

Evolutionary rescue of bacterial populations by heterozygosity on multicopy plasmids

Ian Dewan^{1,2*} and Hildegard Uecker^{1,3}

¹Research Group Stochastic Evolutionary Dynamics, Department of Theoretical Biology, Max Planck Institute for Evolutionary Biology, 24306 Plön, Germany.

²ORCID: 0000-0001-8964-3800.

³ORCID: 0000-0001-9435-2813.

*Corresponding author(s). E-mail(s): dewan@evolbio.mpg.de;

Abstract

Bacterial plasmids and other extrachromosomal DNA elements frequently carry genes with important fitness effects for their hosts. Multicopy plasmids can additionally carry distinct alleles of host-fitness-relevant genes on different plasmid copies, allowing for heterozygosity not possible for loci on haploid chromosomes. Plasmid-mediated heterozygosity may increase the fitness of bacterial cells in circumstances where there is an advantage to having multiple distinct alleles (heterozygote advantage); however, plasmid-mediated heterozygosity is also subject to constant loss due to random segregation of plasmid copies on cell division. We analyze a multitype branching process model to study the evolution and maintenance of plasmid-mediated heterozygosity under a heterozygote advantage. We focus on an evolutionary rescue scenario in which a novel mutant allele on a plasmid must be maintained together with the wild-type allele to allow population persistence. We determine the probability of rescue and derive an analytical expression for the threshold on the fitness of heterozygotes required to overcome segregation and make rescue possible; this threshold decreases with increasing copy number of the plasmid. We further show that the formation of cointegrates from the fusion of plasmid copies increases the probability of rescue. Overall, our results provide a rigorous quantitative assessment of the conditions under which bacterial populations can adapt to multiple stressors through plasmid-mediated heterozygosity.

Keywords: evolutionary rescue, multicopy plasmids, bacterial evolution, plasmid copy number, fitness trade-off

MSC Classification: 92D15 , 60J85

1 Introduction

Many bacteria carry, in addition to the bacterial chromosome, extrachromosomal genetic elements called plasmids, which replicate independently of the chromosome and often exist in multiple copies. The simplest plasmids carry only the core backbone genes required to ensure their stable maintenance in the bacterial host, but many plasmids also carry genes that have important effects on the host phenotype (Garcillán-Barcia et al, 2011). Perhaps the most famous of these are antibiotic resistance genes, which pose a serious threat to the effectiveness of the clinical treatment of bacterial infections (Carattoli, 2013), but phenotypes contributed by plasmid-borne genes also include heavy-metal resistance, virulence, metabolism of novel carbon sources, and symbiotic interactions with other organisms, among many others (see e.g. Portnoy and Martinez, 1985; Silver, 1992; Beijersbergen et al, 1994; Yu et al, 1996; Anda et al, 2015; Wardell et al, 2021). The plasmids carried by a bacterium form an important part of the bacterial genome, and the contribution of plasmids must be considered in explaining and predicting bacterial evolution. In this context, much attention has been given to the ability of many plasmids to transfer horizontally between bacteria, infecting new hosts, in a process called conjugation (Falkow, 1975; Smillie et al, 2010); however, very many plasmids are transmitted only vertically, or depend on other plasmids for horizontal transfer (Coluzzi et al, 2022), and purely vertically transmitted plasmids can play an important role in bacterial adaptation.

The backbone genes of a plasmid are responsible for ensuring the stable maintenance of the plasmid by vertical transmission. This requires ensuring that both daughter cells receive copies of the plasmid at host cell division, and also limiting the possible reduction in host growth imposed by fitness costs of plasmid carriage. To this end, the number of copies of a plasmid per host cell is regulated by the plasmid itself to remain approximately constant; the plasmid copy number is therefore an intrinsic property of the plasmid-host system which plays an important role in the biology of

the plasmid. Although for some plasmids the copy number is kept quite low, at only one or a few copies, many plasmids are found in their hosts in high numbers of copies, often in tens or sometimes even hundreds. Interest in the evolution of these multicopy plasmids and their contribution to bacterial adaptation has recently increased, particularly due to their contribution to antibiotic resistance (San Millan et al, 2009; Gama et al, 2018).

Existence in multiple copies per cell has multiple effects on the contribution of multicopy plasmids to bacterial evolution. The high copy number of the plasmid means the copy numbers of genes carried on the plasmid are also higher than those of genes on the (for many species haploid) bacterial chromosome. The mutational input for these genes is higher because of this; higher gene copy number can also increase the dosage of genes on multicopy plasmids (these effects have been shown in experiments by, e.g., San Millan et al (2016)). The presence of the plasmid in multiple copies also provides the possibility of heterozygosity for loci on a multicopy plasmid, which is not possible for loci on the haploid chromosome.

Once there are multiple variants of the plasmid in the cell, the process of plasmid segregation at host cell division will have important effects on the distribution of plasmid variants among cells in the population (Ilhan et al, 2018). At host cell division, the copies of a plasmid in the host are distributed between the two daughter cells: for low copy number plasmids this is often aided by an active partitioning system to ensure each daughter gets at least one copy of the plasmid, while high copy number plasmids may simply rely on random distribution of the plasmids between daughter cells (Zielenkiewicz and Ceglowski, 2001). Random assortment of the plasmid copies at host cell division (i.e. segregation independent of the allele at the heterozygous locus) tends to eliminate heterozygosity, since any series of cell divisions—even every single cell division—has some nonzero probability of producing a homozygote daughter cell by chance, and the homozygotes are absorbing states of this process. An analogous

situation appears when two different (multicopy) plasmids that share the same replication system reside within the same cell—over the course of time, cell lines carrying either one or the other will emerge. Such segregative loss of one or the other plasmid has already been observed and quantified in early studies on plasmid incompatibility (e.g. Uhlin and Nordström, 1975; Cullum and Broda, 1979). This loss can be counteracted, however, by selection for maintenance of the plasmid variants together (see e.g. Uhlin and Nordström, 1975; Cullum and Broda, 1979; Rodríguez-Beltrán et al, 2018).

Driven by these empirical observations, early models were developed to describe the kinetics of the loss of heterozygosity (or incompatible plasmids) in multicopy plasmids in the absence of selection under varying assumptions about the mechanisms of segregation and replication (Ishii et al, 1978; Novick and Hoppensteadt, 1978; Cullum and Broda, 1979). Recently, interest in the evolution of multicopy plasmids and their hosts has gained momentum (e.g. San Millán et al, 2016; Rodríguez-Beltrán et al, 2018; Santer and Uecker, 2020; Rodríguez-Beltrán et al, 2021; Garoña et al, 2021; Hernández-Beltrán et al, 2022; Santer et al, 2022; Garoña et al, 2023). This includes experimental studies of multicopy plasmids with accompanying models: some of these focus on segregation of plasmid copies as a limitation on the fixation of novel beneficial alleles on plasmids (Ilhan et al, 2018; Garoña et al, 2023), while others focus on the effect of selection for heterozygotes on maintaining plasmid variants together despite segregation (Rodríguez-Beltrán et al, 2018). More recent modelling studies have also examined the effects of assumptions about plasmid replication and segregation on the fate of novel beneficial alleles that appear on (a single copy of) multicopy plasmids (Halleran et al, 2019; Santer and Uecker, 2020).

We here focus on scenarios of heterozygote advantage, where the optimal bacterial fitness is obtained by having multiple variants of a single plasmid carrying distinct alleles. This could be because, for example, the two alleles confer resistance to distinct antibiotics present in the environment—as can be the case with β -lactamases, for which

point mutations can alter the affinity of the enzyme for different β -lactams (Rodriguez-Beltran et al, 2018). The models we present are of an *evolutionary rescue* scenario, of the kind introduced by Gomulkiewicz and Holt (1995), extending previous models of rescue on plasmids (Tazzyman and Bonhoeffer, 2014; Santer and Uecker, 2020). In such a scenario, an environmental change exposes a population to new conditions to which it is maladapted, and the population therefore enters a demographic decline which would lead to extinction. However, if a novel mutation emerges which adapts individuals to the new conditions, and this mutation survives to spread in the population, the population can be rescued from extinction. Since this rescue depends on the mutation occurring and surviving the initial period in which it is rare in the population, rescue is inherently a stochastic process, and the key question is the probability that rescue will occur in a given population. When maintenance of both plasmid variants is required for rescue, rescue does not only require the establishment of the new mutation against the force of genetic drift, but also the persistence of both variants against the force of segregation at all times. It is intuitively clear, that such persistence is only possible if the fitness of heterozygote cells is large enough to bear the production of unfit homozygous cells. Based on a multi-type branching process model, we determine the probability of evolutionary rescue and derive analytical conditions on the plasmid copy number and the fitness of heterozygous cells that need to be fulfilled for maintenance of heterozygosity, and thus rescue, to be possible at all.

2 The model

Consider a demographically stable population of bacteria which carries a plasmid present at a fixed copy number n in each bacterial cell. This population is then exposed to novel environmental conditions to which it is maladapted, and the population begins a demographic decline which will eventually lead to its extinction. Evolution, however, may rescue this population from extinction. A mutation might appear at a locus on

the plasmid which, *if maintained together with the wild-type allele at that locus*, will adapt the population to the novel environmental conditions; that is, mutant *heterozygotes* can survive in the novel environmental conditions, but mutant and wild-type homozygotes are both maladapted. We assume that the population is homozygous for this particular plasmid before the environmental change; that is, that there is only one variant of the plasmid in the population, at least with respect to the locus of interest.

These conditions constitute the evolutionary-rescue-due-to-heterozygote-advantage scenario we wish to describe in the model. As an example of a circumstance in which this scenario might arise, consider a bacterial population resistant to antibiotic A due to a resistance gene located on a plasmids. If the population is exposed to antibiotics A and B simultaneously, then lack of resistance to antibiotic B will cause the population to decline. But if a mutation occurs on a plasmid that converts the A-resistance allele into a B-resistance allele, its host is now resistant to antibiotic B, and (provided that $n > 1$) the remaining wild-type plasmids continue to confer resistance to antibiotic A: therefore the host can survive and grow in the novel conditions. Its descendants are resistant to both antibiotics and the population might be rescued from extinction, provided that both plasmid types are maintained. A descendent cell which loses the mutant and has only wild-type plasmids has no resistance to antibiotic B, while a descendent which has only mutant and no wild-type plasmids has no resistance to antibiotic A; neither will be able to survive. [Rodriguez-Beltran et al \(2018\)](#) showed an example of a very similar situation arising in an evolution experiment. While such a combination of two closely-related antibiotics would not be in a clinical context to treat a patient, our scenario might arise if the two antibiotics co-occur in the environment.

Our model extends that of [Santer and Uecker \(2020\)](#) to the case of heterozygote advantage. Suppose that cells with i mutant plasmids (and therefore $n - i$ wild-type plasmids) reproduce at a rate $\lambda_i^{(n)}$ and die at a rate $\mu_i^{(n)}$. Since the units of time are

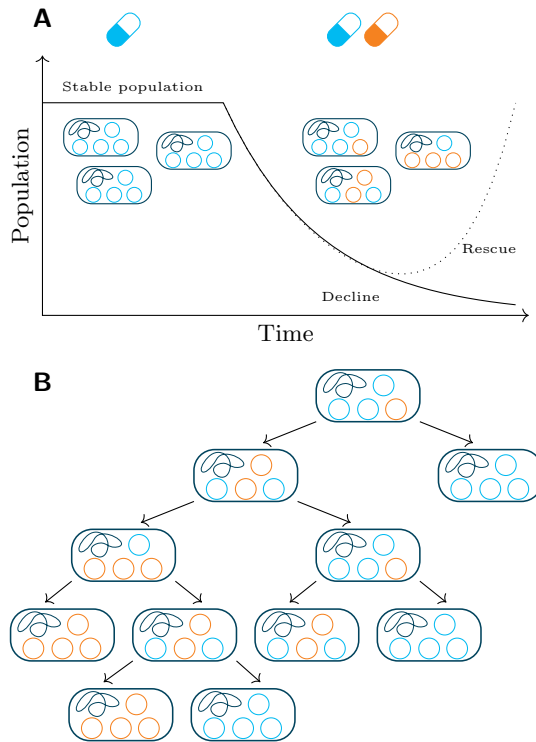


Fig. 1 A model of the rescue of a bacterial population by plasmid-mediated heterozygosity. (A) The evolutionary rescue scenario. An initial stable population carrying a plasmid with an adaptive gene on it is exposed to a sudden environmental change, which leads the population to decline. A mutation on the plasmid (orange plasmids) can rescue the population, but only if it is maintained together with the wild-type (blue) plasmid. (B) The loss of plasmid-mediated heterozygosity due to segregation. From an initial heterozygote cell, random segregation results in the two plasmid types ending up in different, homozygous cells. Based on Figure 1 of [Santer and Uecker \(2020\)](#).

arbitrary, we may fix $\mu_i^{(n)} = 1$. If we then set $\lambda_i^{(n)} = 1 + s(n, i)$, the function $s(n, i)$ gives the net growth rate (or Malthusian fitness) of cells with i mutant plasmids when the plasmid copy number is n . In the scenario we are modelling, where maintenance of both wild-type and mutant plasmids is necessary for the population to persist, we have that $s(n, 0) < 0$ and $s(n, n) < 0$; beyond this, the choice of fitness function is free.

We will consider two different possible fitness functions. The simplest case is

$$s(n, i) = \begin{cases} s_0 & i = 0 \text{ or } i = n, \\ s_{\max} & \text{otherwise,} \end{cases}$$

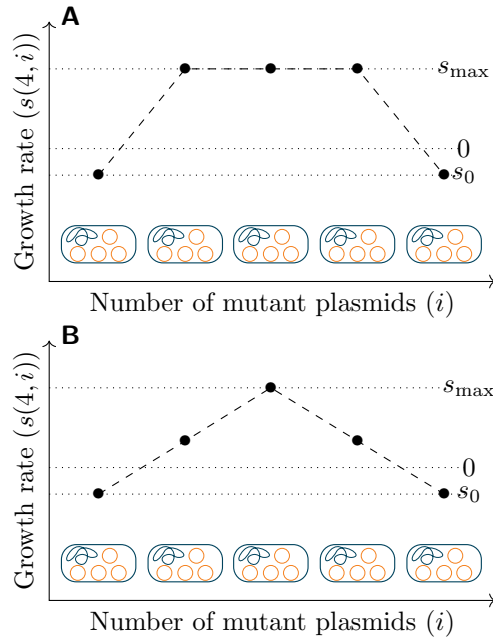


Fig. 2 Examples of the two fitness functions considered, for $n = 4$: (A) the dominant fitness function; (B) the peaked fitness function.

where $s_0 < 0$ and $s_{\max} > 0$. We shall call this the *dominant fitness function* (shown in Figure 2A), since it implies that a single copy of each of the two alleles confers the full fitness effect of that allele. It is also possible to imagine that dominance is intermediate or that there is a gene dosage effect, in which additional copies of a gene increase resistance. Combined with the necessity to maintain both alleles, this might produce a fitness function like

$$s(n, i) = s_{\max} - \frac{2(s_{\max} - s_0)}{n} \left| i - \frac{n}{2} \right|,$$

with $s_0 < 0$ and $s_{\max} > 0$ as before; this function linearly interpolates between a maximum fitness at a exact half-and-half mixture of wild-type and mutant plasmids, and a minimum fitness in the homozygotes (shown in Figure 2B). We shall call this the *peaked fitness function*. If there are gene dosage effects, s_{\max} would increase with

n ; we do not consider this in our numerical examples, but all results for a given n also apply to this scenario.

We have assumed above that every cell has exactly n copies of the plasmid of interest. To maintain this number, we assume the cell replicates its plasmids to a copy number of $2n$ before reproduction, and then segregates equal numbers into each daughter cell. At replication of wild-type plasmids, mutations happen with probability u , but we disregard this for the moment. We denote by $P^{(n)}(i \rightarrow \{j, k\})$ the probability that when a cell with i mutant plasmids (out of n total plasmids) divides it produces daughters with j and k mutant plasmids (again out of n total). We will consider two models of replication, called *regular* and *random replication*, after the taxonomy of Novick and Hoppensteadt (1978).

In regular replication, every plasmid copy in the parent cell is replicated once, and then the copies are segregated randomly, but maintaining the copy number, into each daughter cell. The number of mutant plasmids in the daughter cell then has a hypergeometric distribution, giving a segregation probability function

$$P^{(n)}(i \rightarrow \{j, k\}) = \begin{cases} 0 & j + k \neq 2i, \\ \frac{\binom{2i}{j} \binom{2(n-i)}{n-j}}{\binom{2n}{n}} & j = k = i, \\ 2 \frac{\binom{2i}{j} \binom{2(n-i)}{n-j}}{\binom{2n}{n}} & \text{otherwise.} \end{cases}$$

In random replication, each new plasmid to be replicated is randomly chosen from the pool of the initial plasmids together with the products of previous replications. This proceeds until there are $2n$ total plasmids, when the number of mutants will be distributed according to a Pólya urn scheme. These $2n$ plasmids are then randomly segregated into the daughter cells as before. The number of mutant plasmids in a

daughter cell is then distributed as

$$P^{(n)}(i \rightarrow \{j, k\}) = \begin{cases} \frac{\binom{j+k}{j} \binom{2n-j-k}{n-j}}{\binom{2n}{n}} \cdot \frac{\binom{2n-j-k-1}{n-j-k+i} \binom{j+k-1}{j+k-i}}{\binom{2n-1}{n}} & j = k, \\ 2 \frac{\binom{j+k}{j} \binom{2n-j-k}{n-j}}{\binom{2n}{n}} \cdot \frac{\binom{2n-j-k-1}{n-j-k+i} \binom{j+k-1}{j+k-i}}{\binom{2n-1}{n}} & \text{otherwise.} \end{cases}$$

The first fraction in the expression is the same hypergeometric distribution for segregation as in the regular replication model, representing the probability that the $j + k$ mutant plasmids produced by replication are divided into j in one daughter cell and k in the other. The second fraction is the probability of producing $j + k$ mutant plasmids from i mutant plasmids during replication under the Pólya urn process model: for the derivation see appendix A.

For rescue to happen, a mutation must occur on the plasmid and be maintained together with the wild-type plasmid in the same cell indefinitely. To determine the probability of this occurring in a given population, we split the process into two parts. In the first, we look at the descendants of a single novel mutant, and determine the probability that these descendants will survive indefinitely rather than suffering stochastic loss. This is called *establishment probability* of the mutant. In this part, we ignore additional mutations that the might recurrently generate the mutant from the wild-type plasmid, which is a rare event in a small cell line. In the second part, we can then estimate the number of mutants which will occur in wild-type homozygous cells before the initial population goes extinct, and determine the probability that at least one of them will establish; this is called the *rescue probability*.

We have formulated our birth-death model thus far as a continuous-time multitype branching process, in which cells divide and die independently of each other, excluding in particular resource competition. Of course, a population cannot grow indefinitely. However, we are here considering a population that is, at least initially, declining due to maladaptation. Mutant cells are rare in the early stages of establishment in which we

are interested and therefore independent from each other to a good approximation; this approximation is commonly made even in more complicated models and dates back to the very early calculations of establishment probabilities of beneficial alleles (Haldane, 1927). We further assume that we are far enough away from carrying capacity that we can also neglect competition with the declining wild-type population.

3 The establishment probability

3.1 Analysis

To determine the establishment probability for the descendants of a given bacterial cell, we first determine the extinction probability, the probability that the cell's descendants will at some point go extinct. While our model was formulated in continuous time, we transition to a discrete-time branching process for the analysis of the establishment probability. We can do this because we are only interested in the final outcome of the process—extinction or survival—and not in the timing of events.

A given cell will have either zero daughter cells (if it dies before reproducing), with probability $\mu_i^{(n)} / (\lambda_i^{(n)} + \mu_i^{(n)})$, or two daughter cells (if it reproduces), with probability $\lambda_i^{(n)} / (\lambda_i^{(n)} + \mu_i^{(n)})$. Since the survival of the cell depends on its complement of plasmids, we will need to track the number of mutant and wild-type plasmids in each daughter cell. The extinction probabilities $Q_0^{(n)}, Q_1^{(n)}, \dots, Q_n^{(n)}$ of the descendants of cells with, respectively, $0, 1, \dots, n$ mutant and $n, n-1, \dots, 0$ wild-type plasmids are given by the least fixed point of the generating function of the distribution of daughter cells of a single cell (Sewastjanow, 1975, Folgerung V.1). This function is given by $f(z_0, \dots, z_n) = (f_0(z_0, \dots, z_n), \dots, f_n(z_0, \dots, z_n))$, where each component f_i is given by

$$f_i(z_0, \dots, z_n) = \frac{\mu_i^{(n)}}{\lambda_i^{(n)} + \mu_i^{(n)}} + \frac{\lambda_i^{(n)}}{\lambda_i^{(n)} + \mu_i^{(n)}} \sum_{\{j,k\}} P^{(n)}(i \rightarrow \{j, k\}) z_j z_k, \quad (1)$$

the sum being taken over all unordered pairs $\{j, k\}$ of numbers of plasmids. The first term of this expression corresponds to the probability of immediate cell death, while the sum gives the probabilities of each possible pair of daughter cell types upon reproduction. The fixed point of the generating function is then the solution to system of $n + 1$ equations (1) in $n + 1$ unknowns $z_i = Q_i^{(n)}$, which can be solved numerically. This is also intuitive: The initial cell has either zero or two daughter cells. In the case that the initial cell has zero daughter cells, its descendants go extinct immediately (corresponding to the first term in the fixed point equation); in the case it has two, its descendants go extinct if and only if the descendants of both daughter cells go extinct (corresponding to the second term).

The fixed point gives us the probability that the descendants of a bacterial cell eventually go extinct, but what we really want is the probability that the descendants of a novel mutant never go extinct: the establishment probability. Under the regular replication model, this is very simple to determine: a novel mutant allele appears at first on a single plasmid in a single host cell, so the establishment of the mutant allele occurs only when the descendants of that cell do not go extinct, and

$$P_{\text{est}} = 1 - Q_1^{(n)}. \quad (2)$$

In the random replication model the situation is slightly more complex. When a mutation occurs during plasmid replication, the novel mutant plasmid is available to be selected for replication during the same bacterial generation. If it is replicated there will be multiple mutant plasmids to segregate between daughter cells, and each daughter cell might get one or more mutant plasmids. Thus the establishment process of the mutation gets a head start, possibly starting from multiple mutant plasmids which may possibly be in two separate cells. The establishment probability becomes

a quadratic form in the extinction probabilities of the individual types

$$P_{\text{est}} = 1 - \sum_{0 \leq j, k \leq n} \underbrace{\frac{(2n + (n-1)(j+k))}{n(j+k)(j+k+1)} \frac{\binom{n}{j+k}}{\binom{2n-1}{j+k}}}_{\text{replication}} \underbrace{\frac{\binom{j+k}{j} \binom{2n-j-k}{n-j}}{\binom{2n}{n}}}_{\text{segregation}} Q_j^{(n)} Q_k^{(n)}, \quad (3)$$

where the first factor inside the sum is the probability of having $j+k$ mutant plasmids after replication starting with no mutant plasmids and having one mutation occur during replication (for the derivation of this probability and its assumptions, see appendix A), and the second factor is the probability of the $j+k$ mutant plasmids being segregated into daughter cells with j and k mutant plasmids.

3.2 Results

The establishment probabilities for the regular replication model with a dominant fitness function are shown in Figure 3A. Perhaps unsurprisingly, increasing the fitness of heterozygotes (s_{max}) increases the establishment probability. More interesting is the impact of the copy number of the plasmid. At low copy number, the establishment probability increases with copy number: this is due to the impact of loss of heterozygosity to segregation. Since each cell division has some probability of producing a homozygote daughter cell (the descendants of which are doomed to extinction), heterozygosity is constantly lost to segregation. This loss is most pronounced for small copy numbers, and reduces the establishment and rescue probabilities for small copy numbers. For larger copy numbers, the establishment probability stabilizes—indeed, a close examination shows that it begins to decrease for larger copy numbers.

The establishment probabilities for the regular replication model with a peaked fitness function are shown in Figure 3B. The establishment probabilities are lower than with the dominant fitness function, and decline more precipitously for high copy

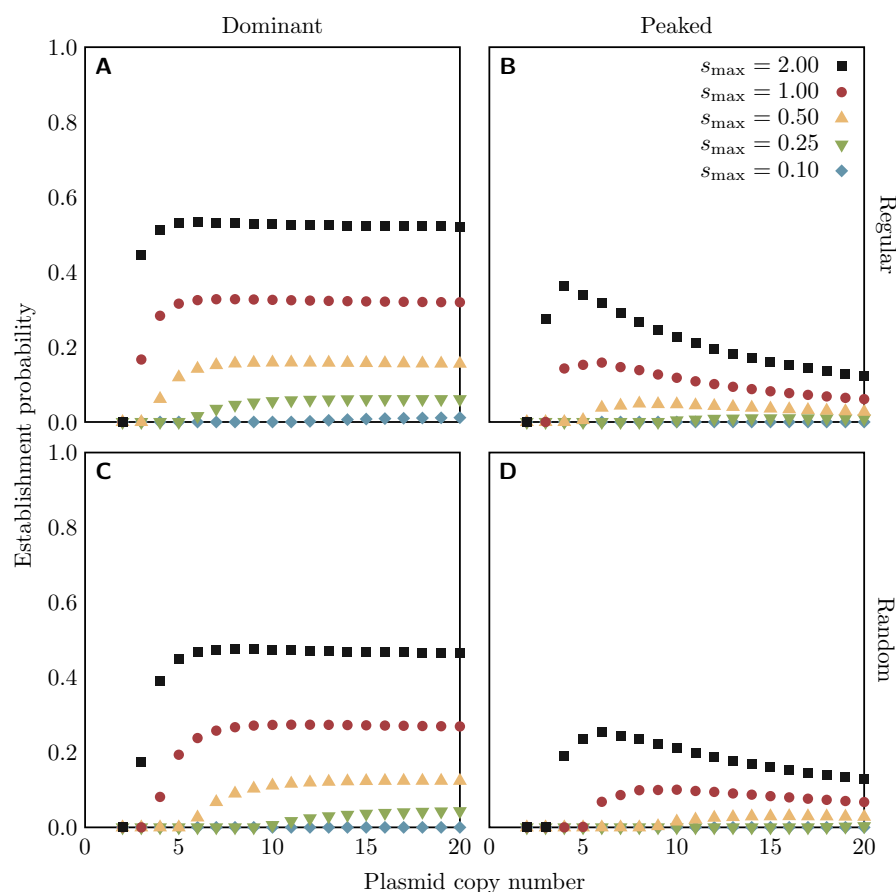


Fig. 3 Establishment probabilities of rescue mutations on a plasmid of a given copy number in the heterozygote advantage scenario, under the regular (top row, from equation (2)) or random (bottom row, from equation (3)) replication assumption. Colours of points indicate the fitness s_{\max} of all heterozygotes (with the dominant fitness function, left column) or the maximum fitness of heterozygotes (with the peaked fitness function, right column). For all cases, $s_0 = -0.1$.

numbers, since the fitness of most cell types, other than the perfectly balanced heterozygotes, has been reduced. In fact, for large enough copy numbers, the fitness of cells with positive but small numbers of one or the other plasmid type becomes negative.

The establishment probabilities for the random replication model are shown in panel C (with the dominant fitness function) and panel D (with the peaked fitness function) of Figure 3. The general trend is the same as under the regular replication model,

but the establishment probabilities are all reduced. This is because random replication increases the rate of loss of heterozygosity to segregation: the “rich-get-richer” behaviour of the Pólya urn scheme means that the proportion of mutant plasmids is on average more unbalanced after replication, and the probability of a homozygote daughter cell is increased.

Note that although the establishment probabilities appear to level off for the random replication model with dominant fitness function, in fact the establishment probability goes to zero for arbitrarily large copy number with random replication. This is because the mutant allele appears initially in a single copy: as the copy number increases, the probability of that mutant plasmid being selected for replication decreases. In the limit, the mutant plasmid is never replicated at all, and eventually the cell containing it dies by chance, and the mutant goes extinct.

4 The rescue probability

4.1 Analysis

Now let us consider the overall probability of rescue for the entire population. If the initial population size N_0 is large, we can treat the initial population deterministically. It has growth rate $\lambda_0^{(n)} - \mu_0^{(n)} = s(n, 0) < 0$, and its size after time t has passed is

$$N(t) = N_0 e^{s(n,0)t}.$$

The total number of reproduction events that occur in this population before it goes extinct is

$$\int_0^\infty \lambda_0^{(n)} N(t) dt = N_0 \frac{\lambda_0^{(n)}}{\mu_0^{(n)} - \lambda_0^{(n)}} = N_0 \frac{1 + s(n, 0)}{|s(n, 0)|},$$

and if mutations occur at a per locus rate u , the expected number of mutations that appear before extinction is

$$unN_0 \frac{1 + s(n, 0)}{|s(n, 0)|}$$

(where n appears because there are n plasmid copies per cell). We assume that, per cell division, at most one mutation occurs. The fraction of these mutations that establish and rescue the population is then just the establishment probability P_{est} . If we assume that mutations occur and establish independently (a safe assumption if they are sufficiently rare), then the number of successfully establishing mutations will be a Poisson random variable of known expected value (being the number of independent events that occur in a given time frame, with a known average rate), and the rescue probability, the probability that there will be at least one mutation that establishes, is

$$P_{\text{rescue}} = 1 - e^{-unN_0 \frac{1+s(n,0)}{|s(n,0)|} P_{\text{est}}}. \quad (4)$$

4.2 Results

The rescue probabilities are shown in Figure 4. The major difference to the trends in the establishment probabilities in Figure 3 is that the rescue probability increases with copy number even when the establishment probability is stable or declining; this is because the increased mutational input with a greater copy number increases the number of mutants that arise before wild-type population extinction.

5 Critical values of n and s_{max}

Examining Figures 3 and 4 shows that not only does the effect of loss of heterozygosity during segregation cause the establishment and rescue probabilities to decline sharply at low copy numbers, but below a certain threshold (which depends on s_{max}) establishment and therefore rescue becomes impossible ($P_{\text{est}} = P_{\text{rescue}} = 0$). We can,

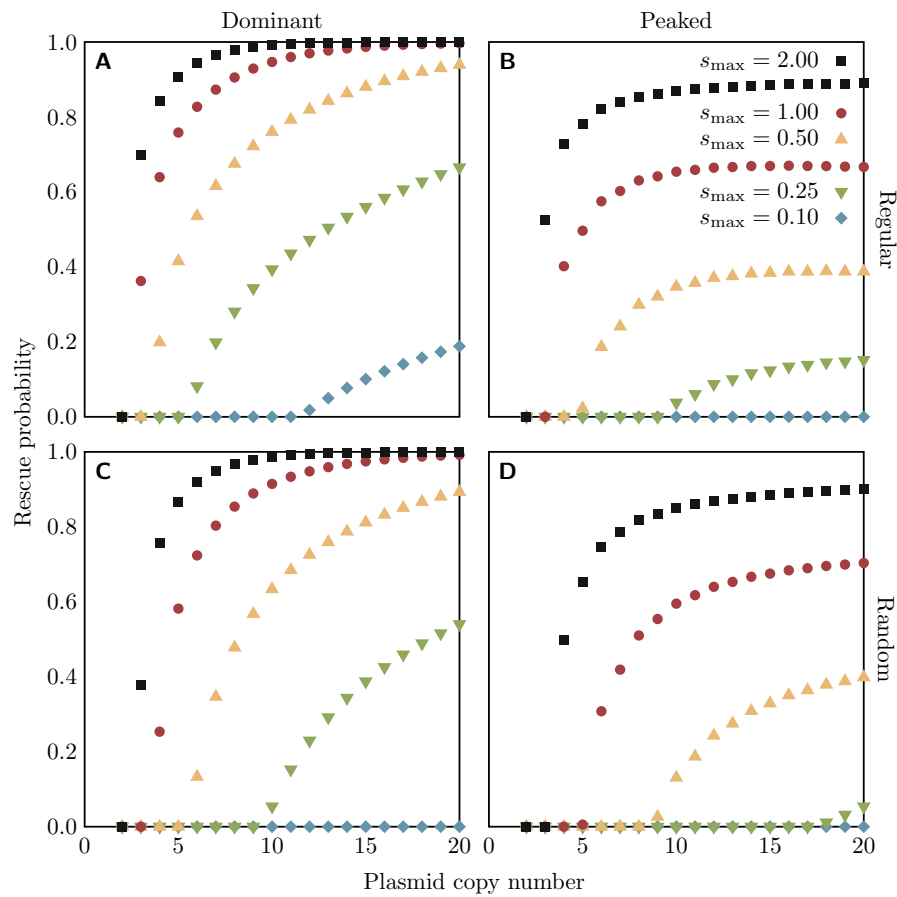


Fig. 4 Rescue probabilities (from equation (4)) of a bacterial population being rescued by a mutation on a plasmid of a given copy number in the heterozygote advantage scenario, under the regular (top row) or random (bottom row) replication assumption. Colours of points indicate the fitness s_{\max} of all heterozygotes (with the dominant fitness function, left column) or the maximum fitness of heterozygotes (with the peaked fitness function, right column). For all cases, $s_0 = -0.1$ and $uN_0 = 0.1$.

in fact, describe analytically the relation between the heterozygote advantage s_{\max} and the copy number n which is required to hold in order for establishment to be possible with the dominant fitness function, in both the regular and random replication models.

Theorem 1. *In the model with regular replication and a dominant fitness function, the establishment probability (and therefore the rescue probability) is nonzero if and*

only if

$$s_{\max} > \frac{2}{2n-3}.$$

Proof. As in the derivation of the establishment probability, we consider the branching process in discrete time. Let M be the matrix the i, j element of which is the expected number of daughter cells with j mutant plasmids from a cell with i mutant plasmids. (A cell has two daughter cells when it divides and zero if it dies.) In the model with regular replication and a dominant fitness function, the elements of this matrix have the form

$$m_{ij} = \frac{\lambda_i^{(n)}}{\lambda_i^{(n)} + \mu_i^{(n)}} \sum_{k=0}^n 2^{\delta_{jk}} P(i \rightarrow \{j, k\}) = 2 \frac{\lambda_i^{(n)}}{\lambda_i^{(n)} + \mu_i^{(n)}} \frac{\binom{2i}{j} \binom{2(n-i)}{n-j}}{\binom{2n}{n}}.$$

By a general result of branching process theory, a discrete-time multitype branching process goes extinct with probability one if and only if it has no final classes and the largest eigenvalue of M is less than or equal to one (Sewastjanow, 1975, Satz V.5). The first condition is a technical one, to exclude the degenerate case where there is no stochasticity in the branching process and the population stays at a fixed size: a final class is a subset of the types in the multitype branching process such that an individual of one of those types will have, with probability one, exactly one offspring and that offspring will be of one of the types in the class. In our case, there can be no final classes because every individual has a nonzero probability of having no offspring. Therefore, there is a nonzero probability of establishment if and only if M has an eigenvalue greater than one.

To find the eigenvalues of M , we factor it as a product $M = 2\Lambda P$, where P is the matrix with i, j element $\binom{2i}{j} \binom{2(n-i)}{n-j} / \binom{2n}{n}$ and Λ a diagonal matrix containing the ratios of birth and death rates, then find the eigenvalues of each factor individually—it will turn out, luckily, that all the eigenvalues of M are products of the eigenvalues of these factors. Finding the eigenvalues of Λ is easy, since it is diagonal. In fact,

because of the dominant fitness function, which assigns the same positive fitness to all heterozygotes and the same negative fitness to all homozygotes, there are only two distinct eigenvalues. The first,

$$\lambda_{\text{hom}} = \frac{\lambda_0^{(n)}}{\lambda_0^{(n)} + \mu_0^{(n)}} = \frac{\lambda_n^{(n)}}{\lambda_n^{(n)} + \mu_n^{(n)}} = \frac{1 + s_0}{2 + s_0},$$

has an eigenspace consisting of all of the “pure homozygote” vectors, that is those with nonzero entries only in the first and last component, while the second,

$$\lambda_{\text{het}} = \frac{\lambda_1^{(n)}}{\lambda_1^{(n)} + \mu_1^{(n)}} = \dots = \frac{\lambda_{n-1}^{(n)}}{\lambda_{n-1}^{(n)} + \mu_{n-1}^{(n)}} = \frac{1 + s_{\text{max}}}{2 + s_{\text{max}}},$$

has an eigenspace consisting of the “pure heterozygote” vectors (those with zero first and last components).

The eigenvalues of P were calculated by [Schensted \(1958\)](#); we present here an elaboration of her argument. Consider the map $T: \mathbb{R}[X] \rightarrow \mathbb{R}^{n+1}$ that takes a polynomial $p(X)$ with real coefficients to the vector

$$Tp = \begin{bmatrix} p(0) \\ p(1) \\ \vdots \\ p(n) \end{bmatrix}.$$

This map is clearly linear, and if we restrict the domain to polynomials of degree at most n , it becomes an isomorphism (since any $n + 1$ points define a polynomial of degree at most n). The sequence of polynomials $\left\{ \binom{X}{0}, \binom{X}{1}, \binom{X}{2}, \dots \right\}$ forms a basis for $\mathbb{R}[X]$; moreover, if we truncate the sequence at $\binom{X}{n}$, we get a basis for the subspace of polynomials of degree at most n , which we can pass through the isomorphism T to obtain a basis for \mathbb{R}^{n+1} .

Looking carefully at the binomial coefficient identity

$$\sum_{\alpha=0}^n \underbrace{\frac{\binom{2i}{\alpha} \binom{2n-2i}{n-\alpha}}{\binom{2n}{n}}}_{P_{i\alpha}} \underbrace{\binom{\alpha}{\beta}}_{T\left(\begin{smallmatrix} X \\ \beta \end{smallmatrix}\right)_{\alpha}} = \frac{\binom{n}{\beta}}{\binom{2n}{\beta}} \underbrace{\binom{2i}{\beta}}_{T\left(\begin{smallmatrix} 2X \\ \beta \end{smallmatrix}\right)_i}$$

(proven in appendix A), we see that it is the i th component of the equation

$$\begin{aligned} PT\left(\begin{smallmatrix} X \\ \beta \end{smallmatrix}\right) &= \frac{\binom{n}{\beta}}{\binom{2n}{\beta}} T\left(\begin{smallmatrix} 2X \\ \beta \end{smallmatrix}\right) \\ &= \frac{\binom{n}{\beta}}{\binom{2n}{\beta}} \sum_{\alpha=0}^{\beta} b_{\beta\alpha} T\left(\begin{smallmatrix} X \\ \alpha \end{smallmatrix}\right), \end{aligned} \tag{5}$$

where the second equality follows from the fact that $\binom{2X}{\beta}$ is a polynomial of degree β , and so must itself be a linear combination of $\binom{X}{\alpha}$ with $\alpha \leq \beta$.

Equation (5) implies that P transformed to the $\{T\left(\begin{smallmatrix} X \\ \beta \end{smallmatrix}\right)\}_{0 \leq \beta \leq n}$ basis is upper triangular, with diagonal elements $b_{ii} \binom{n}{i} / \binom{2n}{i}$. This of course means that the eigenvalues χ_0, \dots, χ_n of P are exactly these diagonal elements. To calculate b_{ii} , we note that the coefficient of X^i in the polynomial $\binom{X}{i}$ is $1/i!$, while in the polynomial $\binom{2X}{i}$ it is $2^i/i!$. None of the $\binom{X}{\alpha}$ for $\alpha < i$ can contribute an X^i term, so it must be the case that $b_{ii} = 2^i$. Thus we have that

$$\chi_i = 2^i \frac{\binom{n}{i}}{\binom{2n}{i}}.$$

The ratio χ_{i+1}/χ_i is less than or equal to 1, showing the sequence of eigenvalues is non-increasing. The first three values are $\chi_0 = 1$, $\chi_1 = 1$, and $\chi_2 = \frac{2(n-1)}{2n-1}$.

Finally we get to the eigenvalues of M itself. In general, of course, it is not the case that the eigenvalues of a product of matrices are the products of the eigenvalues of the factors; however, if a pair of eigenvalues of the two factors share a common eigenspace, then their product is an eigenvalue of the product. Even this condition will not be true for all the pairs of eigenvalues of Λ and P we are interested in, but it *is* true for λ_{hom}

(an eigenvalue of Λ) and 1 (an eigenvalue of P). Recall that the λ_{hom} -eigenspace of Λ consists of the “pure homozygote” vectors: these are also exactly the *left* eigenvectors of P corresponding to the eigenvalue 1 (this can be seen from the fact that the first and last rows of P have a one on the diagonal and zeros everywhere else). Therefore $2\lambda_{\text{hom}}$ is an eigenvalue of M (with multiplicity 2). But since $s_0 < 0$, this eigenvalue is always less than 1, and can never be the eigenvalue we are looking for.

Let $\chi \neq 1$ be another eigenvalue of P and x a corresponding *left* eigenvector. We now show that $2\lambda_{\text{het}}\chi$ is an eigenvalue of M . Unfortunately, it will not be so easy this time: x does not have to lie in the λ_{het} -eigenspace of Λ . However, we can split x into a pure homozygote component x_{hom} , which is the component of x in the λ_{hom} -eigenspace of Λ , and a pure heterozygote component $x - x_{\text{hom}}$, which lies in the λ_{het} -eigenspace. Of these three vectors, x is a (left) eigenvector of P , and $x - x_{\text{hom}}$ is an eigenvector of Λ , but x_{hom} is an eigenvector of *both*: this enables us to construct a vector that will be a left eigenvector of ΛP corresponding to $\lambda_{\text{het}}\chi$, namely $\lambda_{\text{het}}(1 - \chi)x_{\text{hom}} + (\lambda_{\text{hom}} - \chi\lambda_{\text{het}})(x - x_{\text{hom}})$. To see it is an eigenvector, note that

$$\begin{aligned} P^T \Lambda (\lambda_{\text{het}}(1 - \chi)x_{\text{hom}} + (\lambda_{\text{hom}} - \chi\lambda_{\text{het}})(x - x_{\text{hom}})) \\ &= P^T (\lambda_{\text{het}}\lambda_{\text{hom}}(1 - \chi)x_{\text{hom}} + (\lambda_{\text{hom}}\lambda_{\text{het}} - \chi\lambda_{\text{het}}^2)(x - x_{\text{hom}})) \\ &= P^T (\lambda_{\text{het}}\chi(\lambda_{\text{het}} - \lambda_{\text{hom}})x_{\text{hom}} + (\lambda_{\text{hom}}\lambda_{\text{het}} - \chi\lambda_{\text{het}}^2)x) \\ &= \lambda_{\text{het}}\chi(\lambda_{\text{het}} - \lambda_{\text{hom}})x_{\text{hom}} + (\lambda_{\text{hom}}\lambda_{\text{het}}\chi - \chi^2\lambda_{\text{het}}^2)x \\ &= \lambda_{\text{het}}\chi ((\lambda_{\text{het}} - \lambda_{\text{hom}})x_{\text{hom}} + (\lambda_{\text{hom}} - \chi\lambda_{\text{het}})x) \\ &= \lambda_{\text{het}}\chi (\lambda_{\text{het}}(1 - \chi)x_{\text{hom}} + (\lambda_{\text{hom}} - \chi\lambda_{\text{het}})(x - x_{\text{hom}})), \end{aligned}$$

where the first equality uses that x_{hom} and $x - x_{\text{hom}}$ are eigenvectors of Λ and the third uses that x and x_{hom} are left eigenvectors of P .

Since χ_2 is the largest of the remaining eigenvalues of P , to check if there is at least one eigenvalue of M greater than 1 we need only check if $2\lambda_{\text{het}}\chi_2 > 1$. We have that

$$2\lambda_{\text{het}}\chi_2 = 2 \frac{(1 + s_{\text{max}})2(n-1)}{(2 + s_{\text{max}})(2n-1)},$$

which is greater than one if and only if $s_{\text{max}} > \frac{2}{2n-3}$. \square

Theorem 2. *In the model with random replication and a dominant fitness function, the establishment probability (and therefore the rescue probability) is nonzero if and only if*

$$s_{\text{max}} > \frac{4n}{2n^2 - 3n - 1}.$$

In addition, in the long term limit, the heterozygote types are all equally abundant.

Proof. Let M once again be the matrix the i, j element of which is the expected number of daughter cells with j mutant plasmids from a cell with i mutant plasmids. In the model with random replication and a dominant fitness function, the elements of this matrix have the form

$$\begin{aligned} m_{ij} &= \frac{\lambda_i^{(n)}}{\lambda_i^{(n)} + \mu_i^{(n)}} \sum_{k=0}^n 2^{\delta_{jk}} P(i \rightarrow \{j, k\}) \\ &= 2 \frac{\lambda_i^{(n)}}{\lambda_i^{(n)} + \mu_i^{(n)}} \sum_{k=0}^n \frac{\binom{j+k}{j} \binom{2n-j-k}{n-j}}{\binom{2n}{n}} \frac{\binom{2n-j-k-1}{n-j-k+i} \binom{j+k-1}{j+k-i}}{\binom{2n-1}{n}}. \end{aligned}$$

As in the proof for the regular replication case, we use general branching process theory (Sewastjanow, 1975, Satz V.5) to determine that there is a nonzero probability of establishment if and only if M has an eigenvalue greater than one.

We shall call the probability of reproduction for homozygote cells

$$\lambda_{\text{hom}} = \frac{\lambda_0^{(n)}}{\lambda_0^{(n)} + \mu_0^{(n)}} = \frac{\lambda_n^{(n)}}{\lambda_n^{(n)} + \mu_n^{(n)}} = \frac{1 + s_0}{2 + s_0},$$

and that for heterozygote cells

$$\lambda_{\text{het}} = \frac{\lambda_1^{(n)}}{\lambda_1^{(n)} + \mu_1^{(n)}} = \dots = \frac{\lambda_{n-1}^{(n)}}{\lambda_{n-1}^{(n)} + \mu_{n-1}^{(n)}} = \frac{1 + s_{\max}}{2 + s_{\max}}.$$

M has the structure

$$\begin{bmatrix} 2\lambda_{\text{hom}} & 0 & 0 \\ 2\lambda_{\text{hom}}A & 2\lambda_{\text{het}}B & 2\lambda_{\text{hom}}C \\ 0 & 0 & 2\lambda_{\text{hom}} \end{bmatrix},$$

where B is a $(n-1) \times (n-1)$ matrix and A and C are $(n-1)$ component column vectors.

(That the first and last rows are almost all zeros can be seen either by substitution into the expression for m_{ij} above—one of the binomial coefficients will be zero unless $i = j$ —or from the fact that a homozygote cell will only have homozygote daughters.)

Then $2\lambda_{\text{hom}}$ is an eigenvalue of M ; this eigenvalue, however, is always less than one, since $s_0 < 0$. Every other eigenvalue of M will be $2\lambda_{\text{het}}$ times an eigenvalue of B , and its corresponding left eigenvectors are of the form $[x_0 \ v \ x_1]$, where v is a left eigenvector of B .

We now show that $(n-1)(2n+1)/(n+1)(2n-1)$ is such an eigenvalue of B with a corresponding left eigenvector v with all components equal to 1 (the identification of this eigenvalue is due to [Novick and Hoppensteadt \(1978\)](#)). The j th element of $B^T v$ is

$$\sum_{i=1}^{n-1} \sum_{k=0}^n \frac{\binom{j+k}{j} \binom{2n-j-k}{n-j}}{\binom{2n}{n}} \frac{\binom{2n-j-k-1}{n-j-k+i} \binom{j+k-1}{j+k-i}}{\binom{2n-1}{n}},$$

which by equation (A2) in appendix A is equal to $(n-1)(2n+1)/(n+1)(2n-1)$; so $B^T v = (n-1)(2n+1)/(n+1)(2n-1)v$. This must also be the largest eigenvalue of B : B is a nonnegative irreducible matrix (since a cell with $1 \leq i \leq n-1$ mutant plasmids can produce a descendant cell with $1 \leq k \leq n-1$ mutant plasmids after at most $|i-k|$ generations, B^{n-1} is a positive matrix), so by the Perron-Frobenius Theorem, its only eigenvalue with a strictly positive eigenvector is the largest eigenvalue.

Thus $2\lambda_{\text{het}}(n-1)(2n+1)/(n+1)(2n-1)$ is the largest eigenvalue of M (other than possibly $2\lambda_{\text{hom}}$) and its corresponding left eigenvector has all heterozygote types equally abundant. By simple rearrangement, we have that

$$2\lambda_{\text{het}} \frac{(n-1)(2n+1)}{(n+1)(2n-1)} = 2 \frac{(1+s_{\text{max}})(n-1)(2n+1)}{(2+s_{\text{max}})(n+1)(2n-1)}$$

is greater than one if and only if $s_{\text{max}} > \frac{4n}{2n^2-3n-1}$. \square

6 The final distribution

If establishment occurs, we can ask what the relative abundances of cells with different numbers of mutant plasmids will be. The distribution of the number of mutant plasmids in a cell in the infinite time limit can be found from the matrix M the i, j element of which is the expected number of daughters with j mutant plasmids from a cell with i mutant plasmids. If establishment is possible ($P_{\text{est}} > 0$) then M has a dominant eigenvalue, and a corresponding left eigenvector to that eigenvalue gives the distribution of cell types in the long run (Mode, 1971, Theorem 4.1). This eigenvector can be found numerically.

Examples of the final stable distribution of the number of mutant plasmids per cell under the four models considered (regular and random replication and dominant and peaked fitness functions) are shown in Figure 5. Even after establishment, significant numbers of homozygote cells are present in the population, because of their constant replenishment by segregation. Heterozygotes are more abundant under regular replication than random replication and with the dominant fitness function than with the peaked fitness function, because heterozygosity is lost to segregation at a lower rate in those conditions. Note that in the model with random replication and dominant fitness function, all heterozygote cell types are equally abundant in the long-term limit (as shown by Theorem 2).

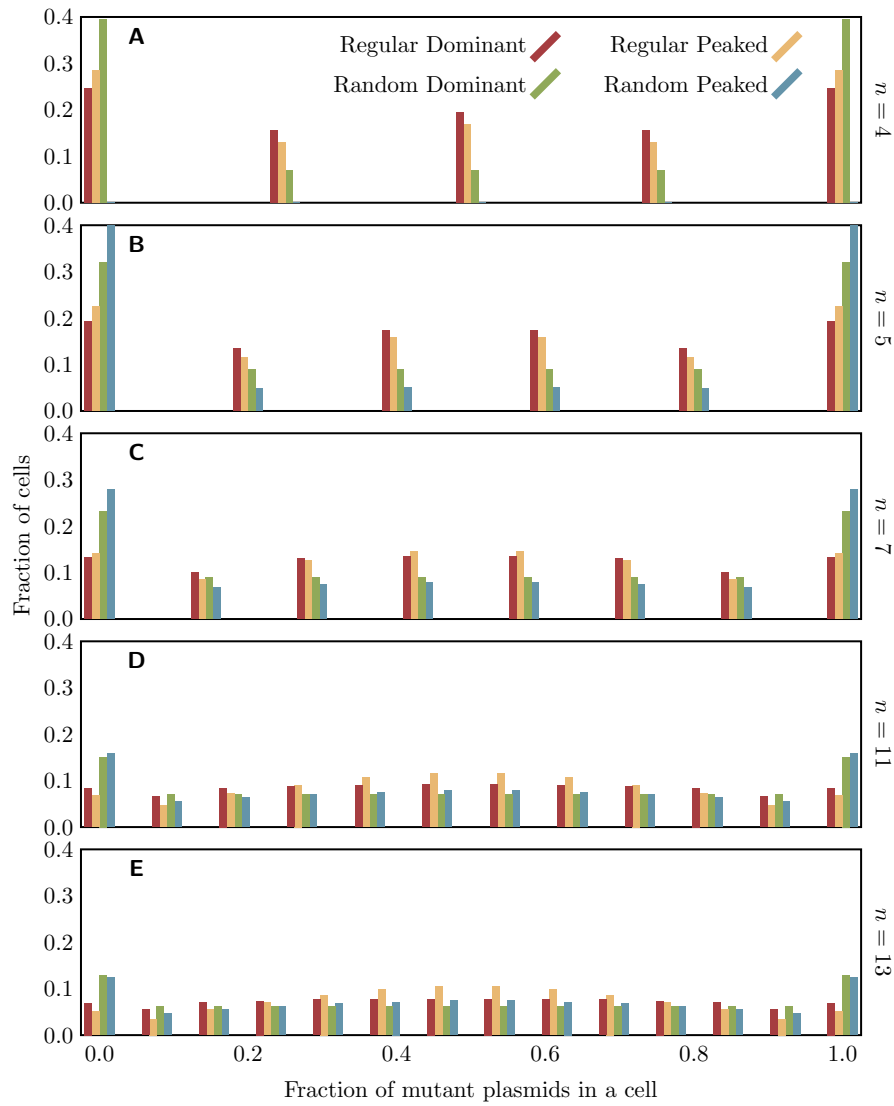


Fig. 5 The final distribution of the proportion of mutant plasmids in a cell for plasmids of given copy numbers, under each of the four combinations of replication assumption and fitness function considered, and with $s_0 = -0.1$ and $s_{\max} = 1.0$. The eigenvector corresponding to the leading eigenvalue of M was calculated numerically; note that for $n = 4$ with random replication and the peaked fitness function, the population goes extinct with probability one, and there is no well-defined final distribution.

7 Cointegration

As we have seen, segregation plays an important role in limiting the possibility of the population being rescued. One process which alters the rate of segregation is the fusion of plasmids to form cointegrates. When plasmids fuse, the resulting multimer has multiple copies of the gene of interest and multiple copies of the plasmid backbone, including multiple copies of the copy number control system (so that an m -mer still “counts against” the copy number m times, even though it is only one molecule). Cointegration can eliminate the possibility of losing heterozygosity to segregation in the descendants of cells in which it occurs (experimentally observed by [Hülter et al, 2020](#); [Garona et al, 2021](#)): if plasmids with distinct alleles fuse, the resulting cointegrate provides the heterozygote phenotype on its own and can no longer be separated by segregation (unless the cointegrate is resolved into independent plasmids again). But conversely, if plasmids with identical alleles fuse, the number of independently segregating units in the cell is reduced, and the rate of loss will increase. We now present a simple extension of the model to include the process of cointegration.

7.1 The cointegration model

In a model which incorporates cointegration, there is a much larger number of cell types possible, since not only are there multiple loci that may have one allele or the other, but also these loci may be distributed among plasmids in a variety of ways. In the following, a cell type will be written in the form $(A)(B)(AB)$, where each letter represents a locus carrying one of the two alleles, A and B, that must be maintained for the population to survive, loci within a pair of parentheses are on the same plasmid, and the total number of loci is fixed at the copy number n . Each cell type \mathbf{a} then has a death rate $\mu_{\mathbf{a}}$, which will be fixed to 1 as before, and a reproduction rate $\lambda_{\mathbf{a}}$, which is determined by the number of mutant and wild-type alleles in the cell (irrespective of their distribution on different plasmids); for the results here, we will use only the

dominant fitness function. Incorporating cointegration into our model requires slight changes to the reproduction probabilities to account for the greater number of cell types, and adds a new process of cointegration which needs to be modelled.

7.1.1 Reproduction

The reproduction probabilities are determined by the regular replication model used previously. The probability that a dividing cell of type \mathbf{a} has daughters of types \mathbf{b} and \mathbf{c} will be denoted $P(\mathbf{a} \rightarrow \{\mathbf{b}, \mathbf{c}\})$. All plasmids are duplicated and then divided into two subsets having the same copy number *of the copy number control system*. This last part is important because now plasmids contribute differently to the total copy number. This means that some assortments of plasmids into daughter cells are no longer possible: for example, an $(\mathbf{AAA})(\mathbf{B})$ cell can only have daughters identical to itself, since a $(\mathbf{B})(\mathbf{B})$ cell has too few and an $(\mathbf{AAA})(\mathbf{AAA})$ cell too many copies of the copy number control system. The probabilities $P(\mathbf{a} \rightarrow \{\mathbf{b}, \mathbf{c}\})$ are calculated by enumeration of the possible pairs of daughter cells and the number of ways to produce each one.

7.1.2 Cointegration

The cointegration rate is divided into two parts: the overall probability of a cointegration event occurring in a given cell, which is the biological part, and the probability, conditional on cointegration occurring, that two particular plasmids fuse, which is purely combinatorial.

The probability, conditional on a fusion occurring, that it takes the cell from type \mathbf{a} to type \mathbf{a}' is denoted $P(\mathbf{a} \Rightarrow \mathbf{a}')$. This probability is the fraction of all pairs of plasmids in the \mathbf{a} -cell that are pairs of the two types that need to fuse for this transition to occur. For example, $P((\mathbf{A})(\mathbf{A})(\mathbf{B}) \Rightarrow (\mathbf{A})(\mathbf{AB})) = \frac{2}{3}$ (there are 3 pairs of plasmids in the mother cell, of which two are $\{(\mathbf{A}), (\mathbf{B})\}$ pairs), and $P((\mathbf{AA})(\mathbf{A})(\mathbf{B}) \Rightarrow (\mathbf{A})(\mathbf{A})(\mathbf{AB})) = 0$.

The probability of cointegration occurring in one generation in cells of type \mathbf{a} is denoted $\kappa_{\mathbf{a}}$. In the model presented here, this will be simply a constant κ , except in cells that have only a single large multimer, where it is 0 because no further cointegration can take place. We assume that at most one cointegration event occurs in each cell cycle.

7.2 Analysis and results

Combining the reproduction and resolution models, the extinction probability $Q_{\mathbf{a}}$ of type \mathbf{a} is given by

$$Q_{\mathbf{a}} = \frac{\mu_{\mathbf{a}}}{\lambda_{\mathbf{a}} + \mu_{\mathbf{a}}} + \frac{\lambda_{\mathbf{a}}}{\lambda_{\mathbf{a}} + \mu_{\mathbf{a}}} \sum_{\{\mathbf{b}, \mathbf{c}\}} \left(\kappa_{\mathbf{a}} \sum_{\mathbf{a}'} P(\mathbf{a} \Rightarrow \mathbf{a}') P(\mathbf{a}' \rightarrow \{\mathbf{b}, \mathbf{c}\}) Q_{\mathbf{b}} Q_{\mathbf{c}} + (1 - \kappa_{\mathbf{a}}) P(\mathbf{a} \rightarrow \{\mathbf{b}, \mathbf{c}\}) Q_{\mathbf{b}} Q_{\mathbf{c}} \right),$$

where the sums are taken over all types or unordered pairs of types. Note that this expression is equally compatible with selection happening before or after the fusion of plasmid copies, since the values of $\lambda_{\mathbf{a}}$ and $\mu_{\mathbf{a}}$ are determined only by the alleles present and not their distribution between plasmid copies. Since we are using the regular replication model and make the assumption that the mutation first appears in a cell with only monomers present, the establishment probability is then given by

$$P_{\text{est}} = 1 - Q_{(\mathbf{A})^{n-1}(\mathbf{B})}.$$

Extinction probabilities for the cointegrate model are shown in Figure 6. An important qualitative difference is that for non-zero κ , the establishment probability no longer exhibits a threshold effect: establishment is theoretically possible for any value of s_{max} .

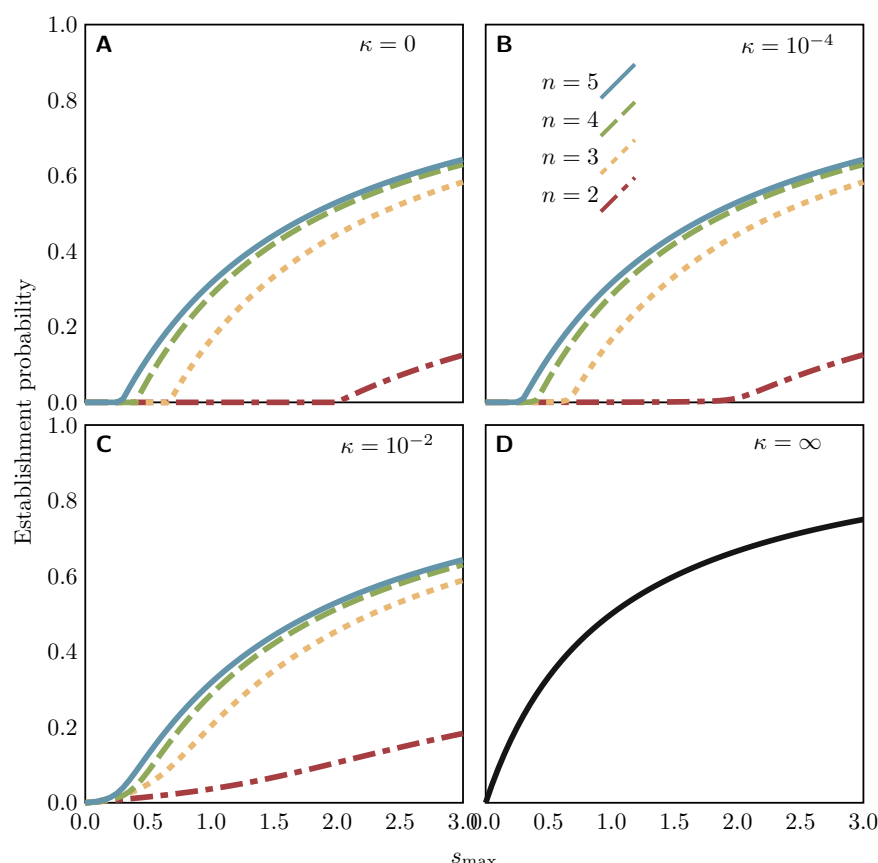


Fig. 6 Establishment probabilities of rescue mutations on possibly-cointegrating plasmids of a given copy number, for varying values of the heterozygote fitness s_{\max} and the cointegration probability κ (for all cases, $s_0 = -0.1$). The $\kappa = \infty$ panel corresponds to a model in which the mutation is immediately followed by the fusion of all plasmids into a single multimer, and thus is the theoretical maximum establishment probability.

8 Discussion

Multicopy plasmids are a frequent component of bacterial genomes, where they may be cryptic passengers or play an important role in bacterial adaptation and evolution. One way in which they can play such a role is through their inherently polyploid nature: loci on a multicopy plasmid can exhibit heterozygosity which is not usually possible for loci on the for many species haploid chromosome. Here we have examined

one scenario in which heterozygosity of plasmids may provide an adaptive advantage to their hosts, by enabling rapid evolution in scenarios of heterozygote advantage.

It is apparent that a key factor in determining the fate of heterozygosity on plasmids is plasmid segregation at host cell division. The two versions of the plasmid originally existing together in a single host are eventually separated by random segregation into distinct hosts. The effects of segregation on the establishment of novel alleles located on multicopy plasmids has already been explored both in models (Santner and Uecker, 2020; Garoña et al, 2023) and experiments (Ilhan et al, 2018; Garoña et al, 2023), showing that segregation reduces the establishment probability of adaptive alleles on multicopy plasmids. Here, we have shown that the effect of segregation is even more drastic in scenarios of heterozygote advantage: the constant loss of heterozygosity to segregation puts strong limits on the probability of successful adaptation unless the fitness of heterozygotes is extremely large. The condition that we derived for the required fitness of heterozygotes more generally holds for the maintenance of heterozygosity on multicopy plasmids, even outside a rescue scenario. While we have phrased our model in terms of wild-type and mutant plasmids (i.e. variants of the same plasmid), the results carry over to the establishment and the maintenance of two different incompatible multicopy plasmids in the same cell, provided they share the same replication and segregation mechanism.

The importance of plasmid-mediated heterozygosity in allowing populations to escape from fitness trade-offs were previously explored by Rodriguez-Beltran et al (2018) in both experiments and models. In a system where different plasmid alleles provide resistance to different antibiotics, they find that fluctuating selection is capable of maintaining heterozygosity, where intermediate antibiotic concentrations maintained heterozygosity the longest (note that unlike in our model, there is always selection for only one plasmid variant at a time). Their model is simpler than ours in one important respect: they bundle all heterozygotes into a single compartment

of their ODE model, ignoring the precise number of mutant and wild-type plasmids. To estimate the rate of segregation, [Rodriguez-Beltran et al \(2018\)](#) assume that all heterozygote cells have half-and-half mutant and wild-type composition, so that the probability of having a homozygote daughter cell is 2^{1-n} . Transferred over to the branching process context, this means that the establishment probability of a mutant on a plasmid in combination treatment is positive if and only if

$$s_{\max} > \frac{1}{2^{n-2} - 1},$$

which is a much weaker condition than our Theorem 1, provided that $n > 3$, since ignoring the unbalanced heterozygotes overestimates the stability of heterozygotes. Application of our model to fluctuating selection remains to be done.

The important role segregation plays in the fate of heterozygosity on multicopy plasmids contrasts strongly with the state of affairs on a haploid bacterial chromosome, where true heterozygosity is impossible, so that our two focal alleles would have to either replace each other completely or be located at different loci and not subject to segregation. Cells in rapidly dividing populations will contain multiple copies of the chromosome and thus be effectively polyploid; although this exposes them to the effects of dominance ([Sun et al, 2018](#)), however, their segregation is not random and heterozygosity cannot be maintained long-term. There is also a strong contrast to the role of segregation in the system one normally thinks of when considering heterozygosity, sexually reproducing diploid organisms. Homologous chromosomes in diploids segregate at meiosis, but heterozygosity is maintained by sexual reproduction; unlike with multicopy plasmids, no selection for heterozygotes is required to maintain heterozygosity at Hardy-Weinberg proportions in the population indefinitely, provided homozygous individuals can thrive. However, if homozygous individuals have fitness less than one as in our rescue scenario and heterozygotes are required for population

persistence, there is also a threshold on the heterozygote fitness below which the population cannot survive. Yet, unlike for multicopy plasmids, this threshold depends on the fitness of homozygotes. The closest analogue in diploids to the role of segregation in multicopy plasmids is inbreeding, which also reduces heterozygosity in the population without altering allele frequencies. Segregation is in a sense a stronger force in reducing heterozygosity than any degree of partial selfing, since the equilibrium heterozygosity is zero in the absence of selection for the heterozygote (unlike partial selfing, which produces a nonzero equilibrium), but less strong than complete selfing; this can be seen by comparing the $n = 2$ case (when segregation has the largest effect on multicopy plasmids, and also the case most directly comparable to diploids), where $2/3$ of the daughter cells of heterozygotes are heterozygotes, to a selfing diploid, where only $1/2$ of the offspring of heterozygotes are heterozygotes. Nevertheless, constraints on heterozygote fitness are stronger for plasmid-mediated rescue of bacterial populations than for rescue of selfing populations: division of a heterozygote cell into two homozygote cells leads to loss of a heterozygote from the population, while offspring reproduction by a selfing heterozygote leaves the heterozygote parent individual intact. This means that the rate of reproduction $1 + s_{\max}$ only needs to be greater than two for the selfing individual rather than three as in the bacterial population (the same holds in a random mating population, in which homozygotes are lethal or infertile). The nonplasmid systems that most closely resemble the genetics of multicopy plasmids are the organelle (mitochondrial and plastid) genomes of eukaryotes (Birky, 1983) and the macronuclei of ciliate protists (Allen and Nanney, 1958)—indeed, it was in the latter context that Schensted (1958) developed a model of chromosomal segregation which is equivalent our regular replication model for plasmids and first determined the eigenvalues of the matrix we refer to as P in the proof of Theorem 1.

We have also seen that the mechanism or mode of replication has a substantial effect on the fate of a novel mutation in heterozygote advantage scenarios. In the regular replication model, there are no additional stochastic effects in replication, while in the random replication model the rich-get-richer property of the replication process is an additional mechanism reducing heterozygosity in the population. Our understanding of the biological mechanisms of plasmid replication would seem to suggest that random replication is a more realistic model of the replication process of most plasmids; however, some empirical studies have found the regular replication model to be a better fit to the data (Garona et al, 2023). Here we have followed Novick and Hoppensteadt (1978) in presenting both models to permit a comparison.

Our model for segregation could, in principle, also be made more realistic. We adopt a fairly strict assumption that every daughter cell receives exactly the same number of copies of the plasmid, meaning that every cell in the model can be assumed to have exactly n total copies. For those low-copy-number plasmids which rely on an active partitioning system to separate plasmids at host cell division, this is a fairly reasonable assumption. But plasmids with higher copy numbers which don't rely on active partitioning cannot guarantee exact (or near exact) partitioning of plasmid copies. For these, a model with random segregation, in which the number of plasmids received by a daughter cell is not fixed (it would be given by a binomial distribution) and in the daughter cell plasmids are replicated from however many copies the cell starts with up to $2n$ for the next host generation, might be more suitable. However, because of clustering and other nonrandom spatial distribution of plasmids within the cell, the random model might not reflect the full biological reality either. Indeed, clustering could make the segregation process not neutral between mutant and wild-type plasmids. The model we have adopted seems a reasonable simplification to make in the pursuit of analytical tractability. Moreover, a combination of modeling and experiments of allele dynamics under directional selection has recently shown that the

adopted model captures the key results obtained from the experiments well (Garona et al, 2023), confirming its usefulness.

Because we are interested in the effect of properties of the plasmid lifecycle on the evolution of the host, we have kept properties such as the plasmid copy number fixed. While this may be a realistic assumption over the short timescale of the rescue process, over the longer term (for example in the rescued population) these properties are also subject to evolution. Given the important role of copy number in the fate of bacterial populations depending on plasmid-mediated heterozygosity, the copy number might well be under selection. Some experimental studies have shown rapid evolution of plasmid copy number in response to selection on plasmid-borne traits in other contexts (Dimitriu et al, 2021; San Millan et al, 2015; San Millan et al, 2016). A mutation in the plasmid replication control system that produces a new plasmid type compatible with the existing plasmid—a plasmid speciation event—would also confer a selective advantage if it allowed the mutant and wild-type alleles to exist on plasmids of distinct compatible types, no longer sharing a common copy number between them, and no longer subject to loss of heterozygosity by segregation. Rapid plasmid speciation has also been observed in experiments (Santos-Lopez et al, 2016).

Another change in the plasmid that would affect the maintenance of plasmid-mediated heterozygosity, which we have modelled, is the formation of plasmid cointegrates (Hülter et al, 2020; Garona et al, 2021). Our results have shown that stable multimerization of plasmids can increase the probability of successful establishment of the novel allele in a heterozygote advantage scenario. Our model of cointegration is deliberately simple and limited, as it is intended as a short demonstration of the importance of cointegration in the heterozygote advantage scenario and not a deep exploration of the cointegration process. Nonetheless, many aspects of the cointegration process not explored here, particularly the precise behaviour of multimers during replication and segregation and the resolution of cointegrates into smaller multimers

or monomers, are of no doubt of importance to the fate of alleles on cointegrates and merit further investigation. We especially have not included consequences of multimerization on the stability of plasmid inheritance, which has been examined in a model by [Summers et al \(1993\)](#). The possibility of evolution of a multi-drug resistance plasmid from two incompatible single-resistance plasmids has been previously demonstrated by [Condit and Levin \(1990\)](#) both experimentally and in a model; it is left open in their model whether the plasmid carrying both resistance genes has been formed through plasmid fusion or exchange of DNA between plasmids. The model by [Condit and Levin \(1990\)](#) does not account for the plasmid copy number. A detailed model that combines elements of their and our models and allows for dependence of recombination on the plasmid copy number or the ratio of variant frequencies could be interesting to further understand the evolution of multidrug resistance on plasmids.

Plasmids are pervasive in natural populations of bacteria, and play an important role in adaptation and evolution of bacteria. Understanding bacterial evolution therefore involves a good understanding of the population genetics of plasmids. The properties of plasmids as independently replicating units within their hosts make them an intriguing genetic system, with complexities not present for haploid or diploid chromosomes. We have explored one aspect of plasmid genetics—heterozygosity on multicopy plasmids—in the context of evolutionary rescue. We have shown that, as intuitively expected, the maintenance of heterozygosity and thus rescue is impossible below a threshold copy number for a given heterozygote fitness. Going beyond this intuition, our concise criterion in terms of the plasmid copy number and the fitness of heterozygous cells quantifies the conditions for the persistence of plasmid-mediated heterozygosity through a heterozygosity advantage, and thus contributes to a population genetics theory of bacterial evolution.

Acknowledgments. Acknowledgments: The authors thank Félix Geoffroy and Mario Santer for helpful discussions.

Statements and Declarations.

- Funding: This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) — project number 400993799 (project 3 within the Research Training Group 2501 “Translational Evolutionary Research”, <https://gepris.dfg.de/gepris/projekt/400993799>). I.D. is a member of the International Max Planck Research School for Evolutionary Biology and gratefully acknowledges the benefits provided by the program.
- Competing interests: The authors have no competing interests to declare.
- Availability of data and materials: No new data was generated in the course of this theoretical work.
- Code availability: The code used to numerically solve the models and produce the figures in the paper is available from as supplementary material.
- Authors’ contributions: ID: Conceptualization, formal analysis, methodology, software, writing—original draft; HU: Conceptualization, methodology, supervision, writing—review & editing.

References

- Allen SL, Nanney DL (1958) An analysis of nuclear differentiation in the selfers of tetrahymena. *The American Naturalist* 92(864):139–160. <https://doi.org/10.1086/282022>
- Anda M, Ohtsubo Y, Okubo T, et al (2015) Bacterial clade with the ribosomal RNA operon on a small plasmid rather than the chromosome. *Proceedings of the National Academy of Sciences* 112(46):14343–14347. <https://doi.org/10.1073/pnas.1514326112>
- Beijersbergen A, Smith SJ, Hooykaas PJ (1994) Localization and topology of VirB proteins of *Agrobacterium tumefaciens*. *Plasmid* 32(2):212–218. [https://doi.org/10.1016/0032-0732\(94\)90038-8](https://doi.org/10.1016/0032-0732(94)90038-8)

[1006/plas.1994.1057](https://doi.org/10.1006/plas.1994.1057)

- Birky CWJr. (1983) Relaxed cellular controls and organelle heredity. *Science* 222(4623):468–475. <https://doi.org/10.1126/science.6353578>
- Carattoli A (2013) Plasmids and the spread of resistance. *International Journal of Medical Microbiology* 303:298–304. <https://doi.org/10.1016/j.ijmm.2013.02.001>
- Coluzzi C, Garcillán-Barcia MP, de la Cruz F, et al (2022) Evolution of plasmid mobility: Origin and fate of conjugative and nonconjugative plasmids. *Molecular Biology and Evolution* 39(6):msac115. <https://doi.org/10.1093/molbev/msac115>
- Condit R, Levin BR (1990) The evolution of plasmids carrying multiple resistance genes: The role of segregation, transposition, and homologous recombination. *The American Naturalist* 135(4):573–596. <https://doi.org/10.1086/285063>
- Cullum J, Broda P (1979) Rate of segregation due to plasmid incompatibility. *Genetical Research* 33(1):61–79. <https://doi.org/10.1017/s0016672300018176>
- Dimitriu T, Matthews AC, Buckling A (2021) Increased copy number couples the evolution of plasmid horizontal transmission and plasmid-encoded antibiotic resistance. *Proceedings of the National Academy of Sciences* 118(31). <https://doi.org/10.1073/pnas.2107818118>
- Eggenberger F, Pólya G (1923) Über die statistik verketteter vorgänge. *Zeitschrift für Angewandte Mathematik und Mechanik* 3(4):279–289. <https://doi.org/10.1002/zamm.19230030407>
- Falkow S (1975) *Infectious Multiple Drug Resistance*. Pion
- Gama JA, Zilhão R, Dionisio F (2018) Impact of plasmid interactions with the chromosome and other plasmids on the spread of antibiotic resistance. *Plasmid* 99:82–88.

<https://doi.org/10.1016/j.plasmid.2018.09.009>

Garcillán-Barcia MP, Alvarado A, de la Cruz F (2011) Identification of bacterial plasmids based on mobility and plasmid population biology. *FEMS Microbiology Reviews* 35(5):936–956. <https://doi.org/10.1111/j.1574-6976.2011.00291.x>

Garoña A, Hülter NF, Picazo DR, et al (2021) Segregational drift constrains the evolutionary rate of prokaryotic plasmids. *Molecular Biology and Evolution* <https://doi.org/10.1093/molbev/msab283>

Garoña A, Santer M, Hülter NF, et al (2023) Segregational drift hinders the evolution of antibiotic resistance on polyploid replicons. *PLOS Genetics* 19(8):e1010829. <https://doi.org/10.1371/journal.pgen.1010829>

Gomulkiewicz R, Holt RD (1995) When does evolution by natural selection prevent extinction? *Evolution* 49(1):201–207. <https://doi.org/10.1111/j.1558-5646.1995.tb05971.x>

Graham RL, Knuth DE, Patashnik O (1994) *Concrete Mathematics*, 2nd edn. Addison-Wesley

Haldane JBS (1927) A mathematical theory of natural and artificial selection, part v: Selection and mutation. *Mathematical Proceedings of the Cambridge Philosophical Society* 23(7):838–844. <https://doi.org/10.1017/s0305004100015644>

Halleran AD, Flores-Bautista E, Murray RM (2019) Quantitative characterization of random partitioning in the evolution of plasmid-encoded traits. *BioRxiv preprint*, <https://doi.org/10.1101/594879>

Hernandez-Beltran JCR, Miró Pina V, Siri-Jégousse A, et al (2022) Segregational instability of multicopy plasmids: A population genetics approach. *Ecology and*

- Evolution 12(12). <https://doi.org/10.1002/ece3.9469>
- Hülter NF, Wein T, Effe J, et al (2020) Intracellular competitions reveal determinants of plasmid evolutionary success. *Frontiers in Microbiology* 11. <https://doi.org/10.3389/fmicb.2020.02062>
- Ilhan J, Kupczok A, Woehle C, et al (2018) Segregational drift and the interplay between plasmid copy number and evolvability. *Molecular Biology and Evolution* 36(3):472–486. <https://doi.org/10.1093/molbev/msy225>
- Ishii K, Hashimoto-Gotoh T, Matsubara K (1978) Random replication and random assortment model for plasmid incompatibility in bacteria. *Plasmid* 1(4):435–445. [https://doi.org/10.1016/0147-619x\(78\)90002-1](https://doi.org/10.1016/0147-619x(78)90002-1)
- Mode CJ (1971) Multitype Branching Processes. *Modern Analytic and Computational Methods in Science and Mathematics*, American Elsevier Publishing Company, Inc.
- Novick RP, Hoppensteadt F (1978) On plasmid incompatibility. *Plasmid* 1(4):421–434. [https://doi.org/10.1016/0147-619x\(78\)90001-x](https://doi.org/10.1016/0147-619x(78)90001-x)
- Portnoy DA, Martinez RJ (1985) Role of a plasmid in the pathogenicity of *Yersinia* species. In: Goebel W (ed) *Genetic Approaches to Microbial Pathogenicity*. No. 118 in *Current Topics in Microbiology and Immunology*, Springer-Verlag, p 29–51, https://doi.org/10.1007/978-3-642-70586-1_3
- Rodriguez-Beltran J, Hernandez-Beltran JCR, DelaFuente J, et al (2018) Multi-copy plasmids allow bacteria to escape from fitness trade-offs during evolutionary innovation. *Nature Ecology & Evolution* 2(5):873–881. <https://doi.org/10.1038/s41559-018-0529-z>

- Rodríguez-Beltrán J, DelaFuente J, León-Sampedro R, et al (2021) Beyond horizontal gene transfer: the role of plasmids in bacterial evolution. *Nature Reviews Microbiology* 19(6):347–359. <https://doi.org/10.1038/s41579-020-00497-1>
- San Millan A, Escudero JA, Gutierrez B, et al (2009) Multiresistance in *Pasteurella multocida* is mediated by coexistence of small plasmids. *Antimicrobial Agents and Chemotherapy* 53(8):3399–3404. <https://doi.org/10.1128/aac.01522-08>
- San Millan A, Santos-Lopez A, Ortega-Huedo R, et al (2015) Small-plasmid-mediated antibiotic resistance is enhanced by increases in plasmid copy number and bacterial fitness. *Antimicrobial Agents and Chemotherapy* 59(6):3335–3341. <https://doi.org/10.1128/aac.00235-15>
- San Millan A, Escudero JA, Gifford DR, et al (2016) Multicopy plasmids potentiate the evolution of antibiotic resistance in bacteria. *Nature Ecology & Evolution* 1(1). <https://doi.org/10.1038/s41559-016-0010>
- Santer M, Uecker H (2020) Evolutionary rescue and drug resistance on multicopy plasmids. *Genetics* 215(3):847–868. <https://doi.org/10.1534/genetics.119.303012>
- Santer M, Kupczok A, Dagan T, et al (2022) Fixation dynamics of beneficial alleles in prokaryotic polyploid chromosomes and plasmids. *Genetics* 222(2). <https://doi.org/10.1093/genetics/iyac121>
- Santos-Lopez A, Bernabe-Balas C, Ares-Arroyo M, et al (2016) A naturally occurring single nucleotide polymorphism in a multicopy plasmid produces a reversible increase in antibiotic resistance. *Antimicrobial Agents and Chemotherapy* 61(2). <https://doi.org/10.1128/aac.01735-16>
- Schensted IV (1958) Model of subnuclear segregation in the macronucleus of ciliates. *The American Naturalist* 92(864):161–170. <https://doi.org/10.1086/282023>

- Sewastjanow BA (1975) Verzweigungsprozesse. R. Oldenbourg Verlag
- Silver S (1992) Plasmid-determined metal resistance mechanisms: Range and overview. Plasmid 27(1):1–3. [https://doi.org/10.1016/0147-619x\(92\)90001-q](https://doi.org/10.1016/0147-619x(92)90001-q)
- Smillie C, Garcillán-Barcia MP, Francia MV, et al (2010) Mobility of plasmids. Microbiology and Molecular Biology Reviews 74(3):434–452. <https://doi.org/10.1128/mmmbr.00020-10>
- Summers DK, Beton CWH, Withers HL (1993) Multicopy plasmid instability: the dimer catastrophe hypothesis. Molecular Microbiology 8(6):1031–1038. <https://doi.org/10.1111/j.1365-2958.1993.tb01648.x>
- Sun L, Alexander HK, Bogos B, et al (2018) Effective polyploidy causes phenotypic delay and influences bacterial evolvability. PLOS Biology 16(2):e2004644. <https://doi.org/10.1371/journal.pbio.2004644>
- Tazzyman SJ, Bonhoeffer S (2014) Plasmids and evolutionary rescue by drug resistance. Evolution 68(7):2066–2078. <https://doi.org/10.1111/evo.12423>
- Uhlen B, Nordström K (1975) Plasmid incompatibility and control of replication: copy mutants of the r-factor r1 in *Escherichia coli* K-12. Journal of Bacteriology 124:641–649. <https://doi.org/10.1128/jb.124.2.641-649.1975>
- Wardell GE, Hynes MF, Young PJ, et al (2021) Why are rhizobial symbiosis genes mobile? Philosophical Transactions of the Royal Society B 377(1842):20200471. <https://doi.org/10.1098/rstb.2020.0471>
- Yu W, Gillies K, Kondo JK, et al (1996) Loss of plasmid-mediated oligopeptide transport system in lactococci: Another reason for slow milk coagulation. Plasmid 35(3):145–155. <https://doi.org/10.1006/plas.1996.0017>

Zielenkiewicz U, Cegłowski P (2001) Mechanisms of plasmid stable maintenance with special focus on plasmid addiction systems. Acta Biochimica Polonica 48(4):1003–1023. https://doi.org/10.18388/abp.2001_3863

Appendix A Mathematical notes

In this appendix, we derive several combinatorial results needed in the main text of the paper.

First, for completeness, we derive the distribution of outcomes of a Pólya urn process (Eggenberger and Pólya, 1923). Suppose we start with an urn of n balls, i blue and $n - i$ white. We pick a ball out of the urn, then replace it together with another ball of the same colour. After repeating this process $m - n$ times, there are m balls in the urn: what is the probability that k of them are blue?

Proposition 1. *Denote the probability described above by $P\binom{n \rightarrow m}{i \rightarrow k}$. Then for $m \geq n > 0$ and $k \geq i \geq 0$,*

$$P\binom{n \rightarrow m}{i \rightarrow k} = \frac{\binom{m-k-1}{(m-k)-(n-i)} \binom{k-1}{k-i}}{\binom{m-1}{m-n}}.$$

Proof. First we prove the special case

$$P\binom{n \rightarrow n}{i \rightarrow k} = \frac{\binom{n-k-1}{i-k} \binom{k-1}{k-i}}{\binom{n-1}{0}};$$

since no balls are drawn in this case, the left-hand side is 1 if $i = k$ and 0 otherwise. The lower indices of the binomial coefficients in the numerator are negatives of each other, so the coefficients can only both be nonzero if $k - i = 0$, and in that case all three coefficients are equal to 1, satisfying the equality.

Next we prove another special case

$$P\left(\begin{matrix} n \rightarrow m \\ i \rightarrow i \end{matrix}\right) = \frac{\binom{m-i-1}{m-n} \binom{i-1}{0}}{\binom{m-1}{m-n}}.$$

Suppose that this equality is established for $P\left(\begin{matrix} n \rightarrow m-1 \\ i \rightarrow i \end{matrix}\right)$. In order to reach m balls in the urn without increasing the number of blue balls, we must first reach $m-1$ balls, then draw a white ball on the final draw. Therefore

$$\begin{aligned} P\left(\begin{matrix} n \rightarrow m \\ i \rightarrow i \end{matrix}\right) &= P\left(\begin{matrix} n \rightarrow m-1 \\ i \rightarrow i \end{matrix}\right) \frac{m-i-1}{m-1} \\ &= \frac{\binom{m-i-2}{m-n-1} \binom{i-1}{0} (m-i-1)}{\binom{m-2}{m-n-1} (m-1)} \\ &= \frac{\binom{m-i-1}{m-n} \binom{i-1}{0} (m-n)}{\binom{m-1}{m-n} (m-n)}, \end{aligned}$$

and by induction on m (with our previous special case as the base case), the equality is established.

Finally, we prove the general case. Suppose we know the equality holds for $P\left(\begin{matrix} n \rightarrow m-1 \\ i \rightarrow k \end{matrix}\right)$ and $P\left(\begin{matrix} n \rightarrow m-1 \\ i \rightarrow k-1 \end{matrix}\right)$. In order for there to be k blue balls in the urn when we get to the point of having m total balls, then either the last ball drawn was white, and there were k blue out of $m-1$ total balls before, or the last ball drawn was blue, and there were $k-1$ blue out of $m-1$ total balls before. Thus

$$\begin{aligned} P\left(\begin{matrix} n \rightarrow m \\ i \rightarrow k \end{matrix}\right) &= P\left(\begin{matrix} n \rightarrow m-1 \\ i \rightarrow k \end{matrix}\right) \frac{m-k-1}{m-1} + P\left(\begin{matrix} n \rightarrow m-1 \\ i \rightarrow k-1 \end{matrix}\right) \frac{k-1}{m-1} \\ &= \frac{\binom{m-k-2}{(m-k-1)-(n-i)} \binom{k-1}{k-i} (m-k-1)}{\binom{m-2}{m-n-1} (m-1)} + \frac{\binom{m-k-1}{(m-k)-(n-i)} \binom{k-2}{k-i-1} (k-1)}{\binom{m-2}{m-n-1} (m-1)} \\ &= \frac{\binom{m-k-1}{(m-k)-(n-i)} \binom{k-1}{k-i}}{\binom{m-1}{m-n}} \left(\frac{((m-k)-(n-i)) + (k-i)}{m-n} \right), \end{aligned}$$

and by induction on m and k , the equality is established. \square

Obviously if we replace blue and white balls in the urn with mutant and wild-type plasmids in a cell, this process describes the replication of plasmids under the random replication model. The probability that starting with a cell with i mutant plasmids there are $j + k$ mutant plasmids after replication is then $P\binom{n \rightarrow 2n}{i \rightarrow j+k}$.

When a mutation occurs during plasmid replication in the random replication model, it is possible for the mutant plasmid to be replicated again the same generation, and so the initial mutant cell can end up with multiple mutant plasmids. Suppose that mutations occur with probability u ; then the probability that a mutation occurs after ℓ plasmid replications, it is the only mutation in that generation, and the cell ends up with i mutant plasmids after all of the replication has occurred but before cell division is

$$A_{\ell i} = (1 - u)^\ell u P\binom{n + \ell + 1 \rightarrow 2n}{1 \rightarrow i} (1 - u)^{n-i-\ell}.$$

We assume that u is small, so that $1 - u \approx 1$; then the probability that a mutation occurs in a given generation is

$$\begin{aligned} \sum_{i=1}^n \sum_{\ell=0}^{n-1} A_{\ell i} &\approx \sum_{i=1}^n \sum_{\ell=0}^{n-1} u P\binom{n + \ell + 1 \rightarrow 2n}{1 \rightarrow i} \\ &= u \sum_{i=1}^n \sum_{\ell=0}^{n-1} \frac{\binom{2n-i-1}{n-\ell-i} \binom{i-1}{i-1}}{\binom{2n-1}{n-\ell-1}} \\ &= u \sum_{\ell=0}^{n-1} \frac{1}{\binom{2n-1}{n-\ell-1}} \sum_{i=1}^n \binom{2n-i-1}{n-\ell-i} \\ &= u \sum_{\ell=0}^{n-1} \frac{1}{\binom{2n-1}{n+\ell}} \sum_{i=0}^{n-1} \binom{2n-i-2}{n+\ell-1} \\ &= u \sum_{\ell=0}^{n-1} \frac{\binom{2n-1}{n+\ell}}{\binom{2n-1}{n+\ell}} \\ &= un. \end{aligned}$$

The probability of ending up with i mutant plasmids conditional on a mutation occurring is then

$$\begin{aligned}
 \frac{\sum_{\ell=0}^{n-1} A_{\ell i}}{\sum_{i=1}^n \sum_{\ell=0}^{n-1} A_{\ell i}} &\approx \frac{\sum_{\ell=0}^{n-1} uP_{1 \rightarrow i}^{(n+\ell+1 \rightarrow 2n)}}{un} \\
 &= \frac{1}{n} \sum_{\ell=0}^{n-1} \frac{\binom{2n-i-1}{n-\ell-i} \binom{i-1}{i-1}}{\binom{2n-1}{n-\ell-1}} \\
 &= \frac{1}{n} \sum_{\ell=0}^{n-1} \frac{n+\ell}{n-\ell} \frac{\binom{n-\ell}{i}}{\binom{2n-1}{i}} \\
 &= \frac{1}{ni \binom{2n-1}{i}} \sum_{\ell=0}^{n-1} (n+\ell) \binom{n-\ell-1}{i-1} \\
 &= \frac{1}{ni \binom{2n-1}{i}} \left(\sum_{\ell=0}^{n-1} 2n \binom{n-\ell-1}{i-1} - \sum_{\ell=0}^{n-1} (n-\ell) \binom{n-\ell-1}{i-1} \right) \\
 &= \frac{1}{ni \binom{2n-1}{i}} \left(2n \sum_{\ell=0}^{n-1} \binom{n-\ell-1}{i-1} - i \sum_{\ell=0}^{n-1} \binom{n-\ell}{i} \right) \\
 &= \frac{1}{ni \binom{2n-1}{i}} \left(2n \binom{n}{i} - i \binom{n+1}{i+1} \right) \\
 &= \frac{1}{ni \binom{2n-1}{i}} \left(2n - \frac{i(n+1)}{i+1} \right) \binom{n}{i} \\
 &= \frac{(2n + (n-1)i)}{ni(i+1)} \frac{\binom{n}{i}}{\binom{2n-1}{i}};
 \end{aligned}$$

by substituting $i = j + k$, we obtain the replication part of equation (3).

Lemma 1. For all integers $n > 0$ and $0 \leq \beta \leq n$,

$$\sum_{\alpha=0}^n \frac{\binom{2i}{\alpha} \binom{2n-2i}{n-\alpha}}{\binom{2n}{n}} \binom{\alpha}{\beta} = \frac{\binom{n}{\beta}}{\binom{2n}{\beta}} \binom{2i}{\beta} \quad (\text{A1})$$

and

$$\sum_{\gamma=1}^{n-1} \sum_{\delta=0}^{2n} \frac{\binom{\delta}{\beta} \binom{2n-\delta}{n-\beta}}{\binom{2n}{n}} \frac{\binom{2n-\delta-1}{n-\delta+\gamma} \binom{\delta-1}{\delta-\gamma}}{\binom{2n-1}{n}} = \frac{(n-1)(2n+1)}{(n+1)(2n-1)}. \quad (\text{A2})$$

Proof. For (A1), we have that

$$\begin{aligned}
 \sum_{\alpha=0}^n \frac{\binom{2i}{\alpha} \binom{2n-2i}{n-\alpha}}{\binom{2n}{n}} \binom{\alpha}{\beta} &= \frac{\binom{2i}{\beta}}{\binom{2n}{n}} \sum_{\alpha=0}^n \binom{2i-\beta}{\alpha-\beta} \binom{2n-2i}{n-\alpha} \\
 &= \frac{\binom{2i}{\beta}}{\binom{2n}{n}} \sum_{\alpha'} \binom{2i-\beta}{\alpha'} \binom{2n-2i}{n-\beta-\alpha'} \\
 &= \binom{2i}{\beta} \binom{2n-\beta}{n-\beta} / \binom{2n}{n} \\
 &= \binom{2i}{\beta} \binom{2n-\beta}{n-\beta} \binom{n}{\beta} / \binom{2n}{n} \binom{n}{\beta} \\
 &= \binom{2i}{\beta} \binom{2n-\beta}{n-\beta} \binom{n}{\beta} / \binom{2n}{\beta} \binom{2n-\beta}{n-\beta} \\
 &= \frac{\binom{n}{\beta}}{\binom{2n}{\beta}} \binom{2i}{\beta},
 \end{aligned}$$

where the first and fifth equalities use trinomial revision and the third Vandermonde convolution. For (A2), we have that

$$\begin{aligned}
 \sum_{\gamma=1}^{n-1} \sum_{\delta=0}^{2n} \frac{\binom{\delta}{\beta} \binom{2n-\delta}{n-\beta}}{\binom{2n}{n}} \frac{\binom{2n-\delta-1}{n-\delta+\gamma} \binom{\delta-1}{\delta-\gamma}}{\binom{2n-1}{n}} \\
 = \frac{1}{\binom{2n}{n} \binom{2n-1}{n}} \sum_{\delta=0}^{2n} \binom{\delta}{\beta} \binom{2n-\delta}{n-\beta} \sum_{\gamma} \binom{2n-\delta-1}{n-\delta+\gamma} \binom{\delta-1}{\delta-\gamma} \\
 = \frac{\binom{2n-2}{n}}{\binom{2n}{n} \binom{2n-1}{n}} \sum_{\delta=0}^{2n} \binom{\delta}{\beta} \binom{2n-\delta}{n-\beta} \\
 = \frac{\binom{2n-2}{n} \binom{2n+1}{n+1}}{\binom{2n}{n} \binom{2n-1}{n}} = \frac{(n-1)(2n+1)}{(2n-1)(n+1)},
 \end{aligned}$$

where the second equality uses Vandermonde convolution (with an implicit change of index to $\delta - \gamma$) and the third uses equation (5.26) of [Graham et al \(1994\)](#). \square