

## THE POPULATION GENETICS OF ADAPTATION: THE ADAPTATION OF DNA SEQUENCES

H. ALLEN ORR

*Department of Biology, University of Rochester, Rochester, New York 14627*

*E-mail: aorr@mail.rochester.edu*

**Abstract.**—I describe several patterns characterizing the genetics of adaptation at the DNA level. Following Gillespie (1983, 1984, 1991), I consider a population presently fixed for the  $i$ th best allele at a locus and study the sequential substitution of favorable mutations that results in fixation of the fittest DNA sequence locally available. Given a wild type sequence that is less than optimal, I derive the fitness rank of the next allele typically fixed by natural selection as well as the mean and variance of the jump in fitness that results when natural selection drives a substitution. Looking over the whole series of substitutions required to reach the best allele, I show that the mean fitness jumps occurring throughout an adaptive walk are constrained to a twofold window of values, assuming only that adaptation begins from a reasonably fit allele. I also show that the first substitution and the substitution of largest effect account for a large share of the total fitness increase during adaptation. I further show that the distribution of selection coefficients fixed throughout such an adaptive walk is exponential (ignoring mutations of small effect), a finding reminiscent of that seen in Fisher's geometric model of adaptation. Last, I show that adaptation by natural selection behaves in several respects as the average of two idealized forms of adaptation, perfect and random.

**Key words.**—Adaptive evolution, adaptive walk, experimental evolution, fitness jump, mutational landscape, natural selection.

Received December 4, 2001. Accepted April 22, 2002.

Evolutionary biologists are nearly unanimous in thinking that adaptation by natural selection explains most phenotypic evolution within species as well as most morphological, physiological, and behavioral differences between species. But until recently, the mathematical theory of population genetics has had surprisingly little to say about adaptation. Instead, population genetics has, for both historical and technical reasons, focused on the fates of neutral and deleterious alleles. The result is a curious disconnect between the verbal theory that sits at the heart of neo-Darwinism and the mathematical content of most evolutionary genetics.

Growing awareness of this gap has encouraged a number of recent forays into the theory of adaptation (Hartl and Taubes 1996, 1998; Barton 1998, 2001; Orr 1998, 1999, 2000; Gerrish 2001; Gillespie 2002). Many of these efforts have focused on a phenotypic model of adaptation due originally to R. A. Fisher. Fisher's (1930) geometric model pictures an organism as a point in a high dimensional space, in which each dimension represents a phenotypic trait value. Because some combination of trait values presumably represents the fittest phenotype in the present environment, some point in this phenotypic space is taken to represent the (local) fitness optimum. Given a recent change in environment—and hence a population that has been thrown off the fitness optimum—the task of adaptation is to move the population toward the new phenotypic optimum. This task is complicated in several ways. First and most important, adaptation must rely on mutations that are random with respect to the organism's needs. Second, a mutation that improves one character might worsen others, a problem of pleiotropy that grows more severe with the dimensionality of the phenotypic space. Last, adaptation is constrained by the laws of population genetics, in particular by probabilities of fixation: Most mutations, even those that are favorable, are accidentally lost when rare (Haldane 1927). Given these difficulties, it is hardly surprising that adaptation in Fisher's model usually requires many substitutions: sub-

stantial progress towards the optimum typically involves an extended adaptive walk (Orr 1998, p. 941).

Study of adaptation in Fisher's model has led to several new results. Perhaps the most interesting is that adaptation is characterized by an approximately exponential distribution of phenotypic effects among the factors fixed. This pattern is robust to a number of model details, for example, the shape of the fitness function, the dimensionality of the phenotypic space, and, most surprisingly, the distribution of mutational effects provided to natural selection (Orr 1998, 1999). Similarly, theory shows that the mean phenotypic effect of favorable mutations fixed at subsequent substitutions  $K = 1, 2, 3, \dots$  falls off as an approximate geometric sequence, that is, decreases by a nearly constant proportion. Last, the average size of the largest factor fixed during a bout of adaptation is larger than suggested by either Fisher's (1930) or Kimura's (1983) earlier analyses. (Mutations typically, however, still travel a short distance to the optimum, especially in high-dimensional organisms. This reflects the "cost of complexity" [Orr 2000].)

Here I turn to a different model, one that considers adaptation in a space of DNA sequences. It may not at first be obvious why one must consider the adaptation of DNA sequences per se, in addition to models like Fisher's. One reason is that Fisher's model is idealized, involving selection on orthogonal characters that experience equivalent (i.e., isotropic) mutational effects. The model thus might serve to train our intuition and to yield heuristic predictions, but it cannot be taken literally. Analysis of adaptation at the DNA level thus represents a step toward biological realism. Second, adaptation in DNA sequence space is characterized by two constraints that do not appear in most phenotypic models, including Fisher's. One is that the space of DNA sequences is discrete. This limits the number of states available at any gene and places a ceiling on the number of substitutions that occur before the fittest possible allele is fixed. The space of

phenotypes in Fisher's model, on the other hand, is continuous. An infinite number of states are available, and an unlimited number of substitutions can occur during adaptation (at least in an indefinitely large population). The typical substitution consequently plays a small part in adaptation. We would like to know whether this is an artifact of the unlimited number of substitutions allowed in Fisher's model. Second, the number of mutant DNA sequences available to natural selection is constrained in a way that is not at first obvious. As Gillespie (1984) pointed out, the mutational landscape has a short horizon. For reasons that will be considered below, natural selection can effectively search only those sequences that are *one* mutational step away from wild type. This fact, which plays a major role in the adaptation of DNA sequences, has no counterpart in Fisher's model. Taking these facts together, it seems possible that the adaptation of DNA sequences might show certain rules or patterns that differ from those seen in more abstract models. We would like to know what these rules or patterns are. Conversely, we would like to know whether any rules or patterns are shared across both DNA-based and more abstract phenotypic models.

Fortunately, analysis of adaptation in a DNA sequence space can build on foundations laid by Gillespie (1983, 1984, 1991). Gillespie considered the fate of a wild-type sequence whose fitness drops following an environmental change and which will, as a result, be replaced by a fitter mutant sequence. Considering the complex scenario in which several favorable mutant sequences might be available, Gillespie described the stochastic process governing the fates of these mutant alleles, calculating waiting times to the fixation of each mutant sequence as well as the probability that a particular mutant will be the next fixed by natural selection. His analysis rests on several important insights (including the use of extreme value theory) that greatly simplify analysis of adaptation through DNA sequence space.

Although the present paper builds on Gillespie's efforts, I focus on a different set of questions. Gillespie was primarily concerned with the statistical properties of molecular evolution, in particular with whether substitutions occur as a Poisson process (as predicted by the neutral theory; Kimura 1983) or in bursts (as would explain the observed overdispersion of the molecular clock). His main interest was therefore in the number of substitutions needed to move from an unfit sequence to the fittest sequence available. He showed that this number is generally small, reflecting a brief burst of substitutions that might explain the overdispersion of the molecular clock.

I am primarily concerned with patterns characterizing the genetics of adaptation. In particular, I focus on problems like the fitness rank of the mutant allele typically chosen by natural selection at the next substitution, the size of the increase in fitness that occurs when that allele is fixed, and the distribution of selection coefficients among mutations fixed during adaptive walks to the best (local) sequence available. I also compare my results with those from Fisher's (1930) model.

#### THE MODEL

##### *The General Approach*

Although several models of adaptive evolution in DNA space have been suggested (reviewed in Gillespie 2002), Gil-

lespie's (1983, 1984, 1991) mutational landscape model seems one of the most realistic. I begin by briefly reviewing this model. I first sketch the overall process described by the model and then more carefully consider a single step in that process.

Gillespie considered a large population in which all alleles behave at any point in time as either definitely favorable or definitely deleterious. It is simplest to think of the population as haploid, although diploidy introduces nothing new as long as all alleles have some heterozygous effect and all effects are taken as heterozygous. Formally, it is assumed that  $2Ns \gg 1$ , where  $N$  is population size and  $s$  is a selection coefficient. It is also assumed that mutation is rare, that is,  $N\mu < 1$ , which is appropriate given a low per nucleotide mutation rate ( $\mu \approx 10^{-9}$ , which is assumed equal at all sites and to all bases). Under these conditions, the population is more or less fixed for a wild-type sequence at any point in time. Although Gillespie (1991) referred to this scenario as "strong selection weak mutation" (SSWM), it is important to note that selection is strong only in relative terms, that is, compared to  $1/N$ . In absolute terms,  $s$  may be small and selection mild.

Gillespie considered a single wild-type sequence that was, until recently, the fittest local allele. By "local," he meant that the wild-type allele was fitter than all of the  $m$  alleles the wild-type can mutate to by a single point mutation ( $m = 3l$ , where  $l$  is sequence length in base pairs). As Gillespie (1984) emphasized, these one-mutational-step neighbors are, to a good approximation, the only sequences that matter. If, for instance, every one-step mutation is deleterious, the wild-type sequence resides at a local optimum and the existence of a fitter sequence *two* mutational steps away is essentially irrelevant. The reason is that, given the low per nucleotide mutation rate, the waiting time to fixation of a double mutant is extremely long—on the order of  $1/\mu^2$  generations—too long to typically matter on the time scale of molecular evolution. (This  $1/\mu^2$  behavior has a simple explanation. One factor of  $\mu$  reflects the low frequency of the deleterious single-mutant sequence at mutation-selection balance [ $=\mu/s_{del}$ ], whereas the other reflects the rate of mutation from the single-mutant to the double-mutant [Gillespie 1984]. Note that in taxa having extremely large population sizes and/or mutation rates, Gillespie's argument may begin to break down; see Discussion.)

The process of interest begins when the environment changes. Throughout, I consider a single environmental change and study the burst of substitutions that occurs in response to this change. If relevant aspects of the environment change on a faster time scale than evolution to the best allele, some of the results (as well as those of Gillespie 1983, 1984, 1991) grow less relevant. (It is worth noting, however, that many of our results will concern only the first step in adaptation and so are less dependent on the time scale of environmental change.) As a result of the environmental change, the wild-type allele drops in fitness. In particular, the wild-type allele drops from being the best to being the  $i$ th best of the collection of  $m + 1$  relevant sequences ( $m$  one-step mutations plus the wild type), where our convention will be to rank the fitness of alleles from the top down: allele 1 is the fittest, allele 2 the next fittest, and so on (see Fig.

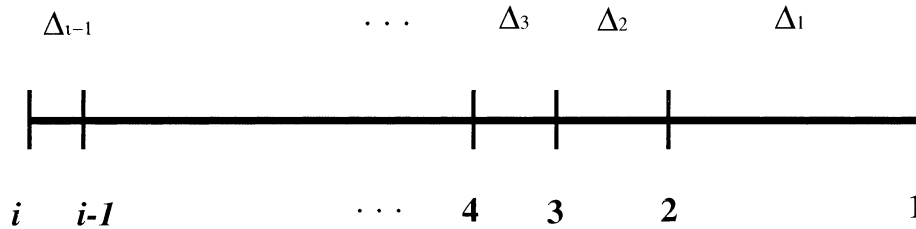


FIG. 1. The scenario being considered. The present wild type is allele  $i$ . The  $i - 1$  alleles to its right are more fit than wild type, and the many alleles to its left (not shown) are less fit than wild type. The fitness spacings between alleles are labeled  $\Delta_1$ ,  $\Delta_2$ ,  $\Delta_3$ .

1). Throughout this analysis, I will make one major assumption: The wild-type allele, while no longer the best, remains near the top in fitness ( $i$  is fairly small). This is almost surely true biologically for two reasons. First, environments are typically autocorrelated through time. The best sequence today will consequently rarely be the worst tomorrow. Second, a large fraction of mutations give rise to unconditionally deleterious or lethal alleles. It seems unlikely that the present wild-type sequence will often fall into the company of such alleles, most of which presumably reflect loss-of-function mutations.

The essence of adaptation in Gillespie's model is simple. Natural selection will move the population from the present wild-type sequence to another sequence that is fitter. Although each of the  $i - 1$  mutant sequences that are fitter than the first wild type enjoys some probability of fixation each time it appears, one will be fixed before the others (see below). This substitution event represents one step in the process we study. At this point, the process described above begins anew. The new wild-type sequence, which retains its absolute fitness,  $W$ , now can produce its own suite of one-step mutations. All of these mutant sequences were inaccessible to the first wild type (as they differ by two mutations). Some of these new mutations may be fitter than the second wild type. If so, the population will again jump through sequence space, moving to yet another wild type. This process continues until evolution arrives at a sequence that is fitter than *all* its one-mutational-step neighbors. Adaptation has, at that point, arrived at a local optimum and is complete.

Formally, one can treat adaptation in DNA sequence space under strong selection and weak mutation as a Markov chain in which the population repeatedly jumps from some point in sequence space to some other, fitter, one. The power of this approach is that, by taking the long view of molecular evolution, we can track adaptation in terms of substitution events, without concerning ourselves with details of allele frequency change as favorable mutations sweep through populations. Such sweeps are effectively instantaneous on the time scale of molecular evolution, a time scale set by the escape of alleles from the boundary layer, that is, by the appearance of lucky mutations that are both favorable and escape accidental loss. My analysis will, however, differ from Gillespie's in one trivial way. I treat adaptation as a discrete time Markov chain in which time is measured in substitution events, whereas Gillespie considered a continuous process in which time was measured in a continuous proxy for generations.

It is important to note that Gillespie's (1983) early work

considered a mutational scheme that differs from that described above. In this work, he assumed that evolution moves through a space of mutually accessible alleles until finally arriving at the best sequence: A wild-type allele can produce  $i - 1$  fitter sequences, each of which can reach all others via a single mutation. I will refer to this model as the "simple" mutational model. This model is obviously artificial. In reality—and as assumed in Gillespie's (1984) later work—each wild-type sequence produces a new suite of  $m$  mutant sequences, all of which were inaccessible to earlier wild-type sequences. Following Gillespie (1984), I refer to this as the "mutational landscape" model. I will be almost exclusively concerned with the mutational landscape model.

#### A Single Step

I now more carefully consider a single step in the process described above. A population is fixed for the  $i$ th fittest allele and natural selection might move it to any one of the  $i - 1$  fitter sequences at the next substitution. Each of these  $i - 1$  alleles enjoys a probability of fixation of  $2s_j$  each time it appears (Haldane 1927), where I assume the absolute strength of selection is sufficiently mild for this approximation to hold. Because mutation is recurrent, each favorable allele would ultimately get fixed in the absence of competition from the other favorable alleles. In reality, however, one of these  $i - 1$  alleles will be fixed before the others. We would like to know the probability,  $P_{ij}$ , that allele  $j = 1$  versus allele  $j = 2$ , etc. is the next fixed. Note that  $P_{ij}$  is nonzero only when  $j < i$ . Evolution cannot, in other words, go backward to a less fit allele nor stay in place. Instead evolution by natural selection always moves forward to a fitter allele when one is available.

$P_{ij}$  obviously depends on the probability of fixation,  $\Pi_j$ : The larger an allele's probability of fixation, the greater the chance the population will jump to it. Gillespie (1983, 1984, 1991) showed how  $P_{ij}$  depends on  $\Pi_j$  when several favorable mutations are available. Because the rate of substitution of allele  $j$  in a haploid population is  $N\mu\Pi_j$ ,  $j$  has an exponentially distributed waiting time to fixation with a mean of  $1/N\mu\Pi_j$  generations. The chance, therefore, that  $j$  is the next allele fixed equals the probability that  $j$ 's waiting time to fixation is smaller than that of all other favorable alleles. A straightforward calculation (Gillespie 1984, 1991) shows that this probability is

$$P_{ij} = \frac{\Pi_j}{\Pi_1 + \Pi_2 + \dots + \Pi_{i-1}}, \quad (1)$$

a result that makes good intuitive sense.

Because an allele's probability of fixation depends on its selective advantage ( $\Pi_j = 2s_j$ ), equation (1) can be written

$$P_{ij} = \frac{s_j}{s_1 + s_2 + \cdots + s_{i-1}}. \quad (2)$$

Note that each of these selection coefficients can, in turn, be written as  $s_j = (W_j - W_i)/W_i$ , where  $W$  is absolute fitness.

Equation (2) will play a key role in this analysis. It tells us how natural selection chooses among the favorable alleles available to it. But equation (2) is only half the machinery required. The other half is provided by extreme value theory. One of the main difficulties in modeling molecular evolution is knowing how to assign selection coefficients to mutant alleles like those in equation (2). As Kimura (1979) and Gillespie (1983) pointed out, one obvious (although not the only) way around this problem is to assign selection coefficients randomly, that is, to draw fitnesses from some probability distribution. But the question is, what distribution? What distribution does nature use when assigning fitnesses to mutant alleles?

Gillespie's (1983) key insight was that we need not rely on meager—indeed nearly nonexistent—data to infer this distribution. Instead we can take advantage of certain limiting results from extreme value theory that describe the tail behavior of almost any distribution. In particular, extreme value theory shows that the distribution of fitnesses of the top several alleles converges to the so-called extreme value distribution *independent* of the distribution used to assign fitnesses to alleles as long as the total number of alleles is large. (Formally, the distribution used to assign fitnesses must belong to the Type III class, which includes all regular distributions, including exponential, gamma, normal, log-normal, and so on. I exclude only exotic distributions that are Cauchy-like and those that are truncated on the right; see Gumbel [1958] and David [1981] for details.)

For our purposes, the most useful limit theorem concerns the spacings in absolute fitness,  $\Delta$ , between the fittest few alleles. As Figure 1 shows,  $\Delta_1$  represents the absolute fitness spacing between the fittest and next-fittest allele,  $\Delta_2$  between the second- and third-fittest alleles, and so on. These extreme spacings show a remarkable behavior: Regardless of the parent distribution used to assign allelic fitnesses, the spacings between the fittest alleles are asymptotically independent exponentially distributed random variables with means

$$E[\Delta_1] = C_m, \quad E[\Delta_2] = \frac{C_m}{2}, \quad E[\Delta_3] = \frac{C_m}{3}, \dots \quad (3)$$

The form of the parent distribution determines only the value of the scaling constant,  $C_m$ . For many distributions,  $C_m$  depends on the number of draws (alleles) from the distribution (see Gumbel 1958, p. 197; hence the subscript). Equation (3) thus shows that the fitness spacings grow smaller as one moves away from the fittest allele, as shown in Figure 1. The distributions represented in (3) were first obtained by J. H. Darwin (1957) and a simple derivation can be found in Gumbel (1958). Thus, assuming only that the present wild-type allele is reasonably fit, we know something about the fitness spacings, and thus selection coefficients, that separate our  $i - 1$  favorable mutant alleles from the wild-type.

A key assumption in Gillespie's (1984) mutational landscape model is that, at each step in adaptation, the fitnesses of new mutant alleles are drawn from the same probability distribution. Because the wild type is constantly increasing in fitness (i.e., moving out along the tail of this distribution), random favorable mutations grow more difficult to come by. Gillespie's model thus captures the biological fact that the population is closing in on a locally optimal sequence.

## RESULTS

### *Preliminary Comments*

I consider six questions about adaptation in a DNA sequence space: (1) What is the probability that a population fixed for the  $i$ th best allele will jump to  $j$ th best allele at the next substitution? (2) What is the average fitness rank of the new allele jumped to by natural selection? (3) How is this jump in fitness rank related to the number of substitutions required to reach the best allele? (4) What is the mean size of the fitness jump that occurs when natural selection drives a substitution? (5) What proportion of the total increase in fitness that occurs during adaptation is due to the first substitution; similarly, what proportion is due to the largest substitution? (6) What is the overall distribution of selection coefficients among mutations fixed during adaptation to the best allele? Most of these questions consider a single step in the adaptive process, but the last two concern an entire adaptive walk. Thus, only these last two questions depend on Gillespie's (1984, 1991) assumption that the distribution of allelic fitnesses stays constant through time.

These questions are, on the whole, simpler than those considered by Gillespie (1983, 1984, 1991) and, not surprisingly, most can be answered analytically. Two questions, however, require computer simulation. Although most of the analytic results are simple, their derivation is often surprisingly tedious (at least in my hands). I have thus taken two steps to make the presentation more intelligible. First, I have placed some of the messier calculations in the Appendices. Second, I have set off important results as numbered conclusions. The reader who prefers to skip the mathematics can get the main findings from these conclusions.

### *Transition Probabilities*

The size of the jump that occurs when natural selection moves a population from allele  $i$  to allele  $j$  is of considerable evolutionary interest. We would like to know if natural selection typically causes a population to leap to the fittest available allele ( $j = 1$ ) or to an allele that represents a small increase in rank (say,  $j = i - 1$ ). More generally, we would like to know the mean probability that evolution jumps from the  $i$ th fittest allele to the  $j$ th fittest allele at the next substitution. I call this probability,  $E[P_{ij}]$ , to emphasize that it represents an average over the distribution of selection coefficients found in nature. The mathematical machinery described in the model section lets us immediately find this quantity. Equations (2) and (3) show that

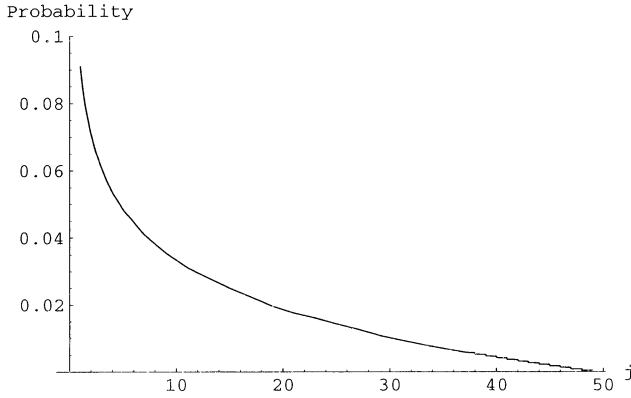


FIG. 2. The mean probability that a population fixed for the  $i$ th fittest allele will jump to the  $j$ th fittest allele at the next substitution (from eq. 5). In the example shown,  $i = 50$  and  $j$  varies from 1 to 49.

$$\begin{aligned}
 E[P_{ij}] &= E\left[\frac{s_j}{s_1 + s_2 + \dots + s_{i-1}}\right] \\
 &= E\left[\frac{\Delta_j + \Delta_{j+1} + \dots + \Delta_{i-1}}{\Delta_1 + 2\Delta_2 + \dots + (i-1)\Delta_{i-1}}\right] \\
 &= E\left[\frac{\frac{Z_j}{j} + \frac{Z_{j+1}}{j+1} + \dots + \frac{Z_{i-1}}{i-1}}{Z_1 + Z_2 + \dots + Z_{i-1}}\right], \quad (4)
 \end{aligned}$$

where, in the last step,  $Z_k = k\Delta_k$  and the  $Z_k$  are thus independent and identically distributed (IID) random variables having the same distribution as  $\Delta_1$ , the spacing between the top two alleles. By symmetry,  $E[Z_k/(Z_1 + Z_2 + \dots + Z_{i-1})] = 1/(i-1)$ , yielding

$$E[P_{ij}] = \frac{1}{i-1} \sum_{k=j}^{i-1} \frac{1}{k}. \quad (5)$$

Equation (5), which is a good probability mass function ( $\sum_{j=1}^{i-1} E[P_{ij}] = 1$ ), is plotted in Figure 2. This figure shows that natural selection is always more likely to jump to the fittest allele ( $j = 1$ ) than the next fittest ( $j = 2$ ) and so on, as one would expect intuitively. As Figure 2 also shows, the probability  $E[P_{ij}]$  decreases nearly logarithmically with  $j$ . In fact equation (5) is roughly  $E[P_{ij}] \approx [1/(i-1)]\ln[(i-1)/j]$ . We will, however, rely on the exact form in (5).

#### The Average Allele Jumped to

Equation (5) shows that the chances that natural selection will move from the present allele  $i$  to the next best allele  $i-1$  is generally small:  $E[P_{i,i-1}] = 1/(i-1)^2$ . Instead natural selection typically makes sizeable jumps in fitness rank, leaping far over the next best allele. But how far? Perhaps the most interesting question is: What allele does natural selection jump to on average?

The answer turns out to be surprisingly simple. It is

$$E[j] = E\left[\sum_{j=1}^{i-1} jP_{ij}\right] = \sum_{j=1}^{i-1} jE[P_{ij}] = \sum_{j=1}^{i-1} j\left(\frac{1}{i-1} \sum_{k=j}^{i-1} \frac{1}{k}\right), \quad (6)$$

which, after some work, reduces to

TABLE 1. Rank of allele substituted in simulations. Distribution of  $S$  shown at left, where  $W = 1 + S$ . Each set of simulations included at least 5000 substitutions (new fitnesses were drawn for each of 5000 runs). All distributions had  $\bar{S} = 0.01$ . For gamma distributed  $S$ , the shape parameter  $\beta = 2$  and the scale parameter  $\alpha = 200$ ; the distribution is thus humped. For normally distributed  $S$ ,  $\sigma = 0.005$ . The second and fourth columns are theoretical predictions, that is, equations (7) and (8), respectively.

Distribution	$E[j]$	$\bar{j}_{obs}$	$\text{Var}[j]$	$\sigma_{obs}^2$
$i = 10$				
Exponential	3.00	3.01	4.22	4.22
Gamma	3.00	3.07	4.22	4.23
Normal	3.00	3.01	4.22	4.18
$i = 50$				
Exponential	13.00	12.99	118.67	118.67
Gamma	13.00	12.76	118.67	115.81
Normal	13.00	13.70	118.67	121.88
$i = 150$				
Exponential	38.00	38.05	1085.33	1088.21
Gamma	38.00	38.62	1085.33	1060.80
Normal	38.00	39.97	1085.33	1095.95

$$E[j] = \frac{i+2}{4}. \quad (7)$$

*Conclusion 1.*—Adaptation by natural selection is characterized by a simple rule that maps the fitness rank of the present wild-type allele onto the mean fitness rank of the mutant allele next fixed:  $E[j] = (i+2)/4$ . Remarkably, this mapping depends only on the fitness rank  $i$  and nothing else: Equation (7) holds regardless of the distribution from which allelic fitnesses are drawn, the length of the sequence under consideration, and so on.

Equation (7) also shows that these jumps in rank are typically large. If, for instance, the population is presently fixed for the  $i = 15$ th best allele, natural selection will on average jump to about the fourth best available allele at the next substitution. In general, the change in rank is  $(2-3i)/4$ , or about 75% of the total gap in rank.

We can also find the variance in the rank of the allele jumped to. A straightforward calculation shows that  $E[j^2] = [i(4i+7)+6]/36$ , and so

$$\text{Var}[j] = \frac{(i-2)(7i+6)}{144}. \quad (8)$$

When  $i = 2$ ,  $\text{Var}[j] = 0$ , as it must, because natural selection can move to only one allele, the fittest. The right side of equation (8) also gives the variance in the change in fitness rank.

Table 1 gives the results of computer simulations that test the accuracy of equations (7) and (8) (these simulations are described in Appendix 1). The analytic theory is quite accurate: equations (7) and (8) nicely predict the mean and variance of the rank of the next allele fixed, regardless of the parent distribution (exponential, gamma, normal) used to assign fitnesses to alleles and regardless of whether adaptation begins from the  $i = 10$ th, 50th, or 150th best allele (where  $m = 3000$  mutations, which corresponds to a sequence of length 1000 bp). Surprisingly, the theory remains reasonably accurate even when starting from alleles that are fairly far

from the top in fitness and where we had no guarantee that extreme value theory would hold.

The above calculations also hint at an unexpected property of adaptation by natural selection. With respect to mean behavior, it falls precisely midway between two idealized and extreme forms of adaptation, perfect and random. In a perfect adaptive process, the best alternative is always chosen, for example, an algorithm calculates the expected return from all possible (local) courses of action and chooses the best. This yields  $E[j] = 1$ . (Kauffman and Levin [1987] referred to this as a “greedy” algorithm.) In a random adaptive process, an alternative is blindly chosen from those that represent improvements, for example, an algorithm can discern better versus worse—but not rank—and so randomly chooses from those (local) alternatives that improve matters. This yields  $E[j] = i/2$ . Curiously, adaptation by natural selection falls exactly halfway between these extremes:  $(1 + i/2)/2 = (i + 2)/4$ .

### The Length of Adaptive Walks

The fact that the mean jump in fitness rank under natural selection falls midway between that seen under perfect and random adaptation suggests another result. The number of substitutions required to reach the fittest allele depends on the size of the jump in rank that occurs in a single substitution: The larger the jump, the fewer the substitutions. It stands to reason then that the mean number of substitutions required to reach the fittest allele under natural selection might also be the average of that required under perfect versus random adaptation. This is in fact true. In a well-known calculation, Gillespie (1983, 1991) showed that, under a simple mutational scheme in which all alleles are mutually accessible (see Model section) and in which adaptation starts at the  $i$ th best allele, a mean of

$$L_{nat\ sel} = \frac{1}{2} + \frac{1}{i} + \frac{1}{2} \sum_{k=2}^{i-1} \frac{k+3}{k(k+1)} \quad (9)$$

substitutions occur before a population arrives at the best allele. Under perfect adaptation, however, a mean of  $L_{perfect} = 1$  substitutions occurs before arriving at the best allele, that is, the population arrives immediately. Similarly, under random adaptation a Markov chain absorption time calculation shows that a mean of  $L_{random} = \sum_{k=1}^{i-1} (1/k)$  substitutions are required. Remarkably, Gillespie’s result is the arithmetic average of these perfect and random absorption times. Equation (9) can, in other words, be written

$$L_{nat\ sel} = \frac{1 + \sum_{k=1}^{i-1} \frac{1}{k}}{2}. \quad (10)$$

I consider the reasons for this behavior in the Discussion. It may be worth noting that equation (10) is to an excellent approximation  $L_{nat\ sel} \approx [1 + \ln(i-1) + \gamma]/2$ , where  $\gamma \approx 0.5772$  is Euler’s constant. The number of substitutions required before absorption at the best allele thus grows logarithmically with  $i$ , as noted by Gillespie (1991).

It must be emphasized that the above absorption time calculations assume a simplified (and artificial) mutational model in which all alleles are mutually accessible. In the more

realistic mutational landscape model with which we are mainly concerned, it does not appear possible to calculate the length of adaptive walks under natural selection analytically (Gillespie 1984, 1991; approximations are, however, possible for random adaptation; Kauffman and Levin 1987). Computer simulations show, however, that this absorption time equals that given by equation (10) above plus  $\sim 0.44$ , for appreciable  $i$ . Simulation work also shows that the above averaging result still holds approximately under the full mutational landscape model. When each new wild-type sequence produces an array of  $m$  new mutant sequences, the mean number of substitutions required until absorption at the locally best allele rises slightly above those given above whether considering naturally selected, random or perfect adaptation. But natural selection remains close to the average of random and perfect adaptation (within 5–10%). For example, when  $i = 10$  and allelic fitnesses are exponentially distributed,  $L_{perfect} = 1.72$ ,  $L_{random} = 3.35$ , and their average is 2.53. Adaptation by natural selection yields  $L_{nat\ sel} = 2.36$  (5000 replicates each, with  $m = 3000$  mutations). Similarly, when  $i = 50$ ,  $L_{perfect} = 1.73$ ,  $L_{random} = 4.99$ , their average is 3.36, and  $L_{nat\ sel} = 3.18$ .

**Conclusion 2.**—Adaptation by natural selection behaves as the average of perfect and random adaptation in at least two respects: (1) the mean fitness rank of the allele jumped to; and (2) the mean number of substitutions required to reach the best allele. The first result is exact under both the simple and mutational landscape models, whereas the second is exact only under the simple mutational model.

### Mean Fitness Jumps

We would like to know the size of the fitness jump that occurs when natural selection drives a substitution. This quantity is of special interest because jumps in fitness are in principle more easily measured than fitness rank. In this section, I derive the mean and variance of this fitness jump. Throughout, the fitness referred to is absolute.

These calculations are more difficult than those above. To see the gist of the calculation of the mean fitness jump, consider the simplest nontrivial case in which  $i = 3$  and evolution can move to two possible favorable alleles. In any particular case, that is, given a particular wild-type sequence and two particular favorable mutant sequences, the mean fitness jump is

$$P_{31}(\Delta_1 + \Delta_2) + P_{32}(\Delta_2), \quad (11)$$

where  $P_{31}$ ,  $P_{32}$ ,  $\Delta_1$ , and  $\Delta_2$  are constants. But averaging over the joint distribution of  $\Delta_1$  and  $\Delta_2$  and expressing the  $P_{3j}$  as functions of these fitness spacings, the expected increase in fitness when starting at  $i = 3$  is

$$E[\Delta W] = E\left[\frac{(\Delta_1 + \Delta_2)^2 + (\Delta_2)^2}{\Delta_1 + 2\Delta_2}\right]. \quad (12)$$

Scaling this argument up to the case in which we start at the  $i$ th fittest sequence and letting  $Z_k = k\Delta_{k2}$  we have

$$E[\Delta W] = E\left\{\left[\left(\frac{Z_1}{1} + \frac{Z_2}{2} + \cdots + \frac{Z_{i-1}}{i-1}\right)^2 + \left(\frac{Z_2}{2} + \cdots + \frac{Z_{i-1}}{i-1}\right)^2 + \cdots + \left(\frac{Z_{i-1}}{i-1}\right)^2\right] \div [Z_1 + Z_2 + \cdots + Z_{i-1}]\right\}. \quad (13)$$

After much rearrangement, it can be shown that equation (13) is equivalent to

$$E[\Delta W] = 2 \left( \sum_{k=1}^{i-2} \frac{k}{k+1} \right) E \left[ \frac{Z_1 Z_2}{Z_1 + Z_2 + \dots + Z_{i-1}} \right] + \sum_{k=1}^{i-1} \frac{1}{k} E \left[ \frac{Z_1^2}{Z_1 + Z_2 + \dots + Z_{i-1}} \right]. \quad (14)$$

Because the  $Z_k$  are IID exponentials, it can be shown that the first expectation equals  $E[\Delta_1]/i$  and the second  $2E[\Delta_1]/i$ . Substituting and simplifying, we arrive at our final answer:

$$E[\Delta W] = \frac{2(i-1)E[\Delta_1]}{i}. \quad (15)$$

This is probably our most surprising result. Despite the fact that a vast array of distributions of allelic fitnesses surely hold in nature and that different sequences come in different lengths, with different starting fitnesses, and so on, the expected fitness jump that occurs when natural selection drives a substitution assumes a very simple form. This form behaves as it should in extreme cases: when  $i = 1$ , for instance,  $E[\Delta W] = 0$  and when  $i = 2$ ,  $E[\Delta W] = E[\Delta_1]$ , as it must. It should be emphasized that equation (15) is not complicated by any change in fitness rank when a new allele produces its own suite of one-mutational-step neighbors. Although fitness rank might change following production of new one-mutational-step neighbors, absolute fitness does not.

While simple, the expression for  $E[\Delta W]$  depends on  $i$  and  $E[\Delta_1]$ , quantities that will not generally be known ( $E[\Delta_1]$  does depend on the form of the parent distribution of fitnesses). Despite this, equation (15) leads to a prediction that, at least in principle, is testable. This prediction hinges on the fact that the mean fitness jump is insensitive to rank:  $E[\Delta W]$  goes as  $(i-1)/i$ . The smallest jump occurs at  $i = 2$ , yielding  $E[\Delta W] = E[\Delta_1]$ ; but for all but very small  $i$  the mean fitness jump is near  $2E[\Delta_1]$ .  $E[\Delta W]$  is, therefore, constrained to a window between  $E[\Delta_1]$  and  $2E[\Delta_1]$ . This fact leads to another. At any given substitution  $K = 1, 2, 3, \dots$  where we order substitutions chronologically, replicate populations might find themselves at different  $i$ , as evolution jumped to different  $i$  at earlier substitutions (at least for  $K > 1$ ). But the mean fitness jump that occurs at substitution  $K$  is a weighted average of the jumps expected at different  $i$ . It follows that the mean fitness jump that occurs at any substitution  $K$  (conditional on a substitution occurring) must *also* be constrained to a window of

$$E[\Delta_1] \leq E[\Delta W]_K \leq 2E[\Delta_1], \quad (16)$$

where the subscript  $K$  emphasizes that we refer to the mean fitness jump occurring at substitution  $K = 1, 2, 3, \dots$  not  $i = 1, 2, 3, \dots$ .

**Conclusion 3.**—*The mean fitness jumps occurring at substitutions  $K = 1, 2, 3, \dots$  (i.e.,  $E[\Delta W]_1, E[\Delta W]_2, E[\Delta W]_3, \dots$ ) will differ by no more than a factor of two under the mutational landscape model. We assume that we begin from a reasonably fit allele and consider only those cases in which a substitution does occur.*

We can also calculate the variance in fitness jump that occurs when natural selection drives a substitution. A lengthy

TABLE 2. Size of fitness jumps in simulations. Distributions as in Table 1. Each set of simulations included at least 5000 substitutions (new fitnesses drawn for each of the 5000 runs).  $E[\Delta W]$  is calculated from equation (15) and  $\text{Var}[\Delta W]$  from equation (17). Theoretical predictions require  $E[\Delta_1]$ , which was found by simulation.

Distribution	$E[\Delta W]$	$\Delta \bar{W}_{obs}$	$\text{Var}[\Delta W]$	$\sigma^2_{\Delta W_{obs}}$
$i = 10$				
Exponential	0.0180	0.0181	$1.67 \times 10^{-4}$	$1.70 \times 10^{-4}$
Gamma	0.0099	0.0098	$5.05 \times 10^{-5}$	$4.90 \times 10^{-5}$
Normal	0.0026	0.0026	$3.52 \times 10^{-6}$	$2.95 \times 10^{-6}$
$i = 50$				
Exponential	0.0196	0.0195	$1.92 \times 10^{-4}$	$1.83 \times 10^{-4}$
Gamma	0.0108	0.0113	$5.82 \times 10^{-5}$	$6.11 \times 10^{-5}$
Normal	0.0028	0.0032	$4.06 \times 10^{-6}$	$4.11 \times 10^{-6}$
$i = 150$				
Exponential	0.0199	0.0198	$1.97 \times 10^{-4}$	$1.93 \times 10^{-4}$
Gamma	0.0109	0.0112	$5.97 \times 10^{-5}$	$5.74 \times 10^{-5}$
Normal	0.0029	0.0037	$4.17 \times 10^{-6}$	$5.11 \times 10^{-6}$

calculation (Appendix 2) shows that  $E[(\Delta W)^2] = 6(i-1)E[\Delta_1]^2/(i+1)$  and so that

$$\text{Var}[(\Delta W)] = 2E[\Delta_1]^2 \left[ \frac{(i^2 + 2)(i-1)}{i^2(i+1)} \right]. \quad (17)$$

As expected,  $\text{Var}[(\Delta W)] = 0$  when  $i = 1$  and  $\text{Var}[(\Delta W)] = E[\Delta_1]^2$  when  $i = 2$ , as it must because  $\Delta_1$  is exponentially distributed.  $\text{Var}[(\Delta W)]$  is also insensitive to  $i$ , and as  $i$  grows, nears  $\text{Var}[(\Delta W)] \approx 2E[\Delta_1]^2$ . Thus, just as the mean fitness jump is constrained to a window between  $E[\Delta_1]$  and  $\sim 2E[\Delta_1]$ , so the variance is constrained to a window between  $E[\Delta_1]^2$  and  $\sim 2E[\Delta_1]^2$ , where I again assume adaptation starts from a reasonably fit allele. The variance-to-mean ratio of fitness jumps thus approaches  $E[\Delta_1]$  as  $i$  grows.

I again used computer simulations to check these results. These simulations, which were identical to those above, show that our expressions for both  $E[\Delta W]$  and  $\text{Var}[(\Delta W)]$  are reasonably accurate. Table 2 gives results for  $i = 10, 50$ , and 150. The theory typically does well, although departures from expectations appear in the normal distribution case at very large  $i$  ( $i = 150$ ). (These departures are not surprising because normal random variables show notoriously slow approach to the extreme value distribution [Gumbel 1958; David 1981, p. 264].) Figure 3 replots the simulation results as the ratio  $E[\Delta W]/E[\Delta_1]$ . This ratio quickly rises to a value of  $\sim 2$ , as predicted. More important, the ratio remains near 2 over a surprisingly broad range of  $i$ —far broader than might be guessed given our use of extreme value theory—although deviations again occur for normally distributed fitnesses at large  $i$ . In general,  $E[\Delta W]$  is constrained to a window between  $E[\Delta_1]$  and  $\sim 2E[\Delta_1]$  as long as adaptation starts from a fairly fit allele, that is,  $i < 50$  or so.

I also used simulations to confirm that the mean sizes of fitness jumps occurring throughout adaptation at substitutions  $K = 1, 2, 3, \dots$  do not differ by more than a factor of about two. This work requires new simulations that follow entire adaptive walks to the best allele (described in Appendix 1). Figure 4 shows the results when the starting  $i = 50$ . The first four substitutions are shown (adaptation typically requires only  $\sim 3.2$  substitutions under these conditions) and means

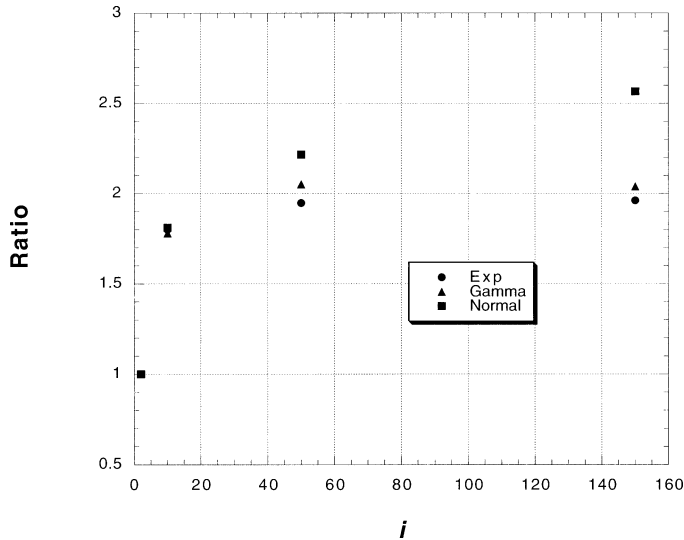


FIG. 3. The ratio  $E[\Delta W]/E[\Delta_1]$  for various distributions of allelic fitnesses and for different starting  $i$  (simulation results). As long as adaptation starts from a reasonably fit allele (small  $i$ ), this ratio is between 1.0 and 2.0. For very large  $i$ , the ratio sometimes exceeds 2.0. Results shown for exponential, gamma, and normal distributions of allelic fitnesses (details of distributions as in Table 1).

are conditional on a substitution occurring (e.g., the mean fitness jump at  $K = 3$  is calculated from runs in which a third substitution *did* occur). Normalizing mean effects at a given substitution by those at substitution 4, it can be seen that, although early substitutions are larger than later ones, their mean effects differ by less than a factor of two. Simulations show this result holds generally as long as adaptation starts from a fit allele.

#### The First and the Largest Substitutions

The above results tell us nothing about the proportion of the total increase in fitness occurring during adaptation that is due to any particular substitution, for example, the first one. This quantity depends on the number of substitutions that actually occur during adaptation. Similarly, the above results tell us nothing about the proportion of the total increase in fitness due to the substitution of *largest* effect (whether at  $K = 1, 2, 3, \dots$ ). These quantities are however, of considerable interest as they can be measured fairly directly.

Although I have not been able to find these proportions analytically, it is easy to find them by simulation. The results are shown in Figure 5. The top plot shows the mean proportion of the total increase in fitness occurring during adaptation due to the first substitution, that is, it provides an estimate of  $E[\Delta W_1/\Delta W_{total}]$ . The bottom plot shows the mean proportion of the total increase in fitness due to the largest substitution, that is, it provides an estimate of  $E[\Delta W/\Delta W_{total}]$ , where the subscript is prefixed to distinguish the order statistic from the jump occurring at  $K = 1$ . In both cases, proportions are plotted as a function of starting rank  $i$ . The first and the largest substitutions clearly account for a large proportion of the overall increase in fitness occurring during adaptive walks to the best sequence.

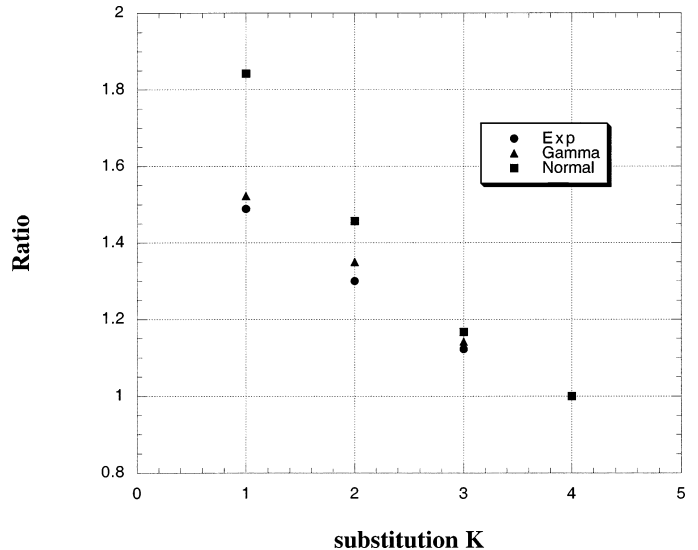


FIG. 4. The mean fitness jumps that occur throughout adaptation differ by no more than a factor of about 2. The data shown are for initial  $i = 50$  (simulation results). For each substitution  $K$ , the mean effect has been standardized by that occurring at  $K = 4$ , which is greater than the number of substitutions on average required to reach the fittest allele. Results shown for exponential, gamma, and normal distributions of allelic fitnesses (details of distributions as in Table 1).

**Conclusion 4.**—*The first substitution accounts on average for at least 30% of the total increase in fitness occurring during adaptation. The largest substitution accounts, on average, for at least 50% of the total increase in fitness. We assume only that we begin from a reasonably fit allele.*

The above results are robust to the form of the underlying distribution of fitnesses (Fig. 5), mean  $s$  (as long as it is reasonably small), and the number of one-mutational-step neighbors  $m$  (as long as it is reasonably large). Indeed the proportion of fitness increase due to the largest substitution may even be robust to the assumption that our starting allele is fairly fit: As the bottom half of Figure 5 shows, this proportion appears to asymptote at  $E[\Delta W/\Delta W_{total}] = 1/2$ . (It is difficult to test this claim rigorously, as computer simulations become prohibitively slow at much larger  $i$ .)

These results likely represent some of our most easily tested predictions. I consider relevant data in the Discussion.

#### Exponential Distribution of Selection Coefficients

A selection coefficient is a normalized fitness jump:  $s = \Delta W/W_i$ , where  $W_i$  is absolute fitness before the substitution. The simulations just described reveal an interesting fact about these selection coefficients: Although  $\Delta W$ , and thus  $s$ , tend to decrease as the population approaches the fittest allele as one would expect intuitively, the ratio of consecutive selection coefficients is roughly constant, whatever the parent distribution of allele fitnesses (not shown). For instance,

$$\frac{E[s]_1}{E[s]_2} \approx \frac{E[s]_2}{E[s]_3}, \quad (18)$$

where the subscripts correspond to substitutions  $K = 1, 2, 3$ . But this means that the mean selection coefficients fixed at  $K = 1, 2, 3, \dots$  form an approximate geometric sequence.

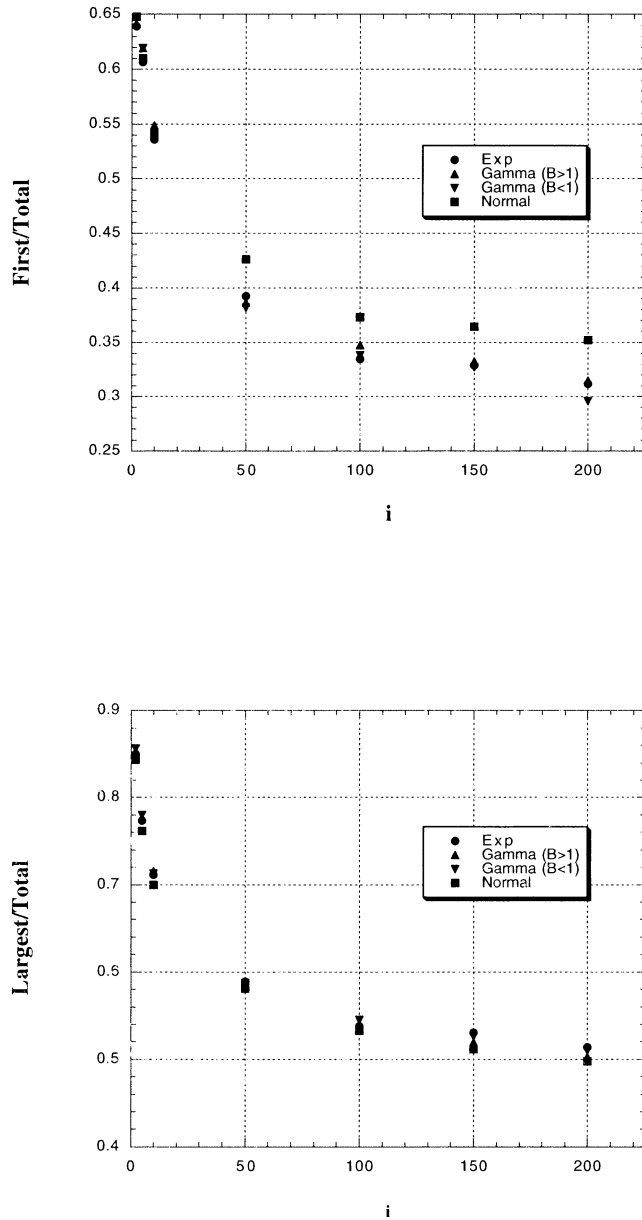


FIG. 5. Top: the mean proportion of the total increase in fitness explained by the first substitution. Bottom: the mean proportion of the total increase in fitness explained by the substitution of largest effect. Results shown for computer simulations using exponential, gamma, and normal distributions of allelic fitnesses. (Details of most distributions as in Table 1. In addition, results from a gamma distribution with shape parameter  $\beta = 1/2$  and  $\bar{s} = 0.01$  are shown as there are technical reasons for thinking  $\beta < 1$  could affect the ratio shown. Fortunately, distributions with  $\beta > 1$  and  $\beta < 1$  behave nearly identically.) A total of 2000 runs were performed at each condition.

This result is interesting for two reasons. The first is that this geometric behavior involves a quantity that is observable. It is, however, important to be clear on just what the theory does and does not predict, a task that is complicated by a somewhat paradoxical result. For, although  $E[\Delta W]$  and  $E[s]$  both decrease with  $K = 1, 2, 3, \dots$ , the expected fitness jump for the *last* substitution often exceeds that for the first. In other words, it is both true that

$$E[\Delta W]_1 > E[\Delta W]_2 > E[\Delta W]_3, \dots \quad (19)$$

and that

$$E[\Delta W]_1 < E[\Delta W]_{last} \quad (20)$$

The paradox is, however, only apparent. The explanation is straightforward: All else being equal, those substitutions with large  $\Delta W$  are more likely to reach all the way to  $j = 1$  and so to be (with high probability) the last substitution. Put differently, although the last substitution sometimes corresponds to substitution  $K = 1$ , sometimes to  $K = 2$ , etc., in each case the substitution is likely to be last only if it is drawn from the right tail of the distribution of  $\Delta W$  for a given  $K$ . The prediction is not, therefore, the vague one that early substitutions are larger than later (where we do not label substitutions by  $K$ ), but the specific one that the mean fitness effect at  $K = 1$  exceeds that at  $K = 2$ , which exceeds that at  $K = 3$ , etc.

The second reason the geometric sequence behavior is interesting is that the same pattern appears in Fisher's (1930) geometric model. In that case, the mean phenotypic sizes of favorable mutations fixed through time form an approximate geometric sequence (Orr 1998, 2000). This behavior appears to have a simple explanation. When looking over many realizations of adaptive walks to the optimum, this geometric behavior gives rise to an exponential distribution of factors fixed (ignoring factors of small effect; Orr 1998, 2000; see Discussion). The fact that  $E[s]$  also shows approximate geometric behavior over  $K = 1, 2, 3, \dots$  suggests that adaptation in DNA sequence space might *also* give rise to an exponential distribution of selection coefficients among factors fixed over many realizations of adaptive walks to the best allele.

This notion is, of course, little more than intuition. But it is easily tested in computer simulations. I thus simulated 10,000 adaptive walks to the locally best allele using exponential, gamma, or normal distributions of allelic fitnesses. In all cases, the distribution of selection coefficients among mutations fixed is exponential, where I again ignore factors of small effect (see Fig. 6, where starting  $i = 2, 5, 50$ , or 150). More exactly, then,  $d \ln \psi(s)/ds$  is a constant over modest to large  $s$ , as indicated by the straight lines in the semi-log plots shown in Figure 6. Many simulations at many different  $i$  confirmed the generality of this result.

*Conclusion 5.—The distribution of selection coefficients fixed over many realizations of adaptive walks to the fittest allele is nearly exponential, where we ignore mutations of small effect.*

## DISCUSSION

Two factors make adaptation a stochastic process. The first is that the course of adaptive evolution depends on the vagaries of which mutations appear when. The second is that the course of adaptive evolution is shaped by probability of fixation: most new favorable mutations suffer a high probability of accidental loss and so rarely contribute to adaptation. When the vagaries of both mutation and probability of fixation are taken into account, we are, fortunately, left with a simple stochastic process: One among the several possible favorable mutations will be the next fixed and adaptation can be viewed as a Markov chain in which different mutant

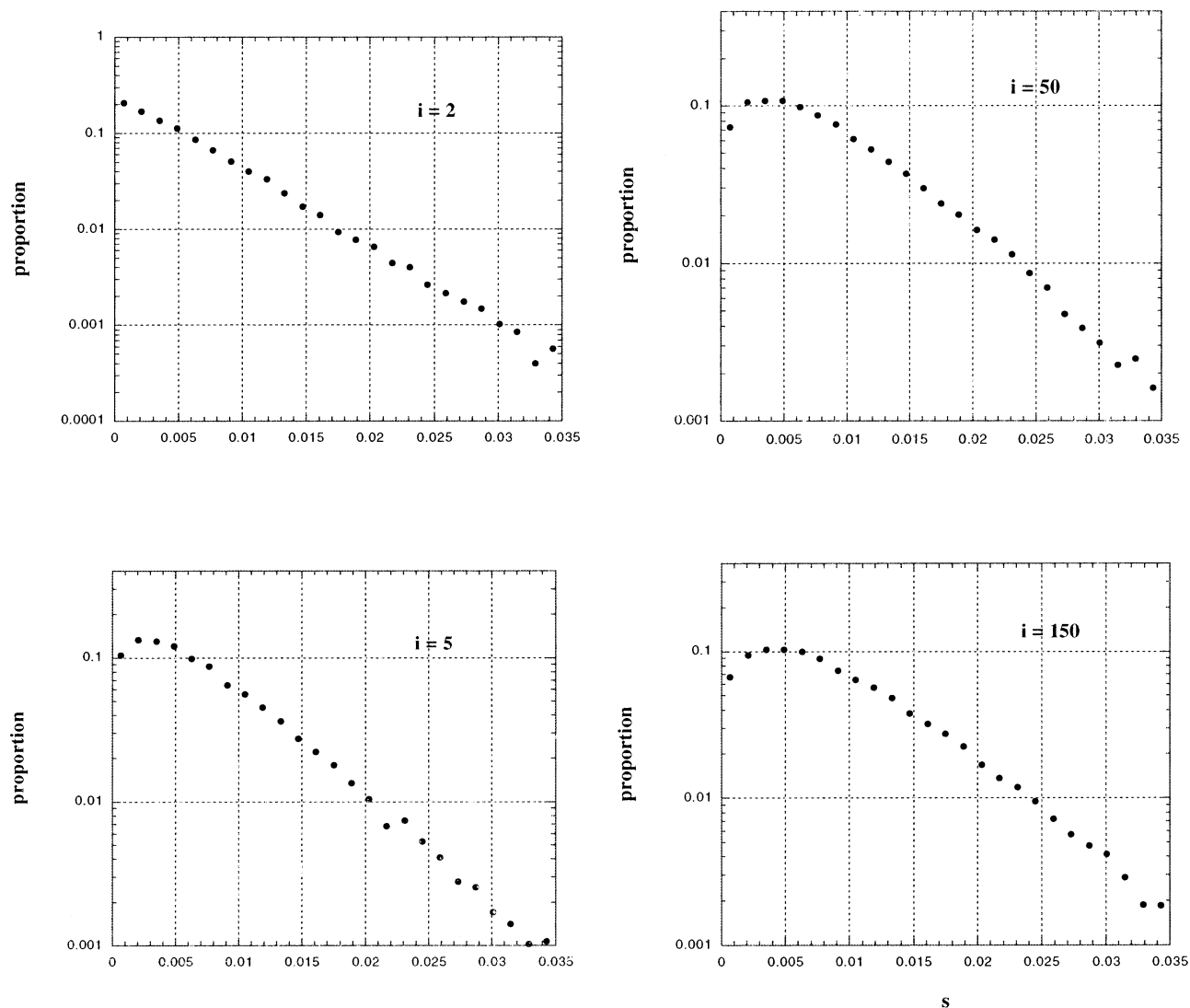


FIG. 6. Selection coefficients fixed throughout adaptation are nearly exponentially distributed (ignoring small  $s$ ). The data shown reflect 10,000 realizations of adaptive walks to the fittest allele when using gamma distributed allelic fitnesses (simulation results; details of gamma distribution as in Table 1; the results are essentially unchanged when the shape parameter  $\beta < 1$ ). The four cases shown correspond to different starting  $i$ , as indicated.

sequences are in turn fixed, until the population arrives at the locally fittest allele. As long as selection is relatively strong and mutation relatively weak—and we consider evolutionary response to a single environmental change—this process is readily quantified (Gillespie 1983, 1984, 1991).

Building on Gillespie's seminal work, I have analyzed this stochastic process. Although I have specified remarkably little about the biological scenario under consideration—no assumptions were made about the distribution used to assign fitnesses to alleles (except that it is Type III); the length of the DNA sequence considered; the number of mutations that are possible; and the local recombinational environment (which is largely irrelevant as beneficial mutations are rare)—a number of nontrivial conclusions were reached. This surprising independence from essentially all biological detail is due to two facts, both identified by Gillespie.

The first is that, at any point in time, we need only consider

the local sequence space, that is, those sequences that are one mutational step away from wild type. The second is that certain limit theorems describe the distribution of fitness spacings between the fittest several alleles regardless of the parent distribution from which fitnesses are drawn. More exactly, we know that the fitness spacings between alleles in the right tail of the distribution of allelic fitnesses—the only alleles we care about—are independent exponential random variables with means given in equation (3). Given these two simplifications, I was able to draw a number of conclusions.

The first is that adaptation is characterized by large jumps in the fitness rank of alleles. In particular, if a population is presently fixed for the  $i$ th fittest allele, it will on average jump to the  $E[j] = (i + 2)/4$ th best allele at the next substitution (eq. 7). This is perhaps our most counterintuitive result. Given only the present wild-type allele's fitness rank, we can predict the rank of the mutant allele fixed by natural selection.

The fact that these jumps in rank are large was implicit in Gillespie's (1983; 1991, p. 243) demonstration that a few substitutions are typically needed to reach the locally best allele. But it was not clear that the mean jump size would assume such a simple form. More important, Gillespie assumed that the population starts from a very fit allele. The present simulations show, however, that the  $E[j] = (i + 2)/4$  result holds for surprisingly large  $i$  (Table 1). As long as the wild-type sequence from which we start resides in the top 5% or so of allelic fitnesses, the theory appears to remain reasonably accurate. It should also be understood that this top 5% refers to the *local* landscape, that is, to those sequences that are one mutational step away from the present allele. In a protein comprising 1000 amino acids (3000 bp), the top 5% of alleles includes 150 sequences. It is difficult to believe that the typical wild-type allele is often less fit than 150 alternative sequences that can be reached by a single base change.

The second main result is that adaptation by natural selection behaves in several respects as the average of two idealized and extreme forms of adaptation. In one, adaptation is perfect, that is, it is characterized by an algorithm that always chooses the best possible alternative. In the other, adaptation is random, that is, it is characterized by an algorithm that randomly chooses from those alternatives that improve matters. Our finding of  $E[j] = (i + 2)/4$  under natural selection falls precisely midway between that seen under perfect ( $E[j] = 1$ ) and random ( $E[j] = i/2$ ) adaptation. In a sense, then, adaptation by natural selection is "not half bad" when compared to a perfect form of adaptation. (This result is exact whether considering the simple or mutational landscape models.) It is worth considering the reasons for this behavior.

It is a common intuition that not-half-bad behavior is an automatic consequence of the fact that probability of fixation is linear in  $s$  ( $\Pi = 2s$ ). This fact, after all, immediately leads to  $P_{ij} = s_j / \sum_{k=1}^{i-1} s_k$  and so guarantees that adaptation is neither deterministic (as when the best alternative is always chosen) nor completely random (as when all alternatives have the same chance of being chosen).  $\Pi = 2s$  is not, however, sufficient to yield not-half-bad behavior. To see this, assume that  $\Pi$  and so  $P_{ij}$  are as above but that fitness spacings are *equal* between all adjacent alleles. In this case, we get  $E[P_{ij}] = 2(i - j)/[i(i - 1)]$ , that is, transition probabilities are linear in  $j$ , and  $E[j] = (i + 1)/3$ . Although intermediate between perfect and random, this mean rank is not midway between the two. This proves that not-half-bad behavior requires that the relationship between transition probability  $E[P_{ij}]$  and  $j$  be steeper than linear (this is a necessary condition). Given  $\Pi = 2s$ , this in turn requires that fitness spacings decrease as one moves to alleles of lesser fitness (as in Fig. 1). Not-half-bad behavior thus depends on the form of extreme spacings between favorable alleles. In sum, when combining probability of fixation, which reflects the inevitable consequences of sampling Mendelian inheritance, and extreme fitness spacings, which reflects the inevitable tail behavior of fitness distributions, we obtain an optimizing algorithm, natural selection, that sits midway between perfect and random in mean efficiency.

The size of the jump in fitness rank is closely connected to the number of substitutions needed to reach the best allele:

If jumps are large, few substitutions are needed. It is not surprising, therefore, that the mean number of substitutions required to reach the best allele under natural selection also falls midway between that required under perfect versus random adaptation. Specifically, we showed that Gillespie's (1983, 1991) well-known calculation of the number of substitutions before absorption at the fittest allele can be written  $L_{nat\ sel} = (1 + \sum_{k=1}^{i-1} 1/k)/2$ , which equals the arithmetic mean of absorption times under perfect and random adaptation (eq. 10). Although these absorption times, like Gillespie's (1983), assume a simple mutational model in which all sequences are mutually accessible, simulations show that a similar, although approximate, result characterizes the full mutational landscape model. By at least two measures, then, adaptation by natural selection behaves (at least approximately) as the average of perfect and random adaptation. It should be noted that, with both of these measures, not-half-bad behavior concerns fitness *rank*, not fitness *per se*, and characterizes mean changes in rank, not higher moments.

The third and perhaps most significant result is that we can write down the mean jump in fitness that occurs when natural selection drives the substitution of a new allele. This expression assumes a surprisingly simple form that is nearly independent of starting rank  $i$  (eq. 15). The size of the mean fitness jump does, however, depend on the mean fitness spacing between the fittest and next-to-fittest allele,  $E[\Delta_1]$ . Neither  $i$  nor  $E[\Delta_1]$  will generally be known. Despite this, the theory allows a prediction that is testable, at least in principle. If, following Gillespie (1984, 1991), we assume that the distribution of allelic fitnesses remains constant throughout a walk to the optimum, the mean fitness jumps that occur throughout adaptation are constrained to a window between  $E[\Delta_1]$  and  $2E[\Delta_1]$  (eq. 16), where we assume adaptation starts from a fairly fit allele. (For much larger  $i$ , we have no guarantee that extreme value theory holds and the mean fitness jump might exceed  $2E[\Delta_1]$ . Indeed such exceedences were seen in computer simulations. But even here the deviations from theory were not qualitative [see Fig. 3]. But because we do not generally know the starting  $i$ , our prediction is perhaps best taken as claiming that the mean fitness jumps will differ by no more than a factor of about two.) The prediction captured in equation (16) does not, of course, mean that the first substitution cannot have 10 times the effect of the second one in any particular bout of adaptation. It can. The prediction concerns the *mean* fitness jumps at the first versus second, etc. substitutions.

The fourth result is that the first substitution and the substitution of largest effect explain a large proportion of the total increase in fitness that occurs during adaptation. Figure 5 shows that the first substitution accounts, on average, for at least 30% of the overall increase in fitness, whereas the largest substitution accounts on average for at least 50%. These findings are independent of, for example, the underlying distribution of fitnesses and the size of the gene. These findings differ from those seen in Fisher's (1930) geometric model of adaptation (Orr 1998). There, it is difficult to specify the proportional contribution of the first or largest substitution: The answer depends on the dimensionality of the organism and the spectrum of mutational effects provided to natural selection, neither of which is known in any actual

case. But roughly speaking, single substitutions in Fisher's model tend to explain little of the overall progress toward the optimum (Orr 1998, eq. 4; Orr 2000, eq. 4). The present results suggest that this finding may be at least partly an artifact of Fisher's model having no necessary last substitution: If adaptive evolution stopped after a small number of substitutions, single steps would explain a considerable portion of total progress. In any case, this is the pattern seen in models of adaptation at the DNA level.

The last result concerns the selection coefficients fixed during adaptation in the mutational landscape model. Computer simulations show that the mean size of selection coefficients fixed at substitutions  $K = 1, 2, 3, \dots$  decreases as an approximate geometric sequence. Early substitutions thus have larger effects than later. This geometric behavior is closely connected to another pattern. When looking over many realizations of adaptive walks to the best allele, the distribution of selection coefficients fixed is exponential, where we ignore factors of small effect. This result is remarkably robust, arising independently of the parent distribution of allele fitnesses and initial  $i$ . This exponential tail behavior is reminiscent of that seen in adaptation in Fisher's geometric model (Orr 1998, 1999). This common result is surprising as Fisher's model and the present one differ in fundamental ways. As emphasized in the introduction, Fisher's is a continuous space model couched in terms of phenotypic effects, whereas the present sequence model is a discrete space model couched in terms of fitness effects. Nonetheless, adaptation in both gives rise to exponential behavior among the factors fixed. The cause of this behavior seems clear. Adaptation in both models involves a repeated rescaling. At each step in adaptation, the population moves closer to the optimum (whether phenotypic or genotypic); thus, at the next substitution, adaptation confronts essentially the same problem as at the previous step, but on a smaller scale. The distribution of factors fixed at earlier versus later substitutions thus shows roughly the same functional form but on a scale that decreases by a constant proportion at each step. Summing over such a self-similar process appears to give rise to a mixed distribution having exponential tail behavior (Orr 1998, 2000; H. A. Orr, unpubl. ms.). It appears then that we have identified a pattern that is robust to the variety of adaptation model considered: Adaptation by natural selection seems characterized by an approximately exponential distribution of factors fixed.

The exponential tail behavior seen here does, however, differ in two technical ways from that seen in Fisher's (1930) model. First, exponential behavior does not emerge in Fisher's model unless the population takes a fairly large number of steps to the optimum. But in the present model exponential behavior emerges even when the initial  $i = 2$  and the population takes few steps to reach a locally optimal allele (Fig. 5). Thus, a caveat that applied to Fisher's model does not apply to the DNA model. Second, the cause of the shoulder of nonexponential effects among small factors differs in the two models. In Fisher's model, this nonexponential region appears to reflect the fact that one does not proceed all the way to the optimum: Because there is no necessary last substitution in Fisher's model, computer simulations must be stopped at some arbitrary point, and simulations show that

the nonexponential region decreases as the population travels 90% to 95% to 98% of the distance to the optimum. In the present model, this argument cannot hold because, in a discrete DNA sequence space, there *are* locally best alleles and true last substitutions. The nonexponential shoulder cannot, then, be an artifact. Despite these technical differences, it remains surprising that such different models yield such similar behavior.

The exponential behavior found here is not, however, the same as that claimed in two previous studies of molecular evolution. Indeed there seems to be considerable confusion about just what is and is not exponentially distributed in these models. In particular, Gillespie (1983) showed that the fitness spacing between the best and next-to-best alleles is exponential (eq. 3 in this paper). Thus, if evolution always involved replacing the second-best with the best allele, the fitness jumps occurring during adaptation would be exponential. But there is no reason to think that adaptation exclusively involves jumps from  $i = 2$  to  $j = 1$  (see also Otto and Jones 2000), and it was not obvious that a process that starts at arbitrary  $i$  and gradually walks to a best local allele would also show exponential behavior. Second, Wahl and Krakauer (2000) claimed that extreme value theory shows that single leaps from the  $i$ th fittest allele to the fittest allele ( $j = 1$ ) yield an exponential distribution of fitness jumps. This is incorrect. From Figure 1 and equation (3) it can be shown that the moment generating function characterizing such leaps is  $M(t) = (i - 1)! \Gamma(1 - t/\lambda) / \Gamma(i - t/\lambda)$ , where  $\lambda = 1/E[\Delta_1]$ . It is not clear what this distribution is, but it is not exponential. (It does, however, reduce to an exponential when  $i = 2$ , as it must.) Exponential behavior characterizes adaptive *walks* to the fittest allele, not single leaps to the fittest allele.

### Conclusions

I have identified several patterns that characterize the adaptation of DNA sequences under certain, hopefully not too restrictive, conditions (e.g., strong selection, weak mutation, random assignment of fitness from a "reasonable" probability distribution). Some of these results are less model dependent than others. In particular, most of our results concern a single step in adaptation and thus do not depend on Gillespie's (1984, 1991) assumption that the distribution of allelic fitnesses stays the same through time. Several results, however, concern entire adaptive walks and so do depend on this assumption. Similarly, whereas some of the present results may be of strictly theoretical interest, others (e.g., large jumps in fitness at the first substitution, the geometric sequence behavior, the largest substitution prediction, and the exponential tail behavior) refer to measurable quantities and are, at least in principle, testable. The most obvious routes to testing these predictions would seem to involve microbial experimental evolution work or comparative sequence analysis (the latter might allow estimation of selection coefficients from polymorphism data, e.g., by measuring the effects of selective sweeps). At present, little relevant data are available from either source, although experimental evolution work provides at least some information.

Such work shows, for instance, that single substitutions

sometimes have large fitness effects, as predicted (Wichman et al. 1999; Bull et al. 2000). Such work also shows that early steps in adaptation typically have larger fitness effects than later ones, also as predicted. While several early experiments suggested this pattern (e.g., Lenski and Travisano 1994; Bull et al. 1997), the best data come from Holder and Bull's (2001) study of two species of DNA bacteriophage. (I ignore RNA phage because their high mutation rates may mean they can substitute two-mutational-step neighbors.) Holder and Bull studied adaptation to high temperature in  $\phi$ X174 and G4. Whole genome sequencing showed that adaptation in  $\phi$ X174 was underlain by five substitutions. The first two substitutions accounted for 80% of the total increase in fitness over the course of the experiment; indeed a single point substitution (the first) appeared to explain about 75% of the ultimate fitness increase. Whereas adaptation in G4 was more complex, early substitutions again tended to be larger than later, with half the total fitness increase due to the first three. Such findings provide qualitative support for the patterns predicted here, but no stronger conclusion seems justified. The problem is that labor-intensive experimental evolution work usually involves less-than-ideal replication; moreover, experimental evolution work usually involves strong selection, whereas the theory presented here assumes small  $s$ . But firmer conclusions would clearly be possible given well-replicated experiments involving weaker selection.

But our most important conclusion does not concern any particular prediction. It instead concerns our more general finding that adaptation is characterized by certain simple patterns. It seems commonly believed that adaptation, unlike neutral or deleterious evolution, is inherently and hopelessly complex, with the basis of adaptive evolution differing in every case as a consequence of differences in the developmental pathway involved, the nature of the environmental challenge, and so on. The present work, building on that of Gillespie (1983, 1984, 1991), suggests that this is not true.

#### ACKNOWLEDGMENTS

I thank N. Barton, A. Betancourt, J. Bull, J. H. Gillespie, J. Huelsenbeck, J. Jaenike, T. Johnson, J. P. Masly, D. Presgraves, L. Rieseberg, G. P. Wagner, and M. Whitlock for helpful comments and discussion. I also thank W. Mosle, L. Orr, and D. Wackeroth for solving the integrals in equation (14). This work was supported by National Institutes of Health grant 2R01 G51932-06A1 and by The David and Lucile Packard Foundation.

#### LITERATURE CITED

- Barton, N. H. 1998. The geometry of natural selection. *Nature* 395: 751–752.  
 ———. 2001. The role of hybridization in evolution. *Mol. Ecol.* 10:551–568.  
 Bull, J. J., M. R. Badgett, H. A. Wichman, J. P. Huelsenbeck, D. M. Hillis, A. Gulati, C. Ho, and I. J. Molineux. 1997. Exceptional convergent evolution in a virus. *Genetics* 147:1497–1507.  
 Bull, J. J., M. R. Badgett, and H. A. Wichman. 2000. Big-benefit mutations in a bacteriophage inhibited with heat. *Mol. Biol. Evol.* 17:942–950.  
 Darwin, J. H. 1957. The difference between consecutive members of a series of random variables arranged in order of size. *Biometrika* 44:211–218.

- David, H. A. 1981. *Order statistics*. John Wiley and Sons, New York.  
 Fisher, R. A. 1930. *The genetical theory of natural selection*. Oxford Univ. Press, Oxford, U.K.  
 Gerrish, P. 2001. The rhythm of microbial adaptation. *Nature* 413: 299–302.  
 Gillespie, J. H. 1983. A simple stochastic gene substitution process. *Theor. Popul. Biol.* 23:202–215.  
 ———. 1984. Molecular evolution over the mutational landscape. *Evolution* 38:1116–1129.  
 ———. 1991. *The causes of molecular evolution*. Oxford Univ. Press, New York.  
 ———. 2002. Why  $k = 4Nus$  is silly. In R. S. Singh and C. B. Krimbas, eds. *Evolutionary genetics: from molecules to morphology*. Vol. III. Cambridge Univ. Press, Cambridge, U.K. *In press*.  
 Gumbel, E. J. 1958. *Statistics of extremes*. Columbia Univ. Press, New York.  
 Haldane, J. B. S. 1927. A mathematical theory of natural and artificial selection. V. Selection and mutation. *Proc. Camb. Philos. Soc.* 28:838–844.  
 Hartl, D., and C. H. Taubes. 1996. Compensatory nearly neutral mutations: selection without adaptation. *J. Theor. Biol.* 182: 303–309.  
 ———. 1998. Towards a theory of evolutionary adaptation. *Genetica* 102/103:525–533.  
 Holder, K. K., and J. J. Bull. 2001. Profiles of adaptation in two similar viruses. *Genetics* 159:1393–1404.  
 Kauffman, S., and S. Levin. 1987. Towards a general theory of adaptive walks on rugged landscapes. *J. Theor. Biol.* 128:11–45.  
 Kimura, M. 1979. Model of selectively neutral mutation in which selective constraint is incorporated. *Proc. Natl. Acad. Sci. USA* 76:3440–3444.  
 ———. 1983. *The neutral theory of molecular evolution*. Cambridge Univ. Press, Cambridge, U.K.  
 Lenski, R. E., and M. Travisano. 1994. Dynamics of adaptation and diversification: a 10,000-generation experiment with bacterial populations. *Proc. Natl. Acad. Sci. USA* 91:6808–6814.  
 Orr, H. A. 1998. The population genetics of adaptation: the distribution of factors fixed during adaptive evolution. *Evolution* 52: 935–949.  
 ———. 1999. The evolutionary genetics of adaptation: a simulation study. *Genet. Res.* 74:207–214.  
 ———. 2000. Adaptation and the cost of complexity. *Evolution* 54:13–20.  
 Otto, S. P., and C. D. Jones. 2000. Detecting the undetected: estimating the total number of loci underlying a trait in QTL analyses. *Genetics* 156:2093–2107.  
 Wahl, L. M., and D. C. Krakauer. 2000. Models of experimental evolution: the role of genetic chance and selective necessity. *Genetics* 156:1437–1448.  
 Wichman, H. A., M. R. Badgett, L. A. Scot, C. M. Boulianne, and J. J. Bull. 1999. Different trajectories of parallel evolution during viral adaptation. *Science* 285:422–424.

Corresponding Editor: M. Whitlock

#### APPENDIX 1

##### Computer Simulations

I performed two kinds of computer simulations. The first were fully stochastic (individual based) but preliminary. In these, the fates of unique recurrent mutant alleles were followed in a haploid population of size  $N$ . These simulations, like those in Gillespie (1983), confirmed the accuracy of his Markov chain approach under strong selection ( $2Ns \gg 1$ ) and weak mutation ( $N\mu < 1$ ); that is, they confirmed that  $P_{ij}$  is as given in equations (1) and (2). These simulations were also used to confirm our main results at small  $i$ :  $E[P_{ij}]$ ,  $E[j]$ ,  $\text{Var}[j]$ ,  $E[\Delta W]$ , and  $\text{Var}[\Delta W]$  approximately behave as predicted by theory. Unfortunately, however, these fully stochastic simulations were very slow and so were limited to small starting  $i$ . (Even at  $i = 3$  and  $i = 4$ , I was restricted to simulating small populations, i.e.,  $N \approx 2000$ –5000.)

I thus turned to a second kind of simulation. These were not fully stochastic, that is, they did not track finite populations. The simulations were therefore much faster, allowing study of large  $i$  and many replicates. These are the simulations presented throughout the text. These simulations assumed two population genetic facts: (1) probability of fixation is  $\Pi = 1 - \exp(-2s)$ ; and (2)  $P_{ij} = \Pi_j / (\Pi_1 + \Pi_2 + \dots + \Pi_{i-1})$ , as in Gillespie (1983). Most of the simulations considered a gene of length 1000 bp. There were thus  $m = 3000$  possible one-step mutations. Fitnesses were assigned to  $m + 1$  alleles as described below. Alleles were then ordered by fitness and one was designated the present wild type, with fitness rank  $i$ . A probability of fixation was calculated for each favorable mutation and the  $P_{ij}$  were directly calculated by equation (1). These  $P_{ij}$ -values were then used in a Monte Carlo scheme to choose the next allele fixed. In simulations of single substitutions, the rank of the allele fixed, the fitness jump, and the selection coefficient involved in its substitution, etc. were recorded and the first run was complete. In each subsequent run of the simulation a new set of  $m + 1$  fitnesses was drawn from the parent distribution of fitnesses, alleles were again ordered by fitness, the  $i$ th fittest was designated the wild type, and the above sequence of events was repeated.

Simulations of entire adaptive walks to the locally best allele were similar to the above, except that they were iterated through several substitutions. A population was allowed to fix a favorable allele (which retains its absolute fitness), produce a new suite of  $m$  mutant alleles drawn from the same distribution of allelic fitnesses, fix a new allele, and so on, until arriving at an allele that was fitter than all  $m$  one-step mutations.

Although these computer simulations are not fully stochastic, they are more exact than the analytic theory in several ways. Most important, the analytic theory assumes that the initial allele  $i$  is near the top in fitness; the simulations assume nothing about extreme value theory or starting fitness rank. I thus used simulations to explore larger  $i$ , where one might expect exponential fitness spacings and the analytic results to begin breaking down. Second, the analytic theory depends on extreme spacings that, although exact for an exponential parent distribution, are only asymptotically correct for all others (Gumbel 1958). I explored the accuracy of this asymptotic theory by studying a variety of parent distributions (exponential, gamma, and normal). Finally, the analytic theory assumes that  $\Pi = 2s$ , whereas the simulations used the more exact  $\Pi = 1 - \exp(-2s)$ , which gives the probability of fixation when  $Ns \gg 1$  (Kimura 1983).

#### Fitness Assignments

Absolute fitnesses were assigned to alleles in the above (not fully stochastic) simulations as follows. Fitness was set as  $W = 1 + S$ , where  $S$  was drawn from a specified probability distribution (exponential, gamma, or normal). Thus  $E[W] = 1 + E[S]$  and  $\text{Var}[W] = \text{Var}[S]$ . This approach allows the production of small fitness differences between alleles, ensuring that selection is mild in absolute terms. Statements in the text that allelic fitnesses or selection coefficients were exponentially, gamma, or normally distributed should thus be taken as shorthand for  $S$  being exponentially, gamma, or normally distributed.

The selection coefficient,  $s$ , separating any two alleles  $m$  and  $n$  is  $s = (W_m - W_n)/W_n = (S_m - S_n)/(1 + S_n)$ . If we let allele  $n$  be a standard allele with  $W_n = 1$  and  $S_n = 0$ , then  $S$  is a traditional selection coefficient,  $s = S_m$ . In practice, however, we calculate a selection coefficient between the present wild-type  $i$  and a favorable

mutant  $j$  where neither necessarily has an absolute fitness of one. Thus,  $s = (S_j - S_i)/(1 + S_i)$ .

#### APPENDIX 2

##### The Variance in Fitness Jumps

Here I sketch the calculation of  $E[(\Delta W)^2]$  and  $\text{Var}[(\Delta W)]$ . By analogy with the argument leading to equation (13) of the text, it is clear that

$$E[(\Delta W)^2] = E \left\{ \left[ \left( \frac{Z_1}{1} + \frac{Z_2}{2} + \dots + \frac{Z_{i-1}}{i-1} \right)^3 + \left( \frac{Z_2}{2} + \dots + \frac{Z_{i-1}}{i-1} \right)^3 + \dots + \left( \frac{Z_{i-1}}{i-1} \right)^3 \right] \div [Z_1 + Z_2 + \dots + Z_{i-1}] \right\}. \quad (\text{A1})$$

After much rearrangement it can be shown that (A1) is equivalent to

$$E[(\Delta W)^2] = c_1 E \left[ \frac{Z_1^3}{Z_1 + Z_2 + \dots + Z_{i-1}} \right] + c_2 E \left[ \frac{Z_1 Z_2 Z_3}{Z_1 + Z_2 + \dots + Z_{i-1}} \right] + c_3 E \left[ \frac{Z_1^2 Z_2}{Z_1 + Z_2 + \dots + Z_{i-1}} \right]. \quad (\text{A2})$$

The constants in (A2) are equal to

$$c_1 = \sum_{k=1}^{i-1} \frac{1}{k^2}, \quad (\text{A3a})$$

$$c_2 = 6 \sum_{k=1}^{i-3} \left( \frac{k}{k+1} \sum_{m=k+2}^{i-1} \frac{1}{m} \right), \quad \text{and} \quad (\text{A3b})$$

$$c_3 = 3 \sum_{k=1}^{i-2} \left( \frac{k}{(k+1)^2} + \frac{1}{k} \sum_{m=k+1}^{i-1} \frac{1}{m} \right). \quad (\text{A3c})$$

Because the  $Z_k$  are IID exponentials, it can be shown that the three expectations in equation (A2) are equal to

$$E \left[ \frac{Z_1^3}{Z_1 + Z_2 + \dots + Z_{i-1}} \right] = \frac{6E[\Delta_1]^2}{i+1}, \quad (\text{A4a})$$

$$E \left[ \frac{Z_1 Z_2 Z_3}{Z_1 + Z_2 + \dots + Z_{i-1}} \right] = \frac{E[\Delta_1]^2}{i+1}, \quad \text{and} \quad (\text{A4b})$$

$$E \left[ \frac{Z_1^2 Z_2}{Z_1 + Z_2 + \dots + Z_{i-1}} \right] = \frac{2E[\Delta_1]^2}{i+1}. \quad (\text{A4c})$$

Substituting (A3) and (A4) into (A2) and extensively rearranging, one finds that the second moment assumes a simple form:

$$E[(\Delta W)^2] = \frac{6(i-1)E[\Delta_1]^2}{i+1}. \quad (\text{A5})$$

It follows that  $\text{Var}[(\Delta W)]$  is given by equation (17) of the text.