# Why Are Phenotypic Mutation Rates Much Higher Than Genotypic Mutation Rates?

**Reinhard Bürger,**[*,†,1] **Martin Willensdorfer**[†] **and Martin A. Nowak**[†]

*\*Department of Mathematics, University of Vienna, 1090 Vienna, Austria and †Program for Evolutionary Dynamics, Department of Mathematics and Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, Massachusetts 02138*

## ABSTRACT

The evolution of genotypic mutation rates has been investigated in numerous theoretical and experimental studies. Mutations, however, occur not only when copying DNA, but also when building the phenotype, especially when translating and transcribing DNA to RNA and protein. Here we study the effect of such phenotypic mutations. We find a maximum phenotypic mutation rate, $u_{max}$, that is compatible with maintaining a certain function of the organism. This may be called a phenotypic error threshold. In particular, we find a minimum phenotypic mutation rate, $u_{min}$, with the property that there is (nearly) no selection pressure to reduce the rate of phenotypic mutations below this value. If there is a cost for lowering the phenotypic mutation rate, then $u_{min}$ is close to the optimum phenotypic mutation rate that maximizes the fitness of the organism. In our model, there is selective pressure to decrease the rate of genotypic mutations to zero, but to decrease the rate of phenotypic mutations only to a positive value. Despite its simplicity, our model can explain part of the huge difference between genotypic and phenotypic mutation rates that is observed in nature. The relevant data are summarized.

THE evolution of mutation rates by natural selection has attracted the attention of evolutionary biologists for many decades (Sturtevant 1937), and a large number of models have been developed to understand various aspects of the evolution of mutation rates (Sniegowski *et al.* 2000). In contrast to the 1930s, a substantial body of empirical data about mutation rates at many levels (per base pair, per gene, or genomic) and for many different organisms is now available (Drake *et al.* 1998). For instance, mutation rates per base pair per replication in microbes with DNA chromosomes range from ~$7 \times 10^{-7}$ down to ~$7 \times 10^{-11}$. There is a strong negative correlation with the genome size, so that the mutation rates per genome differ only by about a factor of two for the organisms cited in Drake *et al.*'s (1998) Table 4. Mutation rates per base pair estimated from specific loci in higher eukaryotes are in the range $2 \times 10^{-10}$ to $5 \times 10^{-11}$ (Drake *et al.*'s 1998 Table 5). Per locus mutation rates also vary widely, even within an organism, with an approximate range from $10^{-4}$ to $10^{-6}$.

In addition to these "genotypic" mutations, organisms are also confronted with what we call "phenotypic" mutations. These are the errors that occur when a DNA-coded gene is transcribed to mRNA and subsequently translated to protein. First measurements of phenotypic mutation rates, in particular, *Escherichia coli* RNA polymerase error rates, were obtained by Springgate and

Loeb (1975). Soon afterward, Edelmann and Gallant (1977) measured the cysteine misincorporation rate for the *E. coli* protein flagellin. These early studies indicated that phenotypic mutation rates are by orders of magnitude larger than genotypic mutation rates. Later studies of *E. coli* (Ellis and Gallant 1982) estimated a global phenotypic error rate of $4.5 \times 10^{-4}$ per codon and confirmed the difference between phenotypic and genotypic mutation rates. Studies in yeast yield similar results (Shaw *et al.* 2002). In contrast to genotypic mutation rates, there does not seem to be a significant difference between eukaryotic and prokaryotic phenotypic mutation rates. Several proofreading and quality control mechanisms exist that increase the accuracy of transcription and translation (Thomas *et al.* 1998; Ibba and Söll 1999; Withey and Friedman 2002). But apparently there is not enough evolutionary pressure to increase the accuracy of the transcription and translation apparatus to DNA replication standards.

Apart from the fact that the huge differences between genotypic and phenotypic mutation rates are puzzling, such high phenotypic mutation rates conceivably could pose a problem because the production of functional protein requires the absence of a deleterious mutation event during transcription and translation. Therefore, cells with a higher phenotypic mutation rate must produce more molecules of this protein than cells with a lower rate. If the production of protein is associated with costs, a selective pressure to reduce the phenotypic mutation rate might be expected. The problem may be exacerbated when more genes have to be transcribed

[1]*Corresponding author:* Department of Mathematics, University of Vienna, Nordbergstrasse 15, A-1090 Vienna, Austria.
E-mail: reinhard.buerger@univie.ac.at

and expressed to increase the "fitness" of a cell, or rather of a single-cell organism, above its current value. In modification of STURTEVANT (1937), who asked "Why does the mutation rate not evolve to zero?" we ask "Why does the phenotypic mutation rate not evolve to lower levels, whereas the genotypic mutation rate has evolved?"

To address this question, we develop a mathematical model of a large population of single-cell organisms with a DNA chromosome, in which genotypic and phenotypic mutations occur. All mutations are assumed to be deleterious. For motivation, let us start by investigating the following simple case. Suppose that a certain gene can perform a phenotypic function if at least $k$ error-free proteins (actually, molecules of the same protein) have been produced. The function leads to a fitness advantage $s$. If the gene is not transcribed and translated, hence the function not performed, the fitness is $f_0$. Each protein molecule that is produced causes costs $c$. The probability that a protein is error free is given by $1 - u$. Thus, $u$ is the phenotypic mutation rate. Then, the expected fitness of an organism that produces $m$ copies of the protein is given by

$$\bar{f} = (f_0 + s)P + f_0(1 - P) - cm = f_0 - cm + sP, \quad (1)$$

where

$$P = \sum_{i=k}^{m} \binom{m}{i}(1 - u)^i u^{m-i}$$

is the probability that generation of $m$ proteins produces at least $k$ error-free proteins.

Taking genotypic mutations of rate $\mu$ per gene into account and assuming that only genes without a mutation can produce functional protein (thus, all mutations considered are detrimental), we need

$$\bar{f}(1 - \mu) > f_0 \quad (2)$$

for the gene to be maintained in the population. This is analogous to the classical error threshold that sets an upper limit on the evolutionary acceptable mutation rate (EIGEN and SCHUSTER 1977; SCHUSTER and FONTANA 1999). Because $P \leq 1$, inequality (2) holds if

$$s > cm + \frac{f_0 \mu}{1 - \mu} \approx cm,$$

where the approximation assumes that $\mu$ is small. We note that $s$, $c$, $\mu$, $u$, $k$, and $m$ are parameters that are constant in the cell population but may depend on the given gene. What varies among cells is the actual number of error-free protein molecules produced.

The mean number of error-free molecules produced is $m(1 - u)$, which needs to be $\geq k$. Therefore, we have

$$s > \frac{c_0}{1 - u}, \quad (3)$$

where $c_0 = ck$ is the minimum cost in terms of protein production that is necessary for performing this phenotypic function. We can rewrite inequality (3) as

$$u < 1 - \frac{c_0}{s}. \quad (4)$$

This condition specifies a (rough estimate for the) phenotypic error threshold. If the phenotypic mutation rate, $u$, exceeds this critical value, then the gene that performs this phenotypic function cannot be maintained in the population by selection alone.

In the following sections, we make this argument more precise by elaborating on a more detailed model that includes an arbitrary number of genes. We show in particular that natural selection leads to phenotypic mutation rates that are much higher than genotypic mutation rates. Specifically, we address (and partially solve) the following questions:

When is it beneficial to transcribe and express a new set of genes that bring about a selective advantage, but production of protein is costly?

What is the optimum number of protein molecules to be produced if $k$ and the other parameters (genotypic and phenotypic mutation rates, fitness advantage, and costs) are given?

Is there an evolutionary explanation for why phenotypic mutation rates are so much higher than genotypic mutation rates?

What are the consequences of costs associated with higher fidelity of protein production?

## THE MODEL

We consider a large population of single-cell organisms (cells, for short) with DNA chromosomes. At each locus $n$ ($1 \leq n \leq L$) under consideration a number $m_n$ of protein molecules is produced. There is no recombination between loci. Errors occur both during DNA replication (*i.e.*, cell division) and during transcription of DNA to RNA and subsequent translation into protein. We call errors of the first kind genotypic mutations and those of the second kind phenotypic mutations. If $\mu_n$ and $u_n$ denote the genotypic and phenotypic mutation rates at locus $n$, respectively, then $Q = \prod_{n=1}^{L}(1 - \mu_n)$ is the probability that DNA is produced without error, and

$$p_{n,i} = \binom{m_n}{i}(1 - u_n)^i u_n^{m_n - i} \quad (5)$$

is the probability that $i$ of the $m_n$ molecules produced by gene $n$ are error free.

We assume that a certain number of mutation-free protein molecules can increase the fitness of a cell because then a beneficial phenotypic function can be performed, but production of protein has costs. The costs per protein molecule produced by gene $n$ are $c_n > 0$. They reduce the fitness of the cell. More precisely, we

assume that every gene, $n$, must produce at least $k_n$ mutation-free copies of the protein so that the fitness of the cell is increased by an amount $s > 0$. If only one of the genes produces less mutation-free protein than required, no such fitness increase occurs. To formalize these assumptions, let $\mathbf{i} = (i_1, \ldots, i_L)$ and denote by $f_\mathbf{i}$ the fitness of a cell that has $i_n$ error-free and $m_n - i_n$ erroneous (protein) molecules produced by gene $n$ ($n = 1, \ldots, L$) as well as an error-free DNA (at all loci). If we denote the total costs of protein production by $c_{\mathrm{tot}} = \sum_{n=1}^{L} c_n m_n$ and assume that the costs reduce the fitness of a cell by an additive amount, we obtain for the (Malthusian) fitnesses

$$f_\mathbf{i} = \begin{cases} f_0 + s - c_{\mathrm{tot}} & \text{if } k_n \le i_n \le m_n \quad (n = 1, \ldots L) \\ f_0 - c_{\mathrm{tot}} & \text{otherwise.} \end{cases} \quad (6)$$

Sometimes, it is convenient to write $\phi_0 = f_0 - c_{\mathrm{tot}}$. This kind of selection involves strong epistasis and is similar to what is called truncation selection in population genetics. Because mutated DNA will always produce mutated RNA, the fitness of cells with mutated DNA is $\phi_0$. Throughout, we require that $\phi_0 \ge 0$. Cells that do not express these genes, because they do not exist or are not activated, have fitness $f_0$.

Let $x_\mathbf{i}$ denote the relative frequency of cells that have error-free DNA and $i_n$ denote error-free protein molecules from locus $n$. Further, let $y$ be the frequency of cells whose DNA carries at least one mutation at one of the $L$ loci. The probability that a cell produces $i_n$ error-free molecules from each locus $n$ is $R_\mathbf{i} = \prod_{n=1}^{L} p_{n,i_n}$. We emphasize that the numbers $\mathbf{i}$ of functional molecules produced by an offspring are independent of the numbers $\mathbf{j}$ produced by its parent. Because we assume that the population size is large enough to ignore stochastic fluctuations, it follows that the dynamics of cell frequencies are given by the system of $\nu = \prod_{n=1}^{L} m_n + 1$ differential equations

$$\dot{x}_\mathbf{i} = QR_\mathbf{i} \sum_\mathbf{j} f_\mathbf{j} x_\mathbf{j} - f x_\mathbf{i}, \quad (7a)$$

$$\dot{y} = (1 - Q) \sum_\mathbf{j} f_\mathbf{j} x_\mathbf{j} + \phi_0 y - f y, \quad (7b)$$

where $f = \sum_\mathbf{j} f_\mathbf{j} x_\mathbf{j} + \phi_0 y$ is the mean fitness. We note that the system of differential equations given by (7a) and (7b) can be written as the replicator equation $\dot{z} = (A - f)z$, where $z = (x, y)$ and $x$ has the components $x_\mathbf{i}$, and the $\nu \times \nu$ matrix $A$ has entries $A_{\mathbf{ij}} = QR_\mathbf{i} f_\mathbf{j}$, $A_{\mathbf{i}\nu} = 0$, $A_{\nu\mathbf{j}} = (1 - Q) f_\mathbf{j}$, and $A_{\nu\nu} = \phi_0$.

Let $\bar{f} = \sum_\mathbf{i} f_\mathbf{i} R_\mathbf{i}$. Then, the matrix $A$ has the eigenvalues

$$\hat{f} = Q\bar{f}, \quad (8)$$

$\phi_0$, and 0, which has multiplicity $\prod_{n=1}^{L} m_n - 1$. The equilibrium solution corresponding to $\hat{f}$, the equilibrium mean fitness, is uniquely determined and locally stable if and only if $\hat{f} > \phi_0$ or, equivalently, if $\bar{f} > \phi_0/Q$. In this case, it attracts all solutions with initial value $y > 0$ because (7a) and (7b) are equivalent to the linear system $\dot{z} = Az$ (THOMPSON and MCBRIDE 1974; BÜRGER 2000). The equilibrium frequencies of cell types are readily shown to be

$$\hat{x}_\mathbf{i} = R_\mathbf{i} \frac{Q\bar{f} - \phi_0}{\bar{f} - \phi_0}, \quad \hat{y} = \frac{(1 - Q)\bar{f}}{\bar{f} - \phi_0}. \quad (9)$$

If $\bar{f} \le \phi_0$, then all solutions converge to $y = 1$. This occurs, for instance, if $s = 0$.

Let $P_n(= P_n(m_n, k_n, u_n)) = \sum_{i_n=k_n}^{m_n} p_{n,i_n}$ denote the probability that in a cell at least $k_n$ molecules produced by gene $n$ are error free. Because these are precisely the cells with fitness $\phi_0 + s$ and $\sum_{\mathbf{i} \ge \mathbf{k}} R_\mathbf{i} = \prod_{n=1}^{L} P_n$, where $\mathbf{i} \ge \mathbf{k}$ means $i_n \ge k_n$ for all $n$, we obtain

$$\bar{f} = (\phi_0 + s) \prod_{n=1}^{L} P_n + \phi_0 \left( 1 - \prod_{n=1}^{L} P_n \right) = \phi_0 + s \prod_{n=1}^{L} P_n. \quad (10)$$

Formula (10) for $\bar{f}$, hence an explicit expression for the mean fitness $\hat{f}$, could have been derived without resorting to the full dynamics (7). However, uniqueness and global stability of the corresponding solution can be inferred only from the complete evolutionary dynamics.

## RESULTS

Our first aim is to determine conditions under which a set of genes, with the properties set out above, increases the fitness of a cell relative to one in which this set of genes does not exist or is not transcribed. Such a cell is assumed to have fitness $f_0$ because no costs from protein production occur. A population of such cells sets the standard, $f_0$, to which $\hat{f}$ has to be compared.

For an analytical treatment we assume that all loci are equivalent; *i.e.*, $c_n \equiv c$, $m_n \equiv m$, $k_n \equiv k$, $u_n \equiv u$, for all $n$. Therefore, we have $c_{\mathrm{tot}} = cLm$ and $\phi_0 = f_0 - cLm$. Furthermore, if we denote by $P \equiv P_n$ the probability that at least $k$ error-free protein molecules per gene are produced, we have

$$P = \sum_{i=k}^{m} p_i, \quad (11)$$

where

$$p_i = \binom{m}{i} (1 - u)^i u^{m-i}. \quad (12)$$

It follows from (10) that

$$\bar{f} = f_0 - cLm + sP^L, \quad (13)$$

which generalizes (1), and from (8) that
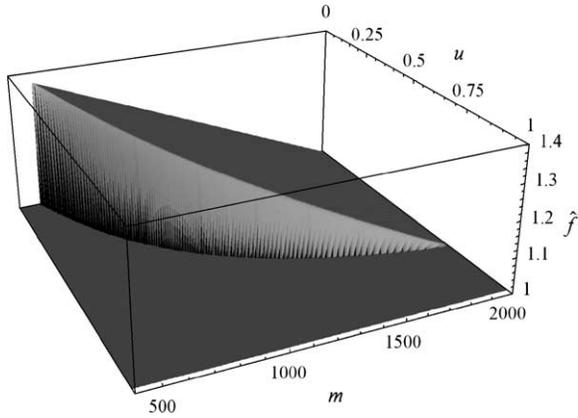
$$\hat{f} = (1 - \mu)^L [f_0 - cLm + sP^L]. \quad (14)$$

FIGURE 1.—Equilibrium mean fitness, $\hat{f}$, as a bivariate function of $m$ ($400 \leq m \leq 2000$) and $u$ ($0 \leq u \leq 1$). The other parameters have the following values: $f_0 = 1$, $\mu = 10^{-4}$, $s = 0.5$, $k = 500$, $c = 0.000025$, $L = 10$. The large $s$ is chosen to obtain a distinctive display of the features of this function.

A population of single-cell organisms that transcribe and translate the set of genes under consideration has an increased equilibrium mean fitness relative to one that does not produce this protein if and only if

$$\hat{f} > f_0. \tag{15}$$

This generalizes (2). Condition (15) imposes a number of restrictions on the parameters of our model, which we explore next. Figures 1 and 2 show how the equilibrium mean fitness $\hat{f}$ depends on various parameters.

Figure 1 displays $\hat{f}$ as a bivariate function of $m$ and $u$. More precisely, it shows $\max(\hat{f}, 1)$ to clearly display the parameter combinations that confer a fitness advantage to the population (in the figures, $f_0 = 1$). Most

distinctive is the steep "wall" surrounding the "mesa-like mountain." It signifies a threshold-like increase of $\hat{f}$ as $m$ increases above a critical value or $u$ decreases below a critical value. For each given $u$, there is a value $m$ that maximizes $\hat{f}$ (see *The optimum number of protein molecules*). The linear decrease of $\hat{f}$ in $m$ is caused by the costs, which are proportional to $cm$. If $m$ is too large [in this case $m \geq 1996$, *cf.* (16)], then $\hat{f} < f_0 = 1$. Moreover, there is a maximum value of $u$, above which $\hat{f} < 1$ for every $m$. Its value is $\sim 0.73$; see (20).

Figure 2 displays $\hat{f}$ as a function of $\log_{10} u$ for several different parameter combinations. The threshold-like dependence on $u$ is distinctive, as is the fact that $\hat{f}$ is effectively constant for values of $u$ below the threshold. This is investigated and explained in *Selection on mutation rates*.

**Error thresholds and other necessary conditions for performing an advantageous function:** From inequality (15), simple conditions for some of the parameters can be derived that must be satisfied so that incorporating or maintaining a set of genes may confer a fitness advantage to the population. Because we always have $P \leq 1$, the following simple upper bound for $m$ is easily deduced from (15) and (14) by setting $P = 1$:

$$m \leq \frac{s - f_0[(1 - \mu)^{-L} - 1]}{cL}. \tag{16}$$

For given $k$, the inequality

$$s \geq f_0[(1 - \mu)^{-L} - 1] + cLk \approx f_0(e^{\mu L} - 1) + cLk \tag{17}$$

must be satisfied independently of $u$ because $m \geq k$. The approximation is valid if $L\mu^2 \ll 1$. For smaller values of $s$, there is no parameter combination that provides a selective advantage to a single-cell population with these
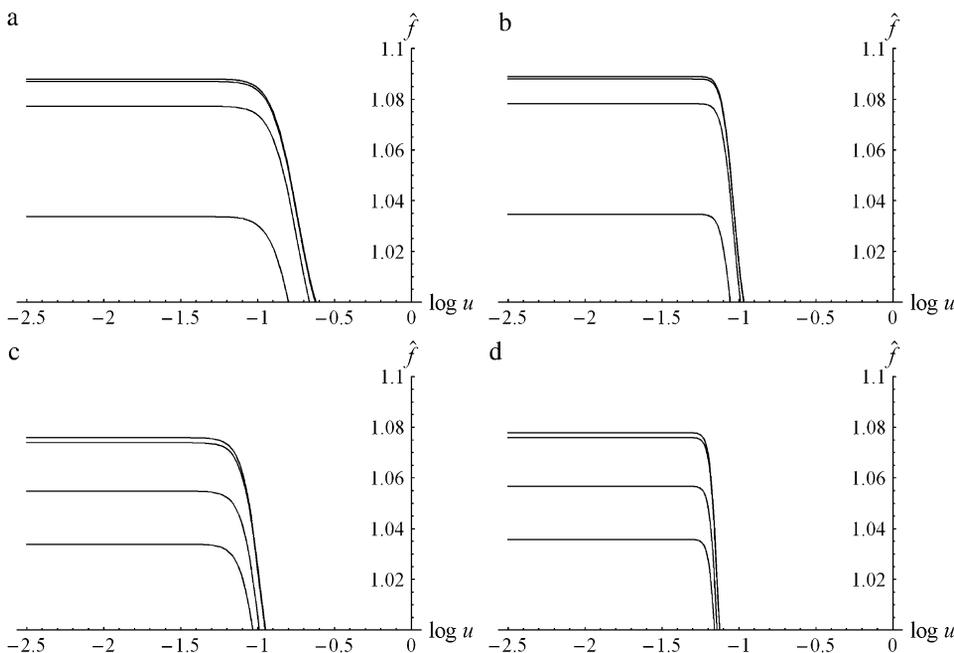


FIGURE 2.—Equilibrium mean fitness, $\hat{f}$, as a function of the (decadic) logarithm of the phenotypic mutation rate, $\log_{10} u$. (a) $L = 1$, $k = 50$, $m = 60$, $c = 2 \times 10^{-4}$; (b) $L = 1$, $k = 500$, $m = 550$, $c = 2 \times 10^{-5}$; (c) $L = 20$, $k = 50$, $m = 60$, $c = 2 \times 10^{-4}$; and (d) $L = 20$, $k = 500$, $m = 550$, $c = 2 \times 10^{-5}$. (a and b) The curves (from top to bottom) are for $\mu = 10^{-4}$, $\mu = 10^{-3}$, $\mu = 10^{-2}$, and $\mu = 5 \times 10^{-2}$. (c and d) The curves are for $\mu = 10^{-5}$, $\mu = 10^{-4}$, $\mu = 10^{-3}$, and $\mu = 2 \times 10^{-3}$. (a–d) $f_0 = 1$ and $s = 0.1$.

genes activated. A simpler, but less accurate, estimate than (17) for the necessary selective advantage is

$$s \geq L(f_0\mu + ck). \tag{18}$$

This complements (3), which has been derived under different assumptions. Condition (17) can be easily reformulated to obtain an upper bound on the tolerable costs of producing a protein molecule.

By simple rearrangement of (14), we obtain from (15) the following upper bound on the genotypic mutation rate per gene:

$$\mu \leq \mu_{ET} = 1 - \left[\frac{f_0}{f_0 - cLm + sP^L}\right]^{1/L}. \tag{19}$$

This may be called a genotypic error threshold. We note that $\mu_{ET} \leq 1 - (1 + s/f_0)^{-1/L}$ and the right-hand side is attained if $u = 0$ (hence $P = 1$) and $c = 0$.

No such simple and precise formula exists for the phenotypic error threshold because $P$ is a complicated function of $u$ and $m$ (but see *Selection on mutation rates*). However, an upper bound on the phenotypic mutation rate $u$ per gene can be derived, above which the mean fitness of a cell population is less than $f_0$ for every $m$. It is given by

$$u_{max} \approx 1 - \frac{cL}{s - f_0\mu}(k + 2 + 2\sqrt{k+1}) \tag{20}$$

(see APPENDIX A) and is more precise than the simple estimate (4) for the phenotypic error threshold derived in the Introduction. Obviously, a set of genes that confers a higher selective advantage can be incorporated and maintained under much higher phenotypic mutation rates and costs than a set of genes conferring a lower advantage. The approximation (20) is very accurate for a single gene ($L = 1$) and is a slight overestimate, on the order of a few percent, if $L > 1$ (results not shown). We remark that our usage of the term "phenotypic error threshold" deviates from that of SCHUSTER and FONTANA (1999).

**The optimum number of protein molecules:** Figure 1 shows that for a given phenotypic mutation rate $u$, there is an optimum number, $m_{opt}$, of molecules to be produced. Indeed, this is intuitive because in the absence of phenotypic mutation, $m = k$ molecules should be produced to minimize the costs. With phenotypic mutation, $k$ or more molecules have to produced to obtain a correctly expressed gene. We can restrict attention to phenotypic mutation rates $u < u_{max}$. Since we have $\hat{f} = Q\bar{f}$ by (8), and because $Q$ is independent of $m$, it is sufficient to find the $m$ that maximizes $\bar{f}$.

Let us first assume $uk \ll 1$. Then the binomial distribution (12) can be approximated by a Poisson distribution with mean $m(1 - u)$, and we obtain from (13)

$$\bar{f} = \begin{cases} f_0 - cLk + s(1 - Lku) + O(k^2u^2), & \text{if } m = k, \\ f_0 - cL(k + 1) + s + O(k^2u^2), & \text{if } m = k + 1. \end{cases} \tag{21}$$
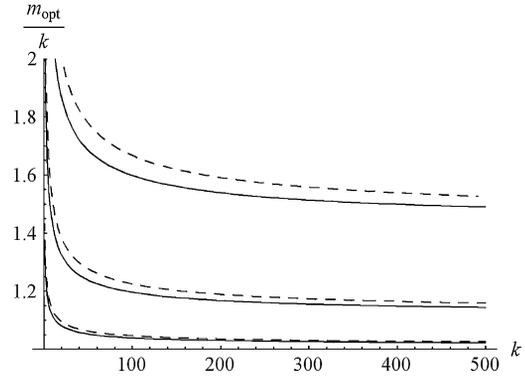


FIGURE 3.—Ratio of optimum number to minimum required number of protein molecules, $m_{opt}/k$ (22), as a function of $k$. Dashed curves, $c/s = 0.0005$; solid curves, $c/s = 0.005$. The three pairs of lines are (from top to bottom) for the following phenotypic mutation rates: $u = 0.3, 0.1$, and $0.01$.

By comparing these two cases, it follows immediately that $\bar{f}$ assumes its maximum at $m = k$ if $uk < c/s$. [Note that (17) implies that $c/s \ll 1$ if $k \gg 1$, so the assumption $uk \ll 1$ is automatically satisfied if a set of genes can be added at all and $k \gg 1$.] Thus, for very small phenotypic mutation rates the fitness is maximized at $m = k$.

If $u$ is sufficiently large, the binomial distribution (12) can be approximated by a normal distribution. [This is accurate for all possible $m$ if $ku(1 - u) \geq 5$.] Then the very accurate approximation

$$m_{opt} = k + \frac{uk + \sqrt{uk}\sqrt{-2\ln\alpha}}{1 - u} \tag{22}$$

for the optimum $m$ is found, where $\alpha = 2(c/s)/\sqrt{\pi uk}$ (APPENDIX B). In particular, the true $m_{opt}$ is always $> k/(1 - u)$. Figure 3 displays $m_{opt}/k$ as a function of $k$ for various parameter combinations. Figure 3 and (22) demonstrate that the optimum number of protein molecules to be produced is only slightly larger than $k$, unless $u$ is very high or $k$ very small. In fact, we have $\lim_{k\to\infty} m_{opt}/k = 1/(1 - u)$; *cf.* the Introduction, before Equation 3. The convergence, however, is slow so that $k$ must be on the order of a few hundred that $1/(1 - u)$ becomes an accurate approximation for $m_{opt}$. It is also important to note that $m_{opt}$ is independent of $L$ and depends only very weakly on $c$ and $s$; larger $c/s$ slightly decreases $m_{opt}$. The mean fitness $\hat{f}$ at $m_{opt}$, however, depends strongly on $L$, $s$, and $c$. Even though the derivation of (22) assumes $ku(1 - u) \geq 5$, (22) remains accurate if $uk < 1$ and correctly predicts that $m_{opt} \to k$ as $u \to 0$. Additional numerical results (not presented) show that the relative error of the approximation (22) rarely exceeds 5% and often is much lower. Finally, we point out that the minimum possible $m$ is only slightly smaller than the optimum $m$ given by (22); this is best seen from Figure 1.

**Selection on mutation rates:** Here, we explore how the equilibrium mean fitness $\hat{f}$ depends on the genotypic

and phenotypic mutation rates. The dependence of $\hat{f}$ on the genotypic mutation rate is very simple because it is proportional to $(1 - \mu)^L$. Therefore, the larger the number of genes involved, the more advantageous is a low genotypic mutation rate. Even for a single gene, there is significant selection for reducing the genotypic mutation rate well below $10^{-2}$ in any cell population of size $10^3$ or higher because selection dominates random genetic drift if population size times selective advantage exceeds ∼10. With $L = 20$ genes, we have $(1 - \mu)^{20} \leq 0.99$ if $\mu \geq 5 \times 10^{-4}$. Thus, in cell populations of size as small as $10^3$ there is already significant selection pressure for reducing the mutation rate below $5 \times 10^{-4}$.

In contrast, for the phenotypic mutation rate, there is hardly any selection pressure to reduce it to such low levels (even in extremely large populations). Indeed, Figure 2 shows that the equilibrium mean fitness becomes nearly independent of $u$ as $u$ gets smaller than ∼$10^{-1}$. There are two reasons for this. First, a reduction in genotypic mutation rate affects fitness in a structurally different way than a reduction in phenotypic mutation rate. In general, the fitness increase caused by a reduction in $\mu$ is only weakly dependent on $s$ because it is proportional to the (typically) much larger term $f_0$. However, fitness changes induced by $u$ are always proportional to $s$ because they enter $\hat{f}$ through changes in $P$, the probability that at least $k$ error-free proteins are produced; see Equation 14. The second reason is that $P^L$ has a sigmoid, often nearly threshold-like, shape. The cumulative binomial density $P$ is extremely close to its maximum value 1 if the mean number of correctly produced molecules, $m(1 - u)$, exceeds $k$ by only a few standard deviations. Hence, if $m$ is sufficiently large so that the binomial distribution can be approximated by a Gaussian, then $\hat{f}$ approaches its maximum value as $u \to 0$ approximately as fast as $e^{-x^2}$ approaches zero as $x \to \infty$.

A simple explicit, but approximate, expression for the minimal mutation rate below which $\hat{f}$ is nearly independent of $u$ can be obtained by approximating the binomial cumulative probability $P$ by a Gaussian. To this aim, let $q$ be a small positive number and let $d$ denote the $(1 - q)^{1/L}$ quantile of the standard normal distribution. Then, we have $1 - q \leq P^L \leq 1$ if $m(1 - u) \geq k + d\sqrt{mu(1 - u)}$. Solving for $u$ yields

$$u_{\min} = \frac{d^2 + 2(m - k) - d\sqrt{d^2 + 4k(1 - k/m)}}{2(d^2 + m)}. \quad (23)$$

We call this the minimum phenotypic mutation rate, because selection will be unable to reduce $u$ below $u_{\min}$, unless the population size is much larger than $(sq)^{-1}$. Notably, $u_{\min}$ is independent of $s$ and $c$, except indirectly if the choice of $q$ is made conditional on $s$. [The reader should not be worried by the fact that the derivation of (23) does not assume $\hat{f} > f_0$. If $\hat{f} < f_0$ for all $u$, for instance, because $s$ is too small or $c$ too large, then there
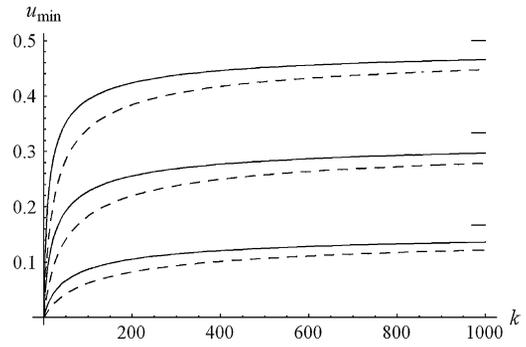


FIGURE 4.—Minimum phenotypic mutation rate, $u_{\min}$, as a function of $k$ for a single locus ($L = 1$). Solid lines, $q = 10^{-3}$; dashed lines, $q = 10^{-6}$. The three pairs of lines differ in $m$; i.e., from top to bottom, $m = 2k$, $m = 1.5k$, and $m = 1.2k$. The three short lines on the right-hand side indicate $\lim_{k \to \infty} u_{\min}$.

will be no selection pressure at all to reduce $u$ because activating the set of genes automatically leads to a fitness disadvantage.]

Figure 4 displays $u_{\min}$ for a single gene as a function of $k$ for several parameter combinations. It shows that $u_{\min}$ is very small only if either $k$ is very small or $m$ is only slightly larger than $k$. The latter *a priori* requires small phenotypic mutation rates. It is also of interest to note that if we assume that $m$ is a fixed multiple of $k$, i.e., $m = ak$, and let $k$ tend to infinity, then $\lim_{k \to \infty} u_{\min} = 1 - 1/a$ and $u_{\min} \leq 1 - 1/a$. Thus, if many more proteins are produced than required ($a$ large), then $u_{\min}$ will be close to 1 if $k$ is large. If, on the other hand, $m$ is not much larger than $k$ [as suggested by expression (22) for the optimal $m$], then $u_{\min}$ will be relatively small. For example, if $a = 1.2$ as in some of the graphs in Figures 4 and 5, then $1 - 1/a = \frac{1}{6}$. These considerations strongly suggest that, on the basis of our model, smaller phenotypic mutation rates are not likely to evolve. It is also notable, although obvious from the derivation, how weakly $u_{\min}$ depends on $q$ and how much larger than $sq$ it is under most conditions. The latter is important because for a single gene the corresponding $\mu_{\min}$ would be $sq$. Hence, $\mu_{\min} \ll u_{\min}$. For $L$ genes, $\mu_{\min}$ would be correspondingly lower, i.e., $\mu_{\min} = 1 - (1 - sq)^{1/L}$.

The above argument, that the selective pressure to reduce the phenotypic mutation rate below $u_{\min}$ is less than $sq$, depends on the assumption that $m$ is given and constant. It does not, however, involve any costs for reducing the mutation rate. Such costs are investigated below. Theoretically, the phenotypic mutation rate could evolve to zero, or at least to much lower levels than given by $u_{\min}$, if $m$ and $u$ could be optimized simultaneously. This can be seen from Figure 1 and would correspond to evolution along the top of the (curved) ridge. Substantial bivariate optimization, and evolution to very low values of $u$, does not appear to be a very likely scenario because it would require extreme fine tuning of $m$ and $u$. If, for given $u$, $m$ is only slightly larger than $m_{\text{opt}} = m_{\text{opt}}(u)$ (∼2% is sufficient), then the selective advantage to reduce $u$ is
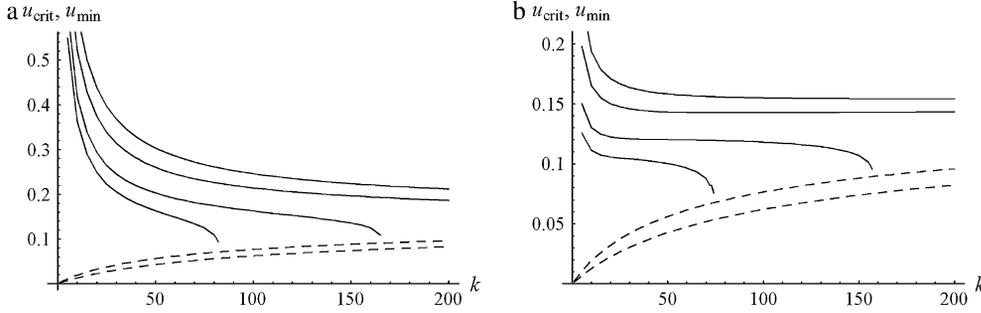
FIGURE 5.—Critical phenotypic mutation rate, $u_{crit}$ (solid lines), and minimum phenotypic mutation rate, $u_{min}$ (dashed lines), as a function of $k$. (a) $L = 1$; (b) $L = 10$. The parameters $\mu = 10^{-4}$, $L = 1$, $c = 10^{-4}$, and $m = 1.2k$ are the same for all curves. For $u_{crit}$ (from top to bottom), $s = 0.5$, $s = 0.1$, $s = 0.02$, and $s = 0.01$. $u_{min}$ is independent of $s$, $\mu$, and $c$. Top dashed line, $q = 10^{-3}$; bottom dashed line, $q = 10^{-5}$. Note that for the corresponding parameter combinations, $u_{crit}$ becomes negative where the lines end.

already vanishingly small (the population is on the gentle slope, which is completely flat in the direction of increasing $u$). If, in contrast, only a few protein molecules less than $m_{opt}$ are produced, then the fitness decrease is substantial (the population "drops off the steep wall"). It seems questionable if mechanisms for the required simultaneous fine tuning of both $m$ and $u$ exist, in particular, because $m_{opt}$ is a population property, not a property of the cell. It would require that a cell knows exactly quite how many error-free molecules it produces.

**The critical mutation rate:** We have already derived an approximation for the phenotypic error threshold, *i.e.*, the maximum mutation rate $u_{max}$ above which the set of genes cannot be maintained. Here we take a closer look at the distinctive threshold-like dependence of $\hat{f}$ on $u$ (Figure 2) and investigate how it depends on $m$ and the other parameters. This threshold-like dependence is a characteristic feature of our model and has a simple explanation. Let us approximate $P$ by the Gaussian cumulative distribution function with mean $m(1 - u)$ and variance $mu(1 - u)$. Then, $P$ switches from a value close to 0 to a value close to 1 near the mean $m(1 - u)$. This transition occurs within about two standard deviations of the mean. For a single gene, this implies that the transition occurs if $m(1 - u) \approx k$, whence $u \approx 1 - k/m$ follows. For the parameter values of Figure 2, a and b, this yields $u \approx 0.17$, or $\log_{10}u \approx -0.78$, a reasonably good approximation to the critical value $u_{crit}$ defined as the solution of $\hat{f} = f_0$. If there is more than one gene, then the transition occurs near $P \approx \left(\frac{1}{2}\right)^{1/L}$ and becomes sharper as $L$ increases. Approximating $P$ by the corresponding Gaussian cumulative distribution, *i.e.*, $P \approx F_G\left((m(1 - u) - k)/\sqrt{2mu(1 - u)}\right)$, where $F_G(x) = \frac{1}{2}(\text{erf}(x) + 1)$, we obtain the critical value $u_{crit}$ by solving

$$\text{erf}\left(\frac{m(1 - u) - k}{\sqrt{2mu(1 - u)}}\right) = 2\left(\frac{1}{2}\right)^{1/L} - 1. \qquad (24)$$

For the parameters of Figure 2c, this yields $\log_{10}u = -0.94$; for Figure 2d, it yields $\log_{10}u = -1.01$. Appar-

ently, both values are a reasonably good approximations for the true $u_{crit}$.

Figure 5 displays the critical phenotypic mutation rate $u_{crit}$ (solid lines), calculated numerically by solving $\hat{f} = f_0$, for four selective coefficients as a function of $k$, and compares it with the minimum phenotypic mutation rate $u_{min}$ (dashed lines), calculated for two choices of $q$. For the two smaller values of $s$, the curves for $u_{crit}$ end when $\hat{f} < f_0 = 1$, *i.e.*, when expression of the set of genes causes a fitness reduction for all $u$.

**The role of costs for reducing the phenotypic mutation rate:** So far, all arguments have assumed that no costs are associated with lower mutation rates. Here, we briefly explore the consequences of such costs. To illustrate the (quite obvious) effects, let us assume that the cost of producing a single protein molecule is $c(1 + \gamma/u)$, where $\gamma \geq 0$.

Figure 6 displays the mean equilibrium fitness $\hat{f}$ as function of $\log_{10}u$ and $m$ with $\gamma = 0.01$. Thus, the costs
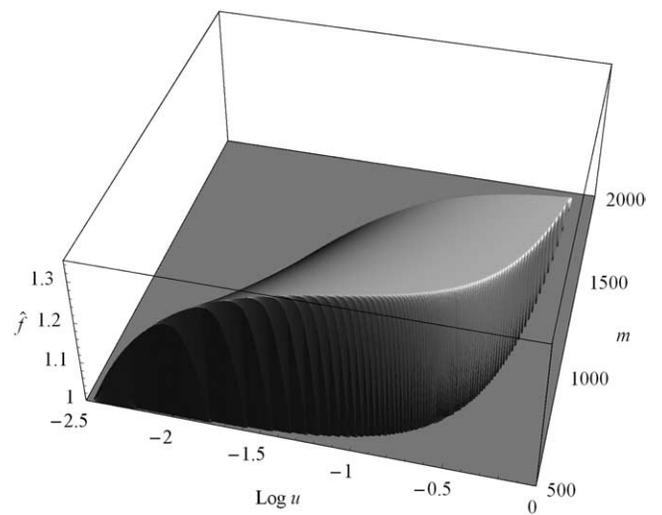


FIGURE 6.—Equilibrium mean fitness, $\hat{f}$, as a bivariate function of $m$ ($500 \leq m \leq 2000$) and $\log_{10}u$ ($-2.5 \leq \log_{10}u \leq 0$) for costly small phenotypic mutation rates. The costs are $c = 0.000025(1 - 0.01/u)$ per molecule. The other parameters are as in Figure 1.

are slowly increasing with decreasing $u$. With $\gamma = 0$ there are no costs for reducing the phenotypic mutation rate and Figure 1 would be obtained (except for the different scaling of the $u$-axis). The fitness optimum is near $(m, u) = (571, 0.0805)$ with $\hat{f} = 1.338$ ($\log_{10}0.0805 = -1.094$). We note that 0.0805 is close to $u_{\min}(m = 571, k = 500, q = 10^{-4}) = 0.0817$.

If the costs increase is even slower, for instance logarithmic, then the optimal $u$ is somewhat smaller (results not shown). In the presence of costs, most of the quantities derived above are much more difficult to compute, and we have not developed an analytical theory that takes into account costs for fidelity.

## DISCUSSION

In this work we argue that there is very little selective pressure to reduce the phenotypic mutation rate $u$ below a minimum mutation rate $u_{\min}$. Usually, transcriptional and translational error rates are measured in number of amino acid substitutions per synthesized amino acid and were estimated to be $4.5 \times 10^{-4}$ (Ellis and Gallant 1982). For this theory, however, we want to know the number of nonfunctional proteins per synthesized protein. This is our unit for the phenotypic mutation rate $u$. For a proper unit conversion we need to know how many amino acid substitutions per protein occur and how many of these substitutions lead to nonfunctional proteins.

The average protein lengths for *E. coli*, *Saccharomyces cerevisiae*, and *Homo sapiens* are 317, 496, and 499 amino acids, respectively. Let us be conservative and use 500 for the average protein length. For a 500-amino-acid long protein and a phenotypic mutation rate of $4.5 \times 10^{-4}$ mistranslations per amino acid, we will have ~0.23 incorrectly synthesized proteins per synthesized protein. However, only a fraction of these 23% will carry amino acid substitutions that render them nonfunctional.

Exhaustive amino acid substitution assays on HIV-1 protease (Loeb *et al.* 1989), T4 lysozyme (Rennell *et al.* 1991), and Lac repressors (Markiewicz *et al.* 1994) showed that 59, 12, and 34% of the examined amino acid substitutions were deleterious (summarized in Saunders and Baker's 2002 Table 1). If we choose 35%, the average of these three values, as the fraction of amino acid substitutions that are deleterious, we have a phenotypic mutation rate of 0.08 deleterious mutations per synthesized protein.

Because a substantial fraction, if not the vast majority, of genotypic mutations are detrimental (Keightley and Eyre-Walker 1999; Keightley and Lynch 2003), deleterious per-locus mutation rates can be expected to be between $10^{-4}$ and $10^{-7}$ (see Introduction), where the upper bound is likely to be an overestimate. In any case, deleterious genotypic mutation rates are several orders of magnitude smaller than phenotypic ones.

In addition to the (deleterious) phenotypic mutation rate we also have to consider the number of protein molecules produced per cell. Early studies showed that only a few hundred proteins account for most of the protein content of a cell and that most of the proteins are present in low copy numbers (O'Farrell 1975). Low copy numbers range from a few proteins per cell to several hundred. For instance, the Lac repressor, a regulatory protein, is thought to be "occurring in about ten copies per gene" (Gilbert and Muller-Hill 1966). *E. coli* DNA photolyase, a DNA repair enzyme, has a copy number of ~10–20 molecules per protein (Harm *et al.* 1968). High copy number proteins can have abundances of many thousand molecules per cell (Gygi *et al.* 1999). A recent study shows that the costs associated with the production of protein may be substantial, and that they increase faster than linear with the amount of protein produced (Dekel and Alon 2005).

Our main results are, first, that there is basically no selective pressure to reduce the phenotypic mutation rate per gene below a minimum value, $u_{\min}$, which is rarely <0.05, and often near 0.1. This compares surprisingly well with the 8% deleterious mutations per synthesized protein calculated above.

In contrast, and despite the simplicity of the model, there is selective pressure to reduce the genotypic mutation rate to much lower levels, one order of magnitude at least. If several genes have to be expressed to increase fitness, the difference becomes larger. Second, for given parameters, there is a critical phenotypic mutation rate, $u_{\text{crit}}$, above which the fitness of the population is actually reduced if the set of genes is expressed. Unless the potential fitness increase, $s$, is very high relative to the costs, $c$, and $k$ very small, $u_{\min}$ is not much smaller than $u_{\text{crit}}$, in particular, if several genes are involved. Both $u_{\min}$ and $u_{\text{crit}}$ depend only very weakly on $c$ and $s$. No simple formulas for $u_{\text{crit}}$ and $u_{\min}$ are available. Their (approximate) calculation involves computation of quantiles of the normal distribution. However, and this is the third result, there is a simple formula for the maximum phenotypic mutation rate, $u_{\max}$, above which there is a fitness disadvantage for expressing the genes under consideration for any number $m$ of actually produced protein molecules. This can be interpreted as a phenotypic error threshold. Unless $kc/s$ is very small or $m$ is only slightly larger than $k$, $u_{\max}$ differs from $u_{\min}$ by less than a factor of 10. Fourth, we show that for all other parameters given there exists an optimum number of protein molecules to be produced, $m_{\text{opt}}$, in the sense that the mean fitness of the population is maximized. We derive a simple and very accurate approximation for $m_{\text{opt}}$. Unless the phenotypic mutation rate is very high or $k$ is small, $m_{\text{opt}}$ is not much larger than $k$ and nearly independent of the selective advantage and the costs. It is independent of the number of loci and of the genotypic mutation rate.

The formal reason for the absence of a selective pressure to reduce the phenotypic mutation rate to such

low levels as that of the genotypic mutation rate is that the two types of mutation rates enter mean fitness, $\hat{f}$ in (14), in qualitatively different ways; this is discussed in *Selection on mutation rates*. A more intuitive reason is that as soon as only a few more than the optimum number, $m_{opt}$, of protein molecules are produced, the selective pressure to reduce the phenotypic mutation rate vanishes because the function can be fulfilled anyway. In such a situation, there is, however, weak selective pressure to reduce the number of actually produced molecules $m$ to $m_{opt}$. In principle, simultaneous evolution of $m$ and $u$ could lead to much lower phenotypic mutation rates. However, as argued in *Selection on mutation rates* this would require extreme fine tuning of these processes (in particular, $m$ has to be adjusted extremely closely to $m_{opt}$) and, thus, seems unlikely. This, together with the role of genetic drift, will be the topic of future investigation.

The above results do not involve any costs for reducing the phenotypic or genotypic mutation rate. If there are costs for reducing the phenotypic mutation rate, the parameter range in which a fitness advantage can be realized by incorporating the set of genes is substantially reduced (or even annihilated if the costs are too high), and fitness is maximized at an intermediate phenotypic mutation rate. Unless the costs are high, this maximum is close to $u_{min}$, as given by (23).

Our model, hence the conclusions, rests on a number of assumptions. We assumed that, if there is more than one locus, all loci are completely equivalent. In reality, this will not be the case because loci can differ in any of the parameters. It appears to be of most interest to study cases in which the number of required error-free protein molecules, $k$, and the actually produced number, $m$, vary among loci. We have not yet studied such a scenario.

Our most critical assumption concerns the dependence of fitness on the number of protein molecules produced. Many fitness functions other than our step-like function (6) are conceivable. For instance, fitness could increase smoothly as the number of error-free proteins increases. We have not studied such a scenario. However, it appears quite reasonable to assume that the performance of a, at least moderately complex, function requires many genes to interact in an appropriate manner. There may be many possibilities of modeling such gene interaction, but none has been studied in the present context.

The following example shows that there are fitness functions that can induce strong selection toward low phenotypic mutation rates. Assume that $k$ error-free proteins are needed to increase fitness by $s$, but that cells that produce one or more erroneous molecules do not have this fitness advantage. Also assume, as in our model, costs $c$ for producing a protein molecule. Then using the previous notation, we have

$$\bar{f} = (f_0 + s)(1 - u)^k + f_0[1 - (1 - u)^k] - ck, \qquad (25)$$

and a (mean) fitness advantage results if $\bar{f} > f_0$. It is trivial to show that this yields the condition

$$u < 1 - \left(\frac{ck}{s}\right)^{1/k} \approx \frac{1}{k}\ln\frac{s}{ck}, \qquad (26)$$

where the approximation requires sufficiently large $k$. Obviously, this is very different from $u_{max}$; *cf.* (20) and (4). If, with this type of fitness function, there are costs associated with the reduction of the phenotypic mutation rate, the incorporation or maintenance of a set of genes that confers a fitness advantage becomes very difficult because the admissible parameter range shrinks dramatically. We do not argue that such a fitness function is realistic in any sense; by contrast, fitness functions like this would make any functional improvement difficult or impossible. Of course, there are many other reasonable fitness functions that could and should be studied; for instance, fitness could be reduced steadily if too many erroneous proteins are produced. Such studies have to be postponed to the future.

## LITERATURE CITED

Bürger, R., 2000   *The Mathematical Theory of Selection, Recombination, and Mutation.* Wiley, Chichester, UK.

Dekel, E., and U. Alon, 2005   Optimality and evolutionary tuning of the expression level in protein. Nature **436:** 588–592.

Drake, J. W., B. Charlesworth, D. Charlesworth and J. F. Crow, 1998   Rates of spontaneous mutation. Genetics **148:** 1667–1686.

Edelmann, P., and J. Gallant, 1977   Mistranslation in *E. coli.* Cell **10:** 131–137.

Eigen, M., and P. Schuster, 1977   The hypercycle: a principle of natural self-organization. A. emergence of the hypercycle. Naturwissenschaften **64:** 541–565.

Ellis, N., and J. Gallant, 1982   An estimate of the global error frequency in translation. Mol. Gen. Genet. **188:** 169–172.

Gilbert, W., and B. Muller-Hill, 1966   Isolation of the Lac repressor. Proc. Natl. Acad. Sci. USA **56:** 1891–1898.

Gygi, S. P., Y. Rochon, B. R. Franza and R. Aebersold, 1999   Correlation between protein and mRNA abundance in yeast. Mol. Cell. Biol. **19:** 1720–1730.

Harm, W., H. Harm and C. S. Rupert, 1968   Analysis of photoenzymatic repair of UV lesions in DNA by single light flashes. II. In vivo studies with *Escherichia coli* cells and bacteriophage. Mutat. Res. **6:** 371–385.

Ibba, M., and D. Söll, 1999   Quality control mechanisms during translation. Science **286:** 1893–1897.

Keightley, P. D., and A. Eyre-Walker, 1999   Terumi Mukai and the riddle of deleterious mutation rates. Genetics **153:** 515–523.

Keightley, P. D., and M. Lynch, 2003   Toward a realistic model of mutations affecting fitness. Evolution **57:** 683–685.

Loeb, D. D., R. Swanstrom, L. Everitt, M. Manchester, S. E. Stamper *et al.*, 1989   Complete mutagenesis of the HIV-1 protease. Nature **340:** 397–400.

Markiewicz, P., L. G. Kleina, C. Cruz, S. Ehret and J. H. Miller, 1994   Genetic studies of the lac repressor. XIV. Analysis of 4000 altered Escherichia coli lac repressors reveals essential and nonessential residues, as well as "spacers" which do not require a specific sequence. J. Mol. Biol. **240:** 421–433.

O'Farrell, P. H., 1975   High resolution two-dimensional electrophoresis of proteins. J. Biol. Chem. **250:** 4007–4021.

Rennell, D., S. E. Bouvier, L. W. Hardy and A. R. Poteete, 1991   Systematic mutation of bacteriophage T4 lysozyme. J. Mol. Biol. **222:** 67–88.

SAUNDERS, C. T., and D. BAKER, 2002   Evaluation of structural and
    evolutionary contributions to deleterious mutation prediction.
    J. Mol. Biol. **322:** 891–901.
SCHUSTER, P., and W. FONTANA, 1999   Chance and necessity in evo-
    lution: lessons from RNA. Physica D **133:** 427–452.
SHAW, R. J., N. D. BONAWITZ and D. REINES, 2002   Use of an in vivo
    reporter assay to test for transcriptional and translational fidelity
    in yeast. J. Biol. Chem. **277:** 24420–24426.
SNIEGOWSKI, P. D., P. J. GERRISH, T. JOHNSON and A. SHAVER,
    2000   The evolution of mutation rates: separating causes from
    consequences. BioEssays **22:** 1057–1066.
SPRINGGATE, C. F., and L. A. LOEB, 1975   On the fidelity of transcrip-
    tion by Escherichia coli ribonucleic acid polymerase. J. Mol. Biol.
    **97:** 577–591.
STURTEVANT, A. H., 1937   Essays on evolution. I. On the effects of
    selection on mutation rate. Q. Rev. Biol. **12:** 467–477.
THOMAS, M. J., A. A. PLATAS and D. K. HAWLEY, 1998   Tran-
    scriptional fidelity and proofreading by RNA polymerase II. Cell
    **93:** 627–637.
THOMPSON, C. J., and J. L. MCBRIDE, 1974   On Eigen's theory of the
    self-organization of matter and the evolution of biological macro-
    molecules. Math. Biosci. **21:** 127–142.
WITHEY, J. H., and D. I. FRIEDMAN, 2002   The biological roles of
    *trans*-translation. Curr. Opin. Microbiol. **5:** 154–159.

## APPENDIX A

Here, we derive an upper bound on the phenotypic mutation rate above which (15) cannot be satisfied. Let us assume $L = 1$. Then a simple calculation reveals that (15) is equivalent to

$$P > \frac{f_0\mu + (1 - \mu)cm}{s(1 - \mu)}. \tag{A1}$$

No general explicit approximation for $P$ is available. However, $P$ can be approximated by the cumulative density function of the normal distribution with mean $m(1 - u)$ and variance $mu(1 - u)$. Therefore, $P = P(m, k, u)$ is close to 1 ($\geq 0.97$) if $m(1 - u) \geq k + 2\sqrt{mu(1 - u)}$ and starts to decline rapidly as $m$ becomes smaller. The inequality $m(1 - u) \geq k + 2\sqrt{mu(1 - u)}$ is satisfied if and only if

$$m \geq m_* = \frac{k + 2u + 2\sqrt{u(k + u)}}{1 - u}. \tag{A2}$$

If we approximate the left-hand side of (A1) by $P = 1$ and the right-hand side by $(cm_* + f_0\mu)/s$ [which is

accurate to order $O(\mu)$], we obtain the desired (approximate) upper bound by solving

$$1 = \frac{cm_* + f_0\mu}{s} \tag{A3}$$

for $u$. This yields

$$u_{\max} \approx \frac{s - f_0\mu - c(k - 2) - 2c\sqrt{k + 1 - ck^2/(s - f_0\mu)}}{s - f_0\mu + 4c}. \tag{A4}$$

Numerical evaluation of the true upper bound shows that this provides an excellent approximation if $k \geq 10$. By ignoring terms of order $c^2$ and higher, we obtain (20). In general, (20) is nearly as good as (A4), but slightly smaller. A similar procedure yields (20) if $L > 1$.

## APPENDIX B

Here, we derive the approximation (22) for the optimum number $m$ of protein molecules to be produced. If $u$ is sufficiently large, *i.e.*, if $mu(1 - u) \geq 5$, then the binomial distribution (12) can be accurately approximated by a normal distribution. By partial differentiation of $\bar{f}$ (13) with respect to $m$ we obtain that the fitness is maximized at the largest solution $m$ of

$$\frac{k + m(1 - u)}{m\sqrt{2mu(1 - u)}} \exp(-A^2)[1 - \mathrm{erf}(A)]^{L-1} = \frac{2^L c\sqrt{\pi}}{s}, \tag{B1}$$

where

$$A = \frac{k - m(1 - u)}{\sqrt{2mu(1 - u)}}. \tag{B2}$$

If $A \leq -2$, which is satisfied if (approximately) $m \geq (k + 2\sqrt{ku})/(1 - u)$, we have $\mathrm{erf}(A) \leq -0.995$, and the terms $[1 - \mathrm{erf}(A)]^{L-1}$ and $2^{L-1}$ cancel.

Because of the rapid decline of $\exp(-A^2)$ for $m > k$, we can approximate $(k + m(1 - u))/m\sqrt{2mu(1 - u)}$ by $(ku)^{-1/2}$ and obtain an excellent approximation for the solution of (B1) by solving $\exp(-A^2) = \alpha$ for $m$, where $\alpha = 2c\sqrt{\pi uk}/s$. Ignoring terms of order $u/k$ and smaller, we arrive at (22).