

**Crossing fitness canyons by a finite population**David B. Saakian,<sup>1,2,3,\*</sup> Alexander S. Bratus,<sup>4</sup> and Chin-Kun Hu<sup>3,5,6,†</sup><sup>1</sup>*Theoretical Physics Research Group, Ton Duc Thang University, Ho Chi Minh City, Vietnam*<sup>2</sup>*Faculty of Applied Sciences, Ton Duc Thang University, Ho Chi Minh City, Vietnam*<sup>3</sup>*Institute of Physics, Academia Sinica, Nankang, Taipei 11529, Taiwan*<sup>4</sup>*Faculty of Computational Mathematics and Cybernetics, Lomonosov Moscow State University, Moscow 119992, Russia*<sup>5</sup>*National Center for Theoretical Sciences, National Tsing Hua University, Hsinchu 30013, Taiwan*<sup>6</sup>*Business School, University of Shanghai for Science and Technology, Shanghai 200093, China*

(Received 8 August 2016; revised manuscript received 2 May 2017; published 9 June 2017)

We consider the Wright-Fisher model of the finite population evolution on a fitness landscape defined in the sequence space by a path of nearly neutral mutations. We study a specific structure of the fitness landscape: One of the intermediate mutations on the mutation path results in either a large fitness value (climbing up a fitness hill) or a low fitness value (crossing a fitness canyon), the rest of the mutations besides the last one are neutral, and the last sequence has much higher fitness than any intermediate sequence. We derive analytical formulas for the first arrival time of the mutant with two point mutations. For the first arrival problem for the further mutants in the case of canyon crossing, we analytically deduce how the mean first arrival time scales with the population size and fitness difference. The location of the canyon on the path of sequences has a crucial role. If the canyon is at the beginning of the path, then it significantly prolongs the first arrival time; otherwise it just slightly changes it. Furthermore, the fitness hill at the beginning of the path strongly prolongs the arrival time period; however, the hill located near the end of the path shortens it. We optimize the first arrival time by applying a nonzero selection to the intermediate sequences. We extend our results and provide a scaling for the valley crossing time via the depth of the canyon and population size in the case of a fitness canyon at the first position. Our approach is useful for understanding some complex evolution systems, e.g., the evolution of cancer.

DOI: [10.1103/PhysRevE.95.062405](https://doi.org/10.1103/PhysRevE.95.062405)**I. INTRODUCTION**

The statistical physics approach to biological evolution has attracted much attention in recent decades. Important problems in such study include the origin of life [1,2], punctuated equilibrium [3,4], the dynamics and phase diagrams of quasispecies models with a mutator gene [5–7], and the effect of finite population size or system size [8–11].

New genes arise due to mutations and evolution works through mutation, selection, and random drift (random sampling due to the finite size of the population) [8–10]. Transitions between adaptive sets of traits may involve multiple mutations separated by neutral or intermediate states with low fitness. This process is also known as crossing a fitness valley and is assumed to be relevant to evolution [12–18]. Here we consider the case of deep valleys, the fitness canyons.

We consider the Wright-Fisher model of the finite population evolution [10,19–22]. In a simplified version of the fitness landscape, only a few  $m + 1$  replicator types are considered and the mathematical problem is to find the first arrival time of new mutants with mutation rates  $\mu$  and the population size  $N$ . We can look at a path of mutations; at the start, we have only replicators without mutations and we are looking for the time when the  $m$ th mutant appears in the finite population. The sequences with less than the maximal number of mutations are referred to as intermediate states. While talking about the fitness difference, we mean the difference of the fitness compared with the first mutant. The crossing a fitness valley

phenomenon plays an important role in the theory of evolution [9,13] and it has been applied to cancer biology [14]. It is especially interesting to calculate the mean first arrival time of double mutants as the simplest nontrivial case and this topic has attracted a great deal of attention in the literature [12,15,16,18,23,24].

The evolution of several genome types has the same level of dynamics complexity as many-body interactions in classical mechanics and few results are available for the simple case of metadynamics (e.g., the succession mutation regime, where the mutations are rare and the majority of individuals belong to a single type) [25]. The problem of finding the first arrival time of multiple mutants can be simplified for the zero selection, which removes most of the interactions in the system, and we are able to solve exactly some aspects of the original  $n$ -body problem. The first exact expression for a neutral version of the double-mutant problem was derived in [12]. In [15,16], exact results for the mean first arrival time of multiple mutants without selection were derived. In [18], a Moran model was considered, which is related to the neutral network fitness models [26–28], and some exact results have been obtained for intermediate neutral mutants (all the sequences along the sequence path have the same fitness).

The more difficult question is how to find the first arrival time for multiple mutants with a nonzero selection for the intermediate number of mutations and the valley crossing time for the whole population, even when the intermediate sequences on the path of sequences are neutral. In [24], the mean first arrival time problem for double mutants was solved exactly for any selection type of the intermediate mutant. In [17], several phases were derived with different scalings for the valley crossing time of the whole population.

\*david.saakian@tdt.edu.vn

†huck@phys.sinica.edu.tw

Taking into account that the solution of the finite population dynamics is too complicated, we will look at a specific case of the simple fitness landscape: All the intermediate mutations are neutral except for the one with either very low (a fitness canyon) or very high (a fitness hill) fitness. The absolute value of the fitness difference  $r$  is much larger than the inverse of the population size  $1/N$ . The population dynamics strongly depends on the location of the hill or canyon on the mutation path. We investigate the mean first arrival time of multiple mutants, which is easier to solve than the problem of the fitness valley crossing by the population. We can assume that the last mutant has much higher fitness than all the intermediate sequences and the mean first arrival time in such a case could be close to the fixation time. We assume that fitness canyons extend the time periods similarly in both cases (the first arrival time of the mutant or the fixation time of the mutant) and our numerics support this point of view for the fitness canyon crossing of the population.

## II. WRIGHT-FISHER MODEL

### A. Known results for the mean first arrival time of neutral mutants

We consider the Wright-Fisher model [9,19–22] for  $m + 1$  sequences. The model studied in [22] corresponds to the case  $m = 1$ . There are unidirectional mutations from the  $l$ th mutant to the  $(l + 1)$ th mutant with a probability  $b_l/N$ . The current state of the system is described by the set of integers  $i_0, \dots, i_m$  that show the numbers of mutants of different types and there is a fixed population size  $\sum_{l=0}^m i_l = N$ . We replace all the population via sampling with the parameters  $\eta_l$  to be defined below. They are chosen to have a proper infinite population limit as a discrete-time Eigen model [21], thus the sampling is related to both selection and mutation. In principle, we can do a sampling just after the selection and before the mutation, but such a model will not have a proper infinite population limit like a discrete-time Eigen model or a continuous-time Crow-Kimura model.

We define the parameters of the sampling  $\eta_l$ ,  $0 \leq l \leq m$ ,

$$\eta_l = \frac{i_l(1 + S_l)(1 - u_l) + i_{l-1}(1 + S_{l-1})u_{l-1}}{\sum_n (1 + S_l)i_n}, \quad (1)$$

where  $1 + S_l$  denotes the fitness (we take  $S_0 = 0$ ) and  $u_l$  are the mutation probabilities. The parameters  $\eta_l$  are slightly modified at the borders. For  $l = 0$  we omit the term  $u_{l-1}$ , while for  $l = m$  we omit the term  $u_l$ . Equation (1) assumes a directed mutation scheme for increasing  $l$ . We consider mainly the scaling

$$u_l = \frac{b_l}{N}. \quad (2)$$

The probability of a transition to a new set of integers (from one generation to the next)  $\hat{i}_0, \dots, \hat{i}_m$  is given by the expression

$$\frac{(\eta_0)^{\hat{i}_0} \dots (\eta_m)^{\hat{i}_m}}{\hat{i}_0! \dots \hat{i}_m!} N!. \quad (3)$$

The procedure to carry out the simulation is described in the Fig. 2 caption with the initial condition that all population are in the wild type.

Consider the mean first arrival time of the  $m$ th mutant, when we have the same mutation rate  $u$  for all the mutants and  $S_l = 0$

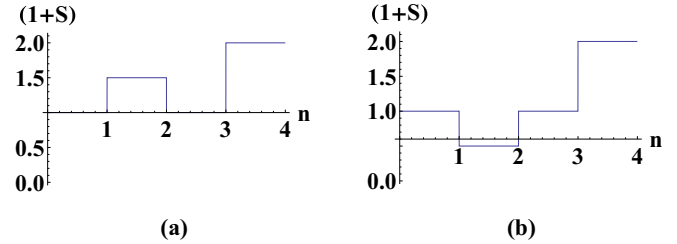


FIG. 1. Fitness along the path of mutations for three intermediate sequences with  $m = 3$ . For the  $l$ th mutant we take  $1 + S_l$  at the point  $l + 0.5$  on our graphics. The fitness of the end-point mutant is always much higher than that of the intermediate mutants. (a) Population climbing the fitness hill. (b) Population crosses the fitness canyon.

(neutral case). In [15,16], the following scaling was derived for the mean first appearance time  $T$  of the first individual with the final mutation in the case of the scaling in Eq. (2),

$$T \sim N^{n(m)}, \quad n(m) = 1 - \frac{1}{2^{m-1}}. \quad (4)$$

Let us give the key idea of their derivation, also following [14]. Consider the Moran model (see [10] and the Appendix), where the mutation events are  $N$  times slower than in the Wright-Fisher model. We look at the first arrival of the  $m$ th mutant from the original population with pure zeroth-type mutants. We replace the model by the branching process with dividing and dying equal rates and some probability rate  $u_l$  for the birth of  $m$ th mutants from the  $(m-1)$ th mutant during the life period. We distinguish the  $m$ -type mutants, born immediately from the  $m - 1$  types (the founder of the families), and those that are offsprings of  $m$ -type mutants. If at some moment of time  $t$  the zeroth class gives  $M_1 \sim Ntu_0$  first type mutants (the founders of the families), every mutant gives its own offsprings. The total number of such progenies is  $\sim (M_1)^2 \sim (Nu_0t)^2$  (this is the key point of the derivation in [16]). These first type of mutants creates the second type of mutants with the rate  $u_1$ . Therefore, at time  $t$  there should be  $(M_1)^2 u_1$  second-type mutants and  $T_2 \sim \frac{1}{Nu_0 \sqrt{u_1}}$ , which gives the scaling in Eq. (4) for our case by Eq. (2).

Consider the time  $t > T_2$ . We obtain the total number of progeny of the second type of mutants as  $M_2 \sim (M_1)^2 u_1$ . The first third-type mutant arises when  $M^2 u_1^2 u_2 \sim 1$ , which gives the result in Eq. (4) for the scaling (2) and  $m = 3$ .

We should verify how the scaling in Eq. (4) is modified in the case of the selection. We will derive the scaling of  $T$  when some  $|S_l|$  are large compared with the inverse of the population size  $1/N$  as well as in the limiting case when  $|S_l| \sim 1$ . We consider the case of the fitness hill and the case of the fitness canyon in Figs. 1(a) and 1(b), respectively.

### B. First arrival time of double mutants

Let us suppose that we have three different types of mutants:  $A$ ,  $B$ , and  $C$ . Type  $A$  can mutate to  $B$  with a probability  $u$  and type  $B$  mutates to  $C$  with probability  $u$  and the first type of mutant has a fitness  $S_1$ .

In [24] we calculated the mean first arrival time of double mutants as an integral of the hypergeometric functions. In the Appendix we derive the following expression for the mean

first arrival time of double mutants:

$$T = 2\sqrt{N} \frac{\Gamma[A] {}_2F_1[A, a, A+1, 1 - \frac{1}{1-k}]}{\Gamma[A+1](1-k)^{ah}}, \quad (5)$$

where  $\Gamma$  is the Euler Gamma function,  ${}_2F_1$  is a hypergeometric function, and

$$\begin{aligned} a &= 2uN, \quad h = \sqrt{S_1^2 + 2u2N}, \\ k &= \frac{1}{2} \left( 1 + \frac{S_1}{\sqrt{S_1^2 + 2u}} \right), \\ A &= \left( 1 + \frac{S_1}{\sqrt{S_1^2 + 2u}} \right) uN. \end{aligned} \quad (6)$$

Consider a large value  $\sqrt{N}S_1 \gg 1$ . We are interested in  $T_1$ , the mean first appearance time of the last mutant with a small difference of the fitness  $S_1\sqrt{N} \sim 1$ , as well as the mean first appearance time of the last mutant  $T_2$  in the case of  $S_1 \sim 11$ . Following [29], we obtain in Eq. (A2) the asymptotic behavior of Eq. (5):

$$T_1 \sim 1/S_1, \quad T_2 \sim 1. \quad (7)$$

Consider the case of a negative  $S_1$ . Following [29], we derive the following asymptotic expressions:

$$T_1 = |S_1|/(Nu^2), \quad T_2 \sim 1/(Nu^2). \quad (8)$$

Figure 2(a) illustrates the accuracy of our asymptotic expressions. We performed the numerical calculations for the valley crossing time by half of the population. Figure 2(b) illustrates that the mean first arrival period grows proportionally to  $c$  defined in the Fig. 2 caption. All the numerical calculations were done by applying the multimodal sampling to solve Eqs. (1) and (3).

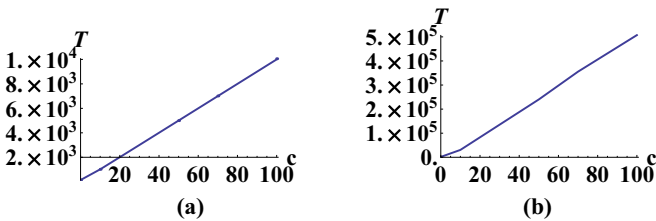


FIG. 2. We generate the random number with a binomial distribution with a probability parameter  $\eta_0$  and a number of trials  $N$  to get  $\hat{i}_0$ . Then we repeat this procedure for the next mutant type with a probability parameter  $\eta_1$  with the number of trials  $N - \hat{i}_0$  to get  $i_1$ . We continue this procedure until the mutant type  $m-1$  to get  $i_{m-1}$  and then we get the number of replicators with the last type  $\hat{i}_m = N - \sum_{l=0}^{m-1} \hat{i}_l$ . (a) The mean first arrival time  $T \equiv T_1$  for a double mutant is  $m = 2$ ,  $N = 10000$ ,  $S_0 = 0$ , and  $S_1 = -c/\sqrt{N}$ . The solid line is given from our analytical formula (5) and the dots are obtained from the numerics. (b) The mean first arrival time  $T \equiv T_1$  for half of the population is  $m = 2$ ,  $N = 10000$ ,  $S_0 = 0$ ,  $S_1 = -c/\sqrt{N}$ , and  $S_2 = 1/\sqrt{N}$ . The solid line shows a linear function that fits the numerical data and the dots correspond to the numerical data.

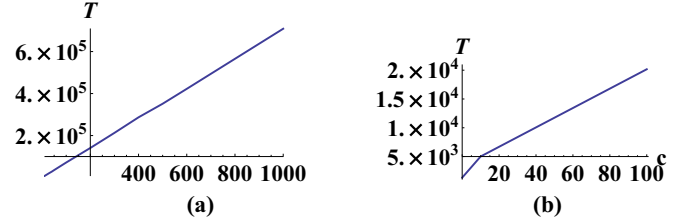


FIG. 3. (a) Mean first arrival time  $T \equiv T_1$  for  $m = 3$ ,  $N = 10000$ ,  $S_0 = 0$ ,  $S_1 = -c/N^{3/4}$ , and  $S_2 = 0$ . The solid line is the fit of numerics for the linear function and the dots correspond to the numerics. (b) Mean first arrival time  $T \equiv T_1$  for  $m = 3$ ,  $N = 10000$ ,  $S_0 = 0$ ,  $S_1 = 0$ , and  $S_2 = -c/N^{1/4}$ . The solid line is the fit of numerics for the linear function and the dots correspond to the numerics.

In the case of a fitness canyon, the mean first arrival time grows proportionally to the fitness depth. We will verify that this effect is also valid for more than two mutants.

Increasing the fitness hill shortens the mean arrival time; this effect works for the general  $m$  case as well, when the hill is near the end of the mutation path. Below we consider only the case when  $b_l = 1$ .

### C. First arrival time on the path with fitness canyon and $m = 3$

Consider the situation when the canyon is located at the first sequence as shown in Fig. 1(a). Using Eq. (4) [15,16], we have for the mean number of the first type of mutants at the time periods like the mean first arrival time period of the neutral case  $\langle i_1 \rangle \sim N^{3/4}$ . From the definition of  $\eta_1$ , we find that the selection term starts to work when  $\langle i_1 \rangle S_1 \sim 1$ , thus  $\hat{S}_1 = N^{-3/4}$  is the typical scale for  $S_1$ . For the smaller values of  $S_1$ , we have the neutral case result by Eq. (4). Then for the large  $S_1 \gg 1/N^{3/4}$  the mean first arrival period becomes  $\frac{|S_1 - S_0|}{\hat{S}_1}$  times larger, similar to the  $m = 2$  case. Eventually, we obtain

$$S_1 \gg 1/N^{3/4}, \quad T_1 \sim N^{3/2}|S_1|, \quad T_2 \sim N^{3/2}. \quad (9)$$

The numerics support these findings [see Fig. 3(a)].

Now let us put the canyon just before the end point of the path. As the mean arrival time scales  $N^{3/4}$  and it can be estimated as  $1/(\langle i_2 \rangle/N)$ , we get  $i_2 \sim N^{1/4}$ . Then we get the scale of the selection at the second sequence  $\hat{S}_2 \sim 1/N^{1/4}$ . The numerics [see Fig. 3(b)] support the following scalings for the mean first arrival times for the given  $S_2$ :

$$|S_2| \gg 1/N^{1/4}, \quad T_1 \sim N|S_2|, \quad T_2 \sim N. \quad (10)$$

Here we denote the mean first arrival time for the case  $T_2$  by  $|S_2| \sim 1$ .

### D. Climbing the fitness hill

Consider now a fitness hill located at the first sequence on the path. With a growing height of the fitness hill, the mean first arrival time  $T$  begins by decreasing and then it grows (see Fig. 4). The maximal value of the acceleration  $A$  (compared with the neutral case) grows as a logarithm of  $N$  (see Fig. 5). For the larger values of  $S_1$ , we found that  $T$  grows, and for  $S_1 \sim 1$  there is an  $N^{1/4}$  time increase of the arrival time  $T$ , which is much smaller than the increase due to the fitness

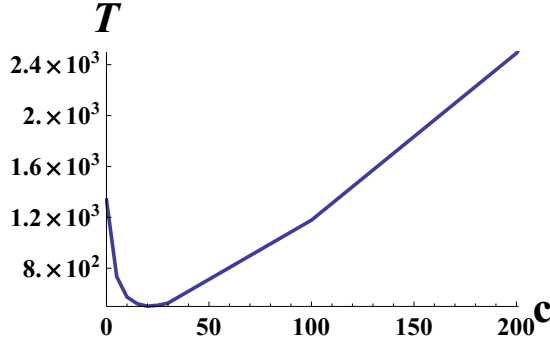


FIG. 4. Mean first arrival time for  $m = 3$ ,  $N = 10\,000$ ,  $S_0 = 0$ ,  $S_1 = c/N^{3/4}$ , and  $S_2 = 0$ .

canyon ( $N^{3/2}$ ). Locating the fitness hill near the end point of the mutation path again leads to the scaling of the  $m - 1 = 2$  case.

#### E. Crossing of a fitness canyon for the general $m$ case

In the case of the fitness canyon at the first sequence, we simply generalize the results of the  $m = 3$  case for the mean first arrival time,

$$T_1 \sim N^{2n(m)}|S_1|, \quad T_2 \sim N^{2n(m)}, \quad (11)$$

where  $n(m)$  is given in Eq. (4). We also investigated the location of the canyon just before the end point of the path. We measured also the fixation time of the  $m$ th mutant. Figure 6 illustrates that this time period grows  $\sim |S_0 - S_1|$ . The numerics are consistent with the last scaling in Eq. (10).

#### F. Climbing of the fitness hill for the general $m$ case

If we put a high fitness hill near the end point of the sequence chain, then the mean first arrival time is shortened to the value of the  $m - 1$  case. The interesting thing is the acceleration of the  $m$ th mutant arrival due to intermediate mutants with a

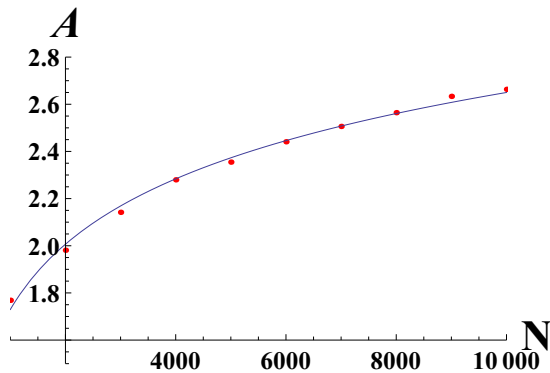


FIG. 5. Here  $A$  is the maximal acceleration of the mean first arrival due to the intermediate high fitness at the position for  $m = 3$ ,  $S_0 = 0$ ,  $S_2 = 0$ , and  $S_1 > 0$  and  $N$  is the population size. We are choosing the value of  $S_1$  giving the maximal acceleration ( $A$  times) compared with the neutral case  $N_1 = 1$ . The solid line is the fit of the data via logarithmic function and dots correspond to data. In addition,  $m = 3$ ,  $N = 10\,000$ ,  $S_0 = 0$ ,  $S_1 = c/N^{3/4}$ , and  $S_2 = 0$ .

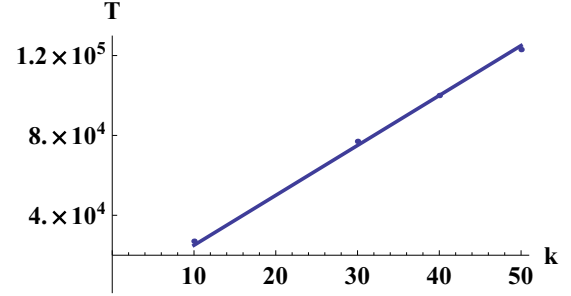


FIG. 6. Mean time of fixation for the third mutant for  $m = 3$ ,  $N = 1000$ ,  $S_0 = 0$ ,  $S_2 = 10$ ,  $S_1 = -c/N^{0.75}$ , and  $S_3 = 10/N$ .

slightly positive difference of  $S_l$ ,  $0 < l < m - 1$ . As a very high value of  $S_l$  prolongs the mean arrival time, there is some optimal set of the fitness values, maximally accelerating the mean first arrival time. It is an interesting problem to estimate how this acceleration scales with  $N$ . We assume a logarithmic acceleration of the first arrival, as in the  $m = 3$  case. A close problem, the optimization of the fitness values along the evolutionary dynamic path, was considered in [30]. They considered the case of very small mutation rates, less than  $1/N^2$ .

A more serious problem is the extension of the mean arrival time due to the fitness hill at the first sequence. In the case of infinite populations, the probability of far mutants has a small factor  $\sim 1/|S_1|^{m-2}$  [31]; then we get an extension of the arrival time, proportional to the inverse of this probability. The numerics for the  $m = 3$  case showed that the arrival time scales about  $|S_1|^{1.3}$ . According to our numerics, the fitness hill extends the arrival time much more than the fitness canyon. For  $S_1 = 1$  the mean first arrival time scales as  $T_2 \sim N^2$ , versus  $N^{7/8}$  for the neutral case for our scaling by Eq. (2).

### III. ALTERNATIVE SCALINGS

While looking at the mean first arrival for the  $m > 2$  case, we used the scaling in Eq. (4). In general, for the neutral case,  $4m - 3$  different scalings were found in [15].

#### A. Rare mutations

The rare mutations correspond to the case

$$u \leq 1/N^2, \quad T_0 \sim 1/u, \quad (12)$$

where  $T_0$  is the mean arrival time of the  $(m + 1)$ th mutant in the neutral case. We numerically found that the existence of the canyon drastically (stronger than linearly in the degree of the height of the canyon) prolongs the mean arrival time in this case for both the mean first arrival and fixation times (see Fig. 7).

#### B. Fast mutations

According to [15], the neutral mutations are fast when

$$u \gg 1/N^{2/m}. \quad (13)$$

In such a case we can apply the infinite population results [see Eq. (A2)]. We verified that for  $m > 2$  the mean first arrival



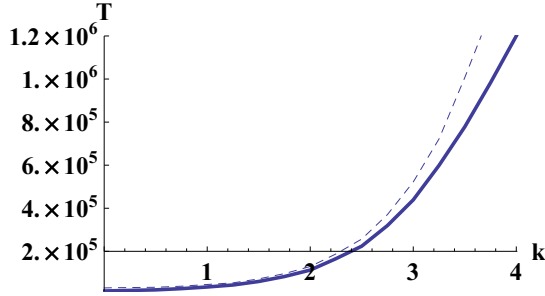


FIG. 7. Mean first arrival time for the third mutant (solid line) and the fixation time (dashed line) versus the canyon depth  $k$  for  $m = 3$ ,  $S_0 = 0$ ,  $S_1 = -k/N$ ,  $S_2 = 0$ ,  $S_3 = 0$ , and  $u = 1/N^2$ .

time grows with the canyon depth less than linearly:

$$T \sim \frac{s^{1/(m-1)}}{(Nu^m)^{1/(m-1)}}. \quad (14)$$

#### IV. DISCUSSION

In this paper we gave a simple formula for the mean first arrival time of double mutants in the case of double mutants with selection, using just a hypergeometric function, versus an integral expression in [24]. We got the asymptotic expression of a large intermediate selection from our exact result for the two-mutant case; we then performed an extensive numerical calculation for the multiple-mutant arrival problem. For the case of multiple sequences  $m > 2$ , we investigated the mean first arrival time for the path of mutations; when all the mutations are neutral besides the single one and the absolute fitness difference of the latter is much larger than the inverse population size, the difference is either a positive number (fitness hill) or negative (the fitness canyon). We investigated the mean first arrival time  $T$ . The location of the fitness canyon (hill) is crucial for the prolongation of  $T$ . When the canyon is located near the original sequence, the  $T$  is rather strongly affected; in the limiting case of fitness difference,  $T$  becomes the square of the neutral case; when the canyon is located near the last sequence,  $T$  is only slightly changed. We gave analytical formulas for the scaling of the first arrival time when the canyon is located either at the first sequence of the mutation path or near the end point.

When there is a high fitness hill  $S_{m-1} \sim 1$  near the end point of the sequence path, the mean first arrival time is shortened until the value of the  $m - 1$  case. The situation with a mean first arrival time is more interesting when there is a fitness hill just after the original sequence. For the small height of the hill,  $T$  decreases several times, while with a further growth of the hill's height  $T$  grows much faster than in the case of the fitness canyon. In the case of the first arrival of the fourth mutant  $T$  is about the square of population size, thus there is an about  $N^{9/8}$  times acceleration compared with the neutral case. This acceleration is even higher for the farther mutants.

We also looked at the optimization of the mean first arrival time, investigated before for the weaker mutation case in [30]. In this work we mainly focused on the first arrival time of a new mutant, while more relevant for virology and oncology is the crossing time by the whole population. We assume that our results are qualitatively correct for the latter case as

well. The fixation time will grow proportionally to the large fitness difference at intermediate sequences, as it is illustrated in Figs. 2(b) and 3(a). Much more dramatic is the hill climbing time, when the hill is located at the beginning of the path. We assume that our findings together with the results known about neutral case [16] give the scaling of the mean first arrival time problem. We have done some numerical calculations for the fixation, deducing how the time scales with the depth of the canyon [Figs. 2(b), 6, and 7]. While the most of our results have been related to the scaling of the mutation rates proportional to the inverse of population size, we also considered the case of rare and fast mutations. We verified that in the case of rare mutations the existence of the canyon drastically prolongs the mean first arrival time, while the latter is only slightly affected by the existence of a canyon in the case of fast mutations. The mean first arrival time, considered in our article, is a much simpler problem than the calculation of the finite population size corrections to the mean fitness [31,32] or the clonal interference phenomenon [33–36].

#### ACKNOWLEDGMENTS

D.B.S. was supported by Grant No. MOST 105-2811-M-001-078 and C.-K.H. was supported by Grant No. MOST 105-2112-M-001-004. This work was initiated with the support of Taiwan-Russia collaborative Grant No. NSC 101-2923-M-001-003-MY3. A.S.B. acknowledges support by Russian Federation Government Subsidy 14.607.21.0091.

#### APPENDIX: MEAN FIRST ARRIVAL TIME FOR A DOUBLE MUTANT

##### 1. Derivation of Eq. (5)

In [24], the diffusion approximation was applied to derive the differential equation for the mean first arrival time of the double mutant. We looked the Wright-Fisher model with the parameters

$$u_0 = b/2N, \quad b_1 = a/2N, \quad S_1 = c/2\sqrt{N}. \quad (A1)$$

We denote by  $T_i$  the mean first arrival time when at the start there are  $i$  B-type alleles. We can write a differential equation for  $T_i$ , following [10]. We introduce

$$x = \frac{i}{\sqrt{N}}, \quad T_i = 2\sqrt{N}y(x). \quad (A2)$$

Then, in [24], the following expression was derived:

$$y(0) = M(A, b, 0) \int_0^\infty \frac{dz e^{kz}}{z M(A, b, z) h \left( \frac{U'(A, b, z)}{U(A, b, z)} - \frac{M'(A, b, z)}{M(A, b, z)} \right)}, \quad (A3)$$

where  $M(A, b, z)$ ,  $U(A, b, z)$  are Kummer functions and the parameters are defined as

$$h = \sqrt{c^2 + 4a}, \quad k = \left( 1 + \frac{c}{\sqrt{c^2 + 4a}} \right) \frac{1}{2}, \quad A = kb. \quad (A4)$$

Consider now the alternative expression using the expression for the Wronskian (13.2.34) in [29]:

$$M(A, b, x)U'(A, b, x) - M'(A, b, x)U(A, b, x) = \frac{e^x \Gamma(b)}{x^b \Gamma(A)}. \quad (\text{A5})$$

We used the relation that the Kummer function  $M$  equals  $\Gamma(b)$  times the Oliver function  $M$ . We obtain

$$y(0) = \frac{\Gamma(A)}{h\Gamma(b)} \int_0^\infty dz \frac{e^{-(1-k)z} U(A, b, z)}{z^b}. \quad (\text{A6})$$

Then, using the known expression for the integral (13.10.7) [taking into account the difference between hypergeometric functions  ${}_2F_1$  and  ${}_2F_1$ ; we use the latter,  $F = F/\Gamma(c)$ ],

$$\begin{aligned} & \int_0^\infty dz t e^{-zt} U(A, B, t) t^{b-1} \\ &= \frac{\Gamma(b)\Gamma(b-B+1)}{\Gamma(c)} \\ & \times z^{-b} {}_2F_1\left(A, b; A+b-B+1; 1-\frac{1}{z}\right), \end{aligned} \quad (\text{A7})$$

we get eventually

$$y(0; a, b, c) = \frac{\Gamma[A] {}_2F_1\left[A, b, A+1, 1-\frac{1}{z}\right]}{\Gamma[A+1] h z^b}. \quad (\text{A8})$$

For the asymptotic behavior of  $F(a; b; c; z)$  as  $z \rightarrow \infty$  with  $a$ ,  $b$ , and  $c$  fixed, we combine (15.2.2) with (15.8.2) or (15.8.8) of [29].

## 2. Relation to the infinite population models

Consider the following system of equations for the growing population with the relative probabilities  $x_1$  and  $x_2$  for the first and second mutants, respectively. We drop the population dilution terms in the equation. First we consider the canyon

case

$$\frac{dx_1}{dt} = u - Sx_1, \quad \frac{dx_2}{dt} = ux. \quad (\text{A9})$$

We have also the initial condition  $x_1(0) = 0$ . We obtain

$$x_1(t) = \frac{u}{S}(1 - e^{-St}). \quad (\text{A10})$$

We calculate the mean first arrival time using the condition  $Nx_2(T) = 1$ , where  $N$  is the population size. Then we obtained

$$x_2 = Nu^2 \left( \frac{1}{S} T - \frac{1}{S^2} (1 - e^{-ST}) \right). \quad (\text{A11})$$

Holding only the first nonzero term, we derive

$$N \frac{u^2}{S} T = 1 \quad (\text{A12})$$

and Eq. (8) for the mean first arrival time for the canyon with the depth  $S$ .

Considering now the case of the fitness hill  $S$ , we get the equation

$$N \left( -\frac{u}{S} T + \frac{u}{S^2} (e^{ST} - 1) \right) = 1. \quad (\text{A13})$$

The principal term is  $\frac{u}{S^2} e^{ST}$ , which gives

$$T \sim \frac{1}{S} \ln \frac{S}{Nu^2}, \quad (\text{A14})$$

which is consistent with Eq. (7).

Consider now the canyon crossing with  $m + 1$  mutant types, where the canyon is again for the first type of mutant case. Now we obtain an equation similar to Eq. (A12),

$$N \frac{u^m T^{(m-1)}}{S} \sim 1, \quad (\text{A15})$$

or Eq. (14).

- 
- [1] M. Eigen, *Naturwissenschaften* **58**, 465 (1971).
  - [2] D. B. Saakian, C. K. Biebricher, and C.-K. Hu, *PLoS One* **6**, e21904 (2011); Z. Kirakosyan, D. B. Saakian, and C.-K. Hu, *J. Phys. Soc. Jpn.* **81**, 114801 (2012).
  - [3] P. Bak and K. Sneppen, *Phys. Rev. Lett.* **71**, 4083 (1993).
  - [4] D. B. Saakian, M. H. Ghazaryan, and C.-K. Hu, *Phys. Rev. E* **90**, 022712 (2014).
  - [5] A. Nagar and K. Jain, *Phys. Rev. Lett.* **102**, 038101 (2009).
  - [6] T. Yakushkina, D. B. Saakian, and C.-K. Hu, *Chin. J. Phys.* **53**, 100904 (2015).
  - [7] D. B. Saakian, T. Yakushkina, and C.-K. Hu, *Sci. Rep.* **6**, 34840 (2016).
  - [8] J. F. Crow and M. Kimura, *An Introduction to Population Genetics Theory* (Harper & Row, New York, 1970).
  - [9] M. Kimura, *J. Genet.* **64**, 7 (1985).
  - [10] W. J. Ewens, *Mathematical Population Genetics* (Springer, New York, 2004).
  - [11] G.-R. Huang, D. B. Saakian, and C.-K. Hu, *J. Phys. Soc. Jpn.* **85**, 044803 (2016).
  - [12] F. B. Christiansen, S. P. Otto, A. Bergman, and M. W. Feldman, *Theor. Popul. Genet.* **53**, 199 (1998).
  - [13] D. Weinreich and L. Chao, *Evolution* **59**, 1175 (2005).
  - [14] Y. Iwasa, F. Michor, and M. A. Nowak, *Genetics* **166**, 1571 (2004).
  - [15] J. Schweinsberg, *Electron. J. Probab.* **13**, 1442 (2008).
  - [16] R. Durrett, D. Schmidt, and J. Schweinsberg, *Ann. Appl. Probab.* **19**, 676 (2009).
  - [17] D. B. Weissman, M. M. Desai, D. S. Fisher, and M. W. Feldman, *Theor. Popul. Biol.* **75**, 286 (2009).
  - [18] J. A. Draghi, T. L. Parsons, G. P. Wagner, and J. B. Plotkin, *Nature (London)* **463**, 353 (2010).
  - [19] R. A. Fisher, *The Genetical Theory of Natural Selection* (Clarendon, Oxford, 1930).
  - [20] S. Wright, *Genetics* **16**, 97 (1931).
  - [21] S. Wright, *Proc. Natl. Acad. Sci. USA* **31**, 382 (1945).
  - [22] D. B. Saakian and C.-K. Hu, *Phys. Rev. E* **94**, 042422 (2016).
  - [23] R. Durrett and D. Schmidt, *Genetics* **180**, 1501 (2005).

- [24] D. B. Saakian, *J. Stat. Mech.* (2015) P05036.
- [25] T. Brotto, G. Bunin, and J. Kurchan, *J. Stat. Mech.* (2016) 033302.
- [26] A. J. Stewart, T. L. Parsons, and J. B. Plotkin, *Evolution* **66**, 1598 (2012).
- [27] P. Schuster, W. Fontana, P. F. Stadler, and I. L. Hofacker, *Proc. Biol. Sci.* **255**, 279 (1994).
- [28] E. V. Nimwegen, J. P. Crutchfield, and M. Huynen, *Proc. Natl. Acad. Sci. USA* **96**, 9716 (1999).
- [29] A. B. O. Daalhuis, in *NIST Handbook of Mathematical Functions*, edited by M. W. J. Olver (Cambridge University Press, Cambridge, 2010), Chap. 13.
- [30] A. Traulsen, Y. Iwasa, and M. A. Nowak, *J. Theor. Biol.* **249**, 617 (2007).
- [31] D. B. Saakian, C. K. Hu, and M. W. Deem, *Europhys. Lett.* **98**, 18001 (2012).
- [32] J.-M. Park, E. Muñoz, and M. W. Deem, *Phys. Rev. E* **81**, 011902 (2010).
- [33] M. M. Desai and D. S. Fisher, *Genetics* **176**, 1759 (2007).
- [34] S. C. Park and J. Krug, *Proc. Natl. Acad. Sci. USA* **104**, 18135 (2007).
- [35] K. Jain, *Genetics* **179**, 2125 (2008).
- [36] I. M. Rouzine, E. Brunet, and C. O. Wilke, *Theor. Popul. Biol.* **73**, 24 (2008).