



**HAL**  
open science

# The mutational landscape of *Bacillus subtilis* conditional hypermutators suggests how proofreading inherently skews polymerase error rates

Ira Tanneur, Etienne Dervyn, Cyprien Guérin, Guillaume Kon Kam King, Matthieu Jules, Pierre Nicolas

## ► To cite this version:

Ira Tanneur, Etienne Dervyn, Cyprien Guérin, Guillaume Kon Kam King, Matthieu Jules, et al.. The mutational landscape of *Bacillus subtilis* conditional hypermutators suggests how proofreading inherently skews polymerase error rates. 2024. hal-04366758

**HAL Id: hal-04366758**

**<https://hal.science/hal-04366758v1>**

Preprint submitted on 29 Dec 2023

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

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



Distributed under a Creative Commons Attribution 4.0 International License

# The mutational landscape of *Bacillus subtilis* conditional hypermutators suggests how proofreading inherently skews polymerase error rates

Ira Tanneur<sup>1,2</sup>, Etienne Dervyn<sup>1</sup>, Cyprien Guérin<sup>2</sup>, Guillaume Kon Kam King<sup>2</sup>, Matthieu Jules<sup>1,‡</sup> and Pierre Nicolas<sup>2,‡</sup>

<sup>1</sup> *Université Paris-Saclay, INRAE, AgroParisTech, Micalis Institute, 78350, Jouy-en-Josas, France.*

<sup>2</sup> *Université Paris-Saclay, INRAE, MaIAGE, 78350, Jouy-en-Josas, France.*

<sup>‡</sup> To whom correspondence should be addressed: [matthieu.jules@inrae.fr](mailto:matthieu.jules@inrae.fr); [pierre.nicolas@inrae.fr](mailto:pierre.nicolas@inrae.fr).

Version: Dec 29, 2023

## Abstract

The sources of spontaneous point mutations in genomes are diverse, including DNA damages and errors introduced by the polymerase during replication. The resulting mutation rate also depends on the counteracting action of DNA repair mechanisms, with mutator phenotypes appearing constantly and playing a role during phases of rapid evolution in nature and in the laboratory. Here, we use the gram-positive model bacterium *Bacillus subtilis* to jointly assess the respective contributions to the mutation rate of DNA polymerase nucleotide selectivity, proofreading, and mismatch repair (MMR). For this purpose, we constructed and analysed several conditional hypermutators with a proofreading-deficient allele of *polC* and/or a deficient allele of *mutL*. By covering a wide range of mutation rates and displaying contrasted mutation profiles, these conditional hypermutators enrich the *B. subtilis* synthetic biology toolbox for directed evolution. Analysis of their mutation profiles in light of several mathematical models exposes the difficulties of interpreting apparent probabilities of error correction that stem from aggregating possibly heterogeneous subclasses of mutations into counts, and from unknowns on the components of the mutation profiles in the presence of the repair systems. Aware of these difficulties, the analysis strongly suggests that proofreading is needed to avoid partial saturation of the MMR in *B. subtilis* and that an inherent effect of proofreading is to skew the net polymerase error rates by reinforcing intrinsic biases of nucleotide selectivity.

## INTRODUCTION

Substitutions, insertions and deletions of single base pairs in the genome can have diverse consequences on encoded molecular functions, from no effect to abrupt change, most often in the direction of deterioration (Eyre-Walker and Keightley 2007). 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. A drastic hypothesis, supported by the comparative study of mutation rates across the tree of life, postulates that they are simply maintained as low as possible; the limit being the “drift-barrier”, at which the strength of selection is matched by the opposite pressure of genetic drift and mutation, presumably biased towards creating weak mutators (Lynch et al. 2016). Mutation rate would thus have nothing to do with the benefit of evolvability for long-term survival, the energetic cost of fidelity, or a biophysical limit. In bacteria, this rate is already typically as low as one mutation in the genome per thousand of generations, but systems with a mutation rate one order of magnitude below those observed in nature were reported to arise under scenarios of artificial evolution (Deatherage et al. 2018; Dervyn et al. 2023).

Mutator phenotypes however occur constantly and their contribution to evolution is difficult to estimate (Taddei et al. 1997; Couce et al. 2017). In asexual populations, where mutator alleles remain linked with the mutations that they generate across successive generations, mutators are advantageous when the potential for fitness improvement is high. In pathogenic bacteria, mutators have for instance been associated with the apparition of complex antibiotic resistances (Dulanto Chiang et al. 2022), rapid evolution within the host during infection (Oliver and Mena 2010), and atypical virulence traits (Rudenko et al. 2020). Mutators also emerge spontaneously during laboratory evolution in response to applied selective pressures (Sniegowski et al. 1997; Swings et al. 2017). To instrumentalize the potential of mutator phenotypes to foster adaptation, conditional systems have been engineered for some organisms as part of the synthetic biology toolbox (Badran and Liu 2015; Sherer and Kuhlman 2020; Molina et al. 2022). 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, primarily arising from DNA lesions caused by endogenous and exogenous agents, from errors introduced by the DNA polymerase during replication and by error-prone polymerases recruited in response to stress (Maki 2002). The resulting mutation rate depends on the intensity of these sources and of the counteracting action of DNA repair mechanisms that work in coordination to achieve transmission of correct genetic material 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 (Foster et al. 2015). The NER (Nucleotide Excision Repair) is necessary for repairing various drug- and UV-induced lesions (*i.e.* bulky lesions), whereas the BER (Base Excision Repair) is essential for repairing lesions caused by a variety of chemical assaults, such as alkylation, oxidation, deamination, etc. (*i.e.* non-bulky lesions). Beyond damages, mutations can also arise from errors made during DNA replication. DNA replication accuracy 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 (Ganai and Johansson 2016).

In bacteria, genome replication is performed by a multiprotein machine classified into the C-family of DNA polymerase holoenzymes, in which the catalytic polymerase  $\alpha$ -subunit exists in two primary forms, DnaE and PolC (Timinskas et al. 2014). 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* (Sanjanwala and Ganesan 1991). In this organism, replication elongation involves two essential polymerases, PolC and DnaE (Dervyn et al. 2001). These enzymes have distinct functions: PolC ensures most of the DNA synthesis, but only DnaE can extend 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 *E. coli*'s  $\alpha$ -subunit replicative machinery has been estimated at approximately  $10^{-6}$  per base pair per generation both *in vitro* (Fujii et al. 1999) and *in vivo* (Niccum et al. 2018). Given this error rate and the size of *E. coli*'s genome, around 5 native replication errors 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 polymerases (Timinskas et al. 2014), including the *B. subtilis* PolC. In contrast, DnaE polymerases do not possess their own exonuclease domain. The proofreading activity of the *E. coli* DNA PolIII holoenzyme containing DnaE relies on an exonuclease domain found in the  $\epsilon$ -subunit. Error made by polymerases devoid of proofreading can also sometimes be corrected by a process known as proofreading *in trans*, or extrinsic proofreading, well described between eukaryotic DNA polymerases (Zhou et al. 2021). In *B. subtilis*, data suggest that PolC exonuclease is able to proofread errors made by the error-prone DnaE polymerase (Bruck et al. 2003; Paschalis et al. 2017).

The MMR is a universal mechanism that is responsible for correcting errors formed during DNA replication and that escaped proofreading. Upon the identification of a replication error, the mismatch sensing protein, MutS, recruits MutL. Most prokaryotic and eukaryotic MutL homologs from human to bacteria possess a highly conserved endonuclease active site that serves to remove mismatches (Pillon et al. 2010; Bolz et al. 2012). *E. coli* has been the primary model for studying MMR, however its MutL does not possess the endonuclease activity, which is encoded in a distinct protein, MutH, which specifically nicks the unmethylated and thus nascent strand bearing the mismatch (Lenhart et al. 2012). 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 (Kadyrov et al. 2006). In bacteria, MMR increases the fidelity of the chromosomal DNA replication pathway approximately 100-fold and MMR is viewed as a system directed to the repair of the most frequent replication errors (Lujan et al. 2012). Mutator phenotypes found in nature are often generated by mutations inactivating the MMR.

Studying organisms, such as *B. subtilis*, with a PolC polymerase and an MMR pathway more widely conserved across biology can offer key insights into the coordinated functioning of these systems within living cells (Klocko et al. 2011). In prolongation of previous works characterising the mutation profiles of MMR-deficient *B. subtilis* strains (Sung et al. 2015; Schroeder et al. 2016),

the main goal of our study was to jointly assess in *B. subtilis* the respective contributions to the mutation rates of nucleotide selectivity, proofreading and MMR and their interdependencies. For this purpose, we constructed and analysed several conditional hypermutators with a proofreading-deficient allele of *polC* and/or a deficient allele of *mutL*. 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* (Schaaper 1988; Niccum et al. 2018), but also that an inherent effect of proofreading is to skew the net polymerase error rates. The conditional hypermutators cover a wide range of mutation rates and display contrasted mutation profiles, enriching the *B. subtilis* synthetic biology toolbox for directed evolution.

## RESULTS

**The mutation rate of *Bacillus subtilis* can be increased up to 6,000 times.** We built five mutant strains with expected hypermutator phenotypes from a *B. subtilis* 168-derived strain. The first two strains are constitutively MMR-deficient as a result of the single deletions of *mutS* and *mutL* (later denoted  $\Delta S$  and  $\Delta L$ ). The three other strains were designed for conditional inactivation of either one or both of these two DNA repair pathways and are thus expected to be inducible hypermutators (**Figures S1 and S2**). In practice, the IPTG-inducible promoter  $P_{hs}$  controls the expression of mutant alleles selected for their reported or suspected ability to competitively displace their functional counterparts. The first allele, denoted here as *mutL\**, has a mutation in the ATP hydrolysis active site of MutL which was described to have a dominant negative effect (Bolz et al. 2012). The second allele, denoted here as *polC\**, encodes a proofreading-deficient variant of PolC with a mutation in its exonuclease domain (Sanjanwala and Ganesan 1991). The last strain, expected to display the highest mutation rate under full induction, expresses these two deficient alleles in a synthetic operon (*mutL\*polC\**). Shorthand notations are used below for the 168-derived parental strain serving as reference and the five mutant strains:  $R^{168}$ ,  $\Delta L$ ,  $\Delta S$ ,  $L^*$ ,  $C^*$ ,  $LC^*$ .

Fluctuation assays were performed (**Figure S3**) to compare the rate of mutation to rifampicin resistance of these strains, in the absence or presence of IPTG. Point estimates and confidence intervals are shown in **Figure 1**, and detailed results in **Table S1**. In the absence of IPTG, the rate of mutation of  $R^{168}$  was estimated at  $9.74 \times 10^{-10}$  per generation. Constitutive inactivation of the MMR in  $\Delta L$  and  $\Delta S$  increased the mutation rate by a factor of approximately 85, with no statistically significant difference between these 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, in a similar fluctuation assay based on rifampicin (Schroeder et al. 2016). Mutation rates without IPTG were slightly higher for the inducible strains ( $L^*$ ,  $C^*$ ,  $LC^*$ ) than for  $R^{168}$  (up to  $1.09 \times 10^{-8}$  for  $LC^*$ ); probably reflecting the basal low-level activity already described for  $P_{hs}$  (Guiziou et al. 2016). From there, the mutation rate of the three inducible strains increased with IPTG concentration, until a plateau reached between 50 and 100  $\mu M$ . At full induction (100  $\mu M$  IPTG), they exhibited clearly distinct mutation rates ranging from  $3.66 \times 10^{-7}$  for  $L^*$  to  $5.78 \times 10^{-6}$  for  $LC^*$ . The mutation rate in  $L^*$  is comparable or slightly higher than in  $\Delta L$  and  $\Delta S$ , while the mutation rate in  $LC^*$  represents an increase by a factor of approximately 6,000 as compared to  $R^{168}$ .

Mutation rates much higher than in the reference strain may induce stress responses, possibly altering the physiology of each mutant differently. Nevertheless, we did not detect any substantial impact on growth in 96-well microtiter plates (**Figure S4**). We also performed transcriptomics experiments on the  $R^{168}$ ,  $L^*$ ,  $C^*$  and  $LC^*$  strains in the presence of 100  $\mu M$  IPTG. The analyses did not reveal any significantly differentially expressed genes between the strains. However, they allowed us to quantify the expression of mutant alleles relative to wild-type alleles for the genes *mutL* and *polC*. Upon induction, the mutant allele accounted for 96-98% of the total mRNA pool for the considered gene (**Table S2**), which is consistent with the results of fluctuation assays giving a mutation rate in  $L^*$  as high as in  $\Delta L$  and  $\Delta S$  (*i.e.* total inactivation of the MMR). We therefore concluded that the mutational profiles of these strains can be compared with each other and 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 experiments give access to the molecular nature of mutations (Lynch et al. 2016), unlike fluctuation assays which only provide a mutation rate aggregated over an array of mutations conferring a screenable phenotype. For each of the five strains ( $R^{168}$ ,  $\Delta L$ ,  $\Delta S$ ,  $L^*$ ,  $C^*$ ,  $LC^*$ ), four independent mutation-accumulation lines (MA-lines) were propagated by repeated cycles (MA-steps, **Figure S5**) of colony sampling, dilution, and plating on LB, in the presence of 100  $\mu$ M IPTG when relevant. By randomly selecting a single colony, each MA step creates a bottleneck in the propagated population. The purpose of these bottlenecks is to limit genetic diversity and thereby maximise random genetic drift and minimise natural selection. The interval between two bottlenecks, or one MA-step, was estimated to represent an average of 25.6 generations. We performed whole-genome sequencing of the endpoint of each line; we also sequenced intermediate time-points (up to 4 for  $LC^*$ ) to detect changes in mutation rates (**Figure S5**). In total, 56 clones isolated after 1 to 37 MA-steps were sequenced. Substitutions, insertions, and deletions were identified, and mutation rates per base pair (bp) and generation (abbreviated  $\text{bp}^{-1}.\text{gen}^{-1}$ ) were estimated in each time interval of each MA-line. To increase statistical power, we also incorporated previously collected data from mutation-accumulation experiments using *B. subtilis* 3610 ( $R^{3610}$ ) and its corresponding  $\Delta\text{mutS}$  mutant strain ( $\Delta S^{3610}$ ) into our analysis (Sung et al. 2015; Sung et al. 2016). The detailed list of all mutations found is provided in **Table S3**.

In the four independent lines of  $R^{168}$ , only one nucleotide substitution was identified after 37 MA-steps, giving an average substitution rate of about  $7 \times 10^{-11} \text{ bp}^{-1}.\text{gen}^{-1}$  (**Table 1**). This was not statistically significantly different from a previous report on *B. subtilis* (Sung et al. 2016), which was recalculated to  $3.4 \times 10^{-10} \text{ bp}^{-1}.\text{gen}^{-1}$  (**Table 1**).

Between 113 and 157 nucleotide substitutions were identified after 21 MA-steps in each of the  $\Delta L$ ,  $\Delta S$  and  $L^*$  strains, resulting in point estimates of the substitution rates between  $1.4 \times 10^{-8}$  and  $1.9 \times 10^{-8} \text{ bp}^{-1}.\text{gen}^{-1}$ , with no statistically significant differences observed between the strains (**Table 1**). The absence of statistical differences between  $\Delta L$ ,  $\Delta S$  and  $L^*$  strains, all derived from *B. subtilis* 168, led us to aggregate the data collected for these three strains under the label  $\text{MMR}^{-168}$  (**Table 1**). For the  $L^*$  strain, sequencing an intermediate time-point located at the end of MA-step 11 did not reveal any difference between the rates of accumulation in the first and in the second part of the evolution (**Figure 2** and **Figure S6**). In the  $\text{MMR}$ -deficient strains, substitutions identified at end points of the MA-lines appeared thus to result from accumulation at a constant rate.

In contrast to  $L^*$  MA-lines,  $C^*$  and  $LC^*$  MA-lines displayed a tendency towards decreasing substitution rates over their evolution, with differences between the two strains in terms of frequency, temporality, and magnitude of decrease (**Figure 2** and **Figure S6**). A statistically significant decrease was detected for a single  $C^*$  MA-line but for all  $LC^*$  MA-lines. Notably, the decrease was only detected during the second half of an evolution of 21 MA-steps for this  $C^*$  MA-line ( $C^*_3$ ), whereas it was detected as early as during the second MA-step for  $LC^*_2$  ( $p$ -value =  $8.4 \times 10^{-3}$ ), and during MA-steps 3-6 for the 3 other  $LC^*$  MA-lines. The magnitude of the decrease was also only a factor  $\sim 2.5$ x for  $C^*_3$  but reached  $\sim 50$ x 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 1,129 identified substitutions occurred in intervals affected by a decrease in the substitution rate (**Table 1**). We therefore decided to retain

only the data corresponding to time intervals before any detected decrease to establish the mutation rates of the strains. This resulted in point estimates for the substitution rates of the  $C^*$  and  $LC^*$  strains of  $5.5 \times 10^{-8}$  and  $5.5 \times 10^{-7}$   $\text{bp}^{-1} \cdot \text{gen}^{-1}$ , respectively (**Table 1**).

Statistically significant decreases were also observed in indel rate (**Figure S7**) and correlated with decreases in substitution rate (**Figure S8**). This is consistent with contributions of both proofreading and MMR in correcting the indel errors introduced by the polymerase activity of PolC (**Supplementary Methods and Results 1.1**) whose rate increases with the length of the homopolymers (**Figure S9**). Insertions and deletions rates computed are given in **Table S4**.

Nonsynonymous mutations were found in the inducible synthetic circuits of 4 out of 5 MA-lines exhibiting a decrease in mutation rate (**Table S5, Figure S10, Supplementary Methods and Results 1.2**). 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 4 times higher than expected in the absence of selection (Chi-squared test with simulated *p*-values, *p*-value= $4.1 \times 10^{-2}$ ). Recent studies have concluded that positive selection is possible in mutation-accumulation experiments despite the extreme bottlenecks imposed on the population (Mahilkar et al. 2022; Wahl and Agashe 2022). Here, the over-representation of mutations in the genetic elements conferring the strongest hypermutator phenotypes indicates adaptive evolution by positive selection to decrease the mutation rate. This finding echoes previous studies that reported changes in the mutation rates in mutation-accumulation experiments (Perfeito et al. 2014; Singh et al. 2017). Additionally, it gives *a posteriori* experimental justification to the choice of restricting the mutation-accumulation experiments on proofreading-deficient *E. coli* to 3 to 6 MA-steps to minimise selection (Niccum et al. 2018).

**Proofreading repairs errors leading to transversions at least as well as those leading to transitions.** The sequence data do not provide information on the DNA strand on which the error that led to the mutation initially 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 prior-free analyses of strand asymmetries in mutational profiles and follows the convention used in cancer research (Tate et al. 2019).

In all strains, a slightly higher number of mutations was found on C than on T bases (**Figure 3**). All strains also exhibited a predominance of transitions over transversions, but the strength of this bias differs between strains (**Figure 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 ones, found in the MMR- strains ( $\Delta S^{3610}$  and MMR-<sup>168</sup> strains). The  $C^*$  strain displays an intermediate proportion (point estimate 0.12). The small number of transversions 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 exhibit an excess of  $C \rightarrow A$  over  $C \rightarrow G$  not seen in other strains. The approximately 10-fold increase in the proportion of transitions when comparing the MMR-<sup>168</sup> strains to the  $R^{3610}$  strain is consistent with previously published results on  $\Delta S^{3610}$  (Sung et al. 2015). Indeed, the MMR has a general tendency to reduce the transition rate much more than the transversion rate across microorganisms (Lujan et al. 2012; Long et al. 2018).



When comparing  $LC^*$  and  $MMR^{-168}$  strains, the inactivation of PolC proofreading activity leads to a slight increase in the proportion of transversions (from 0.01 to 0.04). However, this difference is not statistically significant (**Table 1**). This comparison suggests that proofreading corrects errors leading to both transversions and transitions with at least similar efficiency. Interestingly, the conclusion drawn from the comparison between  $LC^*$  and  $MMR^{-168}$  strains may appear contradictory to what one might infer from the comparison between the  $C^*$  and  $R^{3610}$  strains. When comparing  $C^*$  and  $R^{3610}$ , the inactivation of proofreading significantly reduces the proportion of transversions among substitutions (from 0.25 to 0.12). Even if there is a substantial statistical uncertainty associated with the estimation of the proportion of transversions, particularly in  $MMR$ -deficient contexts, these opposite trends pose a question. Together with the magnitude of fold-change which is larger when inactivating  $MMR$  in presence of proofreading than in its absence (**Figure 3** and **Table 1**), this indeed suggests that  $MMR$  may already be partially inactivated in absence of proofreading ( $C^*$ ), as already described in *E. coli* where the fold-change in substitution rate upon inactivation of the  $MMR$  was also reported to be much lower in cells deficient for PolIII holoenzyme proofreading than in wild-type (Schaaper 1988; Niccum et al. 2018). This was interpreted as resulting from  $MMR$  saturation by the high number of errors introduced during DNA replication; an hypothesis further supported by direct assays of  $MMR$  activity and restoration of the  $MMR$  by overexpression (Schaaper 1988; Schaaper and Radman 1989). We will later formally explore this interpretation using a model-based analysis integrating information on the chromosomal context of mutations (adjacent nucleotides and strand).

**Strand-asymmetry of substitution rate at C:G sites is apparent only after proofreading.** To further characterise the two DNA repair systems, we investigated for each strain how substitution rates were altered by the distance to the origin of replication, the orientation of the strand relative to replication or transcription, the coding or noncoding regions, and the level of transcription. For these analyses, we counted substitutions in the different chromosomal contexts and computed the corresponding “local” substitution rates.

In  $MMR^{-168}$  strains, like in  $\Delta S^{3610}$  and  $R^{3610}$  strains, the substitution rates are significantly higher when C is on the leading than on the lagging strand (**Figure 4A**). This agrees with previous results in *B. subtilis* (Sung et al. 2015), 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  $R^{3610}$  and  $\Delta S^{3610}$  strains. This bias is not present in the  $C^*$  and  $LC^*$  strains, which are proofreading-deficient. Besides the strong asymmetry at C:G sites, we detected a weaker, but statistically significant, replication-oriented asymmetry on 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 (**Figure 4A**). In  $R^{3610}$ , this bias on T:A sites seems less pronounced and is possibly inexistent. In keeping with the conclusions of Schroeder et al. (2016), analysis of localization with respect to transcription, which is most often collinear to replication in *B. subtilis*, did not point to a contribution of transcription-related processes to these asymmetries between strands in any of the hypermutator strains (**Supplementary Methods and Results 1.3, Figure S11AB**).

In wild-type, the replication-oriented asymmetry of substitution rates at G:C sites is a prominent characteristic of the mutational profile. In line with previous analyses of  $MMR$ -deficient strains,

among which  $\Delta S^{3610}$ , this bias was also detected in MMR-<sup>168</sup> strains. Our data reveal its absence in *C\** and *LC\** strains. A similar observation was made in *E. coli* proofreading-deficient strains and led to the interpretation that proofreading is strand-biased and creates this bias (Niccum et al. 2018), but a concurrent explanation will be discussed below. In contrast, the wild-type exhibits little or no asymmetry at A:T sites, whereas such asymmetry exists in the hypermutator strains; error correction systems, and in particular the MMR, tend to reduce or even cancel 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 (**Figure S11C**) and, to a lesser extent, in the half-chromosome near replication origin (**Figure S11D**). A higher substitution rate in non-coding than in coding regions, specific to the wild-type, was also pointed out in *E. coli* and was interpreted as an indication that “MMR preferentially repair coding sequences” (Lee et al. 2012; Foster et al. 2018). Alternatively, trends exclusively observed in the wild-type may correspond to substitutions originating from processes not subject to correction by proofreading and MMR, which could be masked by elevated rates of substitutions in hypermutator strains.

**Polymerase errors still shape the distribution of mutations after proofreading.** In all strains, the substitution rate is heavily influenced by the nucleotide adjacent in the 5' or 3' position to the focal pyrimidine (**Figure 4C** and **Figure 4D**). This observation, previously made in wild-type and MMR-deficient strains (Sung et al. 2015), extends to proofreading-deficient strains. Considering simultaneously the adjacent nucleotides on both sides (**Table S6** and **Figure S12**) and the replication strand, requires binning the counts into 64 replication-stranded triplets. To mitigate the dimensionality issue, exemplified by the absence of observed substitutions for some bins, we adopted a Bayesian estimation framework which incorporates the mean and standard-deviation of log-transformed rates as strain-specific hyperparameters (**Supplementary Methods and Results 1.4**). Information is thereby borrowed from the whole distribution to establish point estimates and credibility intervals of the substitution rate for a given triplet (**Figure 5A**).

Based on these estimates, we measured a very strong correlation between the substitution rates of the 64 replication-stranded triplets in MMR-<sup>168</sup> and  $\Delta S^{3610}$  (Pearson correlation  $r=0.91$  on log-transformed rates), which confirms that our MMR-<sup>168</sup> background is very close to the mutant previously studied (Sung et al. 2015). There is also a strong correlation between the substitution profiles of *LC\** and proofreading-proficient MMR-deficient strains ( $r=0.79$  between *LC\** and  $\Delta S^{3610}$ ). This is in line with the idea that, after PolC proofreading and before correction by the MMR, most of the errors leading to substitutions originate from misincorporations by the PolC polymerase that escaped proofreading. Since MMR removes ~98-99% of these errors (**Table 1**) they also constitute the bulk of errors corrected by the MMR.

We further noticed very substantial correlations between the substitution profiles of the wild-type and those of all hypermutator strains. They were the highest with the MMR-deficient proofreading-proficient strains ( $r=0.72$  with  $\Delta S^{3610}$ ,  $r=0.71$  with MMR-<sup>168</sup>) and remained statistically significant with the proofreading-deficient strains ( $r=0.58$  with *LC\** and  $r=0.57$  with *C\**,  $p$ -values $<10^{-6}$ ). The correlation between triplet substitution rate 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 which escaped the proofreading and the MMR shape substantially the substitution profile of the wild-type.

**Proofreading squares the biases of polymerase error rates.** As measured by the standard deviation and the span of their distributions, point estimates in MMR-<sup>168</sup> and  $\Delta S^{3610}$  are approximately 10-fold more dispersed than in other strains, doubling the dispersion in log-scale (**Figure 5B**). Because point estimates of triplet substitution rates cannot be precisely estimated based on small counts, we sought an alternative approach that could directly quantify dispersion. Using a coincidence-counting method for entropy estimation developed to be robust to sparse sampling (Nemenman et al. 2001), we calculated the Kullback-Leibler (KL) divergence between the unknown underlying distribution of the observed counts and the theoretical distribution assuming a uniform substitution rate (expected count proportional to the number of occurrences of the triplet in the chromosome). The estimates of KL divergences (**Figure 5C**) confirmed both the similar level of dispersion of the substitution rates in the wild-type ( $R^{3610}$ ) and the proofreading-deficient strains, and the comparatively much higher dispersion in MMR-deficient strains (approximately 2.5-fold higher divergence from uniform). In other words, PolC proofreading activity contributes to dispersing the substitution rates, while MMR activity counteracts this dispersion. This observation, which suggests a compensation for the biases of proofreading by opposite biases of MMR correction in *B. subtilis*, may seem surprising. However, it is consistent with the conclusions of a study on *E. coli* PolIII holoenzyme proofreading (Niccum et al. 2018).

After proofreading (but before MMR correction), the substitution profile becomes more dispersed but remains very similar in terms of the direction of biases compared to before proofreading. This observation is puzzling since it means that proofreading amplifies biases which 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 implication of this observation, we can formulate a minimal model in which the proofreading activity reduces the initial probability of error 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 re-incorporation of a nucleotide with the probability of error  $\gamma[i]$  characteristic of the initial polymerase activity. The probability of error after proofreading then writes  $e[i]=\gamma[i](1-d[i])/(1-\gamma[i]d[i])$ , 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 cycle, 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 **Figure 5B**), remains extremely small in comparison to 1. To increase biases to the extent of doubling the dispersion in log-scale (*i.e.*  $e[i]/e[j]=(\gamma[i]/\gamma[j])^2$ ), as approximately observed in our data, it supposes a probability of non-detection and removal in the first step of proofreading  $(1-d[i])$  proportional to  $\gamma[i]$ , the probability of error of the sole polymerase activity.

**A theoretical model suggests saturation of MMR beyond 5 errors per DNA replication cycle.**

Rates of error corrections measured for proofreading and MMR are heavily influenced by the presence or absence of the other system. For instance, proofreading decreases the overall substitution rate by a factor of 162 in the presence of MMR versus only 32 in its absence (resp.  $\mu_{R3610}/\mu_{C^*}$  and  $\mu_{MMR-168}/\mu_{LC^*}$ , denoting by  $\mu$  the considered mutation rate and using the ML estimates

of **Table 1**). Similarly, MMR reduces the substitution rate more significantly when proofreading is present compared to when it is absent. Using Bayesian estimates of substitution rates, none of the triplets or types of mutations (ts, tv, ins, del) clearly contradicts this observation (with a posterior probability of being above the diagonal greater than 5% for all triplets; **Figure 6**). The trend becomes even more pronounced if  $\Delta S^{3610}$  is used instead of MMR-<sup>168</sup> as a representative of the MMR-deficient proofreading-proficient background (**Figure S12**).

To thoroughly examine 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:

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

where  $\gamma[i]$  is the error rate before correction by proofreading or MMR, and  $\theta$  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 subjected to correction by MMR, reducing the number of errors by a factor  $q_{\text{MMR}}[i]$ . In this model, whose assumptions are presented along with an associated Bayesian estimation procedure in **Supplementary Methods and Results 1.5**, the mutation rate in  $C^*$  can be expressed as a function of the rates in the three other backgrounds by identifying  $q_{\text{MMR}}[i]$  to  $\mu_{\text{wt}}[i]/\mu_{\text{MMR}}[i]$  and  $\gamma[i]$  to  $\mu_{\text{LC}^*}[i]$ .

This parameterization was used to estimate the mixture parameter  $\theta$  and to check the agreement of the model to the experimental data (**Figure 7**). In practice, the posterior distribution of  $\theta$  was estimated based either on the rates of substitutions for the 64 replication-oriented triplets (with MMR-<sup>168</sup> or  $\Delta S^{3610}$  as representative of MMR-deficient background) or on the rates of the 4 types of mutations (transition, transversion, insertion, deletion). These three posterior distributions were very similar (**Figure 7A**), the posterior mean for  $\theta$  varying only from 0.071 to 0.084. Observed counts aggregated by type of mutations (**Figure 7C**) and by triplet (**Figures S14 and S15**) fall within the prediction intervals of the model, indicating a good fit to the experimental dataset. Of note, although the fraction of errors introduced by PolC that arise in a context of MMR saturation remains small ( $\theta < 10\%$ ), these errors, which cannot be corrected by the saturated MMR anymore, are responsible for the majority of the mutations observed in strain  $C^*$  (**Figure 7A**).

In a simple mechanistic model, these values of  $\theta$  may correspond to the capacity of the MMR to handle (with a probability of failure  $q_{\text{MMR}}[i]$ ) 4 to 5 errors introduced in the genome per replication cycle (**Figure 7B, Supplementary Methods and Results 1.5**). Alternatively,  $\theta$  values may also be interpreted as reflecting the fraction of errors made by the polymerase before covering a certain distance after a first error.

## DISCUSSION

**Aggregating mutations in counts: a necessary evil.** Our experimental data could, in principle, also be explained by a model which does not involve MMR saturation. This alternative model acknowledges that counting mutations on the genome implies aggregating them by types or contexts, themselves probably encompassing subclasses of mutations originating from errors which are not corrected with the same probability of success. In the extreme scenario where MMR and proofreading correct non-overlapping sets of errors, the effects of separately inactivating the two systems would simply add to each other in cells where both systems are inactivated. This idea has already been mentioned and refuted for mutational signatures of human cancers (Haradhvala et al. 2018). In fact, while changes in mutation rates between strains may not be strictly multiplicative, they are clearly more than additive. This raises interest for a more general model that considers aggregated subclasses of mutations, not distinguished in the counts and arbitrary correction rates. We algebraically explored this model in its simplest form with only two subclasses (**Supplementary Methods and Results 1.6**). Interestingly, as soon as 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, measured in terms of mutation rate fold-change, depends on the presence or absence of the other system.

With more parameters than data points (5 parameters for 4 mutation rates), even the simplest scenario with two subclasses can fit almost any data set, explaining from additive to super-multiplicative effects. The model is therefore difficult to falsify. Nevertheless, we note that the sign of this epistasis depends on whether the systems exhibit similar or opposite specificities: if they do tend to repair different subclasses of mutations, the apparent epistasis will be negative (sub-multiplicative, as it is the case in our data), otherwise the epistasis will be positive (super-multiplicative). If this explanatory model had a prevalent role in the explanation of the apparent epistasis reported here, it is difficult to understand why positive epistasis is not observed in any triplet context or type of mutation. Indeed, explaining negative epistasis with this model would be at odds with the generally admitted idea that MMR correction is mostly coreplicative and targets the same errors as those corrected by proofreading, *i.e.* Watson-Crick mismatches introduced during DNA replication by the DNA polymerase. If aggregation of subclasses does not substantially contribute to the observed apparent epistasis which seems well explained by MMR saturation, it complicates interpretation of the mutation profiles.

### **A cautious interpretation of the apparent efficiency of proofreading and MMR.**

Acknowledging saturation of the MMR, the proofreading-deficient strain  $C^*$  should not be considered as a fully MMR-proficient. Thus, it cannot serve to estimate the efficiency of the proofreading and MMR, *i.e.* as the numerator in a ratio  $\mu_{C^*}/\mu_{LC^*}$  to estimate the probability that an error escapes the proofreading, and as the denominator in a ratio  $\mu_{wt}/\mu_{C^*}$  to estimate the probability that an error escapes correction by the MMR in physiological conditions relevant for the wild-type. Instead, these escape probabilities should be estimated by the ratios  $\mu_{MMR-}/\mu_{LC^*}$  and  $\mu_{wt}/\mu_{MMR-}$ , which are represented in **Figure 8** (MMR- being either MMR-<sup>168</sup> or  $\Delta S^{3610}$ ) for each replication-oriented triplet and type of mutation (ts, tv, ins, del). The mutation rates in  $LC^*$  and wild-type, which

respectively correspond to the rates of errors leading to mutation before and after correction by the combined action of proofreading and MMR are also shown.

Proofreading and MMR escape probabilities ( $\mu_{\text{MMR-}}/\mu_{\text{LC}^*}$  and  $\mu_{\text{wt}}/\mu_{\text{MMR-}}$ ) are correlated to the substitution rates of the triplets (**Figure S16**). In **Figure 8**, triplets are ordered according to wild-type substitution rates (as in **Figure 5A**), which highlights general trends: proofreading escape tends to be higher on triplets where the polymerase initially makes more incorporation errors (**Figure S16**,  $r=0.32$   $p$ -value=0.0097); MMR escape tends to be lower on triplets with high mutation rate in wild-type (**Figure S16**,  $r=-0.45$ ,  $p$ -value=0.00018). Of note, these two correlations are of opposite sign to the artefactual correlations that would be generated by noise in the estimates of  $\mu_{\text{LC}^*}[\text{i}]$  and  $\mu_{\text{wt}}[\text{i}]$ , which also enter in the ratios  $\mu_{\text{MMR-}}[\text{i}]/\mu_{\text{LC}^*}[\text{i}]$  and  $\mu_{\text{wt}}[\text{i}]/\mu_{\text{MMR-}}[\text{i}]$ .

The positive correlation between proofreading escape and mutation rate in  $\text{LC}^*$  matches the observation that proofreading increases the biases of polymerase error rates. Statistical uncertainty on proofreading escape probability makes it difficult to draw conclusions on specific triplets, and few triplets show clear strand asymmetry. It is however interesting to note that this positive correlation is not apparent within pairs of triplets that only differ by the strand relative to DNA replication. In particular,  $\text{GT}\underline{\text{G}}$  and  $\text{AT}\underline{\text{G}}$  are more affected by substitutions in  $\text{LC}^*$  when the focal pyrimidine (T) is on the lagging strand, but no difference is detected in terms of proofreading escape probability between strands. For these triplets, it might be that errors are primarily made by the DnaE polymerase during synthesis of the lagging strand, generating strand asymmetry of substitution rates in  $\text{LC}^*$  without affecting proofreading. Reciprocally, errors affecting  $\text{GCT}$  and  $\text{GCG}$  seem less corrected by proofreading when C is on the leading strand, but the substitution rate in  $\text{LC}^*$  is similar between strands. Deamination of cytosine to uracil in the lagging strand template (*i.e.* the leading strand) due to exposure of single-stranded DNA within the context of the replication fork has been, for a long time, suspected to be the reason behind the near-universal GC-skew between replication strands in prokaryotes (Frank and Lobry 1999) and does represent a substantial source of C to T transitions in wild-type (Bhagwat et al. 2016). Strikingly, MMR-deficient strains show a strong similarly replication-oriented bias for substitutions at C sites (Lee et al. 2012; Sung et al. 2015), which may have the same origin as in wild-type. The hypothesis that cytosine deamination could also create this bias in MMR-deficient strains has already been evoked (Foster et al. 2018). In the absence of proofreading, these errors caused by deamination would be drowned out by polymerase misincorporation errors, and this could explain why substitutions at C sites apparently escape more the proofreading when on the leading strand. This would indeed be consistent with studies reporting a role of the *B. subtilis* MMR in counteracting the effects of base deamination (López-Olmos et al. 2012; Patlan-Vazquez et al. 2022).

The hypothesis of proofreading efficiency connected to initial polymerase incorporation accuracy was proposed based on the analysis of substitution rates by triplet context without distinguishing transversion and transition. Transversions represent only 4% of the substitutions in  $\text{LC}^*$ . The lower proofreading escape probability of errors leading to transversions than to transitions (**Figure 8**) is consistent with the above hypothesis. However, the difference between the proofreading escape probabilities of the two types of errors, which cannot be estimated precisely, seems smaller than expected. A possible explanation would be an overestimation of the proofreading escape for misincorporation errors leading to transversions, as would happen if a substantial fraction of the

rare errors leading to transversions that persist after proofreading in MMR-deficient strains did not stem from polymerase errors and would thus simply not be subjected to proofreading.

Concerning MMR escape, which is not the primary focus of this study, the strong negative correlation with proofreading escape (**Figure S17**,  $r=-0.63$ ,  $p\text{-value}=2.6\times 10^{-8}$ ) could suggest an evolution towards correcting the most frequent errors resulting from DNA replication. This would indeed echo the proposed role of differential MMR efficiency in balancing the DNA replication fidelity between the two strands in eukaryotes (Lujan et al. 2012; Andrianova et al. 2017; Zhou et al. 2021). However, we need to underscore the lack of information on the fraction of the residual errors seen as mutations in wild-type which originate from DNA replication and are subjected to -but missed by- MMR correction. Indeed, the contribution of different sources of spontaneous mutations in wild-type remains poorly understood (Schroeder et al. 2018). Among studies addressing this question, it has been shown that the inactivation of oxidative damage repair pathways increases the mutation rate in *E. coli* under conditions of optimal growth without external stress, suggesting that events escaping correction by oxidative damage repair pathways may contribute to the mutation profile of the wild-type (Foster et al. 2015). Since oxidation damages tend to cause transversions, this hypothesis would indeed be in line with the higher proportion of transversions in wild-type than in any of our hypermutator strains. If a substantial fraction of the substitutions in wild-type originate from other sources of errors not subjected to MMR correction, variations not directly linked to MMR efficiency may enter into the ratio  $\mu_{wt}/\mu_{MMR-}$  used to estimate the MMR escape probability (**Figure 8**), via its numerator (*e.g.*, the number of errors from these other sources) and its denominator (*e.g.*, the number of errors remaining after proofreading). Indeed, variations in the denominator clearly drive the variations of the ratio (**Figure S16**,  $r=-0.94$ ,  $p\text{-value}=2.3e-31$ ) and because the proofreading escape probability ( $\mu_{MMR-}/\mu_{LC*}$ ) is itself positively correlated with the denominator (**Figure S16**,  $r=0.71$ ,  $p\text{-value}=7.3e-11$ ) this could contribute to the negative correlation observed between the MMR escape probability and the proofreading escape probability.

**Conclusion.** Saturation of the MMR in the absence of proofreading, greater dispersion of the 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 shared traits between *B. subtilis* and *E. coli*. Given the considerable divergence between these two organisms in terms of phylogenetic distance and molecular organisation of the DNA polymerase and MMR, these traits are thus probably common to many other organisms. Characterization of the overdispersion of the mutation rates in MMR-deficient proofreading-proficient strains compared to other strains led us to propose the hypothesis that proofreading intrinsically skews DNA polymerase error rates. This could represent a drawback of the principle of proofreading, which relies on the DNA polymerase to detect its own errors, a function of judge and party, leading to introducing the same biases in nucleotide misincorporation and error escape rates.

This study also attempts to examine carefully the consequences of aggregating mutations in counts, and more generally acknowledges the difficulty of interpreting apparent error escape rates. First, unaware aggregation in the analysis of subclasses of mutations resulting from different molecular pathways can create almost any pattern of apparent interactions in the effects of deactivating different error correction systems. Second, interpreting the effect of disabling one error correction system is generally complex due to the uncertainty on the contribution of errors corrected by this system to the mutations observed in its presence. This uncertainty makes it difficult to interpret the appearance of strand bias upon activation of proofreading, which may be caused by a mutation process post-proofreading, such as deamination, rather than by a bias in proofreading. Similarly, the “flattening” of biases in mutation rates upon activation of the MMR is difficult to interpret since it could result from a better efficiency of MMR at correcting the most common errors but also from the contribution of multiple sources of mutations to the profile of the wild-type. The hypothesis of a substantial contribution of multiple sources seems consistent with the drift-barrier hypothesis (Lynch et al. 2016). In wild-type, the rates of the mutations resulting from different molecular pathways may have been pushed down below the same point where they do not encounter significant counter-selection.

The construction of strains with inducible hypermutator phenotypes circumvented the problem of stability which compromised the reproducibility of some studies on *E. coli* proofreading-deficient strains (discussed in (Niccum et al. 2018)). It was also motivated by the interest in tools for synthetic biology applications. Concerning the future use of our inducible systems to accelerate evolution, a proofreading-deficient strain would yield a less biased mutation spectrum than an MMR-deficient strain. If needed, extreme mutation rates can be obtained by inactivating the two reparation systems simultaneously but strong counterselection against mutational load makes that such an induction can be envisioned only over a short period of time.



## MATERIALS AND METHODS

**Media and bacterial strains.** *E. coli* DH5 $\alpha$  was used for plasmid construction and transformation using standard techniques (Sambrook et al. 1989). *B. subtilis* strains used in this study derived from the Master Strain (MS), a prophage-free and *trp*<sup>+</sup> derivative of *B. subtilis* 168 (Dervyn et al. 2023), also denoted here R<sup>168</sup>. Lysogeny broth (LB) was used to grow *E. coli* and *B. subtilis*. Transformation of *B. subtilis* cells was performed using the protocol from (Konkol et al. 2013). When required, 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.** The  $\Delta S$  and  $\Delta L$  mutant strains were generated by transforming the PCR-amplified *kan mutL* and *mutS kan* sequences, using the P1-P2 primer pair along with genomic DNA from  $\Delta\text{mutS}::\text{kan}$  and  $\Delta\text{mutL}::\text{kan}$  mutant strains, respectively. These strains were obtained from the previously published single-gene deletion library of *B. subtilis* (Koo et al. 2017). For the construction of the *L*<sup>\*</sup> strain, the first and second half of the *mutL* gene were PCR-amplified using the P5-P8 and P6-P7 primer pairs (**Figure S1**), respectively, with the P7 and P8 primers both carrying the desired point mutation (as indicated in **Table S7**). The two fragments were then assembled by PCR to lead to the *mutL*(N34H) allele. The backbone of the pDR111 plasmid (kind gift from D. Rüdner), which contains the isopropyl- $\beta$ -D-1-thiogalactopyranoside (IPTG) inducible  $P_{\text{hyperspank}}$  promoter (denoted  $P_{\text{hs}}$ ) and the *spec* gene (conferring resistance to spectinomycin), was PCR-amplified using primers P3 and P4. The 5' extensions of the P5 and P6 primers then allowed for 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  $P_{\text{hs}}$  in a *B. subtilis amyE*-integrative plasmid (**Figure S1**). Similarly, for the construction of the *C*<sup>\*</sup> strain, the *polC* allele found in *B. subtilis mut-1* (Bazill and Gross 1973), characterised by the G430E and S621N mutations (Sanjanwala and Ganesan 1991), was PCR-amplified using the P11-P12 primer pair 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  $P_{\text{hs}}$  in a *B. subtilis amyE*-integrative plasmid (**Figure S2**). For the construction of the *LC*<sup>\*</sup> strain, a *mutL*<sup>\*</sup> *polC*<sup>\*</sup> synthetic operon was generated by assembly of the *mutL*(N34H) allele PCR-amplified from strain *L*<sup>\*</sup> using P5 and P11, and the *polC mut-1* allele PCR-amplified from *C*<sup>\*</sup> using P12 and P10, 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 *mutL*<sup>\*</sup> *polC*<sup>\*</sup> synthetic operon under the control of  $P_{\text{hs}}$  in a *B. subtilis amyE*-integrative plasmid (**Figure S2**). Plasmids were transformed into *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 were deposited on GEO under the accession GSE239804.

**Fluctuation assays.** For each strain to be tested, a single colony was grown in 1 mL LB at 37°C for 90 minutes. This preculture was serially diluted in fresh LB to initiate cultures with a small number of cells  $N_0$ . Cells were then grown during 7.5 h hours to reach saturation. In case the induction with IPTG was to be tested, LB medium with the desired concentration of IPTG was prepared right before use from an IPTG stock concentration of 1 mM. When the volume of culture was 1 mL, cultures were centrifuged before plating so as to keep the cells viable, then 750  $\mu\text{L}$  of supernatant

were removed. The remaining 250  $\mu\text{L}$  were gently vortexed before plating on LB supplemented with rifampicin ( $10 \mu\text{g}\cdot\text{mL}^{-1}$ ). For each assay, a predefined number of cultures (8 for the  $R^{168}$  assay, 4 for  $\Delta S$  and  $\Delta L$  assays, and 3 for all other assays) was not plated on LB medium supplemented with rifampicin, but instead serially diluted and plated on LB agar, in order to determine the final number of cells ( $N_t$ ) in each culture. Fluctuation assays performed on the same day were considered to have the same distribution of the final number of cells. All other cultures were plated on LB agar with rifampicin to obtain the number of Rif-resistant colony-forming units (CFUs). All plates were incubated at  $37^\circ\text{C}$  and scored for CFUs after 24 h of growth. The maximum likelihood estimator (MLE) of the number of mutations by assay ( $m$ ), as well as the confidence interval, were computed under the Luria-Delbrück model, and by taking into account the variation in the final number of cells (Zheng 2016), using the `newton.B0` with default parameters and `confint.B0` functions of R package “Rsalvador” v1.7 (Zheng 2017). For the computation of confidence intervals, the initial guess for the parameter  $m$  was taken as the  $m$  given by the “`newton.B0`” function. We consider here that the mutation probability is constant over the cell cycle, so that the mutation rate per base per generation is the mutation rate per base per cell division (Foster 2006). 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<sup>R</sup> emergence was therefore calculated as  $\mu_{Rif} = m/N_t$ .

**Mutation-accumulation experiments and sequencing.** An isolated colony was collected each day (24 h at  $37^\circ\text{C}$ ), suspended in culture medium + glycerol 20%, and diluted by  $2 \times 10^5$ , a factor that allows for distinguishable colonies, before plating on LB agar (+ 100  $\mu\text{M}$  IPTG for the  $L^*$ ,  $C^*$  and  $LC^*$  strains) for the next MA-step. Counting of the colonies present on the agar plate gave an estimate of the number of bacteria initially present in the diluted colony and thereby of the number of generations per MA-step. Four parallel MA-lines of successive MA-steps were propagated per strain (21 MA-steps for  $\Delta S$ ,  $\Delta L$ ,  $L^*$  and  $C^*$ , 11 for  $LC^*$ ).

For sequencing at intermediate and end points of the MA-lines, 5 to 50% of the picked colony was cultured in LB medium to collect cells. DNA was extracted using the GenElute<sup>TM</sup> Bacterial Genomic DNA Kit (Sigma-Aldrich) following the supplied protocol. The DNA samples corresponding to an intermediate time-point in the four parallel MA-lines for a same strain were pooled in equimolar proportions. Simple and pooled DNA samples were sequenced (150-bp paired-end reads) on an Illumina platform (NovaSeq 6000) to an average depth of  $\sim 300$ . The reads can be accessed in NCBI SRA (BioProject PRJNA995423).

**Detection of mutations.** The reads were aligned to the reference sequence of the *B. subtilis* 168 genome (GenBank: AL009126.3) using the BWA-MEM v0.7.17 (Li and Durbin 2009), after quality control and trimming using sickle v1.33 (command “`sickle pe`” with options “`-t sanger -x -q 20 -l 20`”) (Joshi and Fass 2011). Properly paired reads, selected using “`samtools view -f 3`” (samtools v1.14, Li 2011), were locally realigned around indels using ABRA2 v2.24 (Mose et al. 2019). The number of occurrences of each nucleotide (base read quality  $\geq 35$ ) and indels at each position of the reference in confidently mapped reads (alignment quality  $\geq 50$ ) was counted using “`samtools mpileup`” with options “`-aa -d 5000 -q 50 -Q 35 -x -B`”. These numbers of occurrences were analysed using R.

For each position, we calculated the effective depth (DPeff) as the total number of informative reads. A reference subset of positions common to all samples was determined for the computation of the mutation rates. This reference consisted of positions well covered (DPeff $\geq$ 100) and sequenced on both strands ( $\geq$ 10% of the reads on the less represented strand) in all samples. Most of the regions with a low coverage matched with the regions that were deleted during the construction of the MS / R<sup>168</sup> strain (Dervyn et al. 2023), which lacks 233.4 kb of the chromosome relative to wild-type *B. subtilis* 168, and with the multicopy structural RNAs. Overcovered regions were also eliminated from this reference subset of positions: the region of gene *upp* and downstream (positions 3,788,426 to 3,789,124), repeated due to pop-ins - pop-outs at that locus during the construction of the R<sup>168</sup> strain; the region from position 2,432,478 to 2,433,315, over-covered in *polC*\* samples; the regions of genes *polC* (1,727,133 to 1,731,446) and *mutL* (1,778,337 to 1,780,539) duplicated by the 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) that served as our reference chromosome for the computation of mutation rates.

The distribution of the proportion of non-reference reads in the different samples was graphically examined to establish relevant cut-offs for identifying mutations. A mutation was identified at the end point of a MA-line if, in the corresponding sample, 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. When intermediate time-points were available for this MA-line, the mutation was traced back to the first time-point in which it appeared at a frequency  $\geq$ 5% of the reads in the corresponding pooled sequence sample. Due to the detection during graphical examination of a contamination from other samples, we lowered the cut-off from 75% to 60% for the identification of mutation in the third MA-line of  $\Delta S$  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 prior to the mutation-accumulation experiment and were discarded for calculation of the mutation rates.

**Detection of mutations in the inducible synthetic circuits.** Reads were also aligned to the reference sequences of the inserted regions represented in **Figure S2**. For the mutant alleles *mutL*\* and *polC*\*, reads originating from the native alleles (*mutL* and *polC*) mapped also on the insert and variant calling, in itself, does not distinguish a mutation in the native and in the mutant allele. However, we found that, at positions where bases differed from the reference on these genes in individual samples, the proportion of these alternative bases was bimodally distributed, with two peaks, around 40% and 60% of the reads (**Figure S10**). Given that the characteristic point mutations of both *polC*\* and *mutL*\* accounted for more than 50% of the reads (resp. between 52% and 71% and between 65% and 75%) at their respective positions, we anticipated that the mutations for which  $\sim$ 60% of the reads differed from the reference were on the mutant allele, while the others were on the native allele. This prediction is consistent with the position of the *amyE* locus in which the mutant alleles are inserted, *i.e.* closer to the replication origin of the chromosome than both native alleles, and thus expected to be more abundant in the sample due to ongoing replication. As a verification, we used specific primers to amplify and determine the sequence of either the native or inserted allele. Of the 5 PCR-verified mutations, all of them were found on the predicted allele (**Table S5**).

**Chromosome partitioning to assess the impact of transcription and replication on mutation rate.** To assess the impact 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 ~10% of the genome and contains transcribed untranslated regions (UTRs) we also sought to assess the impact of transcription with more statistical power and in a more accurate way than permitted by the GenBank annotation. For this purpose, two categories of regions of roughly equivalent sizes were also defined based on the transcribed regions identified across 269 samples of a wild-type strain and representative of a wide diversity of growth conditions (Nicolas et al. 2012). These two categories reflected the quantity of transcripts in LB as measured in 9 samples corresponding to the growth in liquid LB (triplicate samples for the exponential, transition and stationary phases) and 2 samples corresponding to 16 hours of growth on LB (non-confluent colonies). The regions of “high” transcription level were those belonging to the top 30% in at least one of these 11 samples. Conversely, the regions of “low” transcription level were those that never belonged to 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 3,622 non-overlapping, transcription-oriented regions covering 84.9% of the reference genome (43.4% for “high”, 41.4% for “low”).

To assess the impact of the DNA replication strand, the leading and lagging strands were defined based on the replication origin (position 1) and the middle of the centromere terminus of replication (position 2,018,289) (Wake 1997). To assess the impact of DNA replication timing, the genome was divided into a replication “first half” corresponding to the 2 Mbp of *B. subtilis* 168 centred on the replication origin (position 1) and a “second half”.

**Mutation rate estimations and comparisons.** To incorporate the list mutations in R<sup>3610</sup> and  $\Delta S^{3610}$  into our analysis, the positions given on the *B. subtilis* NCIB 3610 genome (GenBank: CM000488.1) by (Sung et al. 2015) and (Sung et al. 2016) were transferred to the *B. subtilis* 168 genome by mapping of the 41 bp-long sequence centred on each mutation site. Keeping only the exact and unique matches, more than 99% of these mutations were transferred (**Table S3**), with a perfect collinearity between the positions of the mutations on both reference genomes.

Maximum-likelihood estimates of the mutation rates were obtained as  $\mu = m / (T \times G)$ , where  $m$  is the total number of mutations of a considered type in a considered genotype and genomic context (nucleotide at focal position and adjacent nucleotides, orientation with respect to replication, transcription, ...),  $T$  is the total number of occurrences of the genomic context in the reference sequence, and  $G$  the number of considered generations in MA-lines. Confidence intervals for these point estimates were calculated using the exact method for Poisson counts implemented in R package “epitools” v0.5-10.1, with  $m$  as the count and  $T \times G$  as the time-person at risk.

To assess if a factor impacts the substitution rates, we used Generalised Linear Models (GLMs) for Poisson distributed count data with log-link, combined Analysis of Variance (ANOVA) (R package “stats” v3.6.3) to compare the fit of a GLM that accounts and a GLM that does not account for the considered factor. This statistical comparison was done separately for each genotype.

Markov chain Monte Carlo methods implemented in JAGS (Plummer 2003) accessed through R package “rjags” were used for Bayesian estimation via posterior sampling, in particular for

estimation replication-stranded triplet mutation rates and MMR saturation parameter  $\theta$ . Details on models and algorithms settings are provided in **Supplementary Methods and Results 1.4**.

#### **Mathematical modelling of mutation rates**

Assumptions and Bayesian estimation procedure for the model with saturation of the MMR are presented in **Supplementary Methods and Results 1.5**. Algebraic analysis of the general model with two subclasses of errors and two repair pathways is presented in **Supplementary Methods and Results 1.6**.

## **ACKNOWLEDGEMENTS**

We are grateful to Guillaume Achaz, Clément Nizak and Paulo Tavares for their valuable advice during the course of this work. We thank the INRAE MIGALE bioinformatics facility (<https://doi.org/10.15454/1.5572390655343293E12>) for providing computational resources.

## **AUTHOR CONTRIBUTIONS**

IT, MJ, and PN conceived the project and designed the experimental plan. ED and IT performed the experiments. IT processed the raw data with the help of CG. GKKK, IT, and PN performed the statistical and mathematical analyses. IT, MJ, and PN interpreted the results and wrote the manuscript with the contribution of all authors.

## **FUNDING**

This work was supported by the French National Research Agency (ANR-18-CE43-0002). The PhD fellowship of IT was also funded by the MathNum division of INRAE.

## REFERENCES

- Andrianova MA, Bazykin GA, Nikolaev SI, Seplyarskiy VB. 2017. Human mismatch repair system balances mutation rates between strands by removing more mismatches from the lagging strand. *Genome Res* **27**: 1336-1343.
- Badran AH, Liu DR. 2015. Development of potent in vivo mutagenesis plasmids with broad mutational spectra. *Nat Commun* **6**: 8425.
- Bazill GW, Gross JD. 1973. Mutagenic DNA polymerase in *B. subtilis*. *Nat New Biol* **243**: 241-243.
- Bhagwat AS, Hao W, Townes JP, Lee H, Tang H, Foster PL. 2016. Strand-biased cytosine deamination at the replication fork causes cytosine to thymine mutations in *Escherichia coli*. *Proc Natl Acad Sci U S A* **113**: 2176-2181.
- Bolz NJ, Lenhart JS, Weindorf SC, Simmons LA. 2012. Residues in the N-terminal domain of MutL required for mismatch repair in *Bacillus subtilis*. *J Bacteriol* **194**: 5361-5367.
- Bruck I, Goodman MF, O'Donnell M. 2003. The essential C family DnaE polymerase is error-prone and efficient at lesion bypass. *J Biol Chem* **278**: 44361-44368.
- Couce A, Caudwell LV, Feinauer C, Hindre T, Feugeas JP, Weigt M, Lenski RE, Schneider D, Tenaillon O. 2017. Mutator genomes decay, despite sustained fitness gains, in a long-term experiment with bacteria. *Proc Natl Acad Sci U S A* **114**: E9026-E9035.
- Deatherage DE, Leon D, Rodriguez AE, Omar SK, Barrick JE. 2018. Directed evolution of *Escherichia coli* with lower-than-natural plasmid mutation rates. *Nucleic Acids Res* **46**: 9236-9250.
- Dervyn E, Planson AG, Tanaka K, Chubukov V, Guerin C, Derozier S, Lecointe F, Sauer U, Yoshida KI, Nicolas P et al. 2023. Greedy reduction of *Bacillus subtilis* genome yields emergent phenotypes of high resistance to a DNA damaging agent and low evolvability. *Nucleic Acids Res* doi:10.1093/nar/gkad145.
- Dervyn E, Suski C, Daniel R, Bruand C, Chapuis J, Errington J, Janniere L, Ehrlich SD. 2001. Two essential DNA polymerases at the bacterial replication fork. *Science* **294**: 1716-1719.
- Dulanto Chiang A, Patil PP, Beka L, Youn JH, Launay A, Bonomo RA, Khil PP, Dekker JP. 2022. Hypermutator strains of *Pseudomonas aeruginosa* reveal novel pathways of resistance to combinations of cephalosporin antibiotics and beta-lactamase inhibitors. *PLoS Biol* **20**: e3001878.
- Eyre-Walker A, Keightley PD. 2007. The distribution of fitness effects of new mutations. *Nature Reviews Genetics* **8**: 610-618.
- Foster PL. 2006. Methods for determining spontaneous mutation rates. *Methods Enzymol* **409**: 195-213.
- Foster PL, Lee H, Popodi E, Townes JP, Tang H. 2015. Determinants of spontaneous mutation in the bacterium *Escherichia coli* as revealed by whole-genome sequencing. *Proc Natl Acad Sci U S A* **112**: E5990-5999.
- Foster PL, Niccum BA, Popodi E, Townes JP, Lee H, MohammedIsmail W, Tang H. 2018. Determinants of Base-Pair Substitution Patterns Revealed by Whole-Genome Sequencing of DNA Mismatch Repair Defective *Escherichia coli*. *Genetics* **209**: 1029-1042.
- Frank AC, Lobry JR. 1999. Asymmetric substitution patterns: a review of possible underlying mutational or selective mechanisms. *Gene* **238**: 65-77.

- Fujii S, Akiyama M, Aoki K, Sugaya Y, Higuchi K, Hiraoka M, Miki Y, Saitoh N, Yoshiyama K, Ihara K et al. 1999. DNA replication errors produced by the replicative apparatus of *Escherichia coli*. *J Mol Biol* **289**: 835-850.
- Ganai RA, Johansson E. 2016. DNA Replication-A Matter of Fidelity. *Mol Cell* **62**: 745-755.
- Guiziou S, Sauveplane V, Chang HJ, Clerte C, Declercq N, Jules M, Bonnet J. 2016. A part toolbox to tune genetic expression in *Bacillus subtilis*. *Nucleic Acids Res* **44**: 7495-7508.
- Haradhvala NJ, Kim J, Maruvka YE, Polak P, Rosebrock D, Livitz D, Hess JM, Leshchiner I, Kamburov A, Mouw KW et al. 2018. Distinct mutational signatures characterize concurrent loss of polymerase proofreading and mismatch repair. *Nat Commun* **9**: 1746.
- Joshi NA, Fass JN. 2011. Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files. <https://github.com/najoshi/sickle>.
- Kadyrov FA, Dzantiev L, Constantin N, Modrich P. 2006. Endonucleolytic function of MutLalpha in human mismatch repair. *Cell* **126**: 297-308.
- Klocko AD, Schroeder JW, Walsh BW, Lenhart JS, Evans ML, Simmons LA. 2011. Mismatch repair causes the dynamic release of an essential DNA polymerase from the replication fork. *Mol Microbiol* **82**: 648-663.
- Konkol MA, Blair KM, Kearns DB. 2013. Plasmid-encoded ComI inhibits competence in the ancestral 3610 strain of *Bacillus subtilis*. *J Bacteriol* **195**: 4085-4093.
- Koo BM, Kritikos G, Farelli JD, Todor H, Tong K, Kimsey H, Wapinski I, Galardini M, Cabal A, Peters JM et al. 2017. Construction and Analysis of Two Genome-Scale Deletion Libraries for *Bacillus subtilis*. *Cell Syst* **4**: 291-305.
- Lee H, Popodi E, Tang H, Foster PL. 2012. Rate and molecular spectrum of spontaneous mutations in the bacterium *Escherichia coli* as determined by whole-genome sequencing. *Proc Natl Acad Sci U S A* **109**: E2774-2783.
- Lenhart JS, Schroeder JW, Walsh BW, Simmons LA. 2012. DNA repair and genome maintenance in *Bacillus subtilis*. *Microbiol Mol Biol Rev* **76**: 530-564.
- Li H. 2011. A statistical framework for SP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**: 2987-2993.
- Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**: 1754-1760.
- Long H, Miller SF, Williams E, Lynch M. 2018. Specificity of the DNA Mismatch Repair System (MMR) and Mutagenesis Bias in Bacteria. *Mol Biol Evol* **35**: 2414-2421.
- López-Olmos K, Hernández MP, Contreras-Garduño JA, Robleto EA, Setlow P, Yasbin RE, Pedraza-Reyes M. 2012. Roles of endonuclease V, uracil-DNA glycosylase, and mismatch repair in *Bacillus subtilis* DNA base-deamination-induced mutagenesis. *Journal of Bacteriology* **194**: 243-252.
- Lujan SA, Williams JS, Pursell ZF, Abdulovic-Cui AA, Clark AB, Nick McElhinny SA, Kunkel TA. 2012. Mismatch repair balances leading and lagging strand DNA replication fidelity. *PLoS Genet* **8**: e1003016.
- Lynch M, Ackerman MS, Gout JF, Long H, Sung W, Thomas WK, Foster PL. 2016. Genetic drift, selection and the evolution of the mutation rate. *Nat Rev Genet* **17**: 704-714.



- Mahilkar A, Raj N, Kemkar S, Saini S. 2022. Selection in a growing colony biases results of mutation accumulation experiments. *Sci Rep* **12**: 15470.
- Maki H. 2002. Origins of spontaneous mutations: specificity and directionality of base-substitution, frameshift, and sequence-substitution mutageneses. *Annu Rev Genet* **36**: 279-303.
- Molina RS, Rix G, Mengiste AA, Alvarez B, Seo D, Chen H, Hurtado J, Zhang Q, Donato Garcia-Garcia J, Heins ZJ et al. 2022. In vivo hypermutation and continuous evolution. *Nat Rev Methods Primers* **2**.
- Mose LE, Perou CM, Parker JS. 2019. Improved indel detection in DNA and RNA via realignment with ABRA2. *Bioinformatics* **35**: 2966-2973.
- Nemenman I, Shafee F, Bialek W. 2001. Entropy and inference, revisited. In *Advances in Neural Information Processing Systems*, Vol 14 (ed. TG Dietterich, et al.). The MIT Press.
- Niccum BA, Lee H, MohammedIsmail W, Tang H, Foster PL. 2018. The Spectrum of Replication Errors in the Absence of Error Correction Assayed Across the Whole Genome of *Escherichia coli*. *Genetics* **209**: 1043-1054.
- Nicolas P, Mäder U, Dervyn E, Rochat T, Leduc A, Pigeonneau N, Bidnenko E, Marchadier E, Hoebeke M, Aymerich S et al. 2012. Condition-dependent transcriptome reveals high-level regulatory architecture in *Bacillus subtilis*. *Science* **335**: 1103-1106.
- Oliver A, Mena A. 2010. Bacterial hypermutation in cystic fibrosis, not only for antibiotic resistance. *Clin Microbiol Infect* **16**: 798-808.
- Paschalis V, Le Chatelier E, Green M, Nouri H, Kepes F, Soultanas P, Janniere L. 2017. Interactions of the *Bacillus subtilis* DnaE polymerase with replisomal proteins modulate its activity and fidelity. *Open Biol* **7**.
- Patlan-Vazquez AG, Ayala-Garcia VM, Vallin C, Cortes J, Vasquez-Morales SG, Robleto EA, Nudler E, Pedraza-Reyes M. 2022. Dynamics of Mismatch and Alternative Excision-Dependent Repair in Replicating *Bacillus subtilis* DNA Examined Under Conditions of Neutral Selection. *Front Microbiol* **13**: 866089.
- Perfeito L, Sousa A, Bataillon T, Gordo I. 2014. Rates of fitness decline and rebound suggest pervasive epistasis. *Evolution* **68**: 150-162.
- Pillon MC, Lorenowicz JJ, Uckelmann M, Klocko AD, Mitchell RR, Chung YS, Modrich P, Walker GC, Simmons LA, Friedhoff P et al. 2010. Structure of the endonuclease domain of MutL: unlicensed to cut. *Mol Cell* **39**: 145-151.
- Plummer M. 2003. JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. In Proceedings of the 3rd International Workshop on Distributed Statistical Computing (DSC 2003), March 20-22, Vienna, Austria. ISSN 1609-395X., volume 124, page 125.
- Rudenko O, Engelstadter J, Barnes AC. 2020. Evolutionary epidemiology of *Streptococcus iniae*: Linking mutation rate dynamics with adaptation to novel immunological landscapes. *Infect Genet Evol* **85**: 104435.
- Sambrook J, Fritsch EF and Maniatis T. 1989. Molecular Cloning: a Laboratory Manual. Cold Spring Harbor ed. Cold Spring Harbor Laboratory, NY.
- Sanjanwala B, Ganesan AT. 1991. Genetic structure and domains of DNA polymerase III of *Bacillus subtilis*. *Mol Gen Genet* **226**: 467-472.

- Schaaper RM. 1988. Mechanisms of mutagenesis in the *Escherichia coli* mutator mutD5: role of DNA mismatch repair. *Proc Natl Acad Sci U S A* **85**: 8126-8130.
- Schaaper RM, Radman M. 1989. The extreme mutator effect of *Escherichia coli* mutD5 results from saturation of mismatch repair by excessive DNA replication errors. *EMBO J* **8**: 3511-3516.
- Schroeder JW, Hirst WG, Szewczyk GA, Simmons LA. 2016. The Effect of Local Sequence Context on Mutational Bias of Genes Encoded on the Leading and Lagging Strands. *Curr Biol* **26**: 692-697.
- Schroeder JW, Yeesin P, Simmons LA, Wang JD. 2018. Sources of spontaneous mutagenesis in bacteria. *Crit Rev Biochem Mol Biol* **53**: 29-48.
- Sherer NA, Kuhlman TE. 2020. *Escherichia coli* with a Tunable Point Mutation Rate for Evolution Experiments. *G3 (Bethesda)* **10**: 2671-2681.
- Singh T, Hyun M, Siegowski P. 2017. Evolution of mutation rates in hypermutable populations of *Escherichia coli* propagated at very small effective population size. *Biol Lett* **13**.
- Sniegowski PD, Gerrish PJ, Lenski RE. 1997. Evolution of high mutation rates in experimental populations of *E. coli*. *Nature* **387**: 703-705.
- Sung W, Ackerman MS, Dillon MM, Platt TG, Fuqua C, Cooper VS, Lynch M. 2016. Evolution of the Insertion-Deletion Mutation Rate Across the Tree of Life. *G3 (Bethesda)* **6**: 2583-2591.
- Sung W, Ackerman MS, Gout J-F, Miller SF, Williams E, Foster PL, Lynch M. 2015. Asymmetric Context-Dependent Mutation Patterns Revealed through Mutation–Accumulation Experiments. *Molecular Biology and Evolution* **32**: 1672-1683.
- Swings T, Weytjens B, Schaleck T, Bonte C, Verstraeten N, Michiels J, Marchal K. 2017. Network-Based Identification of Adaptive Pathways in Evolved Ethanol-Tolerant Bacterial Populations. *Mol Biol Evol* **34**: 2927-2943.
- Taddei F, Radman M, Maynard-Smith J, Toupance B, Gouyon PH, Godelle B. 1997. Role of mutator alleles in adaptive evolution. *Nature* **387**: 700-702.
- Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, Boutselakis H, Cole CG, Creatore C, Dawson E et al. 2019. COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res* **47**: D941-D947.
- Timinskas K, Balvociute M, Timinskas A, Venclovas C. 2014. Comprehensive analysis of DNA polymerase III alpha subunits and their homologs in bacterial genomes. *Nucleic Acids Res* **42**: 1393-1413.
- Wahl LM, Agashe D. 2022. Selection bias in mutation accumulation. *Evolution* **76**: 528-540.
- Wake RG. 1997. Replication fork arrest and termination of chromosome replication in *Bacillus subtilis*. *FEMS Microbiol Lett* **153**: 247-254.
- Zheng Q. 2016. A second look at the final number of cells in a fluctuation experiment. *J Theor Biol* **401**: 54-63.
- Zheng Q. 2017. rSalvador: An R Package for the Fluctuation Experiment. *G3 (Bethesda)* **7**: 3849-3856.
- Zhou ZX, Lujan SA, Burkholder AB, St Charles J, Dahl J, Farrell CE, Williams JS, Kunkel TA. 2021. How asymmetric DNA replication achieves symmetrical fidelity. *Nat Struct Mol Biol* **28**: 1020-1028.

## TABLES

**Table 1.** Aggregated numbers of substitutions, substitution rate, and proportion of transversions for each investigated strain.

Strain <sup>a</sup>	Lines	Generation <sub>s<sup>b,c</sup></sub>	Substitutions <sup>c</sup>			Substitution rate [95% CI]	Proportion of transversions [95% CI]
			Total	Ts <sup>d</sup>	Tv <sup>e</sup>		
R <sup>168</sup>	4	3,790	1	1	0	7.0×10 <sup>-11</sup> [0.18-39×10 <sup>-10</sup> ]	0.00 [0.00-0.98]
R <sup>3610</sup>	49	248,920	319	238	81	3.4×10 <sup>-10</sup> [3.0-3.8×10 <sup>-10</sup> ]	0.25 [0.21-0.31]
ΔS <sup>3610</sup>	19	38,000	4,844	4,711	133	3.4×10 <sup>-8</sup> [3.3-3.5×10 <sup>-8</sup> ]	0.03 [0.02-0.03]
ΔL	4	2,151	149	147	2	1.8×10 <sup>-8</sup> [1.5-2.1×10 <sup>-8</sup> ]	0.01 [0.00-0.05]
ΔS	4	2,151	157	155	2	1.9×10 <sup>-8</sup> [1.6-2.2×10 <sup>-8</sup> ]	0.01 [0.00-0.05]
L*	4	2,151	113	111	2	1.4×10 <sup>-8</sup> [1.1-1.7×10 <sup>-8</sup> ]	0.02 [0.00-0.06]
MMR- <sup>168</sup>	12	6,453	419	413	6	1.7×10 <sup>-8</sup> [1.6-1.9×10 <sup>-8</sup> ]	0.01 [0.01-0.03]
C*	4	1,895 (256)	395 (19)	348 (18)	47 (1)	5.5×10 <sup>-8</sup> [5.0-6.1×10 <sup>-8</sup> ]	0.12 [0.09-0.16]
LC*	4	230 (897)	502 (627)	484 (599)	18 (28)	5.5×10 <sup>-7</sup> [5.2-6.3×10 <sup>-7</sup> ]	0.04 [0.02-0.06]

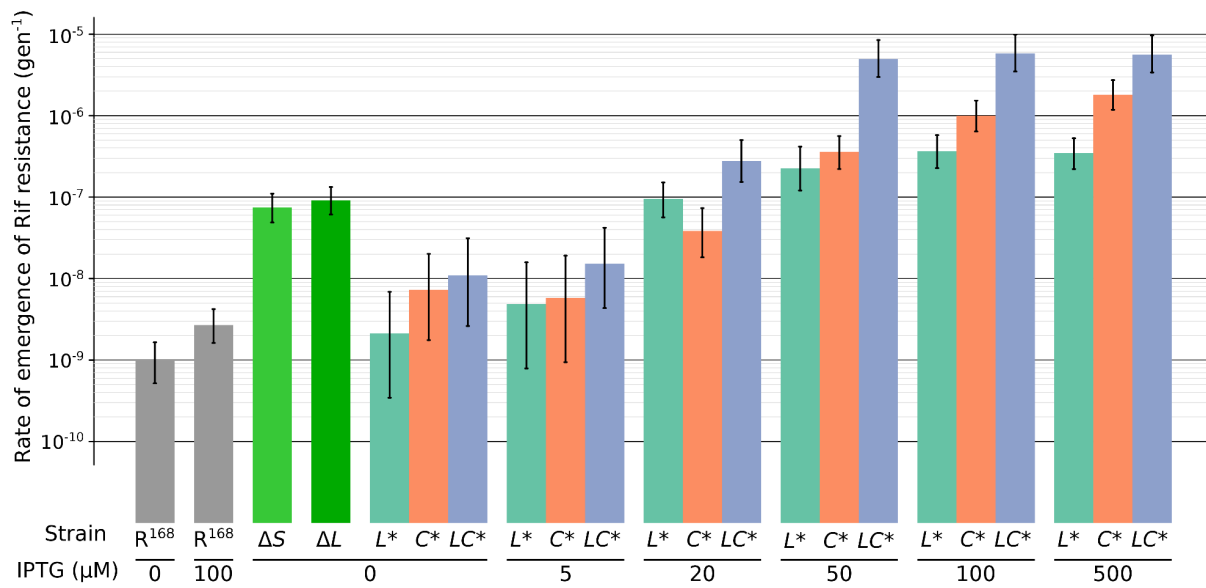
<sup>a</sup> The label MMR-<sup>168</sup> corresponds to the aggregation of the data from the three MMR-deficient strains constructed in this study from R<sup>168</sup> (ΔL, ΔS, L\*). Data for strains R<sup>3610</sup> and ΔS<sup>3610</sup> retrieved from Sung *et al.* (2015) and mapped to R<sup>168</sup> genome sequence.

<sup>b</sup> Conversion between MA-steps and generation based on an estimated average number of generations per step: 25.6 here, 27.53 in Sung *et al.* (2015).

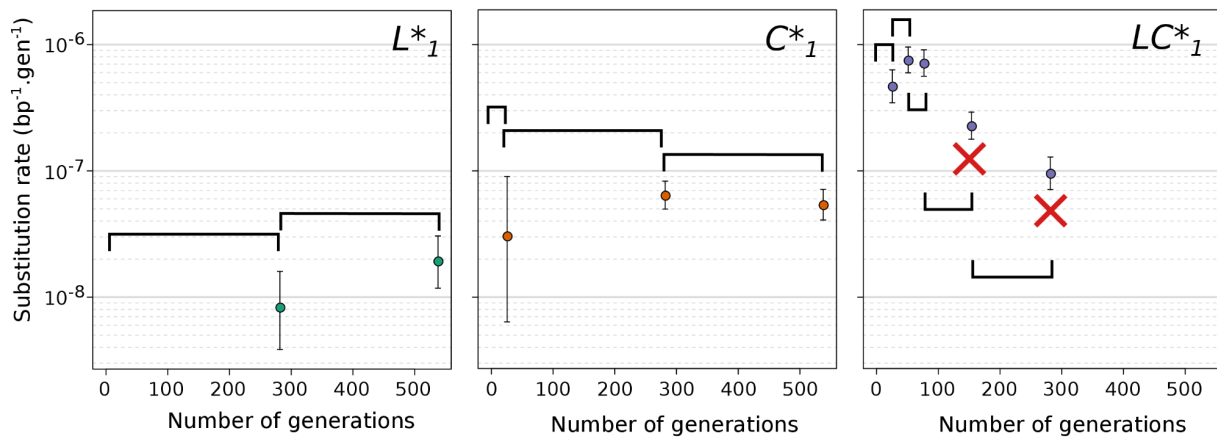
<sup>c</sup> Number of substitutions in the reference subset of positions well covered by the sequencing data (3,795 kbp). Between parentheses: number of generations or substitutions in time intervals with decreased mutation rates (discarded from the analysis).

<sup>d</sup> Number of transversions. <sup>e</sup> Number of transversions.

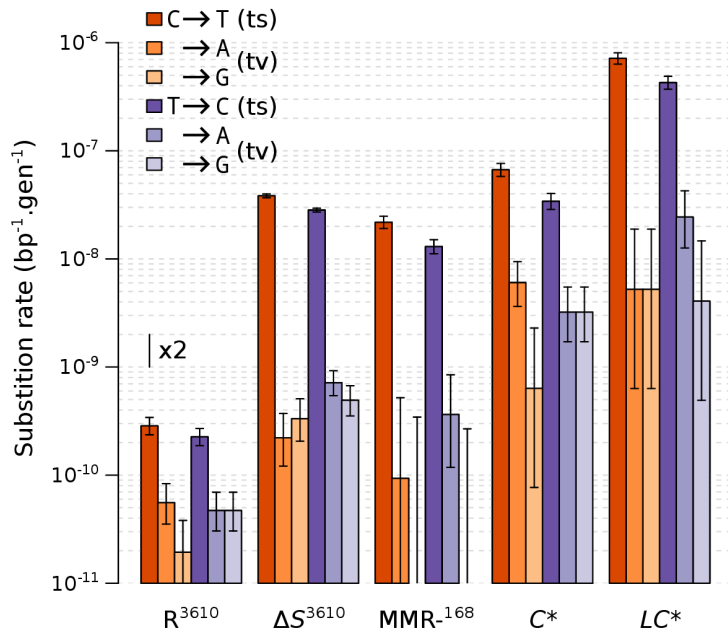
## FIGURES



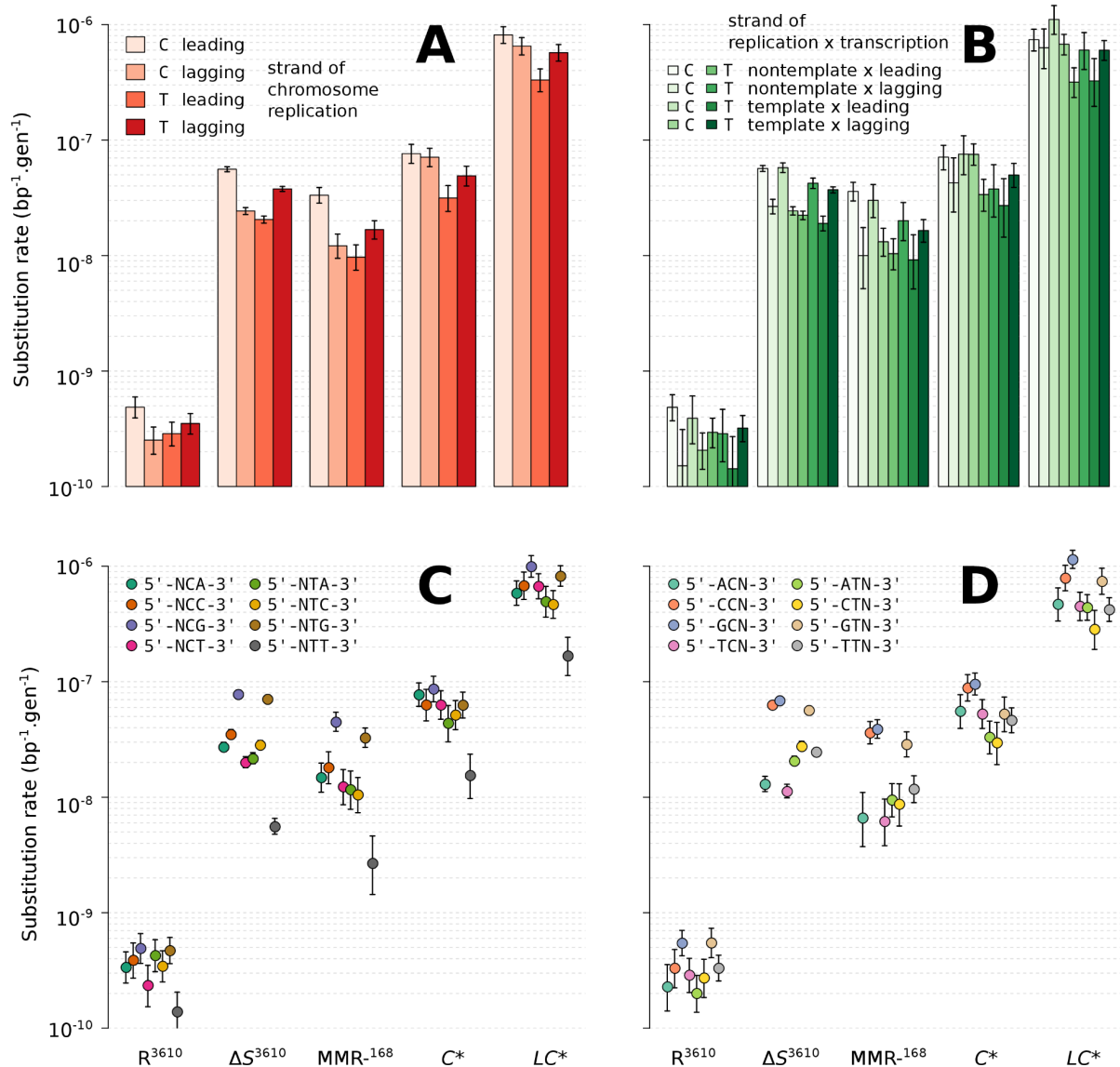
**Figure 1. Mutation rate to rifampicin resistance for increasing IPTG concentration as measured by fluctuation assays.** Each color corresponds to a strain; the vertical bars represent 95% confidence intervals.



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

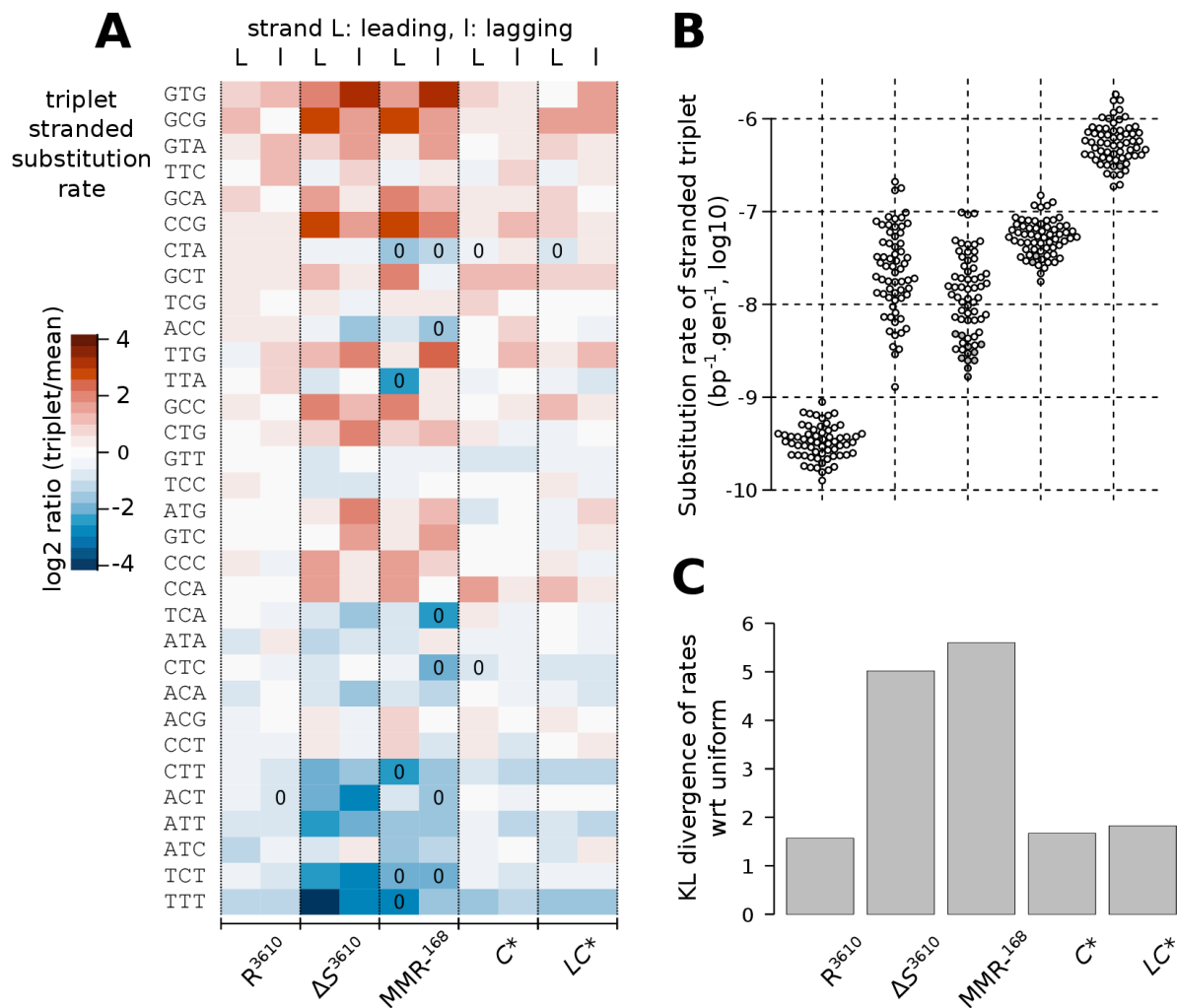


**Figure 3. Substitution rates measured by mutation accumulation experiments.** Rate per position and per generation for each type of substitution. 95% confidence intervals are indicated and were computed using the Poisson distribution. The rates of C → G and T → G mutations for MMR<sup>-168</sup> (aggregation of the ΔS, ΔL, and L\* strains) are not displayed since no mutations of these types occurred in these strains, only the upper limit of the 95% CI is shown.



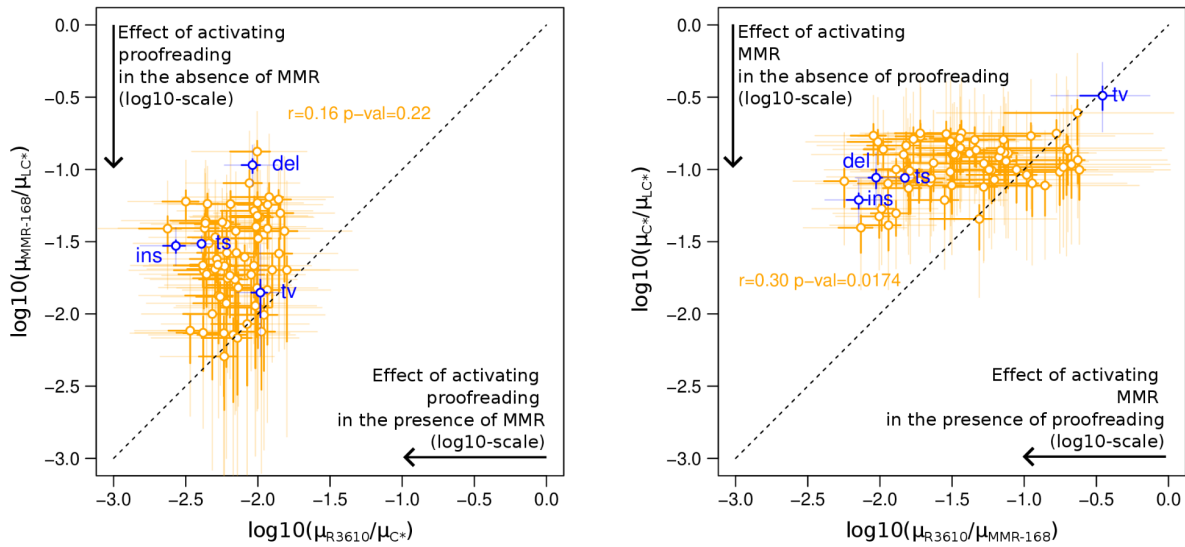
**Figure 4. Impact of replication strand and neighbor nucleotides on substitution rates.**

Substitution rates depending on the different parameters considered. Error bars represent the 95% confidence interval according to the Poisson distribution. **A.** Substitution rate depending on the mutated pyrimidine of the pair and its orientation with regard to the replication strand. **B.** Substitution rate depending on the mutated pyrimidine of the pair and its orientation both with regard to the replication and the transcription strand. **C.** Substitution rate depending on the mutated pyrimidine of the pair and the 3' adjacent nucleotide. **D.** Substitution rate depending on the mutated pyrimidine of the pair and the 5' adjacent nucleotide.

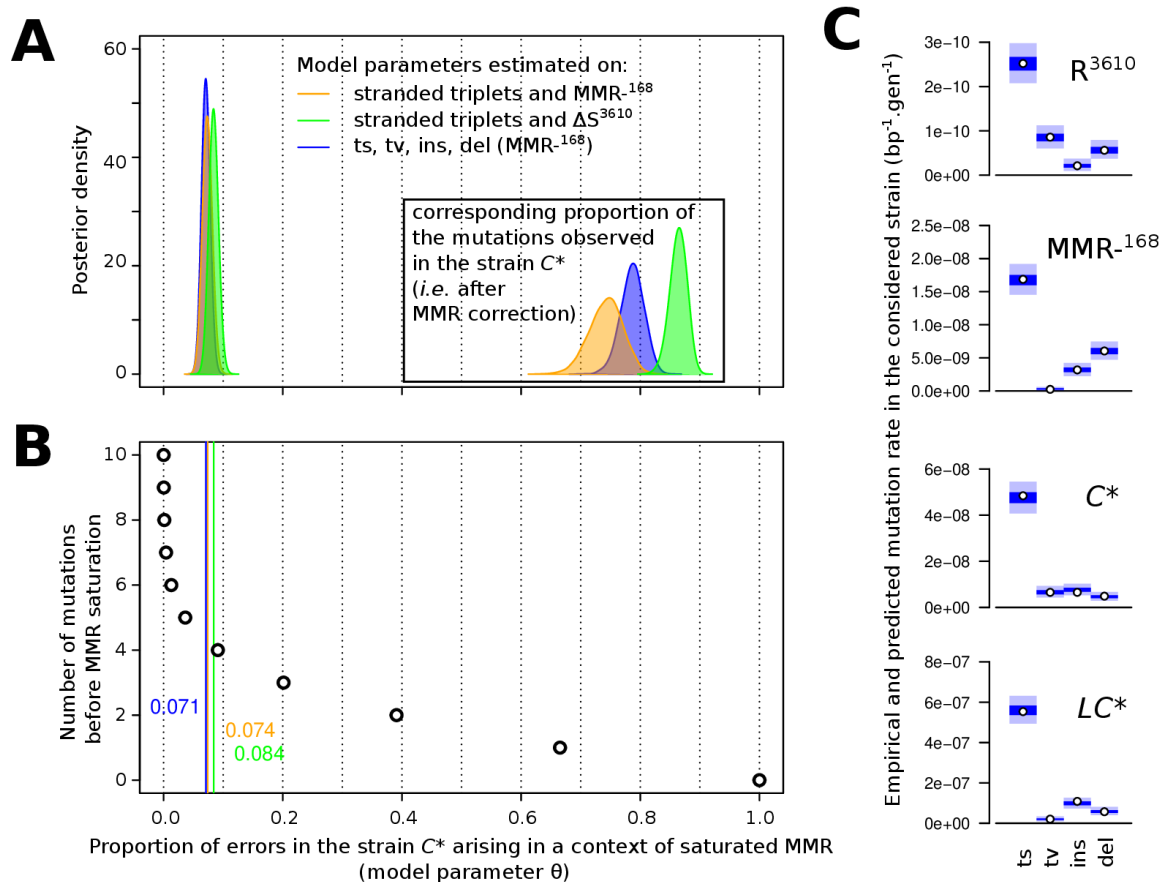


**Figure 5. Comparison of stranded-triplet substitution rates between genotypes.** Each stranded triplet corresponds to the mutated pyrimidine and its 5' and 3' nucleotides, distinguishing pyrimidines on the leading and lagging strand of chromosomal replication.  $\Delta S$ ,  $\Delta L$ , and  $L^*$  (IPTG 100  $\mu$ M) are aggregated as MMR<sup>-168</sup>. A. Heatmap representation of the stranded-triplet substitution rate (log<sub>2</sub> ratio of estimated rate wrt to the mean for each genotype). Rates were estimated with a Bayesian methodology 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 by decreasing order of non-stranded empirical substitution rates. B. Beeswarm representation of the distribution of Bayesian estimates of stranded-triplet substitution rates, the standard deviation of the estimates (in log<sub>10</sub>-scale) is reported above each beeswarm. C. KL divergence with respect to a uniform distribution (same rate for each triplet), derived from robust entropy estimates.

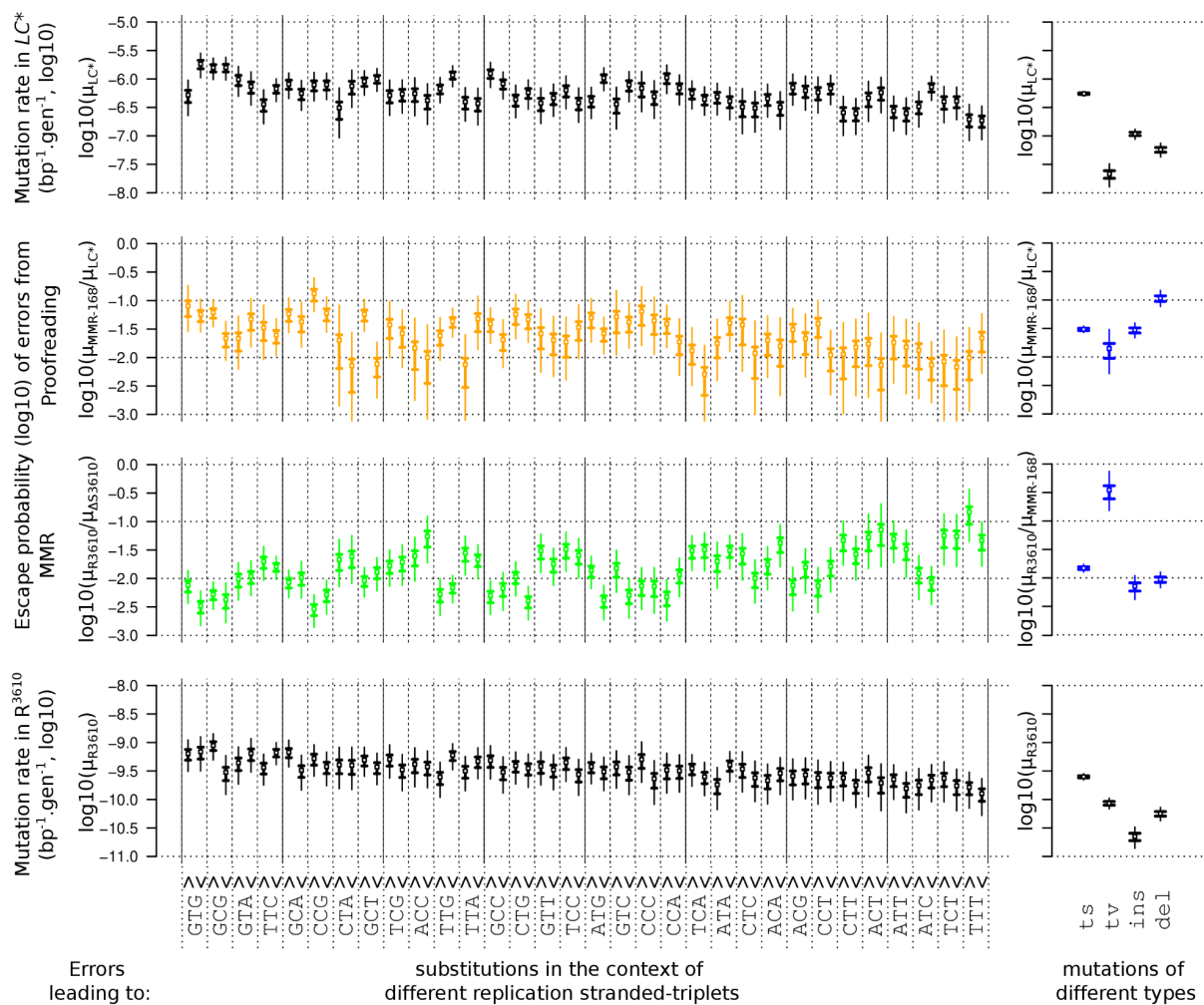




**Figure 6. The effect of activating an error correction system depends on the presence or absence of the other system.** Left plot: comparison of the effect of activating proofreading in the presence or the absence of MMR activity. Right plot: comparison of the effect of activating MMR in the presence or the absence of proofreading activity. 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 (resp. bold and dark vs. thin and light). The data used for MMR- substitution profile is MMR-<sup>168</sup>, *i.e.* the aggregation of  $\Delta S$ ,  $\Delta L$ , and  $L^*$ . The amplitudes of the effects of activating proofreading or MMR are greater in the presence than in the absence of the other system (resp. MMR or proofreading).



**Figure 7. Parameter estimation of the MMR-saturation model and assessment of the fit to experimental data.** **A.** Posterior distribution of the mixing parameter  $\theta$  corresponding to the proportion of errors in strain C\* that arise in a context of saturated MMR. The corresponding proportion of the mutations observed in C\* (i.e. after MMR correction) is shown in the insert plot. **B.** Relationship between the mixing parameter  $\theta$  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 saturation of the MMR. **C.** Assessment of the fit of the MMR-saturation model to experimentally measured rates of transition, transversion, insertion, and deletion. Points represent empirically calculated mutation rates, *i.e.* the number of observed mutations divided by the number of possible sites on the genome and the number of generations. Colored areas represent the distribution of values for empirical rates simulated under the posterior distribution of the model parameters (50% of density in dark area, 95% including also the light area).



**Figure 8. Apparent efficiency of proofreading and MMR across replication-stranded triplets and types of mutations.** Estimated proofreading and MMR escape probabilities are represented (middle plots) along with estimated mutation rates in absence and presence of both correction systems (upper and lower plots). Replication stranded triplets are ordered by decreasing order of non-stranded empirical substitution rates and then pyrimidine on the leading and lagging strands of replication. Bold and thin vertical bars 50% and 95% credibility intervals, respectively.