49**, 1685.
- PENG, C.-K., HAVLIN, S., STANLEY, H. E., and GOLDBERGER, A. L., 1995, *Chaos*, **5**, 82.
- PENG, C.-K., MIETUS, J., HAUSDORFF, J. M., HAVLIN, S., STANLEY, H. E., and GOLDBERGER, A. L., 1993, *Phys. Rev. Lett.*, **70**, 1343.
- PERRIN, J., 1915, *Atoms* (Cambridge University Press).
- PETERSON, I., 1996, *Sci. News*, **150**, 104; 1997, *The Jungles of Randomness* (New York: Viking).
- PRAKASH, S., HAVLIN, S., SCHWARTZ, M., and STANLEY, H. E., 1992, *Phys. Rev. A*, **46**, R1724.
- SCHENKEL, A., ZHANG, J., and ZHANG, Y.-C., 1993, *Fractals*, **1**, 47.
- SHLESINGER, M. F., and WEST, B. J., 1991, *Phys. Rev. Lett.*, **67**, 2106.
- STANLEY, H. E., 1971, *Introduction to Phase Transitions and Critical Phenomena* (Oxford University Press).
- STANLEY, M. H. R., AMARAL, L. A. N., BULDYREV, S. V., HAVLIN, S., LESCHHORN, H., MAASS, P., SALINGER, M. A., and STANLEY, H. E., 1996, *Nature*, **379**, 804.
- SUKI, B., ANDRADE, J. S., JR., COUGHLIN, M., STAMENOVIC, D., STANLEY, H. E., SUJEER, M., and ZAPPERI, S., 1997, *Ann. Biomed. Engng*, **26**, 1.
- SUKI, B., BARABÁSI, A.-L., HANTOS, Z., PETÁK, F., and STANLEY, H. E., 1994, *Nature*, **368**, 615.
- SUJEER, M. K., BULDYREV, S. V., ZAPPERI, S., ANDRADE, J., STANLEY, H. E., and SUKI, B., 1997, *Phys. Rev. E*, **56**, 3385.
- VISWANATHAN, G. M., AFANASYEV, V., BULDYREV, S. V., MURPHY, E. J., PRINCE, P. A., and STANLEY, H. E., 1996, *Nature*, **381**, 413.
- VISWANATHAN, G. M., BULDYREV, S. V., HAVLIN, S., DA LUZ, M. G., RAPOSO, E., and STANLEY, H. E., 1997, *Phys. Rev. Lett.* (submitted).
- WEIBEL, E. R., and GOMEZ, D. M., 1962, *Science*, **137**, 577.
- VESPIGNANI, A., and ZAPPERI, S., 1997, *Phys. Rev. Lett.*, **78**, 4793.
- VOSS, R., 1992, *Phys. Rev. Lett.*, **68**, 3805.

Copyright of Philosophical Magazine B is the property of Taylor & Francis Ltd and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.