

The mutational landscape of *Bacillus subtilis* conditional hypermutators shows how proofreading skews DNA polymerase error rates

Ira Tanneur^{1,2}, Etienne Dervyn¹, Cyprien Guérin², Guillaume Kon Kam King²,
Matthieu Jules^{1,*}, Pierre Nicolas^{2,*}

¹Université Paris-Saclay, INRAE, AgroParisTech, Micalis Institute, 78350 Jouy-en-Josas, France

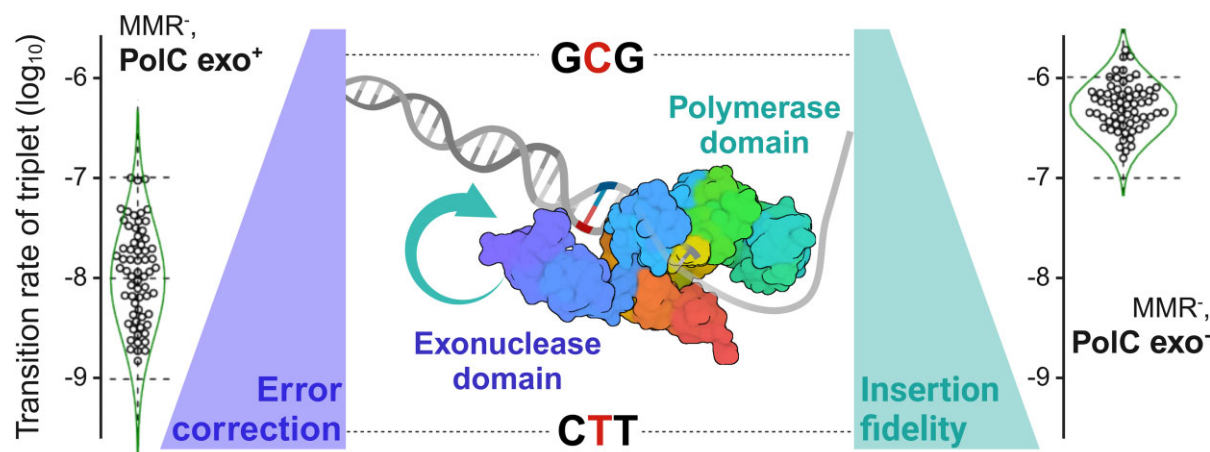
²Université Paris-Saclay, INRAE, MalAGE, 78350 Jouy-en-Josas, France

*To whom correspondence should be addressed: Email: matthieu.jules@inrae.fr
Correspondence may also be addressed to Pierre Nicolas. Email: pierre.nicolas@inrae.fr

Abstract

Polymerase errors during DNA replication are a major source of point mutations in genomes. The spontaneous mutation rate also depends on the counteracting activity of DNA repair mechanisms, with mutator phenotypes appearing constantly and allowing for periods of rapid evolution in nature and in the laboratory. Here, we use the Gram-positive model bacterium *Bacillus subtilis* to disentangle the contributions of DNA polymerase initial nucleotide selectivity, DNA polymerase proofreading, and mismatch repair (MMR) to the mutation rate. To achieve this, we constructed several conditional hypermutators with a proofreading-deficient allele of *polC* and/or a deficient allele of *mutL* and performed mutation accumulation experiments. These conditional hypermutators enrich the *B. subtilis* synthetic biology toolbox for directed evolution. Using mathematical models, we investigated how to interpret the apparent probabilities with which errors escape MMR and proofreading, highlighting the difficulties of working with counts that aggregate potentially heterogeneous mutations and with unknowns about the pathways leading to mutations in the wild-type. Aware of these difficulties, the analysis shows that proofreading prevents partial saturation of the MMR in *B. subtilis* and that an inherent drawback of proofreading is to skew the net polymerase error rates by amplifying intrinsic biases in nucleotide selectivity.

Graphical abstract



Introduction

Substitutions, insertions and deletions of single base pairs in genomes can have diverse consequences on encoded molecular functions, from no effect to abrupt change, most often in the direction of deterioration [1]. As such, point mutations are a constant threat to the integrity of the genetic information, even if they are also essential to adaptive evolution. In the long term, point mutation rates themselves result from an

evolutionary process. An hypothesis, supported by the comparison of mutation rates across the tree of life, postulates that they are simply maintained as low as possible; the limit being the “drift-barrier,” where the strength of selection is matched by the opposite pressure of random genetic drift and mutation, presumably biased toward creating weak mutators [2]. Mutation rate might thus have nothing to do with the benefit of evolvability for long-term survival, the energetic cost of

Received: March 21, 2024. Revised: February 3, 2025. Editorial Decision: February 5, 2025. Accepted: February 20, 2025

© The Author(s) 2025. Published by Oxford University Press on behalf of Nucleic Acids Research.

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.

fidelity, or a biophysical limit. In bacteria, this rate is already typically as low as one mutation introduced in the genome per one thousand generations, but mutation rates one order of magnitude below those observed in nature were reported to arise under scenarios of artificial evolution [3, 4].

However, mutator phenotypes occur everywhere and their contribution to evolution is difficult to estimate [5, 6]. In asexual populations, where mutator alleles and the mutations that they generate remain linked across successive generations, mutators are advantageous when the potential for fitness improvement is high. For example, in pathogenic bacteria, mutators have been associated with complex antibiotic resistance [7], rapid evolution within the host during infection [8], and atypical virulence traits [9]. Mutators also emerge spontaneously during laboratory evolution in response to applied selective pressures [10, 11]. To foster adaptation, conditional mutator circuits and systems have been engineered for some organisms as part of the synthetic biology toolbox [12–14]. Understanding the molecular factors that determine mutation rates is therefore of strong fundamental and applied interest.

The sources of spontaneous mutations in living cells are diverse, including DNA lesions caused by endogenous and exogenous agents, errors introduced by the DNA polymerase during replication, and the activity of error-prone polymerases recruited in response to stress [15]. The resulting mutation rate depends on the intensity of these sources and of the counteracting activity of DNA repair mechanisms, which work in concert to ensure that the correct genetic material is passed on to daughter cells. Two essential mechanisms, conserved from prokaryotes to eukaryotes, ensure accurate repair of both bulky and non-bulky lesions resulting from DNA damage, including those induced by reactive oxygen species, a major source of DNA errors [16]. Nucleotide excision repairs various drug- and UV-induced lesions (i.e. bulky lesions), while base excision repair lesions caused by various chemical assaults, such as alkylation, oxidation, deamination, etc. (i.e. non-bulky lesions). The accuracy of DNA replication depends on three critical mechanisms: the initial selectivity of the DNA polymerase, which is responsible for inserting the correct nucleotide; the proofreading, which removes misincorporated nucleotides through polymerase-associated exonucleases; and the mismatch repair (MMR), which adds a second layer of error correction shortly after replication [17].

In bacteria, genome replication is carried out by a multiprotein machine classified into the C-family of DNA polymerase holoenzymes, in which the catalytic polymerase α -subunit exists in two primary forms, DnaE and PolC [18]. A representative example of DnaE is found in the extensively studied Gram-negative model bacterium *Escherichia coli*. Conversely, PolC is predominant in low-GC Gram-positive bacteria, such as *Bacillus subtilis* [19]. In this organism, replication elongation involves two essential polymerases, PolC and DnaE [20], which have distinct functions: PolC does most of the DNA synthesis, but only DnaE can elongate from RNA primers on the lagging strand before passing the DNA fragment to PolC. In the absence of proofreading and MMR, the error rate of the *E. coli* α -subunit replication machinery has been estimated to be approximately 10^{-6} per base pair per generation both *in vitro* [21] and *in vivo* [22]. Given this error rate and the size of the *E. coli* genome, about 5 mutations are expected to be introduced per generation.

The exonuclease domain essential for proofreading is encoded as an integral part of the vast majority of PolC poly-

merases [18], including *B. subtilis* PolC. In contrast, DnaE polymerases do not have their own exonuclease domain. The proofreading activity of the *E. coli* DNA PolIII holoenzyme containing DnaE is based on an exonuclease domain located in the ϵ -subunit. Errors made by polymerases devoid of proofreading can also sometimes be corrected by a process known as proofreading *in trans*, or extrinsic proofreading, which is well described between eukaryotic DNA polymerases [23]. In *B. subtilis*, data suggest that the PolC exonuclease is able to correct errors made by the error-prone DnaE polymerase [24, 25].

The MMR is a universal mechanism that is responsible for correcting errors that occur during DNA replication and escape proofreading. Upon detection of a replication error, the mismatch-sensing protein, MutS, recruits MutL. Most prokaryotic and eukaryotic MutL homologs, from humans to bacteria, possess a highly conserved endonuclease active site that serves to remove mismatches [26, 27]. *E. coli* has been the primary model for studying MMR, but its MutL lacks the endonuclease activity encoded here in a distinct protein, MutH, which specifically nicks the unmethylated and thus nascent strand [28]. In the absence of MutH and Dam methylation, the process that guides MutL to the nascent strand remains unclear in most prokaryotes and all eukaryotes [29]. In bacteria, MMR increases the fidelity of the chromosomal DNA replication pathway approximately 100-fold, and MMR is considered to be a system directed toward the repair of the most frequent replication errors [30]. Mutator phenotypes found in nature are often caused by mutations that inactivate the MMR.

Studying organisms such as *B. subtilis*, with a PolC polymerase and an MMR pathway widely conserved across biology is important to understand the coordinated functioning of these systems [31]. Extending previous work characterizing the mutation profiles of MMR-deficient *B. subtilis* strains [32, 33], the main goal of our study was to jointly assess the respective contributions of initial nucleotide selectivity, proofreading and MMR to the mutation rates and their interdependencies. For this purpose, we constructed and analyzed several conditional hypermutators. Analysis of the data in light of several mathematical models suggests that proofreading is needed to avoid partial saturation of the MMR, as previously reported in *E. coli* [22, 34], but also that an inherent effect of proofreading is to skew the net polymerase error rate. With their wide range of mutation rates and contrasting mutation profiles, the conditional hypermutators also enrich the *B. subtilis* synthetic biology toolbox for directed evolution.

Materials and methods

Media and bacterial strains

E. coli DH5 α was used for plasmid construction and transformation using standard techniques [35]. The *B. subtilis* strains used in this study were derived from our Master Strain (MS), a prophage-free and *trp*⁺ derivative of *B. subtilis* 168 [4], denoted here R¹⁶⁸. Lysogeny broth (LB) was used to grow *E. coli* and *B. subtilis*. Transformation of *B. subtilis* cells was performed using the protocol of [36]. When relevant, the media were supplemented with the following antibiotics: ampicillin 100 $\mu\text{g.mL}^{-1}$ for *E. coli* and spectinomycin 100 $\mu\text{g.mL}^{-1}$ or kanamycin 5 $\mu\text{g.mL}^{-1}$ for *B. subtilis*.

Construction of hypermutator strains

Genomic DNAs from the $\Delta mutS::kan$ and $\Delta mutL::kan$ mutant strains, previously constructed in [37], served as a template for polymerase chain reaction (PCR) amplification of the *mutS mutL* genome region with the P1-P2 primer pair (Supplementary Table S1). The resulting PCR products were then used to transform MS by homologous recombination, and ΔS and ΔL mutants were selected on kanamycin resistance.

To construct the L^* strain, the first and second halves of the *mutL* gene were PCR-amplified using primer pairs P5-P8 and P6-P7, respectively (Supplementary Fig. S1), where P7 and P8 both carry the desired point mutation (as indicated in Supplementary Table S1). The two fragments were then assembled by PCR, resulting in the *mutL* (N34H) allele. The backbone of the pDR111 plasmid (kind gift of D. Rüdner) containing the isopropyl- β -D-1-thiogalactopyranoside (IPTG) inducible $P_{hyperspank}$ promoter (denoted here P_{hs}) and the *spec* gene (conferring resistance to spectinomycin), was PCR-amplified using P3 and P4. The 5' extensions of P5 and P6 then allowed the assembly of the *mutL* (N34H) allele with the PCR-amplified pDR111 using the HiFi DNA assembly protocol (New England Biolabs, USA). This resulted in cloning the *mutL* (N34H) allele under the control of P_{hs} into a *B. subtilis amyE*-integrative plasmid (Supplementary Fig. S1).

To construct the C^* strain, we used the *polC* allele found in *B. subtilis mut-1* [19, 38] in a version that contains only two point mutations: *mut-1A* (G430E) located within the exonuclease domain of the PolC protein, spanning amino acids 412–617, and *mut-1B* (S621N) which lies just beyond this domain [39, 40]. This allele was PCR-amplified using the primer pair P11-P12 and assembled to the PCR-amplified pDR111 (using P3 and P4) using the HiFi DNA assembly protocol (New England Biolabs, USA). This resulted in cloning the *polC mut-1* allele under the control of P_{hs} into a *B. subtilis amyE*-integrative plasmid (Supplementary Fig. S2).

For the construction of the LC^* strain, a *mutL^* polC^** synthetic operon was generated by assembling the *mutL* (N34H) allele PCR-amplified from strain L^* using P5 and P11, and the *polC mut-1* allele PCR-amplified from C^* using P12 and P10, to the PCR-amplified pDR111 (using P3 and P4) using the HiFi DNA assembly protocol (New England Biolabs, USA). This resulted in the cloning of the *mutL^* polC^** synthetic operon under the control of P_{hs} into a *B. subtilis amyE*-integrative plasmid (Supplementary Fig. S2).

Plasmids were transformed into the *B. subtilis amyE* locus by double recombination events. All strains were verified by sequencing, and transcriptomics experiments were performed to compare global gene expression. The RNA-seq reads and detailed protocols and results have been deposited in GEO.

Fluctuation assays

For each strain to be tested, a single colony was grown in 1 mL LB at 37°C for 90 min. This preculture was serially diluted in fresh LB to start cultures with a small number of cells (N_0). Cells were then grown for 7.5 h to reach saturation. When induction with IPTG was tested, LB medium with the desired concentration of IPTG was prepared from an IPTG stock concentration of 1 mM just before use. When the culture volume was 1 mL, the cultures were centrifuged before plating to retain the cells, and 750 μ L of supernatant was removed. The remaining 250 μ L were gently vortexed before plating onto LB

supplemented with rifampicin (10 μ g.mL⁻¹). For each assay, a number of cultures (eight for the R^{168} assay, four for the ΔS and ΔL assays, and three for all other assays) were not plated on LB medium supplemented with rifampicin, but were serially diluted and plated on LB agar to determine the final number of cells (N_t). Fluctuation assays performed on the same day were considered to have the same distribution of final cell numbers. All other cultures were plated on LB agar with rifampicin to determine the number of Rif-resistant colony-forming units (CFUs). All plates were incubated at 37°C and scored for CFUs after 24 h of growth.

The maximum likelihood estimator of the number of mutations per assay (m) and the confidence interval, were calculated under the Luria-Delbrück model taking into account the variation in the final number of cells [41], using the newton.B0 with default parameters and confint.B0 functions of the R package “Rsalvador” v1.7 [42]. For the calculation of confidence intervals, the initial guess for the parameter m was taken as the m given by the “newton.B0” function. The mutation probability was assumed to be constant over the cell cycle, so that the mutation rate per base per generation is the mutation rate per base per cell division [43]. The final number of cells, N_t , is the result of $N_t - N_0$ cell divisions, i.e. $\sim N_t$ divisions. The rate of Rif^R emergence was therefore calculated as $\mu_{Rif} = m/N_t$.

Mutation accumulation experiments and sequencing

One isolated colony was collected each day (24 h at 37°C), suspended in culture medium + 20% glycerol, and diluted by 2×10^5 , a factor that allows distinguishable colonies, before plating on LB agar (+ 100 μ M IPTG for L^* , C^* and LC^* strains) for the next MA-step. Counting of the colonies rescent on the agar plate provided an estimate of the number of bacteria initially present in the diluted colony and thus the number of generations MA-step. Four parallel MA-lines were propagated per strain (21 consecutive MA-steps for ΔS , ΔL , L^* and C^* , 11 for LC^*).

For sequencing at intermediate and end points of the MA-lines, 5–50% of the picked colony was cultured in LB medium to collect cells. DNA was extracted using the GenElute™ Bacterial Genomic DNA Kit (Sigma-Aldrich) according to the supplied protocol. DNA samples corresponding to an intermediate time-point in the four parallel MA-lines for the same strain were pooled in equimolar proportions. Individual and pooled DNA samples were sequenced (150 bp paired-end reads) on an Illumina platform (NovaSeq 6000) to an average depth of ~ 300 . Reads are available in NCBI SRA.

Detection of mutations

The reads were aligned to the reference sequence of the *B. subtilis* 168 genome (GenBank: AL009126.3) using BWA-MEM v0.7.17 [44], after quality control and trimming using Sickle v1.33 (command “sickle pe” with options “-t sanger -x -q 20 -l 20,” <https://github.com/najoshi/sickle>). Properly paired reads, selected using “samtools view -f 3” (samtools v1.14, [45]), were locally realigned around indels using ABRA2 v2.24 [46]. The number of occurrences of each nucleotide (base read quality ≥ 35) and indel at each position of the reference in confidently mapped reads (alignment quality ≥ 50) was counted using “samtools mpileup” with the options “-aa -d 5000 -q 50 -Q 35 -x -B.” These counts were analyzed in R. The reads were

also aligned to the reference sequences of the inserted regions represented in [Supplementary Fig. S2](#) to detect mutations in inducible synthetic circuits.

For each position, the effective sequencing depth (DPeff) was calculated as the total number of informative reads. For the computation of the mutation rates, a reference subset of positions common to all samples was determined. This reference consisted of positions that were well covered on both strands (DPeff ≥ 100 and $\geq 10\%$ of the reads on the less represented strand) in all samples. Most of the regions with low coverage corresponded to the regions deleted in the construction of the MS / R¹⁶⁸ strain [4], which lacks 233.4 kb of chromosome relative to AL009126.3, and to the multicopy structural RNAs. Over-covered regions were also eliminated from this reference subset of positions and consisted of: the region of gene *upp* and downstream (positions 3,788,426 to 3,789,124), which was repeated due to pop-ins and pop-outs at this locus during the construction of the R¹⁶⁸ strain; the region from position 2,432,478 to 2,433,315, over-covered in *polC** samples; the regions of the genes *polC* (1,727,133 to 1,731,446) and *mutL* (1,778,337 to 1,780,539) duplicated by insertion of the mutant alleles *polC** and *mutL**. This resulted in a reference subset of 3,794,734 positions (out of a total of 4,215,606 bp in AL009126.3), which served as our reference chromosome for the mutation rate calculations.

The distribution of the proportion of non-reference reads in the different samples was examined graphically to establish relevant cut-offs for the identification of mutations. A mutation was identified at the endpoint of an MA-line if a variant accounted for $\geq 75\%$ of the DPeff at a position, with $\geq 10\%$ of the non-reference reads on the less represented strand. If intermediate time-points were available for this MA-line, the mutation was traced back to the first time-point where it occurred at frequency $\geq 5\%$ in the corresponding pooled sequence sample. Due to the detection, during graphical examination, of contamination from other samples, we lowered the cut-off from 75% to 60% for the identification of mutations in the third MA-line of ΔS and from 5% to 2% for the analysis of the pool corresponding to the intermediate time-point for *L** strain MA-lines. Mutations found in all samples or in the four MA-lines of a same strain were interpreted as fixed before the MA experiment and were discarded for the calculation of mutation rates.

Chromosome partitioning to assess the effect of transcription and replication on mutation rate

To assess the effect of transcription, the “gene” features of the GenBank annotation served to define the dichotomy between “template” and “nontemplate” strand as well as between “coding” and “noncoding.” Since the “noncoding” represents only $\sim 10\%$ of the genome and includes transcribed untranslated regions, we also sought to assess the impact of transcription with more statistical power and precision than allowed by the GenBank annotation. For this purpose, two categories of regions of approximately equal size were defined based on the transcribed regions identified in 269 samples of a wild-type strain representative of a wide variety of growth conditions [47]. These two categories reflected the amount of transcripts in LB as measured in nine samples corresponding to growth in liquid LB (triplicate samples for exponential, transition and stationary phases) and two samples corresponding to 16 h of growth on LB agar (non-confluent

colonies). The “high” transcription level regions were those that were in the top 30% in at least one of these 11 samples while the “low” transcription level regions were those that were never in the top 30%. All overlapping regions (i.e. both strands were transcribed) were eliminated, as well as all regions shorter than 100 bp. This resulted in a set of 3622 non-overlapping transcription-oriented regions covering 84.9% of the reference genome (43.4% for “high,” 41.4% for “low”).

To assess the effect of DNA replication asymmetry, the leading and lagging strands were defined based on the origin of replication (position 1) and its most central terminus (position 2,018,289) [48]. To assess the effect of DNA replication timing, the genome was divided into a “first half” corresponding to the 2 Mbp of the chromosome centered on the origin of replication (position 1) and a “second half.”

Mutation rate estimations and comparisons

To include the list of mutations in R³⁶¹⁰ and ΔS^{3610} from [33] and [49] in our analysis, the positions on the *B. subtilis* NCIB 3610 genome (GenBank: CM000488.1) were transferred to the *B. subtilis* 168 genome by mapping the 41-bp long sequence centered on each mutation site. Retaining only exact and unique matches, more than 99% of these mutations were transferred ([Supplementary Table S3](#)), with perfect collinearity between the positions of the mutations on both reference genomes.

Maximum-likelihood estimates of mutation rates were obtained as $\mu = m / (T \times G)$, where m is the total number of mutations of a given type in a given genotype and genomic context (which can be defined by the nucleotide at the focal position and its adjacent nucleotides, as well as the chromosome partitioning with respect to replication and transcription as defined above), T is the total number of occurrences of the genomic context in the reference sequence, and G is the number of generations considered in MA-lines. Confidence intervals for mutation rates and proportions were calculated using the exact methods implemented in R package “epitools” v0.5-10.1 for Poisson and binomial counts, respectively.

To assess whether a factor affects the substitution rates, we used Generalised Linear Models (GLMs) for Poisson distributed count data with log-link, combined with an Analysis of Variance (ANOVA) (R package “stats” v3.6.3) to compare the fit of a GLM including and a GLM excluding the factor of interest. This statistical comparison was performed separately for each genotype.

Markov chain Monte Carlo methods implemented in JAGS [50], accessed via the R package “rjags,” were used for Bayesian estimation via posterior sampling, in particular for the estimation of mutation rates of replication-stranded triplets and the MMR saturation parameter θ . Models and algorithm settings are described in detail in [Supplementary Methods and Results 1.1](#). The code is made available at <https://doi.org/10.5281/zenodo.14850407>.

Mathematical modeling of mutation rates

Assumptions and Bayesian estimation procedure for the model with saturation of the MMR are presented in [Supplementary Methods and Results 1.2](#). Algebraic analysis of the general model with two subclasses of errors and two repair pathways is presented in [Supplementary Methods and Results 1.3](#).

Results and Discussion

The mutation rate of *B. subtilis* can be increased up to 6,000 times

We have constructed five mutant strains with expected hypermutator phenotypes from a strain derived from *B. subtilis* 168. The first two strains are constitutively MMR deficient as a result of single deletions of *mutS* and *mutL*. The other three were engineered for conditional inactivation of either or both of MMR or proofreading and were therefore expected to be inducible hypermutators (Supplementary Figs. S1 and S2). In these three strains, the IPTG-inducible promoter P_{hs} controls the expression of mutant alleles selected for their ability to competitively displace their functional counterparts. The first allele, designated *mutL*^{*}, has a mutation in the region encoding the ATP hydrolysis active site of MutL which has been reported to have a dominant negative effect [26]. The second allele, designated *polC*^{*}, encodes a proofreading deficient variant of PolC with a mutation in its exonuclease domain [19]. The final strain, which was expected to have the highest mutation rate under full induction, expresses these two deficient alleles in a synthetic operon (*mutL*^{*} *polC*^{*}). The following notations are used for the reference parental strain derived from 168 and the five mutant strains: R¹⁶⁸, ΔL , ΔS , *L*^{*}, *C*^{*}, *LC*^{*}.

Fluctuation assays were performed to compare the rate of mutation to rifampicin resistance of these strains. Point estimates and confidence intervals are shown in Fig. 1, and results are detailed in Supplementary Table S2. In the absence of IPTG, the estimated mutation rate of R¹⁶⁸ was 9.74×10^{-10} per generation. Constitutive inactivation of the MMR in ΔL and ΔS increased the mutation rate by a factor of approximately 85, with no statistically significant difference between the two strains. This is close to the factor of about 60 previously obtained for a double deletion of *mutL* and *mutS* in the *B. subtilis* PY79 genetic background with a fluctuation assay also based on rifampicin [32]. Mutation rates in the absence of IPTG were slightly higher for the inducible strains (*L*^{*}, *C*^{*}, *LC*^{*}) than for R¹⁶⁸ (up to 1.09×10^{-8} for *LC*^{*}), likely reflecting the low basal activity already described for P_{hs} [51]. From there, the mutation rate of the three inducible strains increased with IPTG concentration, reaching a plateau between 50 and 100 μ M. At full induction (100 μ M IPTG), the mutation rates ranging from 3.66×10^{-7} for *L*^{*} to 5.78×10^{-6} for *LC*^{*}. The mutation rate in *L*^{*} is comparable to or slightly higher than in ΔL and ΔS , while the mutation rate in *LC*^{*} represents an approximately 6000-fold increase over R¹⁶⁸.

Much higher mutation rates than in the reference strain may induce stress responses, potentially altering the physiology of each mutant differently. However, we did not detect any substantial impact on growth in 96-well microtiter plates (Supplementary Fig. S3). We also performed transcriptomics experiments on the R¹⁶⁸, *L*^{*}, *C*^{*} and *LC*^{*} strains in the presence of 100 μ M IPTG. The analyses did not reveal any significantly differentially expressed genes between the strains. However, they did allow us to quantify the expression of the mutant alleles relative to wild-type alleles for *mutL* and *polC*. Upon induction, the mutant allele accounted for 96-98% of the total mRNA pool for the gene in question (Supplementary Table S3), which is consistent with the results of fluctuation assays that give a mutation rate in *L*^{*} as high as in ΔL and ΔS (i.e. complete inactivation of the MMR). We therefore concluded that the mutation profiles of these strains

can be attributed to the sole inactivation of the two targeted DNA repair pathways.

Highest mutation rates are counter-selected in mutation accumulation experiments

Mutation accumulation (MA) experiments give access to the molecular nature of mutations [2], in contrast to fluctuation assays, which only provide a mutation rate aggregated across a range of mutations that confer a screenable phenotype. For each of the six strains (R¹⁶⁸, ΔL , ΔS , *L*^{*}, *C*^{*}, *LC*^{*}), four independent MA-lines were propagated by repeated cycles of colony picking, dilution, and plating on LB (Supplementary Fig. S4), in the presence of 100 μ M IPTG when relevant. By randomly selecting a single colony, each MA-step creates a bottleneck in the propagated population, the purpose of which is to limit genetic diversity, thereby maximizing random genetic drift and minimizing natural selection. The interval between two bottlenecks, one MA-step, was estimated to be 25.6 generations on average. We performed whole-genome sequencing at the endpoint of each line; we also sequenced intermediate time-points (up to 4 for *LC*^{*}) to detect changes in mutation rates (Supplementary Fig. S5). A total of 24 clones and 32 pooled samples isolated after 1 to 37 MA-steps were sequenced. Point mutations - substitutions, insertions (ins), and deletions (del) - were identified. Mutation rates per base pair (bp) and generation (abbreviated $\text{bp}^{-1}.\text{gen}^{-1}$) were estimated for each time interval of each MA-line. Previously collected data from MA experiments on wild-type and MMR-deficient strains were also incorporated to increase statistical power: *B. subtilis* 3610 (R³⁶¹⁰) and its ΔmutS mutant (MMR⁻³⁶¹⁰), *B. subtilis* PY79 (R^{PY79}) and several MMR-mutants (aggregated as MMR^{-PY79}) [32, 33, 49, 52, 53]. As laboratory descendants of the original Marburg strain, *B. subtilis* 168, 3610 and PY79 are all closely related, but the comparison between our data and 3610 is in principle the most appropriate because the PY79 data were obtained at a slightly lower temperature (30°C). The detailed list of all mutations found is provided in Supplementary Table S4.

In the four independent lines of R¹⁶⁸, only one nucleotide substitution was identified after 37 MA-steps, giving an estimated substitution rate of $7 \times 10^{-11} \text{ bp}^{-1}.\text{gen}^{-1}$ (Table 1). This rate was not statistically significantly different from a previous report on R³⁶¹⁰ and R^{PY79} [32, 49], recalculated in Table 1.

Between 113 and 157 nucleotide substitutions were identified after 21 MA-steps in each of the ΔL , ΔS and *L*^{*} strains, resulting in point estimates of the substitution rates between 1.4×10^{-8} and $1.9 \times 10^{-8} \text{ bp}^{-1}.\text{gen}^{-1}$ (Table 1). These rates were not statistically significant from each other which led us to aggregate the data collected for these three strains under the label MMR⁻¹⁶⁸ (Table 1). Sequencing an intermediate time-point located at the end of MA-step 11 for the *L*^{*} strain did not reveal any difference in the rates of accumulation between the first and second parts of the evolution (Fig. 2 and Supplementary Fig. S6). This suggests that in MMR-deficient strains, substitutions identified at endpoints of the MA-lines result from accumulation at a constant rate.

LC^{*} MA-lines showed a tendency toward decreasing substitution rates during their evolution, with differences in terms of frequency, temporality, and magnitude of decrease (Fig. 2 and Supplementary Fig. S6). A statistically significant decrease

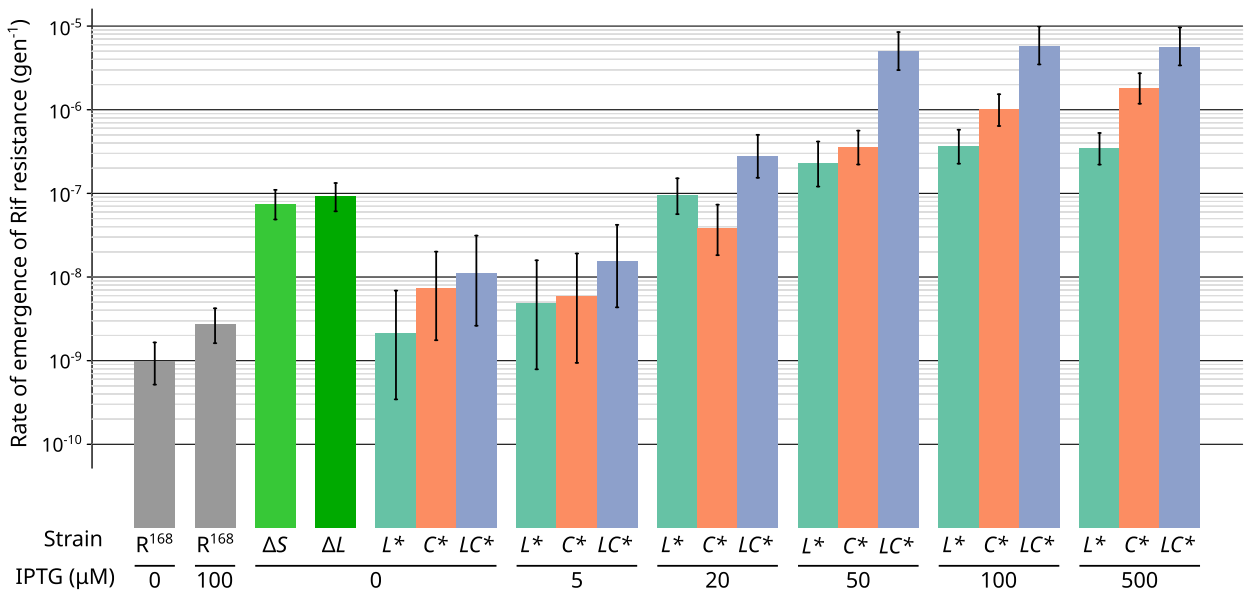


Figure 1. Mutation rate to rifampicin resistance measured by fluctuation assays for increasing IPTG concentration. Each color corresponds to a strain; vertical bars represent 95% confidence intervals.

Table 1. Aggregated number of substitutions, substitution rate, and proportion of transversions for each strain

Strain ^a	Lines	Generations ^{b,c}	Substitutions ^c			Substitution rate (bp ⁻¹ .gen ⁻¹) [95% CI]	Proportion of transversions [95% CI]
			Total	ts ^d	tv ^e		
R ¹⁶⁸	4	3790	1	1	0	7.0e-11 [0.18–39e-11]	0 [0–0.98]
R ³⁶¹⁰	50	251,000	319	238	81	3.3e-10 [3.0–3.7e-10]	0.25 [0.21–0.31]
R ^{PY79}	75	272,700	222	178	44	2.1e-10 [1.9–2.4e-10]	0.20 [0.15–0.26]
ΔS	4	2,151	157	155	2	1.9e-08 [1.6–2.2e-08]	0.01 [0.00–0.05]
ΔL	4	2,151	149	147	2	1.8e-08 [1.5–2.1e-08]	0.01 [0.00–0.05]
L*	4	2,151	113	111	2	1.4e-08 [1.1–1.7e-08]	0.02 [0.00–0.06]
MMR- ¹⁶⁸	12	6,454	419	413	6	1.7e-08 [1.6–1.9e-08]	0.01 [0.01–0.03]
MMR- ³⁶¹⁰	19	38,000	4,844	4,711	133	3.4e-08 [3.3–3.5e-08]	0.03 [0.02–0.03]
MMR- ^{PY79}	118	105,020	5,538	5,432	106	1.4e-08 [1.4–1.4e-08]	0.02 [0.02–0.02]
C*	4	1,895 (256)	395 (19)	348 (18)	47 (1)	5.5e-08 [5.0–6.1e-08]	0.12 [0.09–0.16]
LC*	4	230 (896)	502 (627)	484 (599)	18 (28)	5.7e-07 [5.2–6.3e-07]	0.04 [0.02–0.06]

^aThe label MMR-¹⁶⁸ corresponds to the aggregation of data from the three MMR-deficient strains constructed from R¹⁶⁸ (ΔL, ΔS, L*). Data for strains R³⁶¹⁰ and MMR-³⁶¹⁰ were obtained from [33] and mapped to the R¹⁶⁸ genome sequence. Data for strains R^{PY79} and MMR-^{PY79} were obtained from [32] and mapped to the R¹⁶⁸ genome sequence.

^bConversion between MA-steps and generation based on an estimated average number of generations per step: 25.61 here, 27.53 in [33].

^cNumber of substitutions in the reference subset of positions well covered by the sequencing data (3795 kbp). Between parentheses: number of generations or substitutions in time intervals after a detected decrease in mutation rate (discarded from analysis).

^dNumber of transitions.

^eNumber of transversions.

was also detected for a single C* MA-line (C*₃) and this decrease was only detected during the second half of the 21 MA-steps evolution. In contrast to L* and C*, a decrease was detected in all LC* MA-lines, as early as during the second MA-step for LC*₂ ($P = 8.4 \times 10^{-3}$), and during MA-steps 3-6 for the three other LC* MA-lines. The magnitude of the decrease was a factor approximately 2.5x for C*₃ but reached approximately 50x for LC*₄. Therefore, despite heterogeneity between MA-lines, decreases were globally more frequent, quicker and of larger magnitude for LC* than for C*. Importantly, in LC* MA-lines, 56% of the 1129 identified substitutions occurred in intervals affected by a decrease in the substitution rate (Table 1). To determine the mutation rates of the strains, we only considered data from time intervals be-

fore any detected decrease. This led to point estimates of the substitution rates for the C* and LC* strains of 5.5×10^{-8} and 5.5×10^{-7} bp⁻¹.gen⁻¹, respectively (Table 1).

Decreases in indel rate were also observed and correlated with decreases in substitution rate (Supplementary Fig. S6). This is consistent with both proofreading and MMR contributing to correction of indel errors made by PolC polymerase activity (Supplementary Methods and Results 1.4). The calculated insertions and deletions rates are shown in Supplementary Fig. S7 and given in Supplementary Table S5, the rate of which increases with homopolymer length (Supplementary Fig. S8).

Nonsynonymous mutations were found in the inducible synthetic circuits of 4 out of 5 MA-lines exhibiting a decrease

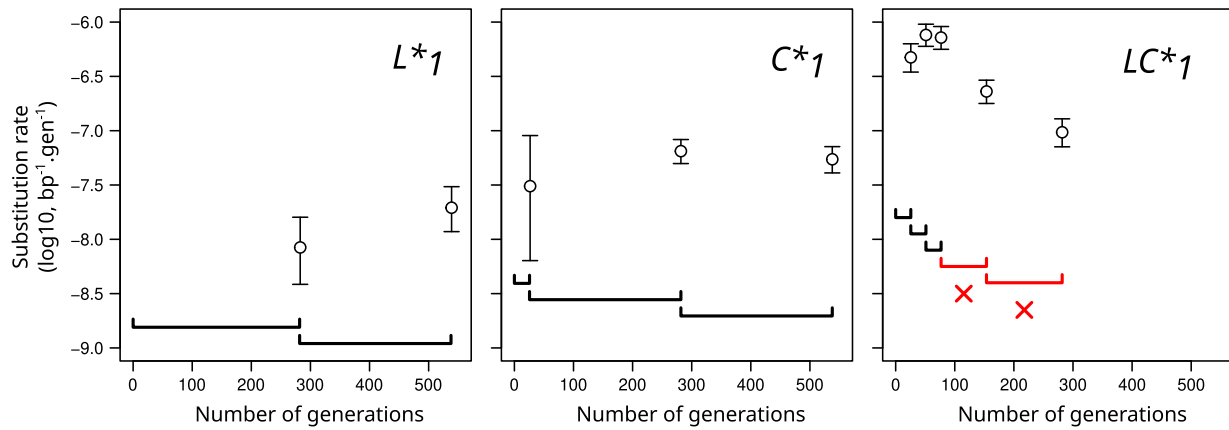


Figure 2. Evolution of substitution rate along MA-lines. Examples are shown for one MA-line from each strain with an inducible mutation rate. The rate is calculated from the number of new substitutions identified within each interval. Sequencing intervals are represented by horizontal brackets and 95% confidence intervals for estimated rates are reported by vertical bars. Red crosses indicate sequencing intervals with significantly decreased mutation rates. These were discarded from downstream analyses.

in mutation rate (Supplementary Table S6, Supplementary Fig. S9, Supplementary Methods and Results 1.5). Given the total number of mutations in the LC^* lines and the size of the *polC* gene, the number of mutations found on the *polC*^{*} allele is four times higher than expected in the absence of selection (chi-squared test with simulated *P*-values, $P = 4.1 \times 10^{-2}$). This finding echoes previous studies reporting changes in mutation rates in MA experiments [54, 55]. Indeed, recent studies have concluded that positive selection is possible in MA experiments despite the extreme bottlenecks imposed on the population [56, 57]. Here, the over-representation of mutations in the genetic elements that confer the strongest hypermutator phenotypes indicates adaptive evolution through positive selection to reduce the mutation rate. Altogether, our results suggest that the ceiling on mutagenesis observed in haploid yeast MMR- and proofreading-deficient mutants has not been reached ([58], Supplementary Methods and Results 1.6). Furthermore, it provides an *a posteriori* experimental justification for the decision to limit MA experiments on proofreading-deficient *E. coli* to 3-6 MA-steps to minimize selection [22].

Proofreading repairs transversion errors at least as well as transition errors

The sequence data do not provide information about the DNA strand on which the error that led to the mutation originally occurred. To record substitutions, we opted for a framework in which the reference base is the pyrimidine (C or T) of the Watson–Crick pair at the genomic position where a mutation is observed. This allows for a prior-free analysis of strand asymmetries in mutation profiles, and follows the convention used in cancer research [59].

In all strains, a slightly higher number of mutations were found on C bases than on T bases (Fig. 3). All strains also showed a predominance of transitions over transversions, but the strength of this bias differed between strains (Fig. 3 and Table 1). The highest proportion of transversions among substitutions, found in R^{3610} (point estimate 0.25), is about 10 times higher than the lowest, found in the MMR- strains (ΔS^{3610} and MMR-¹⁶⁸ strains, point estimates 0.03 and 0.01, respectively). The C^* strain showed an intermediate proportion (point estimate 0.12). The small number of transver-

sions made it difficult to compare those changing the reference pyrimidine (C or T) to an A and to a G. Nevertheless, the C^* and R^{3610} strains may present an excess of C→A over C→G not seen in other strains. The approximately 10-fold increase in the proportion of transitions when comparing the MMR-¹⁶⁸ strains to the R^{3610} strain is consistent with previously published results on ΔS^{3610} [33]. Indeed, there is general tendency across microorganisms for MMR to reduce the transition rate much more than the transversion rate [30, 52].

The substitution rate in C^* divided by the rate in the wild-type (R^{3610} or R^{PY79}) is greater than the rate in LC^* divided by the rate in MMR-¹⁶⁸ (Fig. 3 and Table 1). Thus, the loss of proofreading has less effect in the MMR-¹⁶⁸ strains than in the wild-type. Otherwise, the increase in mutation rate when both repair systems are inactivated compared to the wild-type would be the product of the fold-changes observed when they are inactivated separately (a scenario later referred to as multiplicative fold-changes). This suggests that MMR is already partially defective in the C^* strain. A similar observation was made in *E. coli* upon inactivation of PolIII holoenzyme proofreading [22, 34] and was interpreted as a consequence of MMR saturation due to the high number of errors introduced during DNA replication, a hypothesis further supported by direct assays of MMR activity and restoration by MMR overexpression [34, 60]. To allow comparisons with *E. coli*, we reanalyzed using our analytical pipeline the profiles of mutations obtained for this organism [22, 61] in presence or absence of MMR and proofreading (Supplementary Fig. S10, Supplementary Table S7). Comparing C^* to R^{3610} , inactivation of proofreading significantly reduced the proportion of transversions among substitutions (from 0.25 to 0.12), which may also result from MMR saturation in C^* . Below, we will formally explore the MMR saturation hypothesis using a model-based analysis that integrate information about the chromosomal context of the mutations (adjacent nucleotides and strand).

In the absence of MMR, inactivation of PolC proofreading resulted in a slight increase in the proportion of transversions (from 0.01 in MMR-¹⁶⁸ to 0.04 in LC^*). This difference is not statistically significant (Table 1) but suggests that proofreading corrects errors leading to transversions with at least as much efficiency as those leading to transitions.

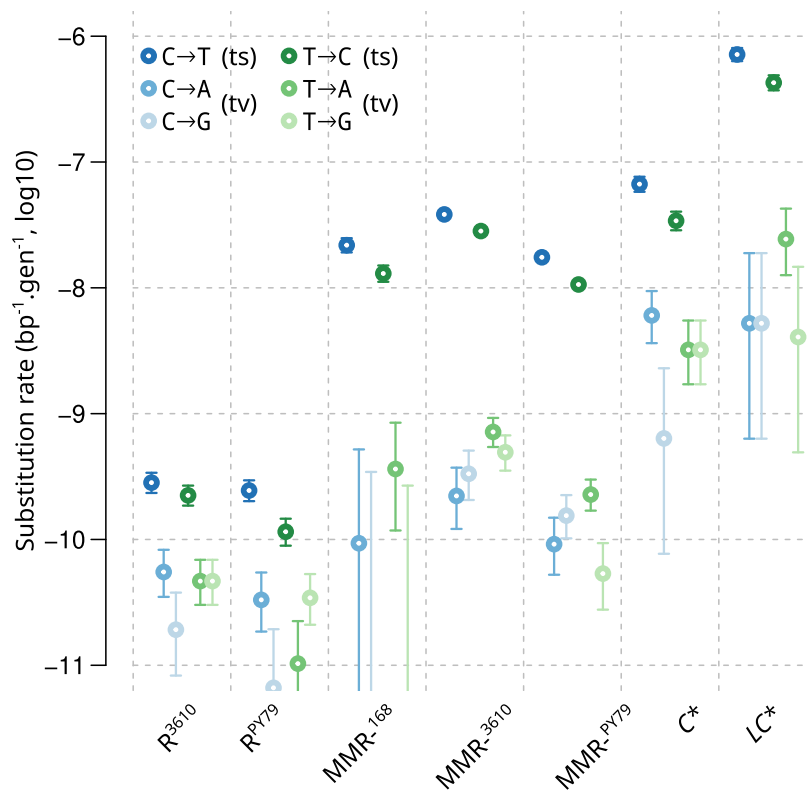


Figure 3. Substitution rates for each type of substitution measured by MA experiments. 95% confidence intervals are reported. No C→G and T→G mutations occurred in MMR-¹⁶⁸, only the upper limit of the 95% CI is shown.

Strand-asymmetry of substitution rate at C:G sites is only visible after proofreading

To further characterize the two DNA repair systems, we examined for each strain how substitution rates varied with distance from the origin of replication, strand orientation relative to replication or transcription, coding or non-coding status, and transcription level. For these analyses, we counted substitutions in the different chromosomal contexts and calculated the corresponding “local” substitution rates.

In MMR-¹⁶⁸ strains, as in MMR-³⁶¹⁰ and R³⁶¹⁰ strains, substitution rates were significantly higher when C is on the leading than on the lagging strand (Fig. 4A). This agrees with the previous analysis of the 3610 dataset [33], where all mutations were recorded as a change on the “strand templating the leading strand” (i.e. the lagging strand), and which showed higher mutation rates for G bases in R³⁶¹⁰ and MMR-³⁶¹⁰ strains. This bias was not present in the C* and LC* strains, which are proofreading deficient. In addition to the strong asymmetry at C:G sites, we detected a weaker but statistically significant replication-oriented asymmetry in substitutions at T:A sites: the substitution rate is higher when T is on the lagging strand, the difference between the two strands being statistically significant for all our hypermutator strains (Fig. 4A). In R³⁶¹⁰, this bias at T:A sites appears to be less pronounced and may even be absent. Consistent with the conclusions of [32], analysis of orientation with respect to transcription, which is most often collinear with replication in *B. subtilis*, did not indicate a contribution of transcription-related processes to these asymmetries between strands in any of the hypermutator strains (Supplementary Methods and

Results 1.7, Supplementary Fig. S11AB). These two replication-oriented biases at G:C and A:T sites appear to be specific to transitions (Fig. 4B).

In wild-type, the replication-oriented asymmetry of substitution rates at G:C sites is a prominent feature of the mutation profile. Consistent with previous analyses of MMR-deficient strains, including MMR-³⁶¹⁰, this bias was also detected in MMR-¹⁶⁸ strains. Our data revealed its absence in C* and LC* strains. A similar observation was made in *E. coli*, leading to the interpretation that proofreading is strand-biased and produces this bias [22], but a concurrent explanation involving strand-biased mutations caused by deamination will be discussed below. In contrast, the wild-type had little or no asymmetry at A:T sites, whereas the hypermutator strains have such asymmetry; error correction systems, and in particular the MMR, tend to eliminate this asymmetry.

In parallel, it is intriguing to observe the presence of two trends detected only in wild-type: a higher substitution rate in non-coding regions (Supplementary Methods and Results 1.7, Supplementary Fig. S11CD) and, to a lesser extent, in the half chromosome far from the replication origin in R³⁶¹⁰ (Supplementary Fig. S11E, ANOVA p-value 0.002). A higher substitution rate in non-coding than coding regions, specific to the wild-type, has also been reported in *E. coli* and was interpreted as an indication that “MMR preferentially repairs coding sequences” [61, 62]. Alternatively, trends observed exclusively in the wild-type may correspond to substitutions originating from processes not subject to correction by proofreading and MMR, which could be masked by increased substitution rates in hypermutator strains.

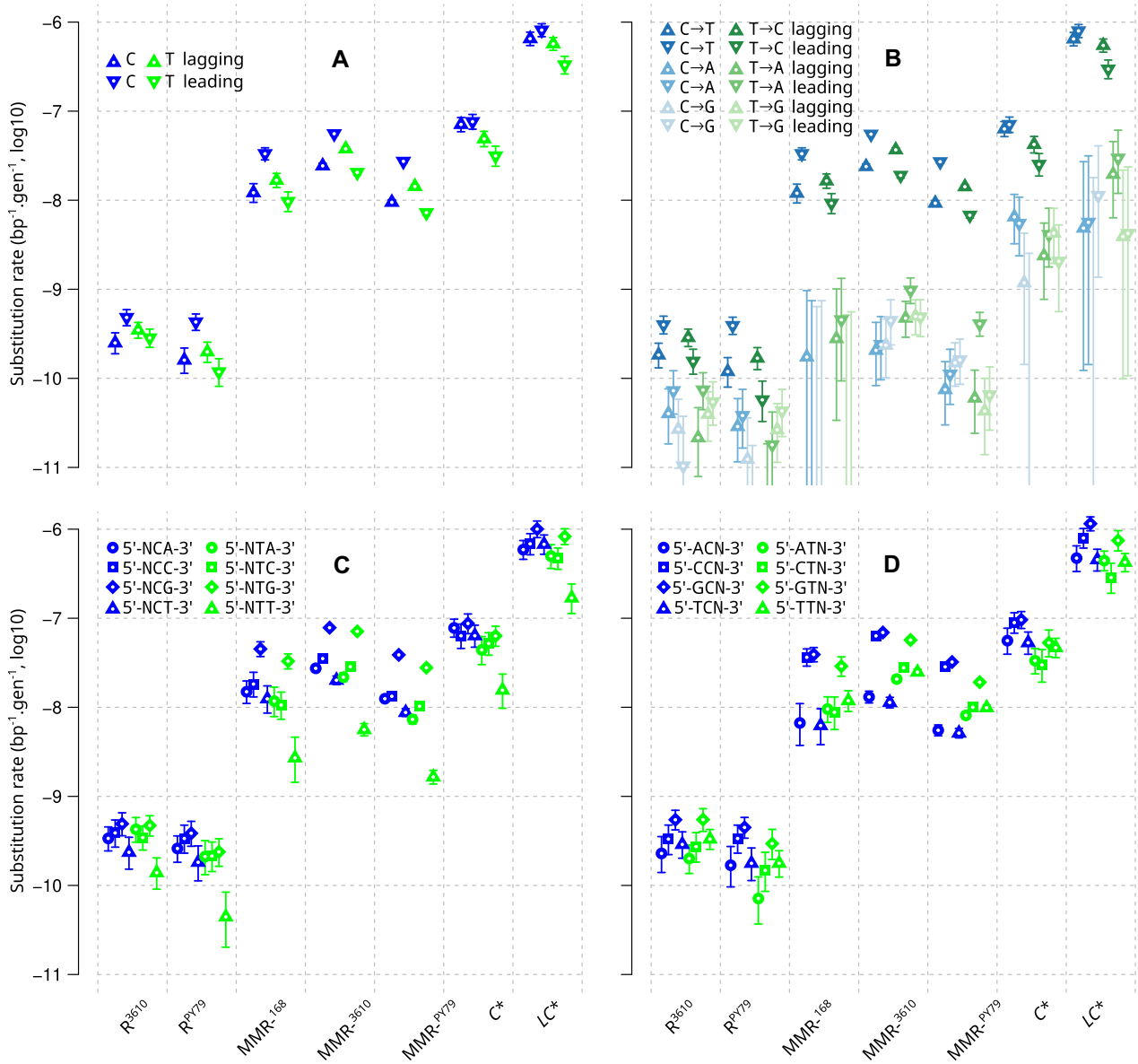


Figure 4. Substitution rate as a function of replication strand and neighboring nucleotides. The pyrimidine of the pair determines the strand of a mutation site. Rates are plotted for each pyrimidine and each genotype. Error bars represent the 95% confidence intervals. **A** and **B**. Effect of orientation with respect to the replication strands (by focal pyrimidine in **A**, by type of substitutions in **B**). **C**. Effect of the 3'-adjacent nucleotide. **D**. Effect of the 5'-adjacent nucleotide. As all the rates reported in this study, those shown in **C** and **D** take into account the number of possible sites in the sequence (see "Materials and methods").

Polymerase initial nucleotide selectivity shapes the distribution of mutations even after proofreading

In all strains, the substitution rate was strongly influenced by the nucleotide adjacent to the focal pyrimidine (Fig. 4C and D). This observation, previously made in wild-type and MMR-deficient strains [33], extends to proofreading-deficient strains. Simultaneously considering the adjacent nucleotides on both sides and the replication strand requires binning the counts into 64 replication-stranded triplets. To mitigate the dimensionality problem, exemplified by the absence of observed substitutions for some bins, we adopted a Bayesian estimation framework that incorporates the mean and standard deviation of log-transformed rates as strain-specific hyperparameters (Supplementary Methods and Results 1.1). Information from the entire distribution is thereby used to establish

point estimates and credibility intervals of the transition rate for a given triplet (Fig. 5A and Supplementary Table S7). The same methodology was applied to the transversions, which unveils distinct profiles that remain difficult to analyze due to the small underlying counts (Supplementary Fig. S12 and Supplementary Table S7).

Based on these estimates, we measured a very strong correlation between the stranded-triplet transition profiles of MMR-168 and the two other MMR- backgrounds (Supplementary Fig. S13, Pearson correlation $r = 0.92$ between log-transformed rates), confirming that our MMR-168 background is very close to the previously studied mutants. There is also a strong correlation between the stranded-triplet transition profiles of LC* and proofreading-proficient MMR-deficient strains (Supplementary Fig. S13, $r = 0.80$ between

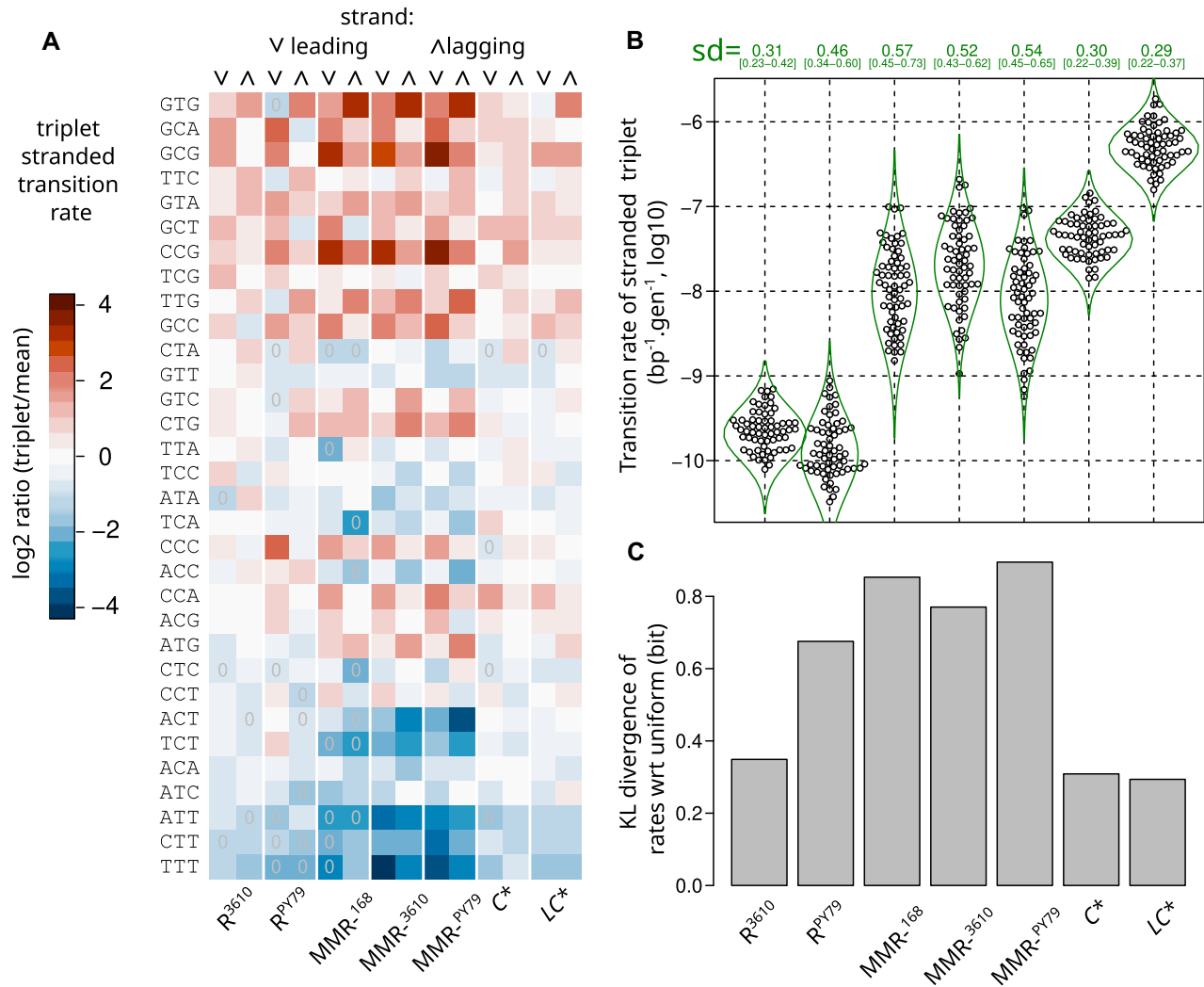


Figure 5. Comparison of stranded-triplet transition profiles between genotypes. Each stranded-triplet corresponds to the pyrimidine of the pair and its 5' and 3' nucleotides, distinguishing pyrimidines on the leading and lagging strands of replication. **A.** Heatmap representation of the stranded-triplet transition profiles (\log_2 ratio of the estimated rate with respect to the mean for the genotype) in the different strains. Rates were estimated using a Bayesian method involving a log-normal prior and hyperparameters. Estimates based on the absence of substitutions in this stranded-triplet context are indicated by "0" in the cells of the heatmaps. Triplets are ordered in decreasing order of transition rate in R^{3610} , averaging the Bayesian estimates of both strands. **B.** Swampplots of the Bayesian estimates of stranded-triplet substitution rates, intended to show individual data points without overlap. The Bayesian estimate of their standard deviation in log10 scale reported above, together with the corresponding probability density function plotted as an envelope. **C.** KL divergences from a uniform distribution (same rate for each triplet), derived from robust entropy estimates. The order of the strains is the same in **B** and **C**.

LC^* and MMR^{-3610}). This is consistent with the idea that, after proofreading and before correction by MMR, most of the errors leading to substitutions are due to misincorporation by the PolC polymerase that escaped proofreading. Since MMR removes ~ 98 – 99% of these errors (Table 1), they also represent the majority of errors corrected by the MMR.

We also found correlations between the wild-type stranded-triplet transition profiles (R^{3610} and R^{PY79}) and those of all hypermutator strains (Supplementary Fig. S13). The correlations are the highest with the proofreading-proficient MMR-deficient strains ($r \geq 0.70$) and remain highly statistically significant with the proofreading-deficient strains ($r \geq 0.67$ with LC^*). In particular, Fig. 4C shows a consistent mutational pattern across all strains, which suggests it likely originates from the initial nucleotide selectivity of the polymerase: NTT substitution rates are significantly lower than those of any

other substitutions, regardless of the genetic background (168, PY79, and 3610) or deficiencies in repair mechanisms (MMR and proofreading). To our knowledge, this specific trend for NTT has not been mentioned in any previous studies and it is not clearly observed in *E. coli* (Supplementary Fig. S10).

The correlation between the stranded-triplet substitution rates profiles in LC^* and the wild-type, whose global substitution rate is 1600 times lower (Table 1), fits remarkably with the working hypothesis that PolC misincorporations that escaped proofreading and MMR substantially shape the substitution profile of the wild-type. Not only is this consistent with the conclusion of the *E. coli* study [22], but there are also similarities between the transition rates measured in the two organisms, the most significant being correlations between the transition profiles of hypermutators (Supplementary Fig. S13). In particular, there is a clear similarity ($r = 0.66$) between

B. subtilis LC^* and the corresponding *E. coli* constitutive proofreading and MMR-deficient mutant (Supplementary Fig. S14). Transitions are thought to result from base pairing errors and base tautomerization [53]. The local sequence context may influence the probability of template tautomerization as well as the probability of incorporating a tautomeric base or a non-cognate base, in conjunction with the structure of the polymerase and the dynamics of the elongating complex.

Proofreading squares the biases of the initial polymerase selectivity

As measured by the estimated span of their distributions, the stranded-triplet transition rates in proofreading-deficient strains (LC^* and C^*) are almost 10-fold more dispersed than in proofreading-proficient MMR- strains, indicating that proofreading approximately doubles the dispersion on a logarithmic scale (Fig. 5B). Because point estimates deduced from small counts are not precise, and to further corroborate the results of the Bayesian standard deviation estimate, we sought a second approach that could directly quantify dispersion. Using an entropy estimation method developed for small counts [63]. We calculated the Kullback-Leibler (KL) divergence between the unknown underlying distribution of observed counts and a uniform substitution rate, where the expected count is proportional to the number of occurrences of the triplet in the chromosome (Supplementary Methods and Results 1.8). Estimates of KL divergences (Fig. 5C) confirmed both the similar level of dispersion of substitution rates in the wild-type R^{3610} and proofreading-deficient strains, and the comparatively much higher dispersion in the MMR-deficient strains (approximately 2.5-fold higher divergence from uniform). In other words, proofreading disperses substitution rates, while MMR activity can counteract this dispersion to a variable extent (much more in R^{3610} than in R^{PY79}). Similar trends are seen in *E. coli* (Supplementary Fig. S15). These observations align with those on *E. coli*, which have been interpreted as a compensation of the proofreading biases of the *E. coli* PolIII holoenzyme by opposite biases of MMR correction [22]; a concurrent hypothesis to explain the apparent compensation of the biases is presented in the discussion.

Interestingly, after proofreading (but before MMR correction), the stranded-triplet substitution rates become more dispersed, but remains very similar in terms of the direction of the biases compared to before proofreading (Supplementary Fig. S13, $r = 0.80$ between LC^* and MMR^{-3610}). This observation is puzzling because it implies that proofreading, while considerably decreasing the mutation rate, increases the dispersion of biases that already arise from the sole polymerase activity, rather than simply masking the initial biases with its own biases of greater amplitude. To better understand the implications of this observation, we can formulate a minimal model in which the proofreading activity reduces the initial error probability of the polymerase, denoted $\gamma[i]$, through a two-step process: first, the detection and removal of a misincorporated nucleotide with probability $d[i]$; second, the reincorporation of a nucleotide with the error probability $\gamma[i]$ characteristic of the initial polymerase activity. The error probability after proofreading writes then $e[i] = \gamma[i] (1 - d[i]) / (1 - \gamma[i]d[i])$, where the term $(1 - \gamma[i]d[i])$ accounts for the possibility of cycling the two-step process if a new incorporation error follows the removal. If $d[i]$ is the same for all i , the possibility of cycling, by itself, already amplifies the initial biases. However, this effect

is negligible as long as $\gamma[i]$, which corresponds to the substitution rate observed in the LC^* strain ($<10^{-5}$ in Fig. 5B), remains extremely small compared to 1. Doubling the biases on a logarithmic scale (i.e. $e[i]/e[j] = (\gamma[i]/\gamma[j])^2$), as almost observed in our data, implies a probability of non-detection and removal in the first step of proofreading $(1 - d[i])$ proportional to $\gamma[i]$, the error probability of the single polymerase activity. The most common errors made during initial incorporation are also the least likely to be corrected during proofreading. As a result, proofreading amplifies the biases of initial polymerase selectivity, approximately raising them to the second power.

A theoretical model suggests that MMR can prevent up to 4 mutations per generation before saturation

The apparent efficiencies of MMR and proofreading are strongly influenced by the presence or absence of the other system. As evoked above, proofreading reduced the total substitution rate by a factor 164 in MMR-proficient cells (μ_{C^*}/μ_{R3610} where μ is the mutation rate under consideration and using the ML estimates from Table 1) but only by a factor 34 in MMR-deficient cells ($\mu_{LC^*}/\mu_{MMR-168}$). Similarly, MMR reduced the substitution rate more in the presence than in the absence of proofreading. According to Bayesian estimates of substitution rates, none of the triplets clearly contradicts this observation (with a posterior probability of being above the diagonal greater than 5% for all triplets; Fig. 6). The trend becomes even more pronounced when MMR^{-3610} or MMR^{-PY79} is used instead of MMR^{-168} to represent the MMR-deficient proofreading-proficient background (Supplementary Fig. S16).

To thoroughly investigate the compatibility of the data collected in *B. subtilis* with a saturation mechanism, we formulated a mathematical model in which the mutation rate, $\mu_{C^*}[i]$, in strain C^* for a given triplet or type of mutation i is determined by the equation:

$$\mu_{C^*}[i] = \gamma[i] (\theta + (1 - \theta) \cdot q_{MMR}[i]),$$

where $\gamma[i]$ is the error rate before correction by proofreading or MMR, and θ is a mixture parameter common to all values of i . It corresponds to the proportion of errors made by the proofreading-deficient PolC that occur in a physiological context of MMR saturation (generating mutations distributed as in LC^*). The complementary proportion $(1 - \theta)$ is subject to correction by MMR, which reduces the number of errors by a factor $q_{MMR}[i]$. In this model, whose assumptions are presented along with an associated Bayesian estimation procedure in Supplementary Methods and Results 1.2, the mutation rate in C^* can be expressed as a function of the rates in the three other backgrounds by identifying $q_{MMR}[i]$ with $\mu_R[i]/\mu_{MMR-168}[i]$ and $\gamma[i]$ with $\mu_{LC^*}[i]$.

This parameterization was used to estimate the mixture parameter θ and to check the agreement of the model with the experimental data (Fig. 7). In practice, the posterior distribution of θ was estimated using either the transition rates for the 64 replication-oriented triplets or the rates of the six types of substitutions (2 transitions and 4 transversions); MMR^{-168} , MMR^{-3610} or MMR^{-PY79} representing the MMR-deficient background. The resulting posterior distributions are very similar (Fig. 7A), with the estimated value of θ varying only between 0.071 and 0.082. The observed counts, aggregated by

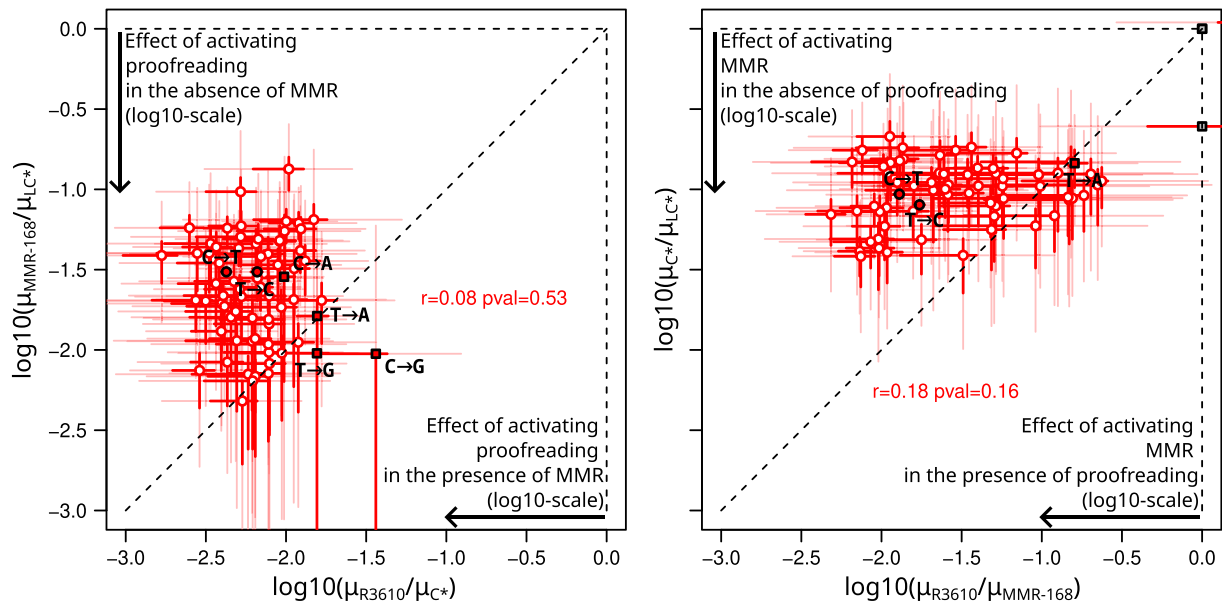


Figure 6. Effects of proofreading and MMR in the presence or absence of the other system. Effects are measured as the log10-ratio of transition rates for stranded triplets (red) and the six types of substitutions (black circles for transitions, black squares for transversions). Left plot: effect of proofreading in the presence or absence of MMR. Right plot: effect of MMR in the presence or absence of proofreading. Around each point, the 50% and 95% marginal credibility intervals on horizontal and vertical axes, computed from the quantiles of the posterior distributions, are represented by segments (bold and dark versus thin and light, respectively). The data used for the reference and MMR- substitution profiles are R³⁶¹⁰ and MMR-¹⁶⁸ (see [Supplementary Fig. S16](#) for similar plots using other reference and MMR- substitution profiles).

substitution type (Fig. 7C) and by triplet ([Supplementary Figs. S17–S19](#)), fall within the prediction intervals of the model, indicating a good fit to the experimental data ([Supplementary Methods and Results 1.1](#)). Notably, although the fraction of errors made by the proofreading-deficient PolC in the context of MMR saturation is small ($\theta < 10\%$), these errors account for the majority ($\geq 75\%$) of the mutations observed in strain C* (Fig. 7A).

In a simple mechanistic model, these values of θ may correspond to the capacity of the MMR to prevent (with a probability of failure $q_{\text{MMR}}[i]$) up to four mutations per generation (Fig. 7B). Alternatively, θ values can be interpreted as reflecting the fraction of errors made by the proofreading-deficient polymerase before it travels a certain distance after a first error ([Supplementary Methods and Results 1.2](#)).

Aggregating mutations in counts: a necessary evil

In principle, our experimental data could also be explained by a model that does not involve MMR saturation. This alternative model recognizes that counting mutations on the genome implies aggregating them by type or context, which are likely to encompass subclasses of mutations arising from errors that are not corrected with the same probability of success. In the extreme scenario where MMR and proofreading correct non-overlapping sets of errors, the increments observed in mutation rates when each system is inactivated would simply add up in cells where both systems are inactivated. This idea has already been mentioned and refuted for mutational signatures in human cancers [64]. In fact, while changes in mutation rates associated with each repair system are not multiplicative, they are clearly more than additive. This raises interest in a more general model that considers aggregated subclasses of mutations in the counts and arbitrary correction rates. We have algebraically explored this model in its simplest form with only

two subclasses ([Supplementary Methods and Results 1.3](#)). Interestingly, when the probability of error correction differs between subclasses for both MMR and proofreading, the aggregation creates apparent epistasis in the sense that the effect of inactivating one system, as measured by the mutation rate fold-change, depends on the presence or absence of the other system.

To explain 4 mutation rates, even the simplest two-subclass scenario has six parameters, consisting of the mutation rate in the absence of any correction and the probability of correction by each system, for both types of error. With more parameters than data points, it can fit almost any data set, explaining additive to super-multiplicative effects. The model is therefore difficult to falsify. We note, however, that the sign of this epistasis depends on whether the systems have similar or opposite specificities: if they tend to repair the same subclasses of errors the apparent epistasis will be positive (super-multiplicative), otherwise the epistasis will be negative (sub-multiplicative, as in our data). Explaining negative epistasis with this model is therefore at odds with the widely accepted idea that MMR correction is mostly coreplicative and targets the same errors as those corrected by proofreading, i.e. Watson–Crick mismatches introduced by DNA polymerase during DNA replication. If subclass aggregation does not contribute significantly to the observed apparent epistasis, which is well explained by MMR saturation, it still complicates the interpretation of the apparent efficiency of repair systems deduced from mutation profiles.

A cautious interpretation of the apparent efficiency of proofreading and MMR

The proofreading-deficient strain C* should not be considered as fully MMR-proficient due to MMR saturation. Therefore, its mutation rate cannot be used to measure the

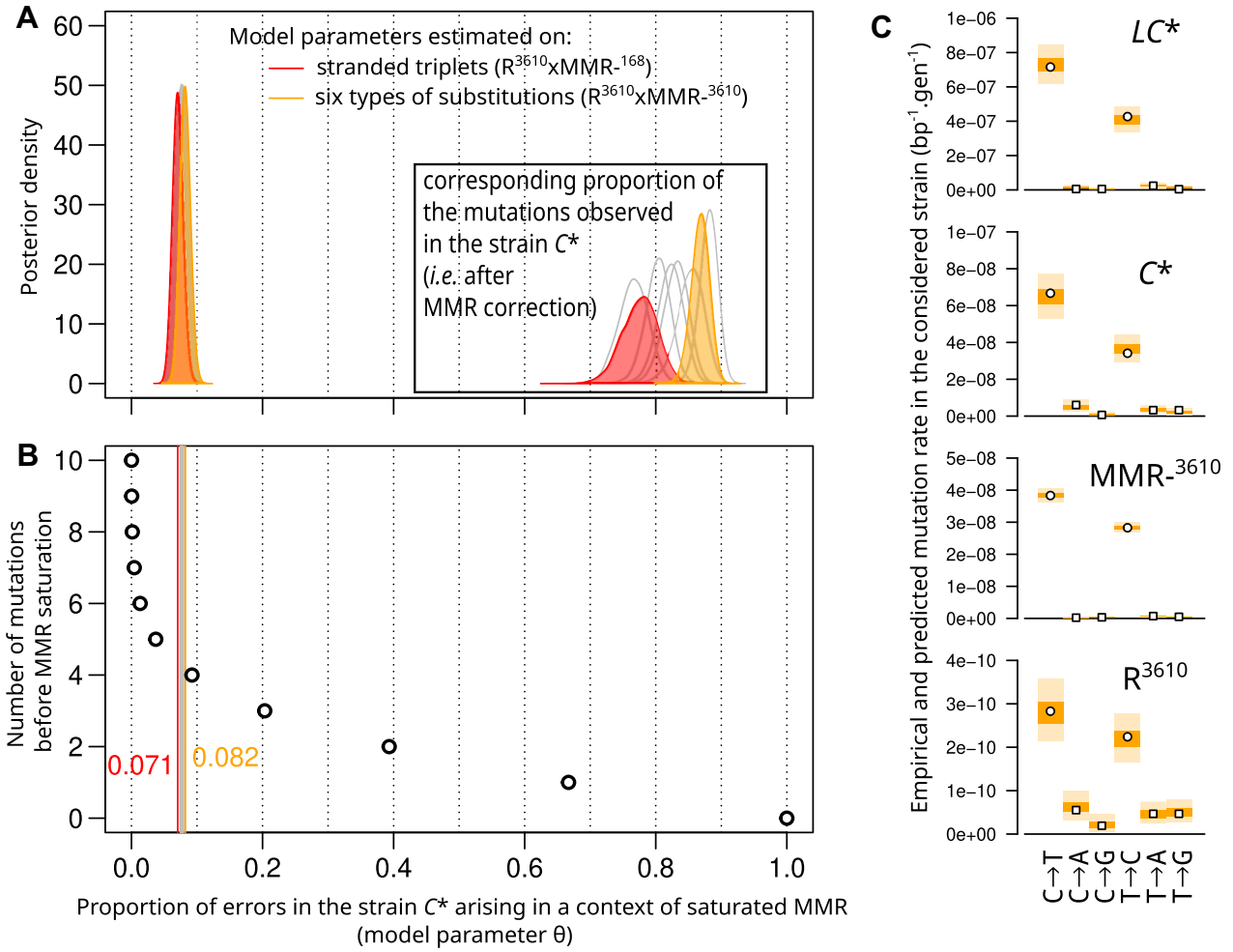


Figure 7. Parameter estimation of the MMR-saturation model and evaluation of the fit to experimental data. **(A)** Posterior distribution of the mixing parameter θ corresponding to the proportion of polymerase errors in strain C^* arising in the context of saturated MMR. The corresponding proportion of mutations observed in C^* (i.e. after MMR correction) is shown in the inset plot. Parameters can be estimated either on data consisting of stranded triplets (transitions) or of the six types of substitutions (two transitions, four transversions) and for four different combinations of reference and MMR profiles ($R^{3610} \times \text{MMR}^{-168}$, $R^{3610} \times \text{MMR}^{-3610}$, $R^{\text{PY79}} \times \text{MMR}^{-168}$, $R^{\text{PY79}} \times \text{MMR}^{-\text{PY79}}$). This gives eight different estimated values of the mixing parameter θ , the posterior distributions associated with the lowest and highest estimates are shown as colored areas; the results corresponding to other combinations of data sets are shown as single lines. **(B)** Relationship between the mixing parameter θ and the number of mutations before MMR saturation in a simplified model of replication: one replication per generation, the first mutations are subjected to MMR correction until MMR saturation. **(C)** Evaluation of the fit of the MMR-saturation model to experimentally measured rates of the six types of mutations. Points represent empirically calculated mutation rates, i.e. the number of observed mutations divided by the number of possible sites in the genome and the number of generations. Colored areas represent the distribution of values for the empirical rates simulated under the posterior distribution of the model parameters (50% of the density in the dark area, 95% including also the light area).

apparent efficiencies of proofreading and MMR under wild-type physiological conditions, i.e. as the numerator in a ratio μ_{C^*}/μ_{LC^*} to estimate the probability that an error escapes proofreading, or as the denominator in a ratio μ_R/μ_{C^*} to estimate the probability that an error escapes MMR correction. Instead, these repair system escape probabilities should be estimated from the ratios $\mu_{\text{MMR-}}/\mu_{LC^*}$ and $\mu_R/\mu_{\text{MMR-}}$ for each replication-oriented triplet and the six substitution types shown in Fig. 8 (where MMR- is either MMR^{-168} or MMR^{-3610} , see [Supplementary Fig. S20](#) for $\text{MMR}^{-\text{PY79}}$). The rates of errors leading to substitutions before and after correction by the combined action of proofreading and MMR (mutation rates in LC^* and R^{3610} , respectively) are also shown.

Proofreading and MMR escape probabilities ($\mu_{\text{MMR-}}/\mu_{LC^*}$ and $\mu_R/\mu_{\text{MMR-}}$) are correlated with triplet substitution rates ([Supplementary Fig. S21](#)). In Fig. 8, the triplets are ordered ac-

cording to wild-type substitution rates (as in Fig. 5A), which highlights general trends: proofreading escape tends to be higher on triplets where the initial polymerase selectivity is the lowest ([Supplementary Fig. S21](#), $r = 0.46$ $P = 0.0001$ when MMR- is MMR^{-3610} or $\text{MMR}^{-\text{PY79}}$); and, in the case of MMR^{-3610} but not $\text{MMR}^{-\text{PY79}}$, MMR escape tends to be lower on triplets with high wild-type mutation rate ([Supplementary Fig. S16](#), $r = -0.39$, $P = 0.002$). Note that these two correlations are of opposite sign to the artifactual correlations that would be generated by noise in the estimates of $\mu_{LC^*}[i]$ and $\mu_R[i]$, which also enter in the ratios $\mu_{\text{MMR-}}[i]/\mu_{LC^*}[i]$ and $\mu_R[i]/\mu_{\text{MMR-}}[i]$.

The positive correlation between proofreading escape and mutation rate in LC^* is consistent with the observation that proofreading increases the bias of polymerase errors. Statistical uncertainty about the probability of proofreading escape

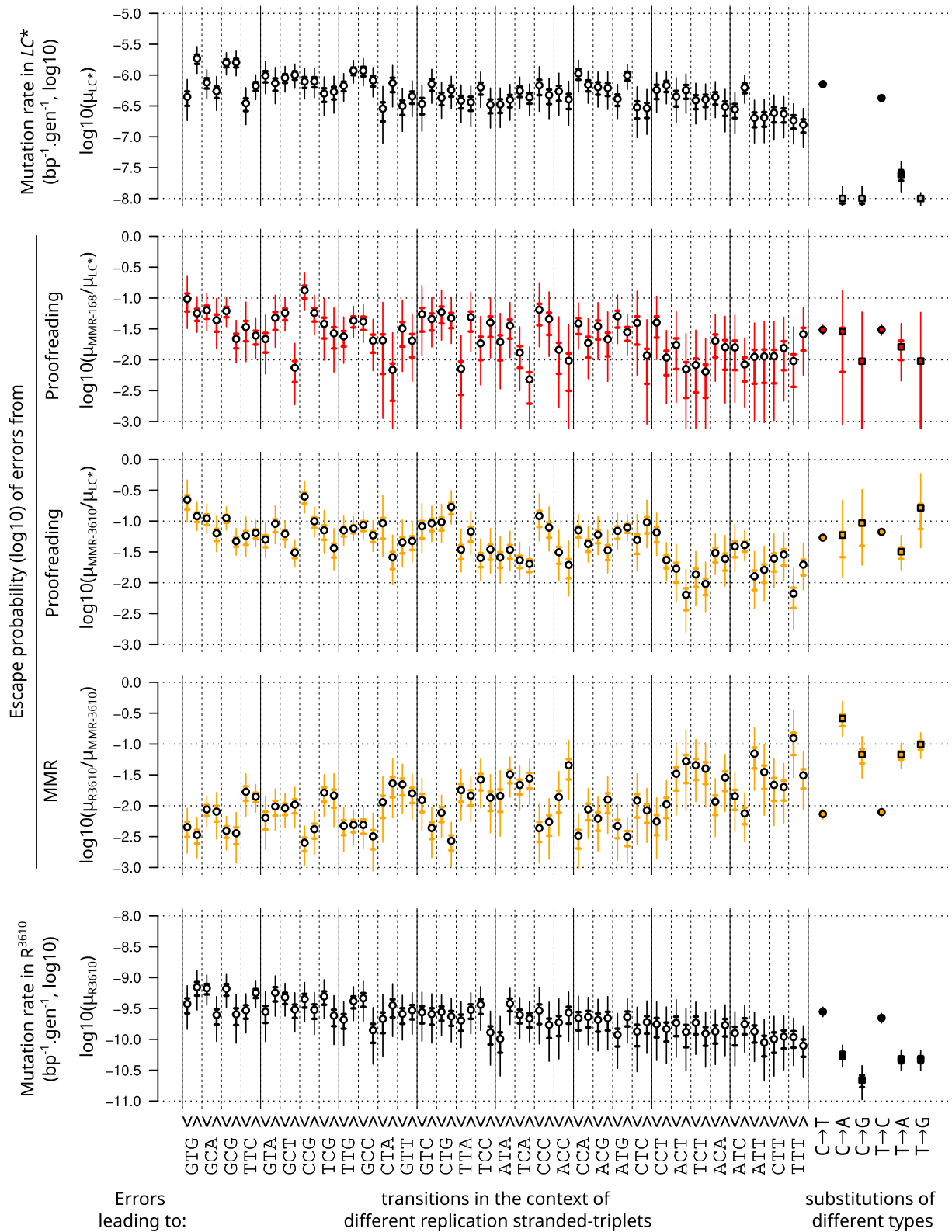


Figure 8. Apparent efficiency of proofreading and MMR across replication-stranded triplets and mutation types. Triplets are ordered in decreasing order of transition rate in R^{3610} and then strand of replication (leading ∇ and lagging \wedge). Estimated mutation rates in the absence or presence of both correction systems are shown in top and bottom plots, respectively. Apparent proofreading escape probabilities are shown either as estimated from MMR^{-168}/LC^* or MMR^{-3610}/LC^* (second and third plots counting from top). Apparent MMR escape probabilities are shown estimated from R^{3610}/MMR^{-3610} (fourth plot). Bold and thin vertical bars represent 50% and 95% credibility intervals, respectively.

makes it difficult to draw conclusions for specific replication-stranded triplets, and few triplets show clear strand asymmetry. It is interesting to note, however, that this positive link is not seen within pairs of triplets that differ only by the strand. In particular, GTG and ATG are more affected by substitutions in LC^* when the focal pyrimidine (T) is on the lagging strand, but no difference in proofreading escape probability between strands is detected. For these triplets, the initial strand asymmetry may originate from incorporation errors made by the DnaE polymerase during synthesis of the lagging strand that are efficiently proofread *in-trans* by PolC. Conversely, errors affecting GCT and GCG appear to be less corrected by proofreading when the C-site is on the leading strand, but the substitution rates in LC^* are similar between strands. In these cases of C-site substitutions on the leading strand, the apparent probability of proofreading escape most likely overestimates the escape of initial polymerase selectivity errors because of the likely contribution of deamination damage to the mutation profile of the proofreading-proficient MMR-deficient strains.

Deamination of cytosine to uracil in the lagging strand template (i.e. the leading strand) due to exposure of single-stranded DNA in the context of the replication fork has long been suspected to be the reason for the near-universal asymmetry between replication strands in prokaryotes [65] and indeed has been shown to be a substantial source of C→T transitions in wild-type [66]. Strikingly, proofreading-proficient MMR-deficient strains exhibit a similarly oriented strong strand asymmetry for substitutions at C sites [33, 62]. It has been hypothesized that cytosine deamination may also cause this asymmetry [61] which is consistent with other studies reporting a role for the *B. subtilis* MMR in counteracting the effects of base deamination [67, 68]. In the absence of proofreading, the high rate of polymerase misincorporation on both strands would mask the deamination-induced asymmetry.

The hypothesis that proofreading efficiency is positively linked to initial polymerase selectivity was proposed above based on the analysis of transition rates by triplet context. Transversions account for only 4% of the substitutions in LC^* . The difference between the apparent proofreading escape probabilities of errors leading to transversions and transitions cannot be precisely estimated, but their relative similarity (Fig. 8) may indicate that the above hypothesis holds only for transitions. Alternatively, the apparent proofreading escape probability may overestimate proofreading escape for initial polymerase selectivity errors that lead to transversions. This would occur if a substantial fraction of the rare transversions in proofreading-proficient MMR-deficient strains did not originate from polymerase errors and thus was simply not subjected to proofreading, as could suggest the differences visible between LC^* and MMR- (MMR-³⁶¹⁰ or MMR-PY79), and also between MMR-³⁶¹⁰ and MMR-PY79, in terms of distribution of transversions across triplets and strand biases (Supplementary Fig. S12).

Regarding MMR escape, which is not the primary focus of this study, its negative correlation with proofreading escape (Supplementary Fig. S22, $r = -0.66$ and $r = -0.57$ when estimated based on 3610 data and PY79 data, respectively) could reflect an evolution of MMR toward correcting the most common DNA replication errors. Indeed, this would echo the proposed role of differential MMR efficiency in balancing DNA

replication fidelity between the two strands in eukaryotes [23, 30, 69].

However, the negative correlation between the apparent probabilities of MMR and proofreading escape could also reflect an overestimation of MMR escape for the sites least affected by DNA replication errors. Indeed, the contribution of different sources of spontaneous mutations in wild-type remains difficult to characterize [53], and thus there is a lack of knowledge about the fraction of the wild-type mutation profile that is not subject to MMR correction. Among the studies addressing this question, it has been shown that inactivation of oxidative damage repair pathways increases the mutation rate in *E. coli* under optimal growth conditions without external stress, suggesting that damage that escapes correction by oxidative damage repair pathways may also contribute to the wild-type mutation profile [16]. Since oxidative damage tends to cause transversions, this hypothesis would be consistent with the higher proportion of transversions in the wild-type than in any of our hypermutator strains. If a substantial fraction of the wild-type (reference) substitution profile μ_R originates from sources not subjected to MMR correction, variations not directly related to MMR efficiency may enter the ratio μ_R/μ_{MMR-} used to estimate the MMR escape probability (Fig. 8), both via its numerator μ_R (e.g. variations in the amount of substitutions originating from these other sources) and via its denominator μ_{MMR-} (e.g. variations in the amount of polymerase errors remaining after proofreading). Indeed, the very strong negative correlation across triplets between μ_R/μ_{MMR-} and μ_{MMR-} (Supplementary Fig. S21, $r = -0.91$ for the 3610 data) suggests that variations in μ_{MMR-} may be largely responsible for variations in the apparent MMR escape probability. Combined with the unsurprising and clear positive correlation between the apparent proofreading escape probability (μ_{MMR-}/μ_{LC^*}) and μ_{MMR-} (Supplementary Fig. S21, $r \geq 0.90$ for MMR-³⁶¹⁰ and MMR-PY79), the substitutions originating from sources not subjected to MMR correction may thus contribute to the observed negative correlation between the apparent probabilities of MMR and proofreading escape.

Conclusion

Saturation of the MMR in the absence of proofreading, greater dispersion of substitution rates in the presence than in the absence of proofreading, and the existence of strand biases that become apparent only in the presence of proofreading appear to be traits shared between *B. subtilis* and *E. coli*. Given the considerable divergence between these two organisms in terms of phylogenetic distance and molecular organization of DNA polymerase and MMR, these traits are likely to be common to many other organisms. The characterization of the overdispersion of mutation rates in MMR-deficient proofreading-proficient strains compared to other strains led us to conclude that proofreading skews DNA polymerase error rates. This could be interpreted as an inherent drawback of the proofreading principle, which relies on the DNA polymerase to detect its own errors. Acting as a judge of its own errors leads to the same biases in initial nucleotide selectivity and error correction efficiency.

This study also examines the consequences of aggregating mutations in counts and, more generally, recognizes the difficulty of interpreting the apparent efficiencies of repair systems.

First, unaware aggregation of subclasses of mutations resulting from different molecular pathways can create almost any pattern of apparent interactions in the effects of disabling different repair systems. Second, interpretation of the effect of disabling one repair system is generally complex because of uncertainty about the contribution of damage or errors subject to correction by that system to the mutations observed in its presence. It is thus difficult to interpret the appearance of strand asymmetry upon proofreading activation, which may be caused by a post-proofreading mutation process, such as deamination, rather than a proofreading bias. Similarly, the “flattening” of mutation rates upon activation of the MMR is difficult to interpret, as it could result from a better efficiency of the MMR in correcting the most common errors, but also from the contribution of multiple sources to the wild-type mutation profiles. Indeed, a substantial contribution of multiple sources seems consistent with the drift-barrier hypothesis [2]: in the wild-type, the rates of mutations resulting from different molecular pathways may have been pushed below the same point where they do not encounter significant counter-selection.

The construction of strains with inducible hypermutator phenotypes addressed the stability issue that affected the reproducibility of some studies on *E. coli* proofreading-deficient strains discussed in [22]. It was also motivated by the interest in tools for synthetic biology applications. Regarding the future use of our inducible systems to accelerate evolution in the laboratory, a proofreading-deficient strain would yield a less skewed mutation spectrum than an MMR-deficient strain. If necessary, extreme mutation rates can be achieved by disabling both repair systems simultaneously, but strong counter-selection against mutational load makes such induction feasible only for a short period of time.

Acknowledgements

We express our gratitude to Guillaume Achaz, Elena Bidnenko, Clément Nizak, and Paulo Tavares for their advice during the course of this work. We also thank the INRAE MIGALE bioinformatics facility (<https://doi.org/10.15454/1.5572390655343293E12>) for providing computational resources and Marina Elez for helpful comments on the manuscript.

Author contributions: I.T., M.J., and P.N. conceived the project and designed the experimental plan. I.T. and E.D. performed the experiments. I.T. processed the raw data with the help of C.G. I.T., and P.N. performed the computational analyses. G.K.K.K. implemented the Bayesian methodology. I.T., M.J., and P.N. interpreted the results and wrote the manuscript with the contribution of all authors.

Supplementary data

Supplementary data is available at NAR online.

Conflict of interest

None declared.

Funding

This work was supported by the French National Research Agency (ANR-18-CE43-0002). The I.T. PhD fellowship was

partially funded by the MathNum division of INRAE. Funding to pay the Open Access publication charges for this article was provided by Institutional and/or grant.

Data availability

The RNA and DNA sequencing data generated in this study have been submitted to the NCBI Gene Expression Omnibus database under accession number GSE239804 and to the NCBI Sequence Read Archive under BioProject accession number PRJNA995423, respectively.

References

1. Eyre-Walker A., Keightley P.D. The distribution of fitness effects of new mutations. *Genome Res* 2007;8:1336–43. <https://doi.org/10.1038/nrg2146>
2. Lynch M., Ackerman M.S., Gout J.F. *et al.* Genetic drift, selection and the evolution of the mutation rate. *Nat Rev Genet* 2016;17:704–14. <https://doi.org/10.1038/nrg.2016.104>
3. Deatherage D.E., Leon D., Rodriguez A.E. *et al.* Directed evolution of *Escherichia coli* with lower-than-natural plasmid mutation rates. *Nucleic Acids Res* 2018;46:9236–50. <https://doi.org/10.1093/nar/gky751>
4. Dervyn E., Planson A.G., Tanaka K. *et al.* Greedy reduction of *Bacillus subtilis* genome yields emergent phenotypes of high resistance to a DNA damaging agent and low evolvability. *Nucleic Acids Res* 2023;51:2974–92. <https://doi.org/10.1093/nar/gkad145>
5. Taddei F., Radman M., Maynard-Smith J. *et al.* Role of mutator alleles in adaptive evolution. *Nature* 1997;387:700–2. <https://doi.org/10.1038/42696>
6. Couce A., Caudwell L.V., Feinauer C. *et al.* Mutator genomes decay, despite sustained fitness gains, in a long-term experiment with bacteria. *Proc Natl Acad Sci USA* 2017;114:E9026–35. <https://doi.org/10.1073/pnas.1705887114>
7. Chiang DulantoPatil A., Beka P.P. *et al.* Hypermutator strains of *Pseudomonas aeruginosa* reveal novel pathways of resistance to combinations of cephalosporin antibiotics and beta-lactamase inhibitors. *PLoS Biol* 2022;20:e3001878. <https://doi.org/10.1371/journal.pbio.3001878>
8. Oliver A., Mena A. Bacterial hypermutation in cystic fibrosis, not only for antibiotic resistance. *Clin Microbiol Infect* 2010;16:798–808. <https://doi.org/10.1111/j.1469-0691.2010.03250.x>
9. Rudenko O., Engelstadter J., Barnes A.C. Evolutionary epidemiology of *Streptococcus iniae*: linking mutation rate dynamics with adaptation to novel immunological landscapes. *Infect Genet Evol* 2020;85:104435. <https://doi.org/10.1016/j.meegid.2020.104435>
10. Sniegowski P.D., Gerrish P.J., Lenski R.E. Evolution of high mutation rates in experimental populations of *E. coli*. *Nature* 1997;387:703–5. <https://doi.org/10.1038/42701>
11. Swings T., Weytjens B., Schalck T. *et al.* Network-based identification of adaptive pathways in evolved ethanol-tolerant bacterial populations. *Mol Biol Evol* 2017;34:2927–43. <https://doi.org/10.1093/molbev/msx228>
12. Badran A.H., Liu D.R. Development of potent in vivo mutagenesis plasmids with broad mutational spectra. *Nat Commun* 2015;6:8425. <https://doi.org/10.1038/ncomms9425>
13. Molina R.S., Rix G., Mengiste A.A. *et al.* In vivo hypermutation and continuous evolution. *Nat Rev Methods Primers* 2022;2:36. <https://doi.org/10.1038/s43586-022-00119-5>
14. Sherer N.A., Kuhlman T.E. *Escherichia coli* with a tunable point mutation rate for evolution experiments. *G3* 2020;10:2671–81. <https://doi.org/10.1534/g3.120.401124>
15. Maki H. Origins of spontaneous mutations: specificity and directionality of base-substitution, frameshift, and sequence-substitution mutageneses. *Annu Rev Genet*

- 2002;36:279–303.
<https://doi.org/10.1146/annurev.genet.36.042602.094806>
16. Foster P.L., Lee H., Popodi E *et al.* Determinants of spontaneous mutation in the bacterium *Escherichia coli* as revealed by whole-genome sequencing. *Proc Natl Acad Sci USA* 2015;112:E5990–5999. <https://doi.org/10.1073/pnas.1512136112>
 17. Ganai R.A., Johansson E. DNA replication—a matter of fidelity. *Mol Cell* 2016;62:745–55.
<https://doi.org/10.1016/j.molcel.2016.05.003>
 18. Timinskas K., Balvociute M., Timinskas A *et al.* Comprehensive analysis of DNA polymerase III alpha subunits and their homologs in bacterial genomes. *Nucleic Acids Res* 2014;42:1393–413.
<https://doi.org/10.1093/nar/gkt900>
 19. Sanjanwala B., Ganesan A.T. Genetic structure and domains of DNA polymerase III of *Bacillus subtilis*. *Mol Gen Genet* 1991;226:467–72. <https://doi.org/10.1007/BF00260660>
 20. Dervyn E., Suski C., Daniel R *et al.* Two essential DNA polymerases at the bacterial replication fork. *Science* 2001;294:1716–9. <https://doi.org/10.1126/science.1066351>
 21. Fujii S., Akiyama M., Aoki K *et al.* DNA replication errors produced by the replicative apparatus of *Escherichia coli*. *J Mol Biol* 1999;289:835–50. <https://doi.org/10.1006/jmbi.1999.2802>
 22. Niccum B.A., Lee H., MohammedIsmail W *et al.* The spectrum of replication errors in the absence of error correction assayed across the whole genome of *Escherichia coli*. *Genetics* 2018;209:1043–54.
<https://doi.org/10.1534/genetics.117.300515>
 23. Zhou Z.X., Lujan S.A., Burkholder A.B *et al.* How asymmetric DNA replication achieves symmetrical fidelity. *Nat Struct Mol Biol* 2021;28:1020–8. <https://doi.org/10.1038/s41594-021-00691-6>
 24. Bruck I., Goodman M.F., O'Donnell M. The essential C family DnaE polymerase is error-prone and efficient at lesion bypass. *J Biol Chem* 2003;278:44361–8.
<https://doi.org/10.1074/jbc.M308307200>
 25. Paschalis V., Le Chatelier E., Green M *et al.* Interactions of the *Bacillus subtilis* DnaE polymerase with replisomal proteins modulate its activity and fidelity. *Open Biol* 2017;7:170146.
<https://doi.org/10.1098/rsob.170146>
 26. Bolz N.J., Lenhart J.S., Weindorf S.C *et al.* Residues in the N-terminal domain of MutL required for mismatch repair in *Bacillus subtilis*. *J Bacteriol* 2012;194:5361–7.
<https://doi.org/10.1128/JB.01142-12>
 27. Pillon M.C., Lorenowicz J.J., Uckelmann M *et al.* Structure of the endonuclease domain of MutL: unlicensed to cut. *Mol Cell* 2010;39:145–51. <https://doi.org/10.1016/j.molcel.2010.06.027>
 28. Lenhart J.S., Schroeder J.W., Walsh B.W *et al.* DNA repair and genome maintenance in *Bacillus subtilis*. *Microbiol Mol Biol Rev* 2012;76:530–64. <https://doi.org/10.1128/MMBR.05020-11>
 29. Kadyrov F.A., Dzantiev L., Constantin N *et al.* Endonucleolytic function of MutLalpha in human mismatch repair. *Cell* 2006;126:297–308. <https://doi.org/10.1016/j.cell.2006.05.039>
 30. Lujan S.A., Williams J.S., Pursell Z.F *et al.* Mismatch repair balances leading and lagging strand DNA replication fidelity. *PLoS Genet* 2012;8:e1003016.
<https://doi.org/10.1371/journal.pgen.1003016>
 31. Klocko A.D., Schroeder J.W., Walsh B.W *et al.* Mismatch repair causes the dynamic release of an essential DNA polymerase from the replication fork. *Mol Microbiol* 2011;82:648–63.
<https://doi.org/10.1111/j.1365-2958.2011.07841.x>
 32. Schroeder J.W., Hirst W.G., Szewczyk G.A *et al.* The effect of local sequence context on mutational bias of genes encoded on the leading and lagging strands. *Curr Biol* 2016;26:692–7.
<https://doi.org/10.1016/j.cub.2016.01.016>
 33. Sung W., Ackerman M.S., Gout J.F *et al.* Asymmetric context-dependent mutation patterns revealed through mutation-accumulation experiments. *Mol Biol Evol* 2015;32:1672–83. <https://doi.org/10.1093/molbev/msv055>
 34. Schaaper R.M. Mechanisms of mutagenesis in the *Escherichia coli* mutator *mutD5*: role of DNA mismatch repair. *Proc Natl Acad Sci USA* 1988;85:8126–30. <https://doi.org/10.1073/pnas.85.21.8126>
 35. Sambrook J., Fritsch E.F., Maniatis T. *Molecular Cloning: A Laboratory Manual*. NY: Cold Spring Harbor ed. Cold Spring Harbor Laboratory, 1989.
 36. Konkol M.A., Blair K.M., Kearns D.B. Plasmid-encoded ComI inhibits competence in the ancestral 3610 strain of *Bacillus subtilis*. *J Bacteriol* 2013;195:4085–93.
<https://doi.org/10.1128/JB.00696-13>
 37. Koo B.M., Kritikos G., Farelli J.D *et al.* Construction and analysis of two genome-scale deletion libraries for *Bacillus subtilis*. *Cell Syst* 2017;4:291–305. <https://doi.org/10.1016/j.cels.2016.12.013>
 38. Bazill G.W., Gross J.D. Mutagenic DNA polymerase in *B. subtilis*. *Nat New Biol* 1973;243:241–3.
<https://doi.org/10.1038/newbio243241a0>
 39. Barnes M.H., Hammond R.A., Kennedy C.C *et al.* Localization of the exonuclease and polymerase domains of *Bacillus subtilis* DNA polymerase III. *Gene* 1992;111:43–9.
[https://doi.org/10.1016/0378-1119\(92\)90601-K](https://doi.org/10.1016/0378-1119(92)90601-K)
 40. Evans R.J., Davies D.R., Bullard J.M *et al.* Structure of PolC reveals unique DNA binding and fidelity determinants. *Proc Natl Acad Sci USA* 2008;105:20695–700.
<https://doi.org/10.1073/pnas.0809989106>
 41. Zheng Q. A second look at the final number of cells in a fluctuation experiment. *J Theor Biol* 2016;401:54–63.
<https://doi.org/10.1016/j.jtbi.2016.04.027>
 42. Zheng Q. rSalvador: an R package for the fluctuation experiment. *G3* 2017;7:3849–56. <https://doi.org/10.1534/g3.117.300120>
 43. Foster P.L. Methods for determining spontaneous mutation rates. *Methods Enzymol* 2006;409:195–213.
[https://doi.org/10.1016/S0076-6879\(05\)09012-9](https://doi.org/10.1016/S0076-6879(05)09012-9)
 44. Li H., Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 2009;25:1754–60.
<https://doi.org/10.1093/bioinformatics/btp324>
 45. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* 2011;27:2987–93. <https://doi.org/10.1093/bioinformatics/btr509>
 46. Mose L.E., Perou C.M., Parker J.S. Improved indel detection in DNA and RNA via realignment with ABRA2. *Bioinformatics* 2019;35:2966–73. <https://doi.org/10.1093/bioinformatics/btz033>
 47. Nicolas P., Mader U., Dervyn E *et al.* Condition-dependent transcriptome reveals high-level regulatory architecture in *Bacillus subtilis*. *Science* 2012;335:1103–6.
<https://doi.org/10.1126/science.1206848>
 48. Wake R.G. Replication fork arrest and termination of chromosome replication in *Bacillus subtilis*. *FEMS Microbiol Lett* 1997;153:247–54.
<https://doi.org/10.1111/j.1574-6968.1997.tb12581.x>
 49. Sung W., Ackerman M.S., Dillon M.M *et al.* Evolution of the insertion-deletion mutation rate across the tree of life. *G3* 2016;6:2583–91. <https://doi.org/10.1534/g3.116.030890>
 50. Plummer M. *3rd International Workshop on Distributed Statistical Computing (DSC 2003)*, Vienna, Austria, 2003, 1–10.
 51. Guiziou S., Sauveplane V., Chang H.J *et al.* A part toolbox to tune genetic expression in *Bacillus subtilis*. *Nucleic Acids Res* 2016;44:7495–508.
 52. Long H., Miller S.F., Williams E *et al.* Specificity of the DNA mismatch repair system (MMR) and mutagenesis bias in bacteria. *Mol Biol Evol* 2018;35:2414–21.
<https://doi.org/10.1093/molbev/msy134>
 53. Schroeder J.W., Yeessin P., Simmons L.A *et al.* Sources of spontaneous mutagenesis in bacteria. *Crit Rev Biochem Mol Biol* 2018;53:29–48. <https://doi.org/10.1080/10409238.2017.1394262>
 54. Perfeito L., Sousa A., Bataillon T *et al.* Rates of fitness decline and rebound suggest pervasive epistasis. *Evolution* 2014;68:150–62.
<https://doi.org/10.1111/evo.12234>
 55. Singh T., Hyun M., Sniegowski P. Evolution of mutation rates in hypermutable populations of *Escherichia coli* propagated at very small effective population size. *Biol Lett* 2017;13:20160849.
<https://doi.org/10.1098/rsbl.2016.0849>

56. Mahilkar A., Raj N., Kemkar S *et al.* Selection in a growing colony biases results of mutation accumulation experiments. *Sci Rep* 2022;12:15470. <https://doi.org/10.1038/s41598-022-19928-5>
57. Wahl L.M., Agashe D. Selection bias in mutation accumulation. *Evolution* 2022;76:528–40. <https://doi.org/10.1111/evo.14430>
58. Herr A.J., Ogawa M., Lawrence N.A *et al.* Mutator suppression and escape from replication error-induced extinction in yeast. *PLoS Genet* 2011;7:e1002282. <https://doi.org/10.1371/journal.pgen.1002282>
59. Tate J.G., Bamford S., Jubb H.C *et al.* COSMIC: the catalogue of somatic mutations in cancer. *Nucleic Acids Res* 2019;47:D941–7. <https://doi.org/10.1093/nar/gky1015>
60. Schaaper R.M., Radman M. The extreme mutator effect of *Escherichia coli* *mutD5* results from saturation of mismatch repair by excessive DNA replication errors. *EMBO J* 1989;8:3511–6. <https://doi.org/10.1002/j.1460-2075.1989.tb08516.x>
61. Foster P.L., Niccum B.A., Popodi E *et al.* Determinants of base-pair substitution patterns revealed by whole-genome sequencing of DNA mismatch repair defective *Escherichia coli*. *Genetics* 2018;209:1029–42. <https://doi.org/10.1534/genetics.118.301237>
62. Lee H., Popodi E., Tang H *et al.* Rate and molecular spectrum of spontaneous mutations in the bacterium *Escherichia coli* as determined by whole-genome sequencing. *Proc Natl Acad Sci USA* 2012;109:E2774–2783. <https://doi.org/10.1073/pnas.1210309109>
63. Nemenman I., Shafee F., Bialek W. In: Dietterich TG, Becker S, Ghahramani Z (eds.), *Advances in Neural Information Processing Systems 14: Proceedings of the 2001 Conference*. The MIT Press, 2001, <https://doi.org/10.7551/mitpress/1120.003.0065>
64. Haradhvala N.J., Kim J., Maruvka Y.E *et al.* Distinct mutational signatures characterize concurrent loss of polymerase proofreading and mismatch repair. *Nat Commun* 2018;9:1746. <https://doi.org/10.1038/s41467-018-04002-4>
65. Frank A.C., Lobry J.R. Asymmetric substitution patterns: A review of possible underlying mutational or selective mechanisms. *Gene* 1999;238:65–77. [https://doi.org/10.1016/S0378-1119\(99\)00297-8](https://doi.org/10.1016/S0378-1119(99)00297-8)
66. Bhagwat A.S., Hao W., Townes J.P *et al.* Strand-biased cytosine deamination at the replication fork causes cytosine to thymine mutations in *Escherichia coli*. *Proc Natl Acad Sci USA* 2016;113:2176–81. <https://doi.org/10.1073/pnas.1522325113>
67. Lopez-Olmos K., Hernandez M.P., Contreras-Garduno J.A *et al.* Roles of endonuclease V, uracil-DNA glycosylase, and mismatch repair in *Bacillus subtilis* DNA base-deamination-induced mutagenesis. *J Bacteriol* 2012;194:243–52. <https://doi.org/10.1128/JB.06082-11>
68. Patlan-Vazquez A.G., Ayala-Garcia V.M., Vallin C *et al.* Dynamics of mismatch and alternative excision-dependent repair in replicating *Bacillus subtilis* DNA examined under conditions of neutral selection. *Front Microbiol* 2022;13:866089. <https://doi.org/10.3389/fmicb.2022.866089>
69. Andrianova M.A., Bazykin G.A., Nikolaev S.I *et al.* Human mismatch repair system balances mutation rates between strands by removing more mismatches from the lagging strand. *Genome Res* 2017;27:1336–43. <https://doi.org/10.1101/gr.219915.116>