

What is complexity?

Christoph Adami^{1,2}

Summary

Arguments for or against a trend in the evolution of complexity are weakened by the lack of an unambiguous definition of complexity. Such definitions abound for both dynamical systems and biological organisms, but have drawbacks of either a conceptual or a practical nature. Physical complexity, a measure based on automata theory and information theory, is a simple and intuitive measure of the amount of information that an organism stores, in its genome, about the environment in which it evolves. It is argued that physical complexity *must* increase in molecular evolution of asexual organisms in a single niche if the environment does not change, due to natural selection. It is possible that complexity decreases in co-evolving systems as well as at high mutation rates, in sexual populations, and in time-dependent landscapes. However, it is reasoned that these factors usually help, rather than hinder, the evolution of complexity, and that a theory of physical complexity for co-evolving species will reveal an overall trend towards higher complexity in biological evolution. *BioEssays* 24:1085–1094, 2002. © 2002 Wiley Periodicals, Inc.

Introduction

Whether or not complexity increases in evolution is one of the central questions of evolutionary biology. Opinions about this subject vary, but generally belong to one of three camps: one suggests that complexity has increased, another claims that there is not enough evidence to argue for or against an increase, and a third denies that “progress characterizes the history of life as a whole, or even represents an orienting force in evolution at all”.⁽¹⁾ Often, these camps disagree not only about the existence of a trend, but also on what type of complexity measure to use, and whether maximum or average complexity is pertinent. Most agree however that nobody

knows precisely what is meant by the word “complexity” when referring to a biological organism. Indeed, while complexity measures abound (many of them invented by physicists, Ref. 2), their relationship to biology is not always clear. I will review here several different kinds of complexity measures (without trying to be exhaustive), and then focus on a recent measure that appears to capture what we intuitively expect from such a measure in biology, and discuss what it implies about a trend in the evolution of complexity.

Complexity is so general a term that it seems to mean something different to everyone. The two main “user-groups” are physicists interested in dynamical systems, and biologists pondering the above-mentioned question of a trend. Perhaps we should expect that a measure exists that is so general that it applies to both biology *and* dynamical systems but, in the meantime, it may be useful to clearly separate the two (at least until we succeed in mapping the behavior of an animal to a dynamical system, a prospect that surely is far off).

In dynamical systems theory, we are interested in the complexity of processes. For example, periodic and random processes are both perceived as simple, with the random processes at “the other end of the scale”, whatever that scale may be. Complex and chaotic processes are deemed to lie somewhere in between. This ordering along a scale supports the general idea of a relationship between structure and complexity, as the consensus is that neither periodic nor random processes possess any structure.

Because all processes can in principle be viewed as computations (and vice versa), complexity measures in dynamical systems theory are usually *computational complexities*. Such constructions allow you to infer the complexity of a sequence of symbols by finding an appropriate *finite state machine* that produces this sequence. (A finite state machine is an abstract automaton that can take on only a finite number of states.) One such measure, called “thermodynamical depth” by Lloyd and Pagels⁽³⁾ attempted to capture “how hard it is to put something together”, but ended up just characterizing the randomness that a process generates.⁽⁴⁾ A measure that satisfies our craving for a “one-humped” criterion (low complexity for both ordered and random systems, with high complexity for those in between) is the *statistical complexity* of Crutchfield and Young.⁽⁵⁾ Their measure has the added bonus of being practical because the statistical complexity can be inferred from observations of the statistics of the sequences that their machine produces. Nevertheless, this measure suffers from a problem that most of the *sequence complexities* (see below)

¹Digital Life Laboratory 136-93, California Institute of Technology, Pasadena.

²Jet Propulsion Laboratory 126-347, California Institute of Technology, Pasadena, CA 91109, USA. E-mail: adami@caltech.edu

Funding agencies: The National Science Foundation Biocomplexity program under contract No. DEB-9981397. Part of this work was carried out at the Jet Propulsion Laboratory, California Institute of Technology, under a contract with the National Aeronautics and Space Administration.

DOI 10.1002/bies.10192

Published online in Wiley InterScience (www.interscience.wiley.com).

have: It characterizes the amount of information necessary to predict the future state of the machine (or the next symbol in a symbolic sequence), but it fails to address their meaning in a complex world.

The complexity of biological organisms cannot as yet be captured by attempting to characterize the dynamics of all their underlying processes. Instead, biological complexity measures refer either to form, function, or the sequence that codes for it. *Structural complexity* is generally what we mean when we consider animals, but this seems to be the hardest measure to define. McShea⁽⁶⁾ has studied several measures of structural complexity, based on number of cell types, different limb-pair types, and even the fractal dimension of sutures in ammonoids, and found some evidence for a trend in these indicators, but nothing as conclusive as one might have anticipated. Bell and Mooers⁽⁷⁾ (see also Ref. 8) found that the diversity of specialized cell types increases with body size, but this correlation is not unexpected, nor does it make a case for a general trend in evolution (some evidence for Cope's rule notwithstanding, Ref. 9). McShea has also made the case for a measure of *functional complexity* of organisms⁽¹⁰⁾ that counts the number of different functions an organism can perform. While such a measure is certainly intuitively satisfying, its drawbacks are obviously the difficulty of defining a range of functions, separating them into distinct non-overlapping ones, and plainly missing some that are not immediately obvious. We should also keep in mind that complexity should not be equated with *evolutionary success*⁽¹¹⁾ a misconception that has led to many controversies. Finally, it seems to be obvious that, during the course of evolution, the number of levels of *nestedness* (number of lower-level entities nested in higher-level ones) has increased. Thus, this nestedness could perhaps be used as a measure of *hierarchical complexity*, but again McShea⁽¹¹⁾ points out the difficulties in using such a measure and documenting a genuine trend.

It is hard to imagine that a *universal* measure for structural or functional complexity can be devised, given that organisms differ so greatly in form and function. However, all these differences are sidestepped when we consider the nucleic acid sequences from whence all creatures derive. Of course, we understand that the difficulty of biology lies precisely in the intricacy of this map from sequence to function. Nevertheless, it is very likely that a properly defined *sequence complexity* should mirror the complexity of the organism that the sequence gives rise to. If this is so (and at this juncture this is only a conjecture) then the problem of defining structural or functional complexity can be demoted to the problem of defining sequence complexity, which is naturally much simpler because sequences are amenable to a mathematical characterization. Many of the complexity measures introduced in Ref. 2 are in fact sequence complexities. Most of them, however, do not appear satisfactory from an intuitive point of view. One of the measures most often put forward as a

candidate, the Kolmogorov complexity (see, e.g., Ref. 2), turns out to be a measure of the *regularity*, rather than complexity, of a sequence. This implies that a random sequence is accorded maximum Kolmogorov complexity, clearly not anything we would be interested in as biologists, because random sequences do not give rise to organisms.

Other sequence complexities, such as Grassberger's *effective measure complexity*,⁽¹²⁾ suffer from the same problem as some of the statistical complexities mentioned above. These measures attempt to characterize short-range and long-range correlations in sequences in such a manner that they optimally predict the *next symbol on the sequence*. But we know that the sequence that codes for a functional protein is completely unrelated to its function, and therefore correlations *among* symbols in a sequence are uninteresting for the purpose of measuring biological complexity. Instead, we should look for a correlation of these symbols with *features of the environment within which this sequence is functional*. In other words, rather than looking for *vertical correlations* between symbols (correlations between symbols along the sequence), we should be looking for *horizontal correlations*, namely correlations between the symbols in a genome and a description of the environment within which that sequence is functional. The *physical complexity* that I introduce below is just such a measure.

Physical complexity⁽¹³⁾ is a measure of sequence complexity that is carefully defined from an automata-theoretic point of view (just as Kolmogorov complexity was), but it has a very simple relationship to information theory, and turns out to be very intuitive. Furthermore, it appears to correspond exactly to what biologists think is increasing when "self-organizing systems organize themselves." Because such a measure can also be applied to sequences of symbols generated by a dynamical system, there is hope that it may bridge the traditional gap between the physical and biological sciences. Rather than starting with the mathematical definition, I will instead describe the intuitive notion, and connect it with the mathematical definition later. The latter is important to clarify the circumstances under which physical complexity can be measured, and to outline the assumptions and approximates going into such an estimate.

In the following, I argue that physical complexity must increase in molecular evolution under certain circumstances,⁽¹⁴⁾ due to the actions of natural selection. This will be illustrated with experiments conducted with digital organisms. Because the circumstances under which the law holds exactly seem so restrictive as to rule out all realistic situations, I discuss how the law of increasing complexity is manifested in real biological systems, and point out the role of co-evolution. Even though the law can be broken (as we know that it must be and has been) we expect it to be responsible for the general trend that has led us from pools of replicating molecules, through prokaryotes, to eukaryotes and multicellular organisms.

Physical complexity

The physical complexity of a sequence refers to the amount of information that is stored in that sequence *about* a particular environment. For a genome, this environment is the one in which it replicates and in which its host lives, a concept roughly equivalent to what we call a *niche*. The definition of physical complexity must be distinguished from *mathematical* (or algorithmic, or Kolmogorov) complexity, which is only concerned with the intrinsic regularity (or, in this case, irregularity) of a sequence. The regularity of a sequence is a reflection of the unchanging laws of mathematics, and not of the physical world in which such a sequence may mean something. Information, on the other hand, is always *about something*. Consequently, a sequence may embody information about one environment (niche) while being essentially random with respect to another. This makes the measure *relative*, or conditional on the environment, and it is precisely this feature that brings a number of important observations that are incompatible with a universal increase in complexity in line with a law of increasing physical complexity.

Randomness is in some ways the “flip side” of information, and is called *entropy* in information theory.⁽¹⁵⁾ Entropy is a measure of potential knowledge, or if applied to a sequence, a measure of how much information a sequence *could* hold, and thus quantifies our uncertainty about the genetic identity of a randomly selected individual from a pool. It is useful to think of sequence entropy as the *length* of a tape, while information is the length of tape containing recordings. Measurement (i.e., recording) turns empty tape into filled tape; entropy into information. As we shall see, this is what happens during adaptation, and it is the force that drives the increase of complexity.

Information is a statistical form of correlation, and thus requires, mathematically and intuitively, a reference to the system that the information is *about*. The sequence on your information-filled tape allows you to make predictions about the state of the system that the sequence is information about. This predictive capability implies that your sequence and the system have something in common, that they are correlated. Your sequence will most likely *not* make predictions about any other system (unless the systems are very similar). If you do not know which system your sequence refers to, then whatever is on it *cannot* be considered information. Instead, it is *potential information* (a.k.a. entropy). This is the fundamental difference between entropy and information, often misrepresented in the literature.⁽¹⁶⁾

Information-theoretic measures of complexity have been considered before, only to be discarded because of erroneous uses of the concept. Most often, *entropy* is used as a candidate for information-theoretic complexity. From the previous discussion, we realize that the entropy of a sequence is the amount of information that it could possibly carry. Of course, this is just the length of the sequence. But it was recognized early on that sequence length is not a good predictor of orga-

nism complexity (the C-paradox), an observation that has discredited information-theoretic approaches to complexity. Physical complexity, a true measure of information, does not suffer from this handicap.

Nonmathematical, intuitive descriptions of complexity often make use of a concept very much akin to the one presented here, namely the idea of *horizontal* correlation. Most often, this is described as “genes embody knowledge about their niche” (Deutsch, Ref. 17) or, as put eloquently by Wilson: “(Organisms) encode the predictable occurrence of nature’s storms in the letters of their genes.⁽¹⁸⁾” This is precisely what physical complexity measures, since physical complexity *is* information about the environment that can be used to make predictions about it. Being able to predict the environment allows an organism to exploit it for survival. In such a manner, physical complexity translates into fitness for the organism. Let us now proceed to the mathematical definition of physical complexity. Such a definition is important because it immediately suggests how complexity can be measured in real adapting populations. I will refer to previous articles^(13,14) for technical points not immediately relevant for the present non-technical discussion.

Technically, physical complexity is defined as the *shared Kolmogorov complexity* between a sequence, and a description of the environment in which that sequence is to be interpreted.⁽¹³⁾ The details of this definition are not important here, in particular because this definition is not practical, since it does not allow the unambiguous determination of sequence complexity from available data. However, it is worth mentioning that it is an instance of *effective complexity*, a concept independently developed by Gell-Mann and Lloyd.⁽¹⁹⁾ When physical complexity is averaged over an *ensemble of sequences*, on the contrary it does become practical, because average mutual (or shared) Kolmogorov complexity is, in the limit of perfect coding, simply equal to the amount of information that the ensemble has about the environment to which it adapts. Perfect coding, in information theory, refers to the limit in which information is coded without loss or waste into a sequence. If this limit is achieved, information is perfectly compressed. Needless to say, this limit is rarely (if ever) achieved in nature, and we will be considering the consequences of imperfect coding (in the form of epistasis) later on.

At this juncture, it is sufficient to think of the physical complexity of a sequence as the amount of information that is coded in the genomes of an adapting population, about the environment to which it is adapting¹. This information is given by the difference between the entropy of the population in the *absence* of selection, and the entropy of the population *given*

¹In the following, my usage of the term “physical complexity of a sequence” should be taken to mean “average physical complexity of an ensemble of sequences” (since that is the only experimentally accessible quantity). We shall not return to the abstract mutual Kolmogorov complexity.

the environment, that is, given the selective forces that the environment engenders. In the section below, I give a technical exposition of the complexity measure. Readers who are satisfied with the intuitive description can skip this section without loss.

Measuring complexity

Because entropies of populations can be measured, the average physical complexity is a practical measure. The entropy of an ensemble (i.e., a population) of sequences X , in which sequences s_i occur with probabilities p_i , is denoted by the symbol $H(X)$ and calculated as

$$H(X) = - \sum_{i=1} p_i \log p_i. \quad (1)$$

The sum in (1) goes over all the different genotypes i in ensemble X . Whether or not selection acts on sequences of the ensemble is crucial for the entropy. When selection does not act, all sequences are equally probable in ensemble X (because in the absence of selection no sequence has an advantage over another). In this case, the probabilities p_i are each equal to the inverse population size, and the entropy takes on its maximal value

$$H_{\max}(X) = - \sum_{i=1}^N (1/N) \log (1/N) = \log N. \quad (2)$$

In an infinite population, the number of all possible genotypes is given by the size of the monomer alphabet, D , to the power of the length of the sequence, L , i.e.,

$$N = D^L. \quad (3)$$

If we agree to take logarithms to the base of the alphabet size, then the *unconditional* entropy of a population of sequences (that is, the entropy in the absence of selection) is just equal to the sequence length:

$$H_{\max}(X) = L. \quad (4)$$

This result is intuitively simple: the amount of information that can potentially be stored in a sequence of length L is just equal to the sequence length.

In the presence of selection, the probabilities of finding particular genotypes i in the population are highly non-uniform: most sequences do not appear (because either they simply never occur, or because their fitness in the particular environment vanishes), while a few sequences are over-represented. As described above, the amount of information that a population X stores about the environment E in which it evolves is then given by the difference:

$$I(X : E) = H_{\max} - H(X|E) = L + \sum_{i=1} p_i \log p_i, \quad (5)$$

Here, I have introduced the standard notation $I(A : B)$ for the entropy shared between A and B (i.e., the information that

A has about B), and the symbol $H(A|B)$ for the *conditional* entropy of A given B . Note that while X in the above formulae represents an ensemble of sequences, E stands for one particular environment, not an ensemble of environments².

Let me re-emphasize at this point that Eq. (5), because it represents the amount of information an ensemble has about its environment in mutation-selection balance, is *the same* as the physical complexity. In order to evaluate it, we must therefore obtain the probabilities p_i .

The probabilities p_i that go into the calculation of the conditional entropy in (5) are in fact *conditional* probabilities, because the probability of finding genotype i in environment E is not equal to the probability of finding the same sequence in, say, environment E' . These probabilities can in principle be estimated by simply counting the abundance of each genotype in the population, n_i , so that

$$p_i \approx \frac{n_i}{N},$$

where N is the population size. Unfortunately, the error committed by approximating the probabilities by the relative abundance gives rise to a sizable error in the entropy of Eq. (1), so large in fact that the estimated entropy is only meaningful for essentially infinite population sizes.^(20,21) Because we need the entropy Eq. (1) in order to estimate the physical complexity, we approximate it instead by summing up the entropy *at every site* along the sequence. This is done by aligning all sequences in the population, and obtaining the substitution probabilities at each site. In this manner, we can obtain the *per-site* entropy

$$H(j) = - \sum_{i=G, C, A, T} p_i(j) \log p_i(j) \quad (6)$$

for site j by compiling the probabilities to find nucleotides i at position j . The entropy Eq. (1) is then approximated by summing over all sites j in the sequence, i.e.,

$$H(X) \approx \sum_{j=1}^L H(j), \quad (7)$$

so that an approximation for the physical complexity of a population of sequences of length L is obtained by inserting Eq. (7) into Eq. (5) above:

$$I(X : E) \approx C_1(X) = L - \sum_{j=1}^L H(j). \quad (8)$$

Here, I defined implicitly $C_1(X)$, the complexity approximated using *single-site* entropies. Technically, this is only a good approximation if there are no correlations *between* sites in a sequence. Such correlations manifest themselves by

²Because E is not an ensemble but a particular *instance*, $I(X : E)$ is strictly speaking a difference of entropies rather than information in the sense of Shannon,⁽¹⁵⁾ but I will use the term information anyway.

epistatic interactions (epistasis) between mutations. It is well known that such epistasis exists (see Ref. 22 for a review), in particular in populations that are not well equilibrated.

Epistasis can be particularly problematic in asexual organisms (and at low mutation rates) because asexuals are at maximal linkage disequilibrium. Therefore, strong epistasis in a gene that could be coded in a much shorter fashion can prevent this compression from happening (perhaps because it would take too many mutations to arrive at a state at which the gene could be compressed). In contrast, recombination can be thought of as a way to *improve* coding efficiency, as it breaks up linkage disequilibrium. In higher organisms, we expect in addition a considerable amount of epistasis *between* genes that are part of a pathway, that regulate each other, or are regulated by the same *cis*-regulatory complex. Because the number of pairs of nucleotides that are correlated in this manner are still expected to be small compared to all the pairs that are independent, we do not expect strong directional effects in epistasis even in such a case. In any case, within each gene, it is possible to correct for directional epistasis (the overall deviation from independence of mutations) if we map out the decrease in fitness of this gene as a function of mutation number (see appendix of Ref. 14, and Ref. 23 for a measurement of directional epistasis in simulated T7 phages). In the following, we are going to assume that epistatic effects are sufficiently weak that the corrections can be ignored.

Natural selection increases physical complexity

Darwinian evolution is often described as a mechanism that increases the fitness of a population. Such a portrayal is problematic because the fitness of a population can depend on many parameters and is difficult to measure. It is probably more appropriate to say that evolution increases the amount of *information* a population harbors about its niche (and therefore, its physical complexity). The only mechanism necessary to guarantee such an increase is natural selection, acting in a single niche, on asexual organisms adapting to a constant unchanging world.

As we saw above, information is revealed, in an ensemble of adapted sequences, as those symbols that are conserved (fixed) under mutational pressure. Imagine then that a beneficial mutation occurs at a variable position. If the selective advantage that it bestows on the organism is sufficient to fix the mutation within the population,⁽²⁴⁾ the amount of information (and hence the complexity) has increased. A beneficial mutation that is lost before fixation does not decrease the amount of information, nor does this happen if a neutral mutation drifts to fixation. A deleterious mutation that occurs at a fixed site could lead to an information decrease, but such a mutation can only drift to fixation in very small populations (Muller's ratchet) or if the mutation rate is so high that the population undergoes a mutational meltdown.

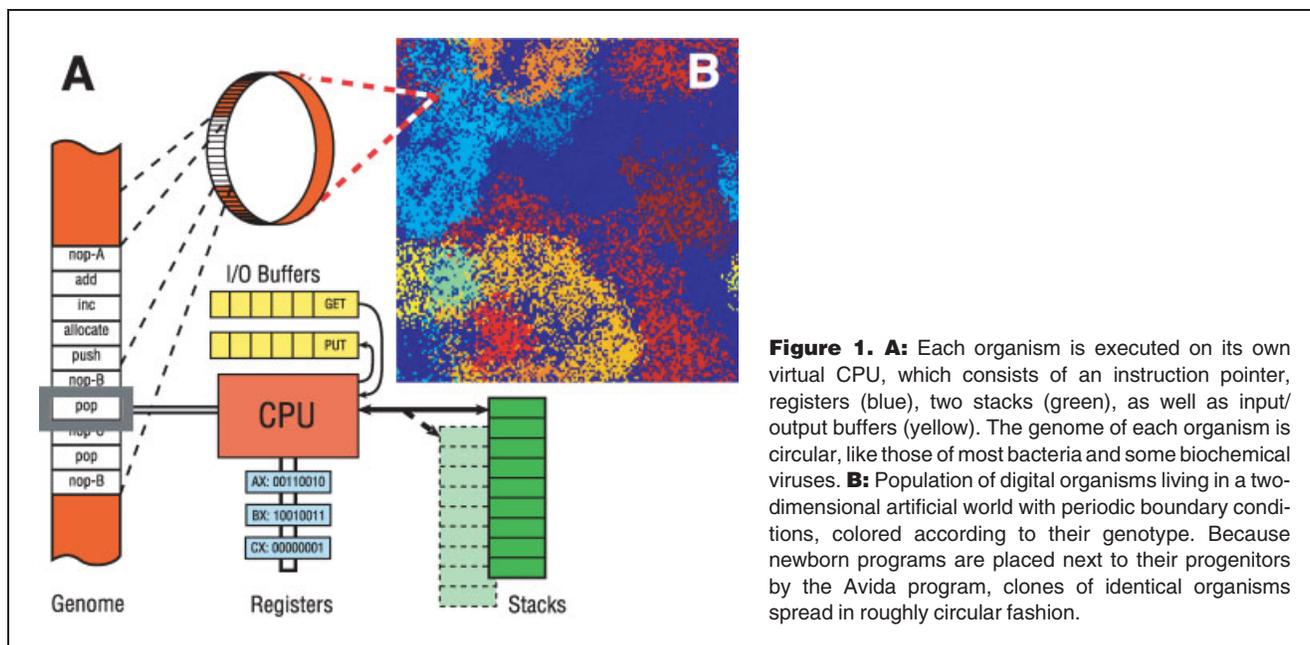
Thus, natural selection can be viewed as a *filter*, a kind of semipermeable membrane that lets information flow into the genome, but prevents it from flowing out. In this respect, the action of natural selection is very much akin to a device known as a Maxwell Demon in physics,⁽²⁵⁾ which implies that natural selection can be perfectly well understood from a thermodynamics perspective as well.

Evolution of complexity in digital organisms

Because evolution is an exceedingly slow process, it is difficult to witness the emergence of novelty and the concomitant increase in complexity in conventional experimental populations of animals, plants, or even bacteria. This obstacle disappears if we have access to a form of life with a much shorter generation time. Digital organisms are just such a form of life: they are computer programs that self-replicate, mutate, and compete for resources.^(26–32) Because digital organisms must copy their entire genome to survive within the computer's memory, and compete for space and computer time with other programs to which they are related by descent, experiments with populations of digital organisms are to be contrasted with more conventional numerical simulations of the evolutionary process. These organisms, because they are defined by the sequence of instructions that constitute their genome, are not simulated. They are physically present in the computer's memory and *live there*. The world to which these creatures adapt, on the other hand, is simulated, which allows the digital experimenter unparalleled precision in the planning, execution and analysis of his experiments. Evolving, self-replicating programs behave just like evolving, self-replicating molecules, and their dynamics are indeed well described by Eigen's⁽³³⁾ theory of macromolecular evolution.⁽³⁴⁾

In creating this virtual world, we do not specify a target sequence that represents the pinnacle of success. Instead, rewards (in the form of bonus execution time for the programs that reap them) are specified for *phenotypes* only, and thus natural selection acts on those. Because the underlying genetic space (the space of computer programs written in this particular language) is so high-dimensional, a large number of genotypes usually map to any particular phenotype, making the identification of a global genotypic optimum practically impossible. Phenotypes in this computational world are computational in nature, as we shall see presently.

In order to survive in their world, digital organisms must replicate fast and use the available resources efficiently. The efficient use of resources concerns chiefly the utilization of the primary energy source for digital organisms: CPU (central processing unit) time. Without CPU time, no digital organism can survive, since they need to copy themselves to survive, and without the code being executed, no copying takes place. Fig. 1 below shows a sketch of the world that is created inside of a standard computer by running the Avida software,⁽²⁵⁾ which is used for all the experiments described here.



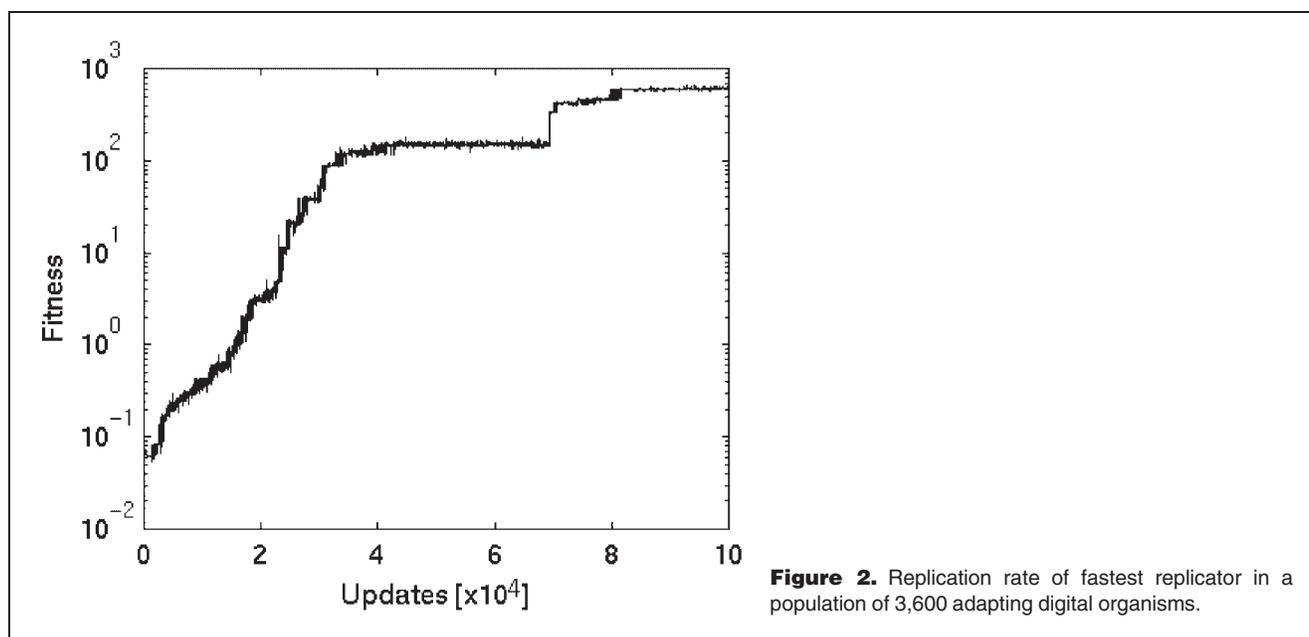
Using random numbers that the organisms can read into their CPU with an appropriate instruction, programs can perform *computations*. Clearly, only very particular sequences of instructions perform meaningful computations on input numbers. In this sense, we can view such a sequence as the equivalent of a nucleotide sequence coding for an enzyme that catalyzes a reaction, involving two input chemicals, producing the energy-rich output chemical. In the evolutionary experiments described below, the rewarded computations are logical operations (such as AND, OR, NOR, etc.) performed on binary input strings. During adaptation, many of these computational reactions evolve among the digital organisms, and are used in a coordinated manner to accelerate their reproduction. In that sense, it can be said that these computational genes play the role of a *computational metabolism*, quite analogous to the enzyme-based biochemical metabolisms. The monomers from which these programs are constructed (the instruction set) are custom-built for the CPU described above. For these experiments,⁽¹⁴⁾ the alphabet has 28 possible instructions, one of which is a logical primitive: NAND (the “not-and” operation).

Consider the behavior of fitness over time (depicted here is the replication rate of the fastest replicator in a population of 3,600 adapting programs whose sequence length is kept fixed at 100, and seeded with a single simple replicator) in Fig. 2. Time is here measured in arbitrary units called “updates”, where one update is the time it takes to execute about 30 instructions for each of the 3,600 programs in the population. One generation corresponds to between 10 and 100 updates in such populations. Note the sudden increase in fitness around update 70,000. At this point in time, a mutation must have created a new genotype much superior to all others.

Following our discussion, we expect this increase in fitness to be associated with an increase in information, so this genotype is a good candidate to inspect for an increase in complexity.

A plot of the approximate complexity (calculated according to Eq. (8)) can be seen in Fig. 3, where it is apparent that the complexity steadily increases, except for a period at the beginning and shortly after each transition. Both observations can easily be explained. During the initial growth of the population, most instructions appear fixed in the population because mutations have not had sufficient time to randomize the non-coding instructions. Evolution may also struggle with a genome (hand-written by the experimenters) that is extremely ill-suited to the environment, but also difficult to re-code. It may simply be badly compressed, and evolution takes a while to find a better way to represent the same information. After each transition, the estimated complexity overshoots its equilibrium value due to the “*hitchhiking*” effect: neutral instructions hitchhiking on beneficial ones appear fixed, until mutations can randomize them again. This is particularly clear in the transition around 70,000 updates in Fig. 3, to which we now turn our attention.

Because of the hitchhiking effect mentioned earlier, the amount of information gained in the transition highlighted in Fig. 3 is not measured very accurately, simply because the time to equilibration (required for an accurate estimate) is longer than the time until the next transition. To get a more accurate estimate of the per-site entropy Eq. (6), we can extract dominating genotypes just before and after the transition. In order to determine whether an instruction is entropy or information, we create all possible one-point mutants of the organisms and obtain their fitness in isolation. In a sense, this is equivalent to building virtual, fully equilibrated populations.



If a mutation does not change the fitness or increases it, it is deemed viable, while all deleterious mutations are classified together with the lethal ones, because they have a low probability of appearing in subsequent generations. After this has been done for each locus, the per-site entropy at locus x_i can be estimated as

$$H(x_i) \approx \log_D(N_{\text{viable}}), \quad (9)$$

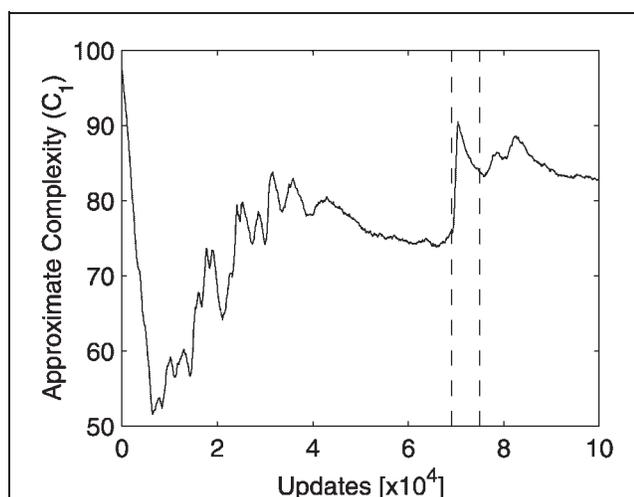


Figure 3. Approximate complexity according to Eq. (8) for a population adapting to a complex world. The dashed lines indicate the times chosen as pre and post-transition, at which the genotypes analyzed in (Fig. 4) were extracted.

where N_{viable} is the number of neutral or beneficial substitutions at that locus. In equation (9), the logarithm is taken to the base of the alphabet size, thus ensuring that our measure for the randomness at each location is normalized to lie between zero and one. If we do this for two organisms before and after the transition, we obtain the per-site entropies of Fig. 4. It is interesting to observe the changes in substitution pattern between these two genomes.

The most radical change seems to have taken place in the region between instructions 66 and 73, where about seven instructions that were moderately variable (in the virtual population) seemed to have turned “cold”, i.e., they have turned vulnerable to mutations. This is precisely the phenomenon of unidirectional information flow pointed out above: entropy is transformed into information. There are other places in the genome where hot instructions turned cold, and vice versa. The net gain in information is about six instructions, which is close to the number that we arrive at if we take into account corrections for epistasis.⁽¹⁴⁾

Causes for complexity declines

In this section, I discuss the mechanisms by which complexity can fail to increase, or even crash. The most obvious origin of a complexity catastrophe is a drastically changing environment. As discussed above, physical complexity is a quantity defined with reference to an environment. If the changes in the environment are fast and extreme, not only will the organism be maladapted to this new environment, but also its measurable physical complexity will have decreased commensurately. High mutation rates can also lead to a loss of complexity, due to the hitchhiking of deleterious mutations on beneficial ones. In small populations, high mutation rates are

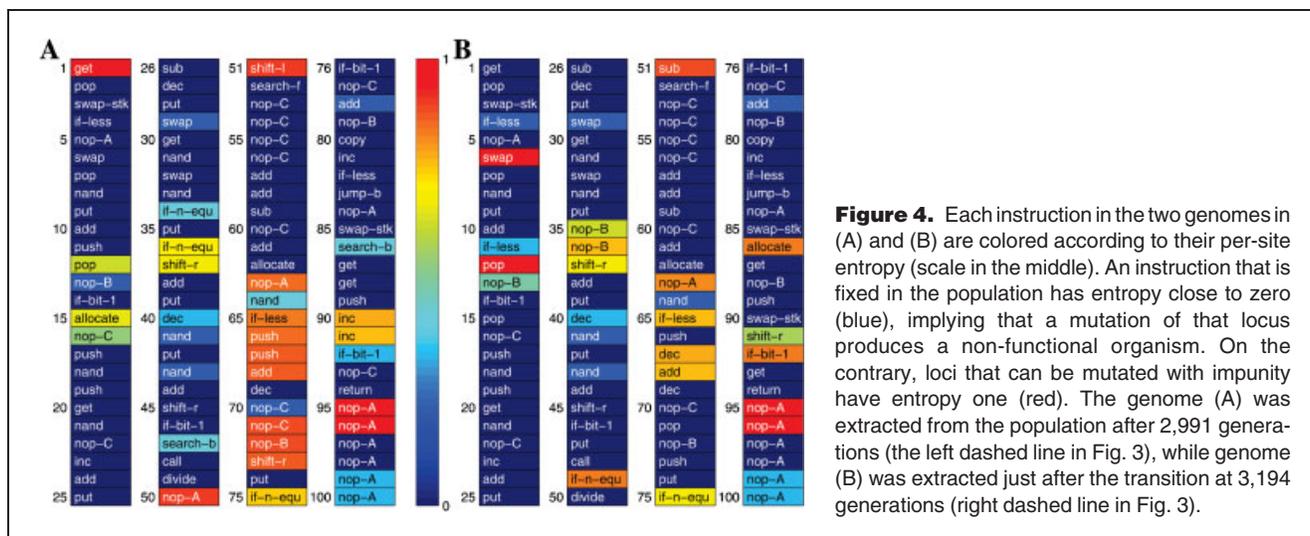


Figure 4. Each instruction in the two genomes in (A) and (B) are colored according to their per-site entropy (scale in the middle). An instruction that is fixed in the population has entropy close to zero (blue), implying that a mutation of that locus produces a non-functional organism. On the contrary, loci that can be mutated with impunity have entropy one (red). The genome (A) was extracted from the population after 2,991 generations (the left dashed line in Fig. 3), while genome (B) was extracted just after the transition at 3,194 generations (right dashed line in Fig. 3).

even more problematic, because the selective filter becomes sloppy and information can leak out. In the extreme case (critically high mutation rates), selection becomes inactive, a phenomenon known as the *error catastrophe* in the molecular evolution literature.⁽³⁵⁾

As is well known, sexual recombination can also lead to an accumulation of deleterious mutations, and thus to a loss of information (and by proxy, complexity). While asexual populations can purge deleterious mutations with certainty (as long as the mutation rate is not too high and the population size too small, as described above) populations of sexual organisms are at risk of gene loss at any mutation rate if deleterious mutations interact antagonistically.⁽³⁶⁾ I believe that there are good reasons to believe that the mechanism of complexity increase that holds for asexual organisms ought to translate unchanged to sexual organisms. First, it is clear that we can treat each tightly linked stretch of DNA, or each single protein, as a symbolic sequence that does not undergo recombination, and within which therefore we expect complexity to increase. Second, the ubiquity of the sexual mode of recombination within eukaryotes implies that selection is not weakened by that mode, perhaps rather to the contrary. Because strength of selection is the ultimate criterion for information maintenance, we do not need to fear a mutational meltdown due to recombination only.

Finally, co-evolution between species occupying different niches is a special case of a changing environment (for each of the interacting species), and thus opens up the possibility of escaping the inexorable growth of complexity promised by perfect selection. In this case, however, there are good reasons to assume that, for the most part, co-evolution will aid, rather than hinder, the evolution of complexity, because co-evolution is a slow rather than drastic environmental change, creating new niches that provide new opportunities for

adaptation. I discuss complexity growth in ecosystems briefly below.

Evolution of ecosystem complexity

With the present tools we cannot, strictly speaking, make any prediction about a trend in the complexity of entire ecosystems of interacting niches, since the concept of physical complexity only makes sense within an organism's own niche. An increase in complexity can only be observed in any particular niche, for the amount of time that this niche exists unchanged. Furthermore, the complexity of an organism can never exceed the potential complexity of the niche. Because niches do change, and because many niches of differing *potential information* coexist at the same time, we cannot expect that a trend in one niche will persist forever, nor that the same trend will be observable in all currently existing niches. In one niche, for example, its inhabitants may have incorporated all of its potential information into their genome (such as some prokaryotes), while another niche may just have been invaded so that its inhabitants show rapid gene turnover. The coexistence of niches with different entropy (different potential complexity) explains the coexistence of organisms with differing complexity in our ecosystems today, and should not be viewed as an argument against a trend.

Should we not expect an *overall* trend if evolution produces more and more diverse niches with more and more potential information? This question addresses the issue of co-evolution, and whether this process indeed produces niches with more and more entropy (which could then host, in turn, organisms with more and more complexity). This question is complicated by the fact that co-evolution necessarily produces *changes* in an organism's niche, which can reduce an organism's complexity. In general, a change in niche will almost

always produce an instant decrease in physical complexity, because only in the most rare circumstances will the change be exactly right to convert an entropic sequence into an informational one. However, if the change in the niche makes it richer (i.e., produces features that are awaiting discovery), then following the initial decline in complexity the species can enter a period of adaptation that can take it into realms of complexity hitherto unattainable. Also, if a species invades a new niche that leads to the loss of a previously functional gene (either through mutation accumulation or antagonistic pleiotropy, Ref. 37), most likely a species would still exist that did not undergo the change, so that the total complexity of the two species would be constant or increasing.

Thus, we have to look at the process of co-evolution and its capacity to create more complicated environments as the possible unifying process that could give rise to an overall trend. Unfortunately, the mathematics of information in co-evolving environments appears as yet too daunting to make a prediction about whether this is the case or not. It seems plausible to me, but it is clear that counterexamples can be manufactured where co-evolution gives rise to catastrophic extinctions, which reduce the environment's complexity and, necessarily, the physical complexity of its inhabitants at the same time. In such a formalism, the total complexity of an ecosystem would have to be defined as the mutual entropy of all organisms, about each other and the world they live in. This is an information-theoretic formula that is not difficult to write down, but the associated quantity promises to be much more difficult to measure.

Conclusions

In order to be able to speak about complexity, we must define it. In this review, I have presented a mathematical definition of sequence complexity that has a very intuitive interpretation for biological genomes, as the amount of information that a population stores about the environment in which it lives. With this definition, we can address the issue of a trend in the evolution of complexity. By recognizing that natural selection in a niche is equivalent to a filter that allows increases in information but not decreases, it is possible to show that, within that niche, physical complexity must increase if the environment does not change.

While natural selection can fail to maintain the acquired information, it is highly likely that the mechanism of interacting niches in an ecosystem will ultimately lead not only to a trend within each niche, but also to a trend in the overall (total) complexity of an ecosystem. Physical complexity increases if selection is efficient, and decreases if it fails. Still, this measure of complexity does *not* translate to *adaptation*. An organism well-adapted to a simple niche can have a lower physical complexity than an organism badly adapted to a complicated niche. Thus, adaptation reflects only the *degree* to which the potential complexity of the niche is reflected in the physical

complexity of the organism, and certainly does not allow complexity comparisons across niches.

Acknowledgments

I thank Charles Ofria and Travis Collier for collaboration in the experimental work reported here, as well as Richard Lenski and Claus Wilke for valuable discussions. I am further indebted to Murray Gell-Mann for discussions on complexity, and for pointing out the relation between physical and effective complexity. Thanks are also due to Allan Drummond for comments on the manuscript.

References

- Gould SJ. Full House. New York: Harmony Books. 1996. p 3.
- Badii R, Politi A. Complexity—Hierarchical Structures and Scaling in Physics. Cambridge: Cambridge University Press. 1997.
- Lloyd S, Pagels H. Complexity as thermodynamic depth. *Ann Phys* 1986; 188:186–213.
- Crutchfield JP, Shalizi CR. Thermodynamic depth of causal states: Objective complexity via minimal representations. *Phys Rev E* 1999;59: 275–283.
- Crutchfield JP, Young K. Inferring statistical complexity. *Phys Rev Lett* 1989;63:105–108.
- McShea DW. Metazoan complexity and evolution: Is there a trend? *Evolution* 1996;50:477–492.
- Bell G, Mooers AO. Size and complexity among multicellular organisms. *Biol J Linn Soc* 1997;60:345–363.
- Bonner JT. The Evolution of Complexity. Princeton: Princeton University Press. 1988.
- Alroy J. Cope's rule and the dynamics of body mass evolution in North American Fossil Mammals. *Science* 1998;280:731–734.
- McShea DW. Functional complexity in organisms: Parts as proxies. *Biology and Philosophy* 2000;15:641–668.
- McShea DW. The hierarchical structure of organisms: A scale and documentation of a trend in the maximum. *Paleobiology* 2001;27:405–423.
- Grassberger P. Toward a quantitative theory of self-generated complexity. *Int J Theor Phys* 1986;25:907–928.
- Adami C, Cerf NJ. Physical complexity of symbolic sequences. *Physica D* 2000;137:62–69.
- Adami C, Ofria C, Collier TC. Evolution of biological complexity. *Proc Natl Acad Sci USA* 2000;97:4463–4468.
- Shannon CE, Weaver W. The Mathematical Theory of Communication. Urbana: University of Illinois Press. 1949.
- Adami C. Information theory in molecular biology. 2002, in review.
- Deutsch D. The Fabric of Reality. New York: The Penguin Press. 1997. p 179.
- Wilson EO. The Diversity of Life. Cambridge: Harvard University Press. 1992. p 9.
- Gell-Mann M, Lloyd S. Information measures, effective complexity, and total information. *Complexity* 1996;2:44–52.
- Basharin GP. On a statistical estimate for the entropy of a sequence of independent random variables. *Theory Probability Appl* 1959;4:333–336.
- Schneider TD, Stormo GD, Gold L, Ehrenfeucht A. Information content of binding sites on nucleotide sequences. *J Mol Biol* 1986;188:415–431.
- Wolf J, Brodie E, Wade M. Epistasis and the Evolutionary Process. Oxford: Oxford University Press. 2000.
- You L, Yin J. Dependence of epistasis on environment and mutation severity as revealed by in silico mutagenesis of phage T7. *Genetics* 2002;160:1273–1281.
- Haldane JBS. A mathematical theory of natural and artificial selection. V: Selection and mutation. *Proc Camb Phil Soc* 1927;23:838–844.
- Adami C. Introduction to Artificial Life. New York: Springer Verlag. 1998.
- Ray TS. An approach to the synthesis of life. In: Langton CG, Taylor C, Farmer JD, Rasmussen S. editors. *Proc Artificial Life II*. Redwood City: Addison Wesley. 1991.

27. Adami C. Learning and complexity in genetic auto-adaptive systems. *Physica D* 1995;80:154–170.
28. Lenski RE, Ofria C, Collier TC, Adami C. Genome complexity, robustness, and genetic interactions in digital organisms. *Nature* 1999;400:661–663.
29. Wagenaar D, Adami C. Influence of chance, history, and adaptation on evolution in *Digitalia*. In: Bedau MA, McCaskill JS, Packard NH, Rasmussen S, editors. *Proc Artificial Life VII*. Cambridge: MIT Press. 2000. p 216–220.
30. Yedid G, Bell G. Microevolution in an electronic microcosm. *Am Nat* 2001;157:465–487.
31. Wilke CO, Wang JL, Ofria C, Lenski RE, Adami C. Evolution of digital organisms at high mutation rate leads to survival of the flattest. *Nature* 2001;412:331–333.
32. Wilke CO, Adami C. The biology of digital organisms. *Trends Ecol Evol* 2002;17 (to be published).
33. Eigen M. Self-organization of matter and the evolution of biological macromolecules. *Naturwissenschaften* 1971;58:465–523.
34. Wilke CO, Adami C. Evolution of mutational robustness. *Mutation Research* 2003; (to be published).
35. Eigen M. Natural selection: a phase transition? *Biophys Chem* 2000;85:101–123.
36. Kondrashov AS. Deleterious mutations and the evolution of sexual reproduction. *Nature* 1988;336:435–440.
37. Cooper VS, Lenski RE. The population genetics of ecological specialization in evolving *E. coli* populations. *Nature* 2000;407:736–739.