



**HAL**  
open science

# Genome Streamlining: Effect of Mutation Rate and Population Size on Genome Size Reduction

Juliette Luiselli, Jonathan Rouzaud-Cornabas, Nicolas Lartillot, Guillaume Beslon

► **To cite this version:**

Juliette Luiselli, Jonathan Rouzaud-Cornabas, Nicolas Lartillot, Guillaume Beslon. Genome Streamlining: Effect of Mutation Rate and Population Size on Genome Size Reduction. *Genome Biology and Evolution*, 2024, 16, 10.1093/gbe/evae250 . hal-04905734

**HAL Id: hal-04905734**

**<https://hal.science/hal-04905734v1>**

Submitted on 22 Jan 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Genome Streamlining: Effect of Mutation Rate and Population Size on Genome Size Reduction

Juliette Luiselli <sup>1,2</sup>, Jonathan Rouzaud-Cornabas <sup>1,2</sup>, Nicolas Lartillot <sup>3</sup>, Guillaume Beslon <sup>1,2,\*</sup>

<sup>1</sup>INSA-Lyon, CNRS, Université Claude Bernard Lyon 1, ECL, Université Lumière Lyon 2, LIRIS UMR5205, Lyon 69621, France

<sup>2</sup>Beagle Team, Inria Lyon La Doua, Villeurbanne, France

<sup>3</sup>Laboratoire de Biométrie et de Biologie Évolutive UMR CNRS 5558, Université Claude Bernard Lyon 1, Université Lyon 1, Villeurbanne, France

\*Corresponding author: E-mail: guillaume.beslon@insa-lyon.fr.

Accepted: November 14, 2024

## Abstract

Genome streamlining, *i.e.* genome size reduction, is observed in bacteria with very different life traits, including endosymbiotic bacteria and several marine bacteria, raising the question of its evolutionary origin. None of the hypotheses proposed in the literature is firmly established, mainly due to the many confounding factors related to the diverse habitats of species with streamlined genomes. Computational models may help overcome these difficulties and rigorously test hypotheses. In this work, we used Aevol, a platform designed to study the evolution of genome architecture, to test 2 main hypotheses: that an increase in population size ( $N$ ) or mutation rate ( $\mu$ ) could cause genome reduction. In our experiments, both conditions lead to streamlining but have very different resulting genome structures. Under increased population sizes, genomes lose a significant fraction of noncoding sequences but maintain their coding size, resulting in densely packed genomes (akin to streamlined marine bacteria genomes). By contrast, under an increased mutation rate, genomes lose both coding and non-coding sequences (akin to endosymbiotic bacteria genomes). Hence, both factors lead to an overall reduction in genome size, but the coding density of the genome appears to be determined by  $N \times \mu$ . Thus, a broad range of genome size and density can be achieved by different combinations of  $N$  and  $\mu$ . Our results suggest that genome size and coding density are determined by the interplay between selection for phenotypic adaptation and selection for robustness.

**Key words:** genome architecture, genome evolution, genome streamlining, mutation rate, modeling, population size.

## Significance

Many bacterial species show reduced genomes. However, the diversity of these species and their life traits makes it difficult to identify the mechanisms that led to this reduction. Indeed, no unifying hypothesis accounts for the whole diversity of genome size reduction. Here, we used simulations to systematically explore the effect of population size and mutation rate on genome size. Our results show that the interaction between these 2 factors tightly determines the size, but also the density of genomes, making it possible to account for the whole diversity of reduced genomes by acting on these 2 parameters only. Our results suggest a theoretical model in which genome reduction is driven by a robustness/fitness trade-off.

## Introduction

Genome size was one of the first studied genome characteristics (Leth Bak et al. 1969; Bachmann 1972), yet its

dynamic and causal factors are still poorly understood. Genome size is hugely variable across life: from less than  $10^4$  base pairs (bp) for viruses (Gago et al. 2009), to more

© The Author(s) 2024. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

than  $10^{11}$  bp for some plants (Pellicier et al. 2010). It does not correlate reliably with the number of genes or other variables throughout the different branches of life (Barow and Meister 2002; Westoby et al. 2021).

The observed range of genome sizes is more restricted when studying only bacterial organisms (Westoby et al. 2021), ranging from  $10^5$  bp for intracellular endosymbiotic bacteria (Chong et al. 2019) to  $10^7$  bp for some myxobacteria (Schneiker et al. 2007). Bacterial genomes are mostly dense, and within this domain of life, genome size is loosely correlated with the number of coding genes (Konstantinidis and Tiedje 2004; Almpanis et al. 2018). However, the precise determinants of bacterial genome size are still unknown, as it is still impossible to accurately predict the total genome size from the number of coding genes or from other genomic characteristics (Petrov 2001; Barow and Meister 2002; Choi et al. 2020). Part of the determinants of genome size are likely to be highly lineage-specific and linked to the ecological or evolutionary history of the lineages (Martinez-Gutierrez and Aylward 2022). Nevertheless, it has been argued that at least a part of the observed variation may be due to universal mechanisms, linked to population genetics and molecular evolutionary processes (Lynch and Conery 2003; Lynch 2007). In particular, it has been suggested that population genetics mechanisms could explain the reductive evolution observed in several bacterial strains (Lynch 2006). However, among the shortest bacterial genomes, one can find 2 types of bacteria which have very different ecological environments and evolutionary history: endosymbionts such as *Buchnera aphidicola* (Moran and Mira 2001) and free-living marine bacteria such as *Prochlorococcus marinus* (Dufresne et al. 2005) or *Pelagibacter ubique* (Giovannoni et al. 2005). Strikingly, both types of bacteria lie at the 2 extremes of bacterial population sizes, questioning the mechanisms that led to genome reduction (Batut et al. 2014; Martínez-Cano et al. 2015; Wernegreen 2015).

*Buchnera aphidicola*, and endosymbionts more generally, are characterized by very small effective population sizes ( $N_e$ ) and high mutation rates ( $\mu$ ). Endosymbiosis also generally entails the introduction to a new stable environment and very close interactions with the host (Moran 1996; Mira and Moran 2002). These many complex factors result in decaying genomes, smaller in total size and with fewer coding genes than those of average bacteria (Heddi et al. 1998). Endosymbionts have typically lost both coding and noncoding genomic content (Moran and Mira 2001; Wernegreen 2002), maintaining a coding fraction on the order of 85% (van Ham et al. 2003), which is quite typical for bacteria (Kuo et al. 2009).

In sharp contrast, free-living marine bacteria such as *Prochlorococcus marinus* or *Pelagibacter ubique* also have reduced genomes (Giovannoni et al. 2005; Batut et al. 2014), but are believed to have very large effective

population sizes (Marais et al. 2008; Flombaum et al. 2013; Giovannoni et al. 2014), although that is an ongoing debate (Chen et al. 2022; Filatov and Kirkpatrick 2024). Noticeably, in their case, genome size reduction is primarily contributed by the loss of noncoding sequences rather than coding sequences (Giovannoni et al. 2005; Batut et al. 2014). This phenomenon is called streamlining and could indicate a very effective selection (Wolf and Koonin 2013; Giovannoni et al. 2014). Many hypotheses have been proposed to account for genome size reduction and the associated changes in genome architecture in such free-living organisms: adaptation to a nutrient-poor environment or to other abiotic factors, the Black Queen hypothesis, or high mutation rates (Koskiniemi et al. 2012; Morris et al. 2012; Batut et al. 2014; Ngugi et al. 2023).

Both endosymbionts and free-living marine bacteria thus show a marked reduction in genome size, linked to an increase in mutation rate (Bourguignon et al. 2020) but, strikingly, also linked to either an increase or a decrease in effective population size  $N_e$ . Indeed, while some observations link the decrease in genome size to the increase in random drift (Moran 2002; Nilsson et al. 2005; Kuo et al. 2009), this is not consensual among the scientific community since a long-term reduction in  $N_e$  is also thought to increase genome complexity and genome size: the increase in genetic drift would cause the fixation of slightly deleterious duplications, which would be more frequent than slightly deleterious deletions (Lynch and Conery 2003; Lefebvre et al. 2017). The balance between insertion and deletion rates and spectra may also play a role in genome size evolution (Petrov 2002) and deletion biases in particular are believed to contribute to the small genome size of prokaryotes (Bingham and Ratcliff 2024). Overall, this suggests that a specific study focusing on the interaction between various mutational biases, variations in mutation rate and variations in effective population size is needed.

In this study, we focus on determining the impact of both an increased mutation rate and a change in population size on genome size evolution. However, mutation rates and population sizes are difficult to estimate. The effective population size is also highly variable through time, such that it is not totally obvious which long-term average is relevant at the macroevolutionary scale (Brevet and Lartillot 2021; Müller et al. 2022). For that reason, many comparative analyses have relied on somewhat indirect proxies, such as life-history traits (Popadin et al. 2007; Romiguier et al. 2012; Figuet et al. 2016). However, the precise quantitative relation between these proxies and effective population size is difficult to assess. Moreover, the very different living conditions and potential mutational biases of the bacterial species that have undergone genome reduction introduce many confounding factors. To avoid these pitfalls, we choose to turn to simulation, which allows us to control all the parameters (population size, mutation rate, and

mutational biases) and the magnitude of their variation. It also ensures that no other factor than the ones investigated will impact the phenomenon under study. Hence, we can gain a theoretical understanding of the relationship between the different factors at stake and genome size reduction.

In silico experimental evolution provides tools to study genomic architecture in detail (Adami 2006; Hindré et al. 2012; Batut et al. 2013). For our study, we need a framework that provides coding and noncoding genomic compartments which can vary independently, and with arbitrary underlying mutational biases for the deletion/insertion balance. Then, running simulations in a perfectly controlled environment covering a broad range of population sizes  $N$  and mutation rates  $\mu$  makes it possible to investigate the conditions and mechanisms leading to genome size reduction. We will hence use Aevol, a simulation platform that provides an explicit genomic structure where both the coding and noncoding genome can evolve freely. Aevol emulates the evolution of bacteria and enables replicated and controlled in silico evolution experiments with known and fixed parameters (Knibbe et al. 2007; Banse et al. 2023). It provides an ideal tool to uncover links between genome size and either population size or mutation rate, as the experimenter perfectly controls these parameters. Throughout the experiments, fitness, genome size, and amounts of coding and noncoding bases are monitored to study the evolution of genome architecture and the response of genome size to changes in  $\mu$  and  $N$ .

Our results show that both an increase in  $N$  or  $\mu$  lead to genome size reduction, regardless of the underlying mutational bias. However, both conditions lead to very different genome structures, as a high  $\mu$  reduces both the coding and noncoding compartments while a high  $N$  reduces only the noncoding compartment. Surprisingly, they both lead to a similar coding proportion when increased by the same factor, such that  $N \times \mu$  appears as a key compound parameter determining this proportion. To understand this result, we measured both the phenotypical adaptation and the replicative robustness of the genomes, *i.e.* their capacity to transmit faithfully their phenotypes to their offspring. Indeed, while the per-base mutation rate is constant within each of our experiments, the genome-wide mutation rate varies with genome size, and the impact of the mutations depends on the genome structure and the type of mutation. Therefore, replicative robustness is tightly linked with genome size and coding proportion. We show that the observed variations in genome size and structure are due to the interaction between selection for phenotypical adaptation to the environment and selection for robustness.

## Results

We perform our experiments using Aevol, a forward-in-time evolutionary simulator (Knibbe et al. 2007; Banse et al.

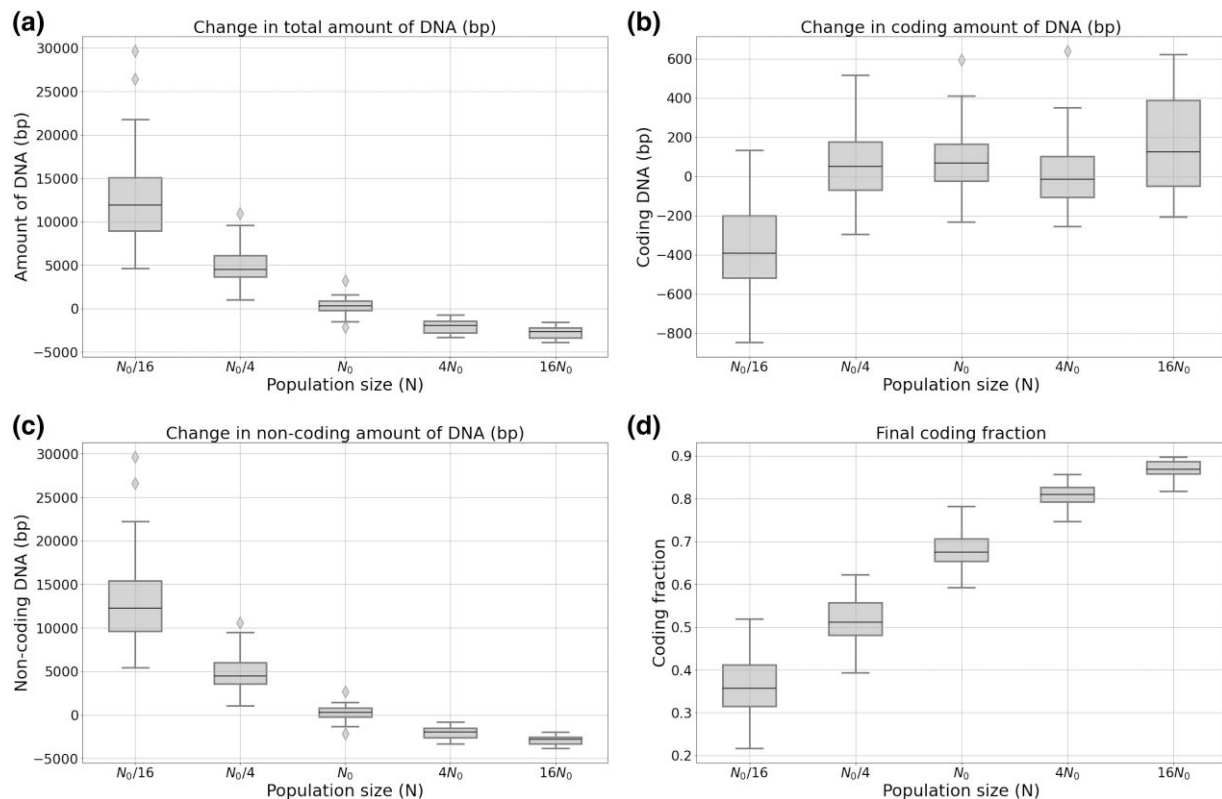
2023). Aevol is an individual-based model which includes an explicit population and in which every organism owns a double-stranded genome. It uses an explicit genome decoding algorithm directly inspired by the central dogma of molecular biology to compute the phenotype, and thus the fitness, of each individual based on its genomic sequence. As Aevol also includes a large variety of mutational operators (including substitutions, InDels, and chromosomal rearrangements), this nonparametric genotype-to-phenotype map allows for changes in the genome architecture (genome size, coding density, overlapping genes or operons, etc.), without assuming a predefined distribution of fitness effects. Indeed, in the model, it is possible to reach similar fitnesses in many ways, by adjusting the number of genes, their loci, their lengths, or the intergenic distances, hence the total amount of noncoding DNA. In Aevol, genes are typically created by duplication-divergence (Kalhor et al. 2024), but they can also be deleted, and some may emerge *de novo*. Hence, the impact of a given mutation highly depends on the preexisting genome structure, which can in turn be indirectly selected (Knibbe et al. 2007). Aevol therefore allows studying changes in size and structure of genomes in response to changes in population size and mutation rates.

Our experiments start from 5 “Wild-Type” (WT) lines, each having evolved for 10 million generations within a population of 1,024 individuals and a mutation rate of  $10^{-6}$  mutations per base pair for each mutation type: substitutions, small insertions, small deletions, duplications, deletions, translocations, and inversions. There is no underlying mutational bias: the insertion and deletion of bases are equally probable. The 5 WT lines display stable genome structures (with small random variations, as exemplified by cases  $N_0$  and  $\mu_0$  on Figs. 1 and 2) although they still slowly gain fitness by fixing rare favorable mutations (see case  $N_0$  on Fig. 5a). Their fitness and genomic characteristics are displayed in Section 4.2, Table 1. In our experiments, these WT lines are used as founders of new populations, which are confronted with new evolutionary conditions for 2 million generations. In parallel, these same WT lines were evolved in the same conditions they first evolved in, providing perfect control experiments. We compare the fitness, genome size, and genome structure of populations that evolved in new conditions with those of the control populations. Finally, we repeat part of these experiments with WT lines that evolved with either an insertion or a deletion bias to understand how an underlying mutational bias might impact our findings.

### Genome Size Evolution Following a Change in Population Size and Mutation Rate

#### *Change in Population Size*

In the absence of mutational bias, increasing the population size by a factor of 4 or 16 results in a reduction in



**Fig. 1.** Total a), coding b) and non-coding c) genome size variation, and final coding fraction d), after 2 million generations. For each of the 5 WT, 10 replicas were performed under a constant mutation rate ( $\mu_0 = 10^{-6}$  per base pair for each type of mutation) with 5 different population sizes ( $N_0 = 1,024$  being the control population size).

the total genome size (see Fig. 1a). Yet, this change does not impact the coding and noncoding parts of the genome proportionally: while the size of the coding compartment is barely affected (see Fig. 1b), the noncoding genome size is greatly reduced (see Fig. 1c). As a result, the coding proportion of the genome increases (see Fig. 1d). Conversely, reducing the population size by a factor of 4 or 16 increases the total genome size (Fig. 1a) by increasing greatly the noncoding genome size (Fig. 1c). In the extreme condition  $N_0/16$ , the coding genome size is also slightly reduced (Fig. 1b). As a result, the coding fraction of the genome is drastically reduced (Fig. 1d).

### Change in Mutation Rate

In the absence of mutational bias, increasing the mutation rate drastically reduces the total genome size (see Fig. 2a). Thus, at first sight, population size and mutation rate seem to have a similar effect on genome evolution. However, in the details, the effect of these 2 variables on genome structure appears to differ, as the reduction now occurs in both the coding and non-coding genomic compartments (see Fig. 2b and c). Both are nevertheless not proportionally affected by the decrease in mutation rate, which affects more

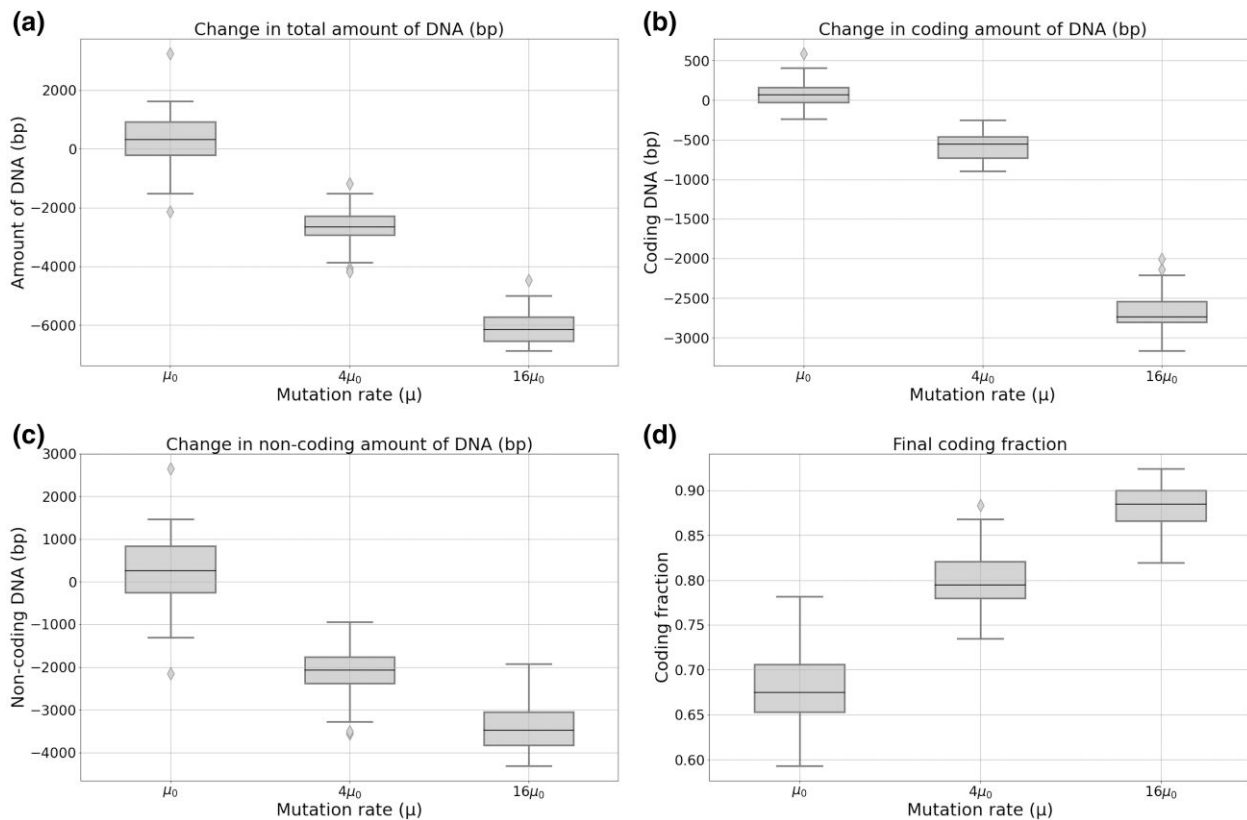
strongly the noncoding part of the genome, such that the final coding fraction of the genome increases with  $\mu$  (see Fig. 2d). Altogether, these results show that streamlined genomes, denser and shorter than their ancestors, can result from either an increase in population size or in mutation rate.

Notably, and despite the very different dynamics displayed in the 2 experiments, a 4-fold increase in  $N$  or in  $\mu$  results in the same final coding proportion of approximately 80%. The same is true for a 16-fold increase (88%). To further investigate this result, we conducted additional experiments to observe the combined effects of a simultaneous modification in both  $N$  and  $\mu$ .

### Linked Effect of Population Sizes and Mutation Rates

Figure 3 shows the variation in the total amount of DNA, coding size, and noncoding size, as well as the variation in coding fraction for several combinations of changes in  $N$  and  $\mu$  (note that, in the panels of Fig. 3, the bottom line and the central column, respectively, correspond to the values presented in Figs. 1 and 2).

Overall, as  $N$  increases, the total amount of DNA decreases, whatever the value of  $\mu$  (see Fig. 3a). A higher  $\mu$



**Fig. 2.** Total a), coding b) and non-coding c) genome size variation, and final coding fraction d), after 2 million generations. For each of the 5 WT, 10 replicas were performed under a constant population size ( $N_0 = 1,024$  individuals) with 3 different mutation rates: the control  $\mu_0 = 10^{-6}$  mutations per base pair for each type of mutation,  $4 \times \mu_0$  and  $16 \times \mu_0$ .

**Table 1.** Characteristics of the 5 WT at the start of our experiments

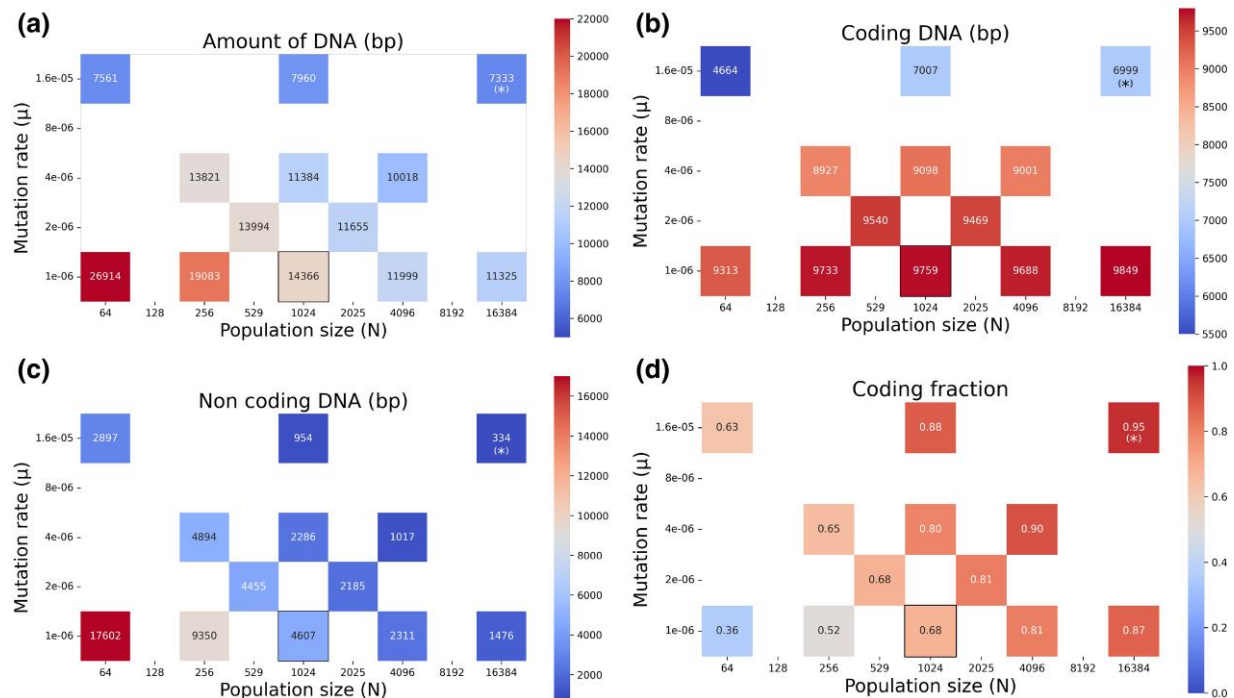
WT id	Fitness (arbitrary unit)	Total genome size (bp)	Coding size (bp)	Non-coding size (bp)	Coding fraction
1	0.014903	13,599	9,395	4,204	0.69
2	0.103795	13,660	8,828	4,832	0.65
3	0.128472	14,171	9,507	4,664	0.67
4	0.035369	14,507	10,003	4,504	0.69
5	0.029588	14,290	10,644	3,646	0.74
Average	0.0624254	14,045.5	9,675.4	4,370	0.69

also leads to a reduction in the total genome size, whatever the value of  $N$ . However, the effect of population size and mutation rate differ when considering the coding size of the genome: specifically, the coding size increases with  $N$  but decreases with  $\mu$  (see Fig. 3b). This is countered by the change in the noncoding size of the genomes (see Fig. 3c), which strongly decreases with both  $N$  and  $\mu$  and drives the overall change in genome size.

The interplay between  $N$  and  $\mu$  results in a surprisingly constant coding fraction across the different constant

values of  $N \times \mu$  (see Fig. 3d). Indeed, we observe that under constant  $N \times \mu$ , and although these 2 factors taken individually have changed in different proportions, the coding fraction remains constant: 80% when  $N_0 \times \mu_0$  is multiplied by 4 compared with the control conditions, and 88% when  $N_0 \times \mu_0$  is multiplied by 16 (see Fig. 3d). Although the coding fraction does slightly vary (from 68% to 63%) for the most extreme tested configuration ( $N_0/16$  and  $16\mu_0$ ), the diagonal of constant  $N_0 \times \mu_0$  also displays an almost constant coding fraction (Fig. 3d).

However, strikingly, the total genome size as well as the coding and noncoding genome sizes vary greatly, even for similar coding densities (Fig. 3b, c, and d). For densities of 63% and 65%, the total amount of DNA can be almost halved (from 13,821 bp to 7,561 bp) by going from  $N_0/4$  and  $4\mu_0$  to  $N_0/16$  and  $16\mu_0$  on the same diagonal of constant  $N \times \mu$ . Conversely, we can reach similar values of genome size (11,300 bp) despite important differences in the coding percentage (80% when  $\mu$  is multiplied by 4, and 87% when  $N$  is multiplied by 16). Altogether, these results show that a large range of genome sizes and structures (here corresponding to coding densities) can result from a combined variation in both the population size  $N$  and the mutation rate  $\mu$ .



**Fig. 3.** Amount of DNA a), coding size b), noncoding size c) and coding fraction d) for the different combinations of  $\mu$  and  $N$  tested, after 2 million generations. For each of the 5 WT, 10 replicas were performed for each tested set of conditions. Control conditions ( $N = 1, 1024$  and  $\mu = 1.10^{-6}$ ) are outlined in black. For the combination of both the highest mutation rate and the largest population size, only the median was tested due to computational limitations, which is indicated by a (\*).

### Mutational Biases Change the Equilibrium Genome Size, but not the Role of $N$ and $\mu$

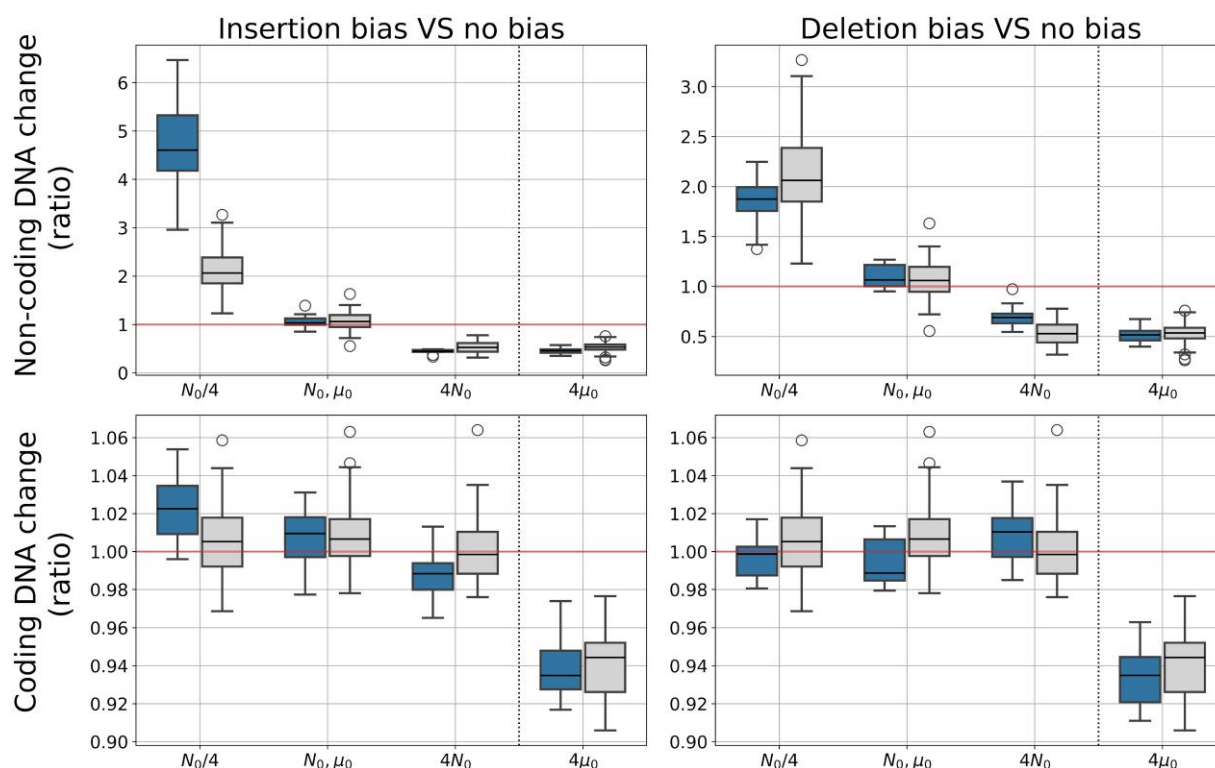
As genome sizes are generally thought to be heavily impacted by mutational biases, we control whether the effect of population size and mutation rate we observed is affected by either a deletion or an insertion bias. To this end, we evolved 5 WT organisms with either an insertion bias (twice as many duplications than large deletions), or a deletion bias (twice as many large deletions than duplications). The rates of all other types of mutations, as well as the sum of all mutation rates, are the same as in the previous experiments. As expected, the equilibrium genome sizes and coding proportions of these WT is affected by the balance between large deletions and duplications, with an average genome size of 11,623 bp in the presence of a deletion bias and 16,350 in the presence of a duplication bias (instead of 14,046 bp without any bias). The coding proportion is also affected: 0.78 and 0.61, respectively, instead of 0.69. This shows that the genome size and structure are, as expected, strongly influenced by the underlying mutation biases (Kuo and Ochman 2009).

We then confronted the median (in terms of genome size) WT of each condition to changes in population size (multiplied or divided by 4) or mutation rate (multiplied by 4) for 10 replicas. Similarly to what is observed without bias, an increase in  $N$  reduces the non-coding genome

size only, while an increase in  $\mu$  reduces both the coding and noncoding genome (see Fig. 4). Notably, a decrease in  $N$  increases the noncoding genome size even in the case of a deletion bias, although an insertion bias greatly amplifies this effect. As a result, and despite the strong mutational biases, we observe that multiplying either the population size or the mutation rate by the same factor leads to a genome compaction in similar proportions (the final coding fraction being 0.85 vs. 0.88 in the case of the deletion bias, and 0.78 vs. 0.77 in case of the insertion bias, respectively). Therefore, although mutational biases influence the equilibrium genome sizes and structures, they do not fundamentally change how the genomes react to variations in population size or mutation rate. In other words, our simulations show that mutational biases only determine the equilibrium set point around which population size and the overall mutation rate then modulate the genome size and structure. Similar experiments were run with biases in InDels and are presented in [supplementary material S2, Supplementary Material](#) online.

### Robustness Selection as the Explanatory Mechanism

We observed that 2 distinct processes, triggered by an increase in either population size or mutation rate, can lead to genome size reduction in our experiments. However, both have different effects on coding and noncoding



**Fig. 4.** Change in coding and noncoding genome sizes in reaction to changes in  $N$  or  $\mu$  for the different mutational biases. Blue boxes (on the left of each condition) show the results with a mutational bias (left: insertion bias, right: deletion bias), and gray boxes (on the right of each condition) show the results without mutational bias. Depicted values are the ratio of the coding/noncoding sizes at the final generation over the value at generation 0.

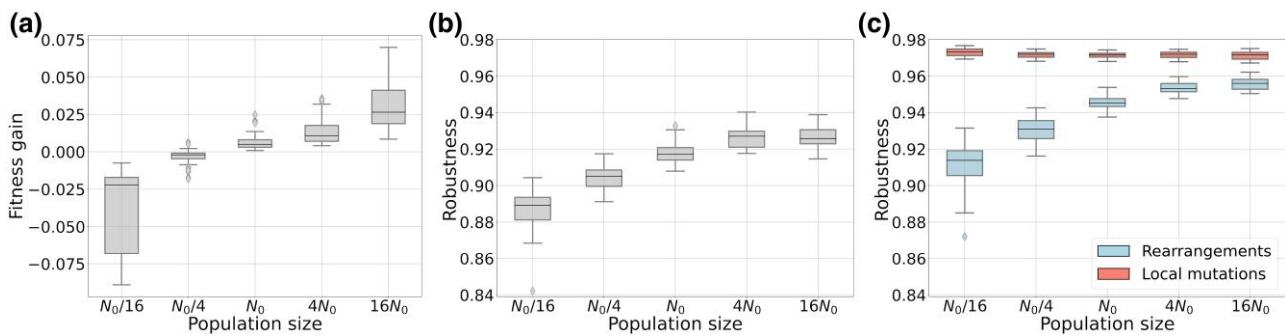
sequences: while an increased  $\mu$  reduces both the coding and noncoding genome sizes, increasing  $N$  reduces only the noncoding genome size.

We propose that these observations can be explained by an interplay between selection for phenotypic adaptation to the environment (hereafter called *direct selection*), and selection for replicative robustness (hereafter referred to as *indirect selection*). More specifically, we define the replicative robustness of an individual as its ability to transmit its fitness to its offspring. It hence corresponds to the proportion of offspring that did not acquire new deleterious mutations. This depends both on the number of mutations occurring at replication (which in turn depends on genome size) and on the probability for a given mutation to be deleterious (usually called mutational robustness Wilke and Adami 2003), which depends on the intertwining between the kind of mutation and the genomic architecture. In our case, WT organisms are very well adapted to their environment, thus most mutations will be deleterious if they affect the coding part of the genome. This is particularly true for chromosomal rearrangements, which can affect large genomic segments (Knibbe et al. 2007; Banse et al. 2023). Conversely, beneficial mutations are extremely rare. We therefore approximate the robustness of our organisms by measuring the proportion of their offspring that have

the exact same fitness, *i.e.* that underwent no mutations or only neutral mutations.

A more robust individual has more chances to pass on its genomic information accurately than a less robust one, thus enabling its lineage to better maintain its fitness in the long term and to outcompete other lineages in which deleterious mutations would accumulate at a higher rate. This results in an indirect selection for replicative robustness. We recall that replicative robustness depends both on the probability for a given mutation to be neutral (hence on the fraction of noncoding sequences in the genome) and on the mean number of mutations undergone by the genome at each generation (hence on the genome-wide mutation rate). Here, while the per base mutation rate is constant within each experiment, the total amount of DNA, and hence the genome-wide mutation rate, varies and can thus be indirectly selected. By contrast, direct selection depends only on the content of the coding compartment, the size of which is likely to be positively correlated with the level of phenotypical adaptation (at least in our model). As a result, indirect selection for robustness favors shorter genomes with a lower coding fraction, while direct selection for phenotypical adaptation maintains or even increases the coding size of the genome.





**Fig. 5.** Fitness gain (a) and Robustness (b: overall and c: by mutation type) at the end of the simulations, for different population sizes  $N$  and without mutational biases. Robustness is defined as the proportion of neutral offspring. The mutation rate is fixed to  $10^{-6}$  per base pair for each type of mutation.

The efficacy of both direct and indirect selection increases with population size, since some deleterious mutations that were quasi-neutral for a low  $N$  can become effectively counter-selected in the context of a high  $N$ , changing the balance of beneficial vs deleterious fixed mutations. To quantify this effect, we measured the robustness of the individuals at time 2,000,000 in the simulations without mutational biases. Figure 5a and b shows that the increase in selection efficacy induced by the increase in population size indeed induces both an increase in fitness (due to direct selection) and an increase in replicative robustness (due to indirect selection). In terms of genomic structure, a more efficient direct selection (*i.e.* a weaker random drift) is thus expected to increase the coding genome size, and a more efficient indirect selection is expected to decrease the overall genome size. The combination of both these effects leads to a decrease in the noncoding genome size, and maintenance of the coding genome size, as exemplified by Fig. 1b and c. Conversely when the population size is reduced, the increased drift leads to the loss of coding sequences and inflation of the noncoding compartment (Fig. 1b and c). This reorganization of the genome structure is associated with a loss in robustness (Fig. 5b).

In Aevol, genomes undergo different types of mutations that can be roughly grouped into local mutations (substitutions, InDels) and chromosomal rearrangements (duplications, deletions, inversions, translocations). Both kinds of events don't have the same effect on robustness. Figure 5c shows the change in robustness induced by the different types of events. It shows that the loss and gain in robustness are driven by chromosomal rearrangements. In contrast, local mutations (substitutions and InDels) do not have a significant effect on robustness.

In the case of an increased mutation rate, things are very different: a sudden increase in  $\mu$  results in an immediate drop in robustness at the beginning of the experiments (Fig. 6a). As the proportion of offspring that bears mutations rises with  $\mu$ , we go from an initial robustness of 92% for  $\mu_0$ , to 71% for  $4\mu_0$ , and only 26% for  $16\mu_0$ . In these new conditions,

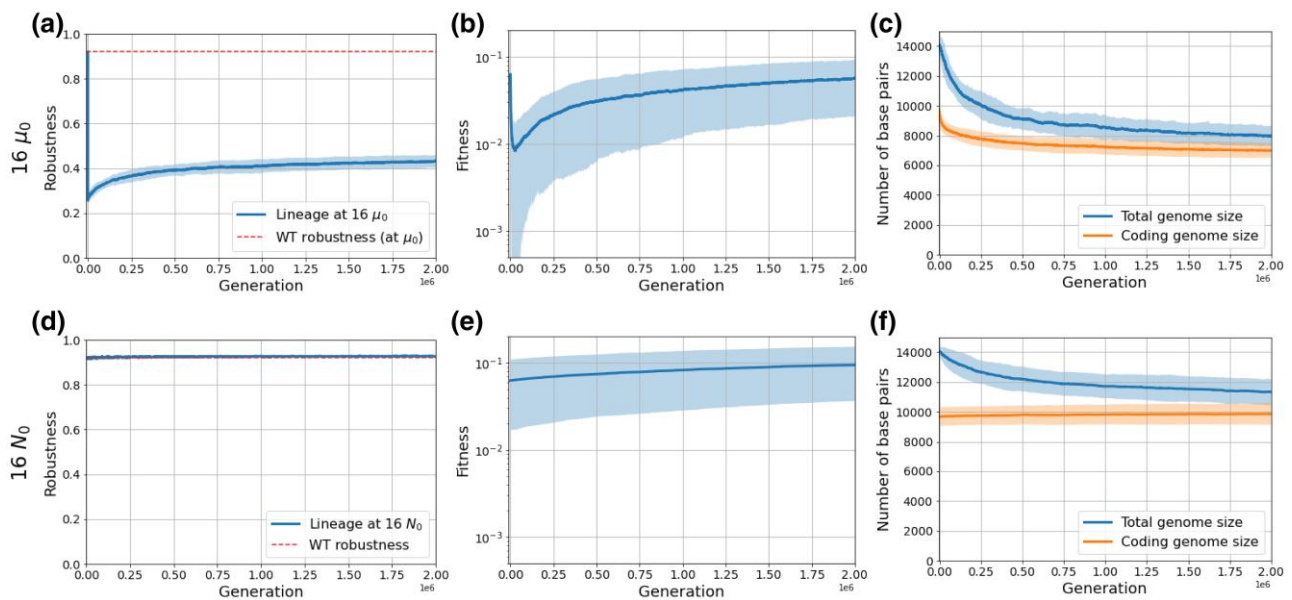
organisms are no longer able to transmit their genome to the next generation without deleterious mutations, and thus the indirect selection for robustness becomes temporarily stronger than the direct selection for phenotypical adaptation. Indeed, features that would not be accurately inherited cannot be selected. This indirect selection for robustness leads to the fixation of mutations that drastically decrease genome size, even at the cost of a loss of fitness for the individuals (see Fig. 6b): the only lineages that survive in the long term are those that have undergone a decrease in genome size, allowing them to reduce their per-genome mutation rate, thus regaining some robustness (see Fig. 6c). Once the robustness has increased sufficiently, direct selection for phenotypical adaptation can resume and the fitness starts to increase again (see Fig. 6b). Interestingly, organisms manage here to continue to lose some coding base pairs while increasing their fitness, probably thanks to global genome restructuring allowing for a more compact encoding of the phenotype, for example, through overlapping genes. This dynamic is very different from when  $N$  is increased (and so the initial robustness is unaffected), as shown by Fig. 6d, e, and f.

Notably, robustness does not reach values as high as that observed before the increase in mutation rate and stays below 50%. Indeed, the genome size could not be divided by 16 while keeping a good enough phenotypical adaptation, and the selection for phenotypical adaptation becomes stronger than the selection for robustness as soon as some organisms can pass on their genomic information reliably enough.

The interplay between direct and indirect selection can therefore explain both types of genome size reduction: affecting both coding and noncoding compartments (although not proportionally) when caused by an increased mutation rate, and restricted to the noncoding compartment when caused by an increased population size.

## Discussion

We found that, in our experiments, genome size reduction can be caused by an increase in population size, mutation



**Fig. 6.** Robustness, fitness, and genome architecture across generations for  $\mu = 1.6 \times 10^{-5}$  ( $16\mu_0$ ) per base pair for each mutation type and  $N = 1,024$  ( $N_0$ ) (top row, panels a, b, and c) and  $N = 16,384$  ( $16N_0$ ) and  $\mu = 1 \times 10^{-6}$  ( $\mu_0$ ) per base pair for each mutation type (bottom row, panels d, e, and f). Lines represent the mean values across the 50 simulations, and the shaded areas represent the standard deviations.

rate, or both, even in case of mutational biases. These 2 factors can nevertheless be distinguished, as they have different effects on the coding and noncoding sequences of the genome. Their combination in various proportions can create a broad range of alternative patterns of genome size and coding density. In particular, by playing independently on mutation rate and population size, our model can reproduce the 2 extreme but different cases of genome size reduction that are seen in some endosymbionts and cyanobacteria. As an example, *Prochlorococcus marinus* is known to have lost both some parts of its coding and noncoding genome, although in different proportion such that its coding density has increased (Dufresne et al. 2005; Batut et al. 2014; Giovannoni et al. 2014). In our model, this would correspond to a population undergoing an increase in population size and a slight increase in mutation rate, which is coherent with the scientific literature on *Prochlorococcus marinus* (Hu and Blanchard 2008; Marais et al. 2008), although the large effective population size of this species has been recently debated (Chen et al. 2022; Filatov and Kirkpatrick 2024). On the other hand, *Buchnera aphidicola* has conserved its coding proportion but greatly reduced its total genome size (Moran and Mira 2001), which could be explained in our model by an increase in mutation rate and a decrease in population size, in similar proportions. This suggests that indirect selection for shorter genomes through robustness selection could be a key factor playing on genome evolution (Wilke et al. 2001; Gabzi et al. 2022), and especially on the evolution of genome size and structure.

Our observations confirm those made by Lynch and Conery (2003), namely that an increased genetic drift, here associated with a decreased population size, increases the genome size. Our results also point toward an equilibrium genome size: a sufficient number of genes makes it possible to fine-tune the phenotype to the environment, but the genome also has to be short enough to prevent the degeneration caused by an excess of chromosomal rearrangements (Knibbe et al. 2007; LaBar and Adami 2020). Increasing the mutation rate or the population size displaces this equilibrium toward shorter genomes, either through a more efficient genome purification of noncoding sequences (when increasing  $N$ ) or a loss of both coding and noncoding sequences to recover a minimal level of robustness (when increasing  $\mu$ ). Of course, mutational biases (regarding the balance between insertions and duplications versus deletions) also play an important role in determining the equilibrium genome size. In particular, deletion biases have been suggested as one main reason explaining why bacterial genomes remain small (Mira et al. 2001). However, we show here that, because of the indirect selection for robustness, a deletion bias is not needed to prevent a runaway inflation in the size of genomes. Instead, selection for robustness provides a counteracting force that increases with genome size, eventually offsetting any underlying bias in favor of insertions or duplications. Importantly, this indirect selection was not postulated in the model but emerged spontaneously in the simulations.

We propose an evolutionary mechanism consisting of a trade-off between direct selection for phenotypical

adaptation and indirect selection for replicative robustness. In this respect, mutations appear to be a weak selective force, as pointed out by Lynch and Walsh (2007). However, the emphasis was previously on the mutational targets contributed by genomic features, such as introns. Here, we emphasize another aspect, which seems to have been overseen thus far: any nonfunctional DNA represents an additional target for initiating macroscopic mutational events that can eventually impact the coding genome. This mechanism requires no additional hypotheses and is very general. It should therefore be pervasive in the living world.

Sung et al. (2012) have observed that, in real populations, the mutation rate scales negatively with both the population size and the amount of coding DNA. They propose that this is a consequence of selection for lower per-base mutation rates induced by the amount of coding DNA. Here, thanks to the use of fixed mutation rates, we have shown that the mutation rate can select the amount of DNA, including both the coding and noncoding compartments. This points towards the per-genome mutation rate being the relevant value, which can evolve due to changes in genome size and per-base mutation rate. This calls for further experiments in which both the genome size and the per-base mutation rate would be allowed to evolve, to study their relative speed of adaptation and their contribution to the variation of the per-genome mutation rate.

Although our main focus was on the final equilibrium reached by the populations after a change in  $N$  or  $\mu$ , our observations are broader than the end equilibrium as we can observe the temporal dynamics (Fig. 6 and [supplementary S3–S15, Supplementary Material](#) online). In particular, we observe that, when the mutation rate increases strongly, the fitness immediately drops drastically (Fig. 6b). This can be related to an error-threshold crossing mechanism (Eigen 1971; Takeuchi and Hogeweg 2007; de Boer and Hogeweg 2010): individuals can no longer pass on to their descendants all the information contained in their genome. They therefore lose fitness, and the lineage that survives in the long term is the one where genomes greatly reduced in size in the early phase of the experiment, thus reducing the number of mutations per replication event and finally reaching a point at which the information can be passed on reliably. The detailed aspects of these temporal dynamics could be the focus of future work. Indeed, it has been shown that genome reduction in endosymbionts occurred very quickly after the endosymbiosis became effective (Moran 2003; Wernegreen 2015), which is also what we observed in our data (Fig. 6).

In our experiments,  $N \times \mu$  stands out as a determining factor of some (although not all) aspects of genome structure, as isoclines of identical  $N \times \mu$  values display similar coding densities, even in the case of reduced genomes or mutational biases. Understanding this invariant is one of

the most exciting perspectives opened by our work. Its importance has already been highlighted by Lynch et al. (2006) in organelles, but our results suggest that this joined factor of drift and mutational pressure is a determinant of genome evolution throughout the tree of life. Notably, there is a small variation in coding fraction along  $N13:17\mu$  isoclines, which could be due to our use here of population size ( $N$ ) instead of effective population size ( $N_e$ ). Indeed, in our setup, the competition is local and thus  $N_e$  is slightly greater than  $N$ , but this relationship is not linear (see [supplementary material S1, Supplementary Material](#) online). Further versions of the model could rely on various measures of the effective population size to reach more accurate predictions, but we believe that our results can be interpreted nonetheless, as changes in population size and in effective population size are very similar over the range of population sizes tested here (see [supplementary material S1, Supplementary Material](#) online).

In order to allow for a fair quantitative comparison between the effect of mutation rates and population size, the amplitudes of the variations applied to the 2 parameters were similar in our experiments. In biological species, the range of variation in mutation rates is much narrower than the range of variation in effective population size, as shown by Lynch et al. (2023). Hence, given our explanatory mechanism, the observed range of variations in genome size is likely to be driven mainly by changes in  $N$ . However, our results show that  $\mu$  and  $N$  do not play an identical role. Indeed, variations in  $N$  change solely the noncoding size of the genome, while the variation in  $\mu$  impacts both the coding and the noncoding sizes. Therefore, even a small variation in  $\mu$  compared with a variation in  $N$  could be significant in determining genome architecture trajectories. This highlights that the correlation of  $N$  and genome size is not enough to understand genome evolution and that  $\mu$ , as well as any underlying mutational bias, also needs to be taken into account as a determining factor.

In this paper, we specifically focused on the effect of the variation in population size and mutation rates on genome size. Of course, it does not imply that the mechanism we identified is the only one, and various additional ones can also impact genome size evolution. For instance, there can be a limitation in available resources for nucleotide production, constraining the total genome size (Ngugi et al. 2023). In the case of endosymbiosis, exchanges can also happen between the host and the endosymbiont genomes, hence contributing to its streamlining (Bock 2017). Recombination could also further complicate the picture by adding a new type of mutation with unexpected interactions. More importantly, mobile genetic elements, and transposable elements (TE) in particular, are often proposed as one of the main drivers of genome expansion (Marino et al. 2024), especially in populations with small effective population sizes that could not eliminate them efficiently

due to the low selective pressure (Lynch and Conery 2003). TE invasions have been shown to increase dramatically genome size in eukaryotes (Kidwell 2002; Oggenfuss et al. 2021), although van Dijk et al. (2022) have demonstrated that they can also lead to streamlining in prokaryotes because genome reduction prevents their invasion. We did not test their impact here, but our results show that the effect of the variations in population size and mutation rate is conserved, even in case of a strong insertion bias (Fig. 4 and [supplementary figure S2, Supplementary Material](#) online). This enables us to conjecture that mobile elements would change the equilibrium genome size (as observed in our simulations, Figs. 4 and [supplementary figure S2, Supplementary Material](#) online), and probably drastically increase the variance of observed sizes, but that they are unlikely to change the response of genome size evolution to changes in  $\mu$  or  $N$ . This remains however to be tested.

To conclude, our experiments show that genome size reduction can occur in 2 very different conditions for bacteria. On the one hand, a very large population size promotes a more efficient selection in the face of random drift, which in turn enhances the robustness of genomes by decreasing their noncoding load. This corresponds to streamlining and leads to genomes with a high coding density. On the other hand, a higher mutation rate results in an instantaneous decrease in the robustness of genomes in the entire population, making the selection for robustness transiently stronger than the selection for phenotypical adaptation. The genome then shrinks rapidly, with both coding and noncoding sequences being discarded until a new robustness equilibrium is reached, all this at a substantial initial cost in phenotypical adaptation. This corresponds to a decaying genome and is compatible with empirical observations in endosymbiotic bacteria (Moran 2003). Strikingly, this remains true even in the presence of a mutational bias. Although the model that we propose here, of a balance between selection for robustness and selection for phenotypical adaptation, can explain the tendencies we observe and the final genome structures in our populations, further work is needed to understand the transient regimes and the mechanisms behind the constant coding fraction along the  $N \times \mu$  isoclines.

## Materials and Methods

### The Aevol Framework

Aevol (Knibbe et al. 2007; Banse et al. 2023) is an individual-based forward-in-time simulation software that has been specifically designed to study the evolution of genome structure. It emulates a population that is composed of a fixed number of individuals on a grid (Fig. 7a). Each individual owns a double-stranded circular genomic sequence, composed of 0s and 1s. To compute

the phenotype, sequences on the genome are recognized as promoters and mark the start of transcription, which stops when a sequence able to form a hairpin structure is encountered. On RNAs, Shine-Dalgarno-like sequences followed by a START codon mark the beginning of translation. The RNA sequence is then read 3 bases at a time until a STOP codon is encountered on the same reading frame. An artificial genetic code allows for each sequence of codons to be converted into a mathematical function, and the sum of all functions encoded on the genome defines the phenotype of the individual (Fig. 7b). The distance between this function and a target function, which represents the ideal phenotype in the specified environment, gives the fitness of the individual with a scaling factor  $k$  that tunes the strength of the selection. A detailed explanation can be found on the dedicated website [www.aevol.fr](http://www.aevol.fr).

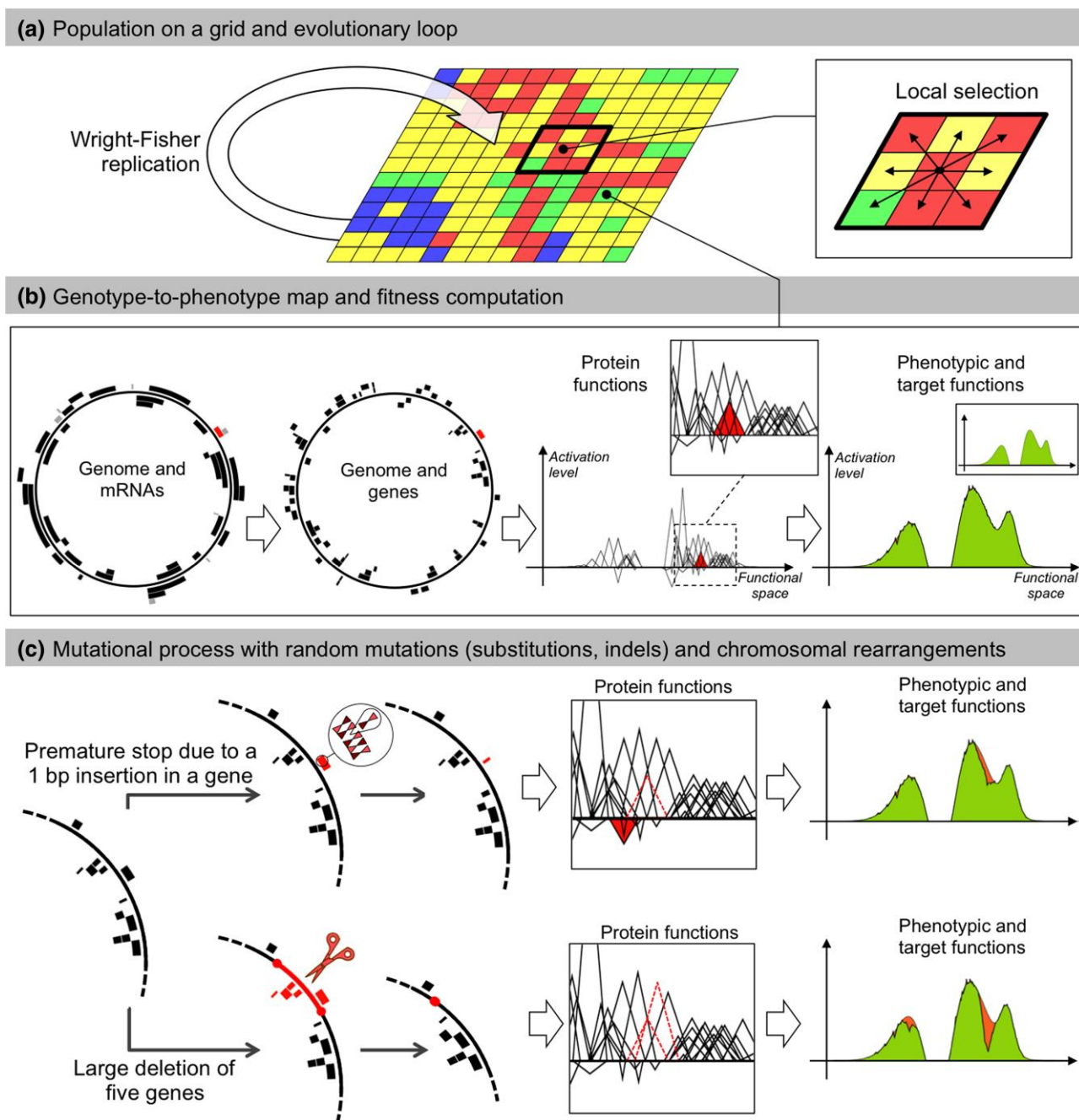
All individuals are replaced at each generation following a spatialized Wright–Fisher model. The number of descendants of each individual depends on its fitness difference with its neighbors. At each reproduction event, point mutations or genomic rearrangements can occur (Fig. 7c). They create diversity in the genomes, hence in the phenotypes, and allow the genome size and structure to change. These changes can be neutral or not, depending on whether mutations alter coding and/or noncoding sequences. These changes do not have a predefined effect on the fitness of the offspring as their genomes will be decoded thereafter, thus the model does not impose an a priori genome structure and allows us to study the evolution of genome architecture in various experimental conditions.

The mutation rate (in  $bp^{-1}$ ) is set for each type of mutation independently. When all mutation rates are equal, there is an equal probability of losing or gaining base pairs. The size distribution of InDels is uniform in  $[1, 6]$ , and the size distribution of large deletions and duplications is uniform in  $[1, L]$  (with  $L$  the genome length).

## Experimental Design

### Wild Types

In order to observe changes in genome architecture induced by changes in the population size and/or mutation rates, we begin our experiments from pre-evolved organisms, which are called “WT”. Having already evolved for millions of generations under constant conditions, WTs are very stable in genome structure and well adapted to their environment (although the fitness never stops increasing). Five different WTs were used for our experiments, all having evolved for 10 million generations at the basal conditions of  $N_0 = 1,024$  individuals and a mutation rate of  $\mu_0 = 10^{-6}$  mutations per base pair per generation for each type of mutations (point mutations, small insertions, small deletions, inversions, duplications, large deletions, and translocations). Importantly, in this experiment, all



**Fig. 7.** The Aevol model. a) Individuals are distributed on a grid. At each generation, the whole population replicates according to a Wright–Fisher replication model, in which selection operates locally within a  $3 \times 3$  neighborhood. b) Each grid cell contains a single organism described by its genome. Genomes are decoded through a genotype-to-phenotype map with 4 main steps (transcription, translation, computation of protein functions, and computation of the phenotype). Here, for illustration purposes, a random gene and the corresponding mRNA are colored in red. The red triangle represents the function of this gene in the mathematical world of the model. The phenotypic function is calculated by summing all protein functions. The phenotype is then compared with a predefined target (in green) to compute the fitness. The individual presented here has evolved in the model during 500,000 generations. c) Individuals may undergo mutations during replication. Two example mutations are shown: A small insertion (top) and a large deletion (bottom). Top: A 1 bp insertion occurs within a gene. It causes a frameshift, creating a premature stop codon. The ancestral function of the gene is lost (dashed triangle) and the truncated protein has a deleterious effect (red triangle). This leads to a greater divergence between the phenotype and the target (orange area on the phenotype). Bottom: The deletion removes 5 genes. The functions of 2 of them can be seen in the box (dotted triangles). This results in a large discrepancy between the phenotype and the target (orange area on the phenotype).

**Table 2.** Experimental conditions tested

Population size	Mutation rate (per base pair, per mutation type)	$N \times \mu$ product
64 ( $N_0/16$ )	$10^{-6}$ ( $\mu_0$ )	$1/16N_0 \times \mu_0$
256 ( $N_0/4$ )	$10^{-6}$ ( $\mu_0$ )	$1/4N_0 \times \mu_0$
<b>1024 (<math>N_0</math>)</b>	<b><math>10^{-6}</math> (<math>\mu_0</math>)</b>	<b><math>N_0 \times \mu_0</math></b>
529 ( $\approx N_0/2$ )	$2 \times 10^{-6}$ ( $2 \times \mu_0$ )	$\approx N_0 \times \mu_0$
256 ( $N_0/4$ )	$4 \times 10^{-6}$ ( $4 \times \mu_0$ )	$N_0 \times \mu_0$
64 ( $N_0/16$ )	$16 \times 10^{-6}$ ( $16 \times \mu_0$ )	$N_0 \times \mu_0$
2,025 ( $\approx 2 \times N_0$ )	$2 \times 10^{-6}$ ( $2 \times \mu_0$ )	$\approx 4N_0 \times \mu_0$
4,096 ( $4 \times N_0$ )	$10^{-6}$ ( $\mu_0$ )	$4N_0 \times \mu_0$
1,024 ( $N_0$ )	$4 \times 10^{-6}$ ( $4 \times \mu_0$ )	$4N_0 \times \mu_0$
4,096 ( $4 \times N_0$ )	$4 \times 10^{-6}$ ( $4 \times \mu_0$ )	$16N_0 \times \mu_0$
16,384 ( $16 \times N_0$ )	$10^{-6}$ ( $\mu_0$ )	$16N_0 \times \mu_0$
1,024 ( $N_0$ )	$16 \times 10^{-6}$ ( $16 \times \mu_0$ )	$16N_0 \times \mu_0$
16,384 ( $16 \times N_0$ )	$16 \times 10^{-6}$ ( $16 \times \mu_0$ )	$256N_0 \times \mu_0$

The control condition is in bold. Note that, as the simulations take place on a squared grid, population sizes could not be exactly divided or multiplied by 2.

types of mutations are equally probable: there is no mutational bias toward the insertion or deletion of base pairs. Bacterial populations are very large and cannot be directly modeled owing to computational load. We hence limit the population sizes in our experiments, but compensate by increasing the mutation rates such that the  $N \times \mu$  parameter is of the same order of magnitude as for real bacterial populations. Finally, to limit the effect of drift, we used a selection strength  $k = 1,000$ , which is relatively high and guarantees an efficient selection. The fitnesses and genome structures of the WTs are listed in Table 1.

### Experimental Conditions

A range of population sizes increases or decreases and mutation rates increases, as well as some combinations of both, are tested. All conditions are listed in Table 2 below. For each combination of conditions, 10 replications of each of the 5 WTs are run. Initial populations are always clonal: all individuals are identical to the specific WT used for the run.

### Data Analyses

To analyze the simulations, we reconstruct the ancestral lineages of the final populations. To this end, simulations are run for 2,100,000 generations, and we identify all the ancestors of a random individual of the final population. We then study the data from generation 0 to generation 2,000,000 and ignore the last 100,000 to ensure that the final population has coalesced and that we study the lineage of the whole final population.

On this lineage, we retrieve the fitness, coding, and non-coding genome size at each generation, as well as the replicative robustness every 1,000 generations. The replicative robustness is measured as the proportion of the offspring of an individual that has the exact same fitness as its parent,

*i.e.* that underwent no mutation at all, or only purely neutral mutations. To estimate replicative robustness for a given individual of the lineage, we generate 10,000 offsprings and compare them to their parent.

To compare experimental conditions, we retrieve the individual at generation 2,000,000 in each lineage. This individual is the common ancestor of the final population (at generation 2, 100, 000), thus ensuring that its genome structure has been conserved by evolution. A visualization of the temporal lineage data (fitness, coding fraction and total, coding, and noncoding genome sizes) for the 50 replicas of each experimental condition is provided in [supplementary S3 \(Figures S3–S15\), Supplementary Material online](#).

### Effect of Mutational Biases

As it is often assumed that mutational biases—toward deletions for bacteria and toward insertions for eukaryotes—are very important for genome size evolution (Petrov 2002), we also tried to confront our experiments to the impact of mutational biases. We tested 4 mutational biases: twice as many large deletions than duplications, twice as many small deletions than small insertions, twice as many duplications than large deletions, and twice as many small insertions than small deletions. In all cases, the sum of all mutation rates is conserved, such that the overall mutational pressure is the same as in the previous experiments.

For each mutational condition, 5 WT evolved for 10,000,000 generations. Then, the median-sized WT of each mutational condition was extracted and confronted with either an increase or decrease in population size ( $4 \times N_0$ ,  $N_0/4$ ) or an increase in all mutation rates proportionally ( $4 \times \mu_0$ —note that, in case of bias,  $\mu_0$  may be different for the different types of mutation) for 2,100,000 generations. By extracting the ancestor of the lineage at generation 2,000,000, we could compare these experiments to the control conditions (where the population size and mutation rates remained stable for 2,100,000 generations).

### Supplementary Material

Supplementary material is available at *Genome Biology and Evolution* online.

### Acknowledgments

The authors would like to thank Laurent Duret and David P. Parsons for fruitful comments on the manuscript.

### Funding

The authors acknowledge the support of the French Agence Nationale de la Recherche (ANR), under grant no. ANR-20-CE02-0008 (NeGA project). J.L., G.B., and N.L. would like to thank the Rhône-Alpes Institute for

Complex Systems (IXXI) for funding. All authors thank the Grid'5000 testbed, supported by a scientific interest group hosted by Inria and including CNRS, RENATER and several Universities as well as other organizations (see <https://www.grid5000.fr>), for computational support.

## Conflict of interest

The authors declare no competing interests.

## Data Availability

The code of Aevol is available on GitLab at <https://gitlab.inria.fr/aeval/aeval>. WTs sequences to reproduce the experiments, as well as the full lineages data and robustness data, are available on Zenodo: <https://doi.org/10.5281/zenodo.10669479>.

## References

- Adami C. Digital genetics: unravelling the genetic basis of evolution. *Nat Rev Genet.* 2006;7(2):109–118. <https://doi.org/10.1038/nrg1771>.
- Almpanis A, Swain M, Gatherer D, McEwan N. Correlation between bacterial G+C content, genome size and the G+C content of associated plasmids and bacteriophages. *Microb Genom.* 2018;4(4):e000168. <https://doi.org/10.1099/mgen.0.000168>.
- Bachmann K. Genome size in mammals. *Chromosoma.* 1972;37(1):85–93. <https://doi.org/10.1007/BF00329560>.
- Banse P, Luiselli J, Parsons DP, Grohens T, Foley M, Trujillo L, Rouzaud-Cornabas J, Knibbe C, Beslon G. Forward-in-time simulation of chromosomal rearrangements: the invisible backbone that sustains long-term adaptation. *Mol Ecol.* 2023;00:1–14. <https://doi.org/10.1111/mec.17234>.
- Barow M, Meister A. Lack of correlation between AT frequency and genome size in higher plants and the effect of nonrandomness of base sequences on dye binding. *Cytometry.* 2002;47(1):1–7. <https://doi.org/10.1002/cyto.v47:1>.
- Batut B, Knibbe C, Marais G, Daubin V. Reductive genome evolution at both ends of the bacterial population size spectrum. *Nat Rev Microbiol.* 2014;12(12):841–850. <https://doi.org/10.1038/nrmicro3331>.
- Batut B, Parsons DP, Fischer S, Beslon G, Knibbe C. In silico experimental evolution: a tool to test evolutionary scenarios. *BMC Bioinformatics.* 2013;14(S15):S11. <https://doi.org/10.1186/1471-2105-14-S15-S11>.
- Bingham EP, Ratcliff WC. A nonadaptive explanation for macroevolutionary patterns in the evolution of complex multicellularity. *Proc Natl Acad Sci U S A.* 2024;121(7):e2319840121. <https://doi.org/10.1073/pnas.2319840121>.
- Bock R. Witnessing genome evolution: experimental reconstruction of endosymbiotic and horizontal gene transfer. *Annu Rev Genet.* 2017;51(1):1–22. <https://doi.org/10.1146/genet.2017.51.issue-1>.
- Bourguignon T, Kinjo Y, Villa-Martin P, Coleman NV, Tang Q, Arab DA, Wang Z, Tokuda G, Hongoh Y, Ohkuma M, et al. Increased mutation rate is linked to genome reduction in prokaryotes. *Curr Biol.* 2020;30(19):3848–3855. <https://doi.org/10.1016/j.cub.2020.07.034>.
- Brevet M, Lartillot N. Reconstructing the history of variation in effective population size along phylogenies. *Genome Biol Evol.* 2021;13(8):evab150. <https://doi.org/10.1093/gbe/evab150>.
- Chen Z, Wang X, Song Y, Zeng Q, Zhang Y, Luo H. *Prochlorococcus* have low global mutation rate and small effective population size. *Nat Ecol Evol.* 2022;6(2):183–194. <https://doi.org/10.1038/s41559-021-01591-0>.
- Choi I-Y, Kwon E-C, Kim N-S. The C- and G-value paradox with polyploidy, repeatomes, introns, phenomes and cell economy. *Genes Genomics.* 2020;42(7):699–714. <https://doi.org/10.1007/s13258-020-00941-9>.
- Chong RA, Park H, Moran NA. Genome evolution of the obligate endosymbiont *Buchnera aphidicola*. *Mol Biol Evol.* 2019;36(7):1481–1489. <https://doi.org/10.1093/molbev/msz082>.
- de Boer FK, Hogeweg P. Eco-evolutionary dynamics, coding structure and the information threshold. *BMC Evol Biol.* 2010;10(1):361. <https://doi.org/10.1186/1471-2148-10-361>.
- Dufresne A, Garczarek L, Partensky F. Accelerated evolution associated with genome reduction in a free-living prokaryote. *Genome Biol.* 2005;6(2):R14. <https://doi.org/10.1186/gb-2005-6-2-r14>.
- Eigen M. Selforganization of matter and the evolution of biological macromolecules. *Naturwissenschaften.* 1971;58(10):465–523. <https://doi.org/10.1007/BF00623322>.
- Figuet E, Nabholz B, Bonneau M, Mas Carrio E, Nadachowska-Brzyska K, Ellegren H, Galtier N. Life history traits, protein evolution, and the nearly neutral theory in amniotes. *Mol Biol Evol.* 2016;33(6):1517–1527. <https://doi.org/10.1093/molbev/msw033>.
- Filatov DA, Kirkpatrick M. How does evolution work in superabundant microbes? *Trends Microbiol.* 2024;32(9):836–846. <https://doi.org/10.1016/j.tim.2024.01.009>.
- Flombaum P, Gallegos JL, Gordillo RA, Rincón J, Zabala LL, Jiao N, Karl DM, Li WKW, Lomas MW, Veneziano D, et al. Present and future global distributions of the marine cyanobacteria *Prochlorococcus* and *Synechococcus*. *Proc Natl Acad Sci U S A.* 2013;110(24):9824–9829. <https://doi.org/10.1073/pnas.1307701110>.
- Gabzi T, Pilpel Y, Friedlander T. Fitness landscape analysis of a tRNA gene reveals that the wild type allele is sub-optimal, yet mutationally robust. *Mol Biol Evol.* 2022;39(9):msac178. <https://doi.org/10.1093/molbev/msac178>.
- Gago S, Elena SF, Flores R, Sanjuán R. Extremely high mutation rate of a hammerhead viroid. *Science.* 2009;323(5919):1308–1308. <https://doi.org/10.1126/science.1169202>.
- Giovanoni SJ, Cameron Thrash J, Temperton B. Implications of streamlining theory for microbial ecology. *ISME J.* 2014;8(8):1553–1565. <https://doi.org/10.1038/ismej.2014.60>.
- Giovanoni SJ, Tripp HJ, Givan S, Podar M, Vergin KL, Baptista D, Bibbs L, Eads J, Richardson TH, Noordewier M, et al. Genome streamlining in a cosmopolitan oceanic bacterium. *Science.* 2005;309(5738):1242–1245. <https://doi.org/10.1126/science.1114057>.
- Heddi A, Charles H, Khatchadourian C, Bonnot G, Nardon P. Molecular characterization of the principal symbiotic bacteria of the weevil *Sitophilus oryzae*: a peculiar G + C content of an endocytobiotic DNA. *J Mol Evol.* 1998;47(1):52–61. <https://doi.org/10.1007/PL00006362>.
- Hindré T, Knibbe C, Beslon G, Schneider D. New insights into bacterial adaptation through in vivo and in silico experimental evolution. *Nat Rev Microbiol.* 2012;10(5):352–365. <https://doi.org/10.1038/nrmicro2750>.
- Hu J, Blanchard JL. Environmental sequence data from the sargasso sea reveal that the characteristics of genome reduction in *prochlorococcus* are not a harbinger for an escalation in genetic drift. *Mol Biol Evol.* 2008;26(1):5–13. <https://doi.org/10.1093/molbev/msn217>.
- Kalhor R, Beslon G, Lafond M, Scornavacca C. A rigorous framework to classify the postduplication fate of paralogous genes. *J Comput Biol.* 2024;31(9):815–833. <https://doi.org/10.1089/cmb.2023.0331>.
- Kidwell MG. Transposable elements and the evolution of genome size in eukaryotes. *Genetica.* 2002;115(1):49–63. <https://doi.org/10.1023/A:1016072014259>.

- Knibbe C, Coulon A, Mazet O, Fayard J-M, Beslon G. A long-term evolutionary pressure on the amount of noncoding DNA. *Mol Biol Evol.* 2007;24(10):2344–2353. <https://doi.org/10.1093/molbev/msm165>.
- Konstantinidis KT, Tiedje JM. Trends between gene content and genome size in prokaryotic species with larger genomes. *Proc Natl Acad Sci U S A.* 2004;101(9):3160–3165. <https://doi.org/10.1073/pnas.0308653100>.
- Koskiniemi S, Sun S, Berg OG, Andersson DI. Selection-driven gene loss in bacteria. *PLoS Genet.* 2012;8(6):e1002787. <https://doi.org/10.1371/journal.pgen.1002787>.
- Kuo C-H, Moran NA, Ochman H. The consequences of genetic drift for bacterial genome complexity. *Genome Res.* 2009;19(8):1450–1454. <https://doi.org/10.1101/gr.091785.109>.
- Kuo C-H, Ochman H. Deletional bias across the three domains of life. *Genome Biol Evol.* 2009;1:145–152. <https://doi.org/10.1093/gbe/evp016>.
- LaBar T, Adami C. Genome size and the extinction of small populations. In: Cham: Springer International Publishing; 2020. p. 167–183.
- Lefebvre T, Morvan C, Malard F, François C, Konecny-Dupré L, Guéguen L, Weiss-Gayet M, Seguin-Orlando A, Ermini L, Sarkissian CD, et al. Less effective selection leads to larger genomes. *Genome Res.* 2017;27(6):1016–1028. <https://doi.org/10.1101/gr.212589.116>.
- Leth Bak A, Black FT, Christiansen C, Freundt EA. Genome size of mycoplasma DNA. *Nature.* 1969;224(5225):1209–1210. <https://doi.org/10.1038/2241209a0>.
- Lynch M. Streamlining and simplification of microbial genome architecture. *Annu Rev Microbiol.* 2006;60(1):327–349. <https://doi.org/10.1146/micro.2006.60.issue-1>.
- Lynch M. The frailty of adaptive hypotheses for the origins of organismal complexity. *Proc Natl Acad Sci U S A.* 2007;104(suppl\_1):8597–8604. <https://doi.org/10.1073/pnas.0702207104>.
- Lynch M, Ali F, Lin T, Wang Y, Ni J, Long H. The divergence of mutation rates and spectra across the tree of life. *EMBO Rep.* 2023;24(10):e57561. <https://doi.org/10.15252/embr.202357561>.
- Lynch M, Conery JS. The origins of genome complexity. *Science.* 2003;302(5649):1401–1404. <https://doi.org/10.1126/science.1089370>.
- Lynch M, Koskella B, Schaack S. Mutation pressure and the evolution of organelle genomic architecture. *Science.* 2006;311(5768):1727–1730. <https://doi.org/10.1126/science.1118884>.
- Lynch M, Walsh B. The origins of genome architecture. vol. 98. Sunderland (MA): Sinauer Associates; 2007.
- Marais GAB, Calteau A, Tenailon O. Mutation rate and genome reduction in endosymbiotic and free-living bacteria. *Genetica.* 2008;134(2):205–210. <https://doi.org/10.1007/s10709-007-9226-6>.
- Marino A, Debaecker G, Fiston-Lavier A-S, Haudry A, Nabholz B. Effective population size does not explain long-term variation in genome size and transposable element content in animals. 2024.
- Martinez-Gutierrez CA, Aylward FO. Genome size distributions in bacteria and archaea are strongly linked to evolutionary history at broad phylogenetic scales. *PLoS Genet.* 2022;18(5):e1010220. <https://doi.org/10.1371/journal.pgen.1010220>.
- Martinez-Cano DJ, Reyes-Prieto M, Martínez-Romero E, Partida-Martínez LP, Latorre A, Moya A, Delaye L. Evolution of small prokaryotic genomes. *Front Microbiol.* 2015;5:742. <https://doi.org/10.3389/fmicb.2014.00742>.
- Mira A, Moran N. Estimating population size and transmission bottlenecks in maternally transmitted endosymbiotic bacteria. *Microb Ecol.* 2002;44(2):137–143. <https://doi.org/10.1007/s00248-002-0012-9>.
- Mira A, Ochman H, Moran NA. Deletional bias and the evolution of bacterial genomes. *Trends Genet.* 2001;17(10):589–596. [https://doi.org/10.1016/S0168-9525\(01\)02447-7](https://doi.org/10.1016/S0168-9525(01)02447-7).
- Moran NA. Accelerated evolution and Muller's ratchet in endosymbiotic bacteria. *Proc Natl Acad Sci U S A.* 1996;93(7):2873–2878. <https://doi.org/10.1073/pnas.93.7.2873>.
- Moran NA. Microbial minimalism: genome reduction in bacterial pathogens. *Cell.* 2002;108(5):583–586. [https://doi.org/10.1016/S0092-8674\(02\)00665-7](https://doi.org/10.1016/S0092-8674(02)00665-7).
- Moran NA. Tracing the evolution of gene loss in obligate bacterial symbionts. *Curr Opin Microbiol.* 2003;6(5):512–518. <https://doi.org/10.1016/j.mib.2003.08.001>.
- Moran NA, Mira A. The process of genome shrinkage in the obligate symbiont *Buchnera aphidicola*. *Genome Biol.* 2001;2(12):1–12. <https://doi.org/10.1186/gb-2001-2-12-research0054>.
- Morris JJ, Lenski RE, Zinser ER. The black queen hypothesis: evolution of dependencies through adaptive gene loss. *mBio.* 2012;3(2):e00036–12. <https://doi.org/10.1128/mBio.00036-12>.
- Müller R, Kaj I, Mugal CF. A nearly neutral model of molecular signatures of natural selection after change in population size. *Genome Biol Evol.* 2022;14(5):evac058. <https://doi.org/10.1093/gbe/evac058>.
- Ngugi DK, Acinas SG, Sánchez P, Gasol JM, Agusti S, Karl DM, Duarte CM. Abiotic selection of microbial genome size in the global ocean. *Nat Commun.* 2023;14(1):1384. <https://doi.org/10.1038/s41467-023-36988-x>.
- Nilsson AI, Koskiniemi S, Eriksson S, Kugelberg E, Hinton JCD, Andersson DI. Bacterial genome size reduction by experimental evolution. *Proc Natl Acad Sci U S A.* 2005;102(34):12112–12116. <https://doi.org/10.1073/pnas.0503654102>.
- Oggenfuss U, Badet T, Wicker T, Hartmann FE, Singh NK, Abraham L, Karisto P, Vonlanthen T, Mundt C, McDonald BA, et al. A population-level invasion by transposable elements triggers genome expansion in a fungal pathogen. *Elife.* 2021;10:e69249. <https://doi.org/10.7554/eLife.69249>.
- Pellicier J, Fay MF, Leitch IJ. The largest eukaryotic genome of them all? *Bot J Linn Soc.* 2010;164(1):10–15. <https://doi.org/10.1111/boj.2010.164.issue-1>.
- Petrov DA. Evolution of genome size: new approaches to an old problem. *Trends Genet.* 2001;17(1):23–28. [https://doi.org/10.1016/S0168-9525\(00\)02157-0](https://doi.org/10.1016/S0168-9525(00)02157-0).
- Petrov DA. Mutational equilibrium model of genome size evolution. *Theor Popul Biol.* 2002;61(4):531–544. <https://doi.org/10.1006/tpb.2002.1605>.
- Popadin K, Polishchuk LV, Mamirova L, Knorre D, Gunbin K. Accumulation of slightly deleterious mutations in mitochondrial protein-coding genes of large versus small mammals. *Proc Natl Acad Sci U S A.* 2007;104(33):13390–13395. <https://doi.org/10.1073/pnas.0701256104>.
- Romiguier J, Ranwez V, Douzery E, Galtier N. Genomic evidence for large, long-lived ancestors to placental mammals. *Mol Biol Evol.* 2012;30(1):5–13. <https://doi.org/10.1093/molbev/mss211>.
- Schneider S, Perlova O, Kaiser O, Gerth K, Alici A, Altmeyer MO, Bartels D, Bekel T, Beyer S, Bode E, et al. Complete genome sequence of the myxobacterium *Sorangium cellulosum*. *Nat Biotechnol.* 2007;25(11):1281–1289. <https://doi.org/10.1038/nbt1354>.
- Sung W, Ackerman MS, Miller SF, Doak TG, Lynch M. Drift-barrier hypothesis and mutation-rate evolution. *Proc Natl Acad Sci U S A.* 2012;109(45):18488–18492. <https://doi.org/10.1073/pnas.1216223109>.
- Takeuchi N, Hogeweg P. Error-threshold exists in fitness landscapes with lethal mutants. *BMC Evol Biol.* 2007;7(1):15. <https://doi.org/10.1186/1471-2148-7-15>.
- van Dijk B, Bertels F, Stolk L, Takeuchi N, Rainey PB. Transposable elements promote the evolution of genome streamlining. *Philos Trans R Soc B.* 2022;377(1842):20200477. <https://doi.org/10.1098/rstb.2020.0477>.



- van Ham R, Kamerbeek J, Palacios C, Rausell C, Abascal F, Bastolla U, Fernández J, Jiménez L, Postigo M, Silva F, et al. Reductive genome evolution in *Buchnera aphidicola*. *Proc Natl Acad Sci U S A*. 2003;100(2):581–586. <https://doi.org/10.1073/pnas.0235981100>.
- Wernegreen JJ. Genome evolution in bacterial endosymbionts of insects. *Nat Rev Genet*. 2002;3(11):850–861. <https://doi.org/10.1038/nrg931>.
- Wernegreen JJ. Endosymbiont evolution: predictions from theory and surprises from genomes. *Ann N Y Acad Sci*. 2015;1360(1):16–35. <https://doi.org/10.1111/nyas.2015.1360.issue-1>.
- Westoby M, Nielsen DA, Gillings MR, Litchman E, Madin JS, Paulsen IT, Tetu SG. Cell size, genome size, and maximum growth rate are near-independent dimensions of ecological variation across bacteria and archaea. *Ecol Evol*. 2021;11(9):3956–3976. <https://doi.org/10.1002/ece3.v11.9>.
- Wilke CO, Adami C. Evolution of mutational robustness. *Mutat Res*. 2003;522(1-2):3–11. [https://doi.org/10.1016/S0027-5107\(02\)00307-X](https://doi.org/10.1016/S0027-5107(02)00307-X).
- Wilke CO, Wang JL, Ofria C, Lenski RE, Adami C. Evolution of digital organisms at high mutation rates leads to survival of the flattest. *Nature*. 2001;412(6844):331–333. <https://doi.org/10.1038/35085569>.
- Wolf YI, Koonin EV. Genome reduction as the dominant mode of evolution. *Bioessays*. 2013;35(9):829–837. <https://doi.org/10.1002/bies.v35.9>.
- Associate editor:** Michael Lynch