

Novel genes arise from genomic deletions across the bacterial tree of life

1 Arya Kaul,^{1,2,*} Fernando Rossine,¹ Karel Břinda,² and Michael Baym^{1**}

2 ¹Departments of Biomedical Informatics and Microbiology, and Laboratory of Systems Pharmacology, Harvard Medical School,
3 Boston, MA, 02115, USA

4 ²Inria, Irista, Univ. Rennes, 35042 Rennes, France

5 *Correspondence: arya_kaul@g.harvard.edu

6 **Correspondence: baym@hms.harvard.edu

7

8 **ABSTRACT**

9 Bacteria are hosts to enormous genic diversity. How new genes emerge, functionalize, and
10 spread remain longstanding questions. Here, we explore a mechanism by which adaptive
11 deletions fuse distant gene fragments. Unlike other gene birth mechanisms that begin with
12 rare, neutral mutations, these “deletion-born fusions” reach high frequency by hitch-hiking on
13 the deletion. The deletion-driven proliferation of the fusion prolongs the mutational supply
14 within these genes before loss, providing additional opportunities for neofunctionalization. We
15 document one such gene fixing and expressing in a long-term *E. coli* evolution experiment, and
16 identify additional fusion events in the *Mycobacterium tuberculosis-bovis* split. Finally, we
17 develop a scalable systematic screen to detect these genes in all 2.4 million public single-isolate
18 genomes and identify deletion-born fusions across the bacterial tree of life. These findings
19 challenge the notion that deletions are solely destructive and highlight their role as potential
20 catalysts for evolutionary innovation.

21

22 **KEYWORDS**

23 structural variation, bacterial gene birth, genomic deletions

24

25 **INTRODUCTION**

26 Bacteria are the most genetically diverse domain of life on Earth.^{1,2} This genetic diversity also
27 translates to a profound diversity in their protein-coding genic repertoire. New genes and gene
28 families continually arise across the bacterial tree of life; pangenome studies of single species
29 can identify tens of thousands of protein families, and metagenomic sequencing of
30 environments like the human gut has revealed millions of protein-coding genes of unknown
31 function.³⁻⁷ This diversity raises the question of where new bacterial genes and their
32 corresponding functions come from.

33 Different mechanisms for bacterial gene birth have been proposed and each uniquely
34 constrains the generation of functional novelty. On the one hand, duplication and subsequent
35 diversification ensures that novel genes derive from functional templates,⁸⁻¹⁰ but the
36 preexisting functionality and domain structure constrain their evolutionary potential. On the
37 other hand, overprinting,^{11,12} the expression of a novel gene overlapping a previously existing
38 sequence but in a different reading frame, generates protein products without homology to
39 existing proteins or prior functional constraints. Yet new proteins created from random
40 peptides often suffer from excessive hydrophobicity, aggregation propensity, and intrinsic
41 disorder.^{13,14} Considering these tradeoffs, gene fusion, the stitching together of entire fragments
42 of distinct genes, provides a gene birth mechanism that balances innovation with functionality

43 by bringing together existing protein domains into new contexts,¹⁵ though such fusions must
44 still avoid being lost to selection or genetic drift.

45 Random loss of newborn genes is compounded a deletional bias that threatens to purge them
46 from bacterial genomes before they can acquire beneficial function. Compact bacterial
47 genomes, often less than fifteen percent non-coding,¹⁶ reflect a pervasive bias favoring loss of
48 non-essential DNA.¹⁷⁻¹⁹ This process, seen in natural,²⁰ clinical,^{21,22} and experimental
49 settings,^{23,24} is typically framed as a loss of genetic potential. These deletions may arise
50 through diverse mechanisms, including homologous recombination,²⁵ replication-associated
51 slippage,^{26,27} erroneous repair,²⁸ or site-specific recombinases.²⁹ Regardless of mechanism,
52 deletions are a hallmark of bacterial genome evolution and can be positively selected by
53 reducing the metabolic cost of replicating DNA or by tuning gene interaction networks.^{16,19,30}

54 Deletions can also generate new arrangements of existing material. In a recent analysis of the
55 Lenski Long-Term Evolution Experiment (LTEE),³¹ uz-Zaman et al. found that deletions moving
56 regulatory elements contributed the largest share to the transcription and translation of
57 previously non-coding DNA.³² In other instances, deletions generated functional gene fusions
58 that led to new anti-phage defense functions,³³ to phenotypes with increased colony spread,³⁴
59 and to potentially adaptive gene products in *Mycobacterium tuberculosis* lineages.³⁵ Despite the
60 importance of gene fusions in generating novel phenotypes and individual cases linking
61 deletions to the emergence of fusion genes, it remains unknown whether this represents a
62 widespread mechanism for bacterial genome innovation.

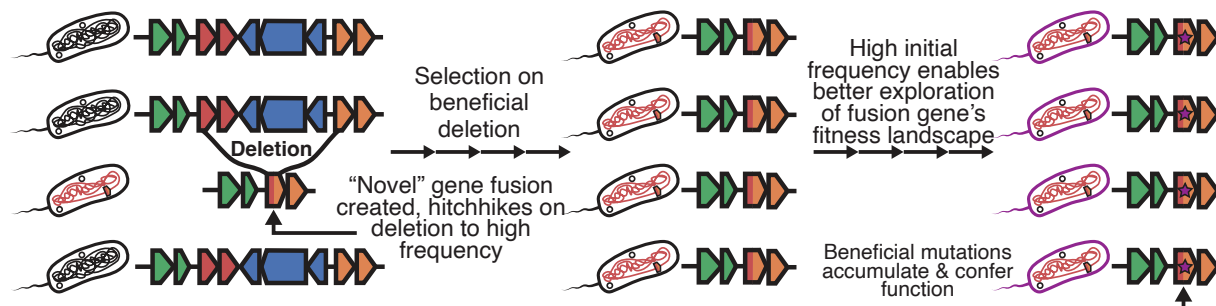
63 Here, we explore a model of bacterial gene birth in which a deletion results in the fusion of the
64 start of one open reading frame (ORF) and the end of another to create a novel ORF. Next, we
65 identify the creation and maintenance of these deletion-born fusions in two densely sampled
66 contexts: in the Lenski LTEE, and during mycobacterial speciation. Finally, we develop a
67 computational technique to efficiently query putative deletion-born fusions across multi-million
68 bacterial isolate genome collections and find evidence for this mechanism of gene birth across
69 the bacterial tree of life. Our findings reframe the bacterial deletional bias as not merely
70 destructive, but as a possible creative force.

71 **RESULTS**

72 ***Fusion genes can spread by hitchhiking on the fitness benefit of their causal deletion***

73 Deletions in a bacterial genome result in the fusing of two distal sections of genetic material
74 spanning the deletion junction. Either one or both deletion boundaries can fall inside an
75 existing gene, thus generating a chimeric ORF at the junction. This is made more likely by the
76 gene-dense architecture of bacterial genomes.³⁶ Both the novel ORF and the removal of
77 intervening material can cause changes in relative fitness.

78 If a deletion confers a fitness benefit, any fusion ORF created as a by-product can increase in
79 frequency through hitchhiking. Although most nascent genes are expected to be neutral or
80 deleterious regardless of the mechanism by which they originated,³⁷ deletion-born fusions may
81 possess an advantage over other novel genes. Unlike mechanisms that depend on rare, initially
82 neutral genes gradually drifting to appreciable frequency, deletion-born fusions can be quickly
83 driven to elevated frequencies as direct correlates of positively selected structural changes. This
84 hitchhiking may prolong their residence time in the population, increasing the likelihood that
85 subsequent mutations will convert them into beneficial, functional alleles before being lost to
86 drift or purifying selection (**Figure 1**). Moreover, because these fusions are assembled from pre-
87 existing coding material, they are more likely to fold into active catalytic, structural, or
88 regulatory domains. Overall, we expect that deletion-born fusion genes should be more likely to
89 be functionalized than novel genes that emerge through different mechanisms.



90
91
92
93
94

Figure 1. Proposed model of deletion-born fusion genes.

(1) An initial deletion occurs in a subset of the population and spontaneously creates a deletion-born fusion gene, (2) The deletion is beneficial and rises in frequency with the fusion gene hitchhiking to high frequency, (3) fusion gene at high initial frequency explores fitness landscape and functionalizes

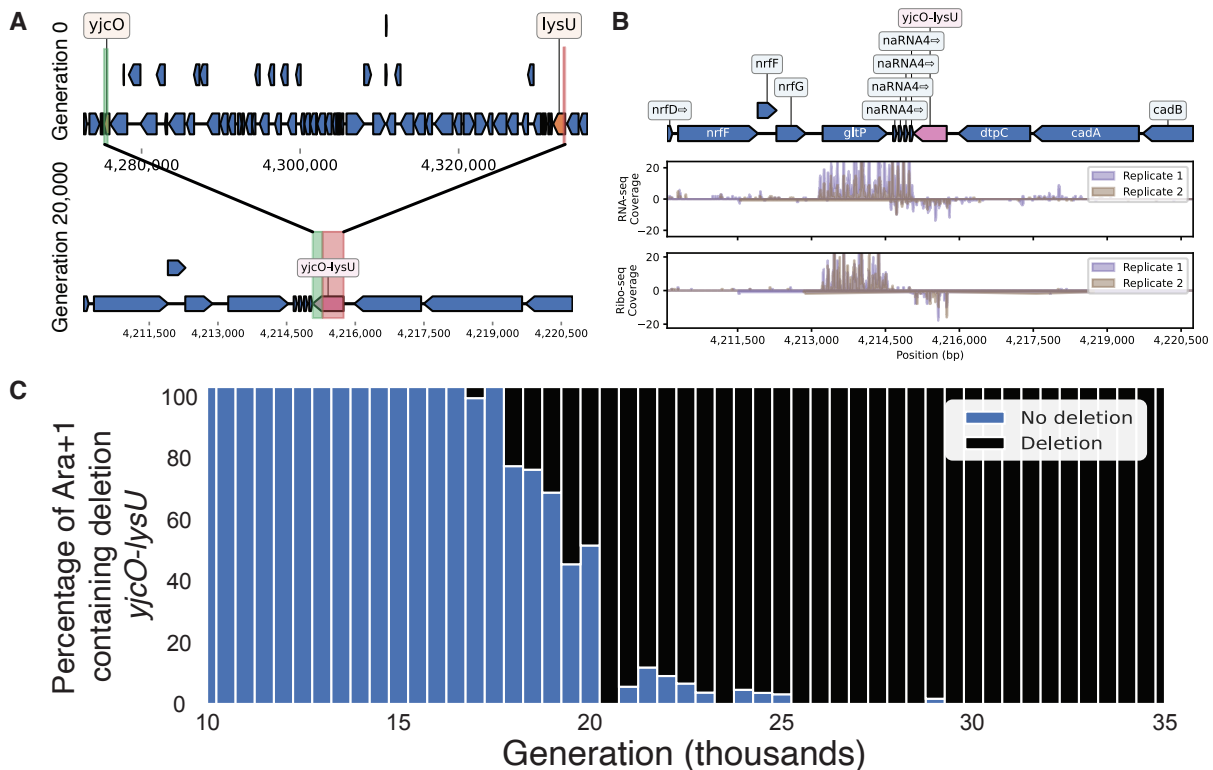
95 To formalize this verbal model, we constructed a stochastic simulation that incorporates the
96 role of genetic hitchhiking and allowed us to contrast the functionalization likelihood of novel
97 genes that arose through different mechanisms. We found that starting novel gene frequency is
98 the dominant driver of neofunctionalization, supporting the hypothesis that hitchhiking on a
99 deletion is a means by which novel ORFs can functionalize (**Supplementary Figure 1**). We
100 constructed a forward-time Wright-Fisher simulation with 10^6 members, each with one of two
101 states (no novel gene, and non-functional novel gene present with cost c). We assumed that
102 each novel gene is initially non-functional, could be purged each generation with probability
103 p_{purge} , and could functionalize with a probability μ . We also assumed that deletion-born fusion
104 genes quickly increased in frequency to some frequency p_{init} (which should be reflective of the
105 fitness advantage of the genomic deletion that gave rise to the gene) while genes born from
106 other mechanisms always had $p_{\text{init}} = 10^{-6}$ (a lone member of the population begins with the
107 gene). We demonstrated that across values for p_{purge} and μ the greatest determinant of at least
108 one member of the population functionalizing the gene before it is purged from the population
109 is p_{init} , the initial frequency of the gene. When $p_{\text{init}} = 10^{-6}$ the gene never functionalized
110 consistently; however, as p_{init} increased, the probability of functionalizing increased
111 substantially. This functionalization probability still fell below many values of μ ; especially if
112 the fusion is deleterious instead of neutral. It was only when p_{init} approached 1 that
113 neofunctionalization occurred consistently under a sweep of parameters. These results
114 suggested that young deletion-born fusion genes should be present in populations that are
115 adapting to new environments, when genomic deletions are typically the most advantageous
116 (i.e., p_{init} is the highest).

117 **A novel deletion-born fusion fixes in the Lenski LTEE**

118 To investigate the emergence of deletion-born fusions, we analyzed data from the Lenski LTEE.
119 Briefly, the LTEE tracks the evolution of 12 replicate populations of *Escherichia coli* in minimal
120 media since 1988, spanning over 82,000 generations as of 2025.^{38,39} By clustering predicted
121 ORFs from sequenced clonal isolates against the ancestral genome, we identified a novel fusion
122 gene created by a 57.5 kb deletion in the Ara+1 lineage. This deletion fused *yjcO* (a gene of
123 unknown function) to *lysU* (lysyl-tRNA synthetase) (**Figure 2A**).

124 The fusion of *yjcO* and *lysU* is expressed and potentially folds but is likely non-enzymatic and
125 functionally inert. Re-analysis of existing RNA-seq and ribosome profiling data from clonal
126 isolates in Ara+1 revealed that the *yjcO-lysU* fusion is both expressed and translated at
127 generation 50,000 (**Figure 2B**).⁴⁰ Domain annotation shows that the fusion retained multiple
128 Sel1-like repeats from *yjcO*, but no catalytic domains were preserved from *lysU*. Across all
129 clonal isolates sequenced after its appearance, the fusion gene shows no nucleotide

130 substitutions. Based on the gene length and known LTEE mutation rates, we estimate the
 131 probability of at least one mutation occurring after its appearance to be <1%, consistent with
 132 the observed absence of variation and suggesting weak or no selection on the fusion.



133
 134 **Figure 2. A 57.5 kb deletion creates a fusion gene in the Lenski LTEE Ara+1 lineage.**
 135 (A) Top: Genomic locus surrounding *yjcO-lysU* in REL606, Bottom: The same genomic locus of the clonal isolate from
 136 Generation 20,000 with the deletion and resulting *yjcO-lysU* fusion gene.
 137 (B) Genomic locus of the *yjcO-lysU* region with coverage plots corresponding to RNA-seq and Ribo-profiling data from
 138 two clonal isolates picked at generation 50,000 in Ara+1.
 139 (C) Metagenomic sequencing results comparing the fraction of reads supporting the deletion or not between Generation
 140 10,000 and generation 35,000 in the Ara+1 lineage.
 141 Both the fusion and the deletion it arose from are first detected in metagenomic sequencing
 142 data at around generation 19,500 and appear to fix within ~500 generations (Figure 2C). The
 143 speed of this selective sweep indicates that either the deletion itself or the resulting *yjcO-lysU*
 144 fusion likely conferred an approximate 6.5% fitness advantage under LTEE conditions
 145 (Supplementary Note).
 146 Supporting the hypothesis that the selective advantage lies in the deletion itself rather than the
 147 fusion gene, a parallel 43.4 kb deletion independently arose and fixed in the Ara-3 lineage at
 148 the same genomic locus, but did not produce a novel fusion (Supplementary Figure 2).
 149 Additional analyses of transposon sequencing data confirmed that disrupting the fusion gene
 150 had no deleterious effect,⁴¹ reinforcing the interpretation that the deletion, not the gene, is the
 151 beneficial variant.
 152 Genomic regions flanking deletions ≥ 1 kb exhibited significantly greater and more variable
 153 changes in expression and translation than randomly sampled regions (Supplementary Figure

154 **3**). These effects decayed with increasing window size, consistent with deletions inducing
155 localized promoter capture and the disruption of transcriptional context. This supports a model
156 in which deletions remodel local expression and translational landscapes, occasionally
157 generating novel transcribed or translated products.

158 These results indicate that, even in the relatively constant conditions of the LTEE, new fusion
159 genes can arise from large deletions and quickly increase in frequency in the population.

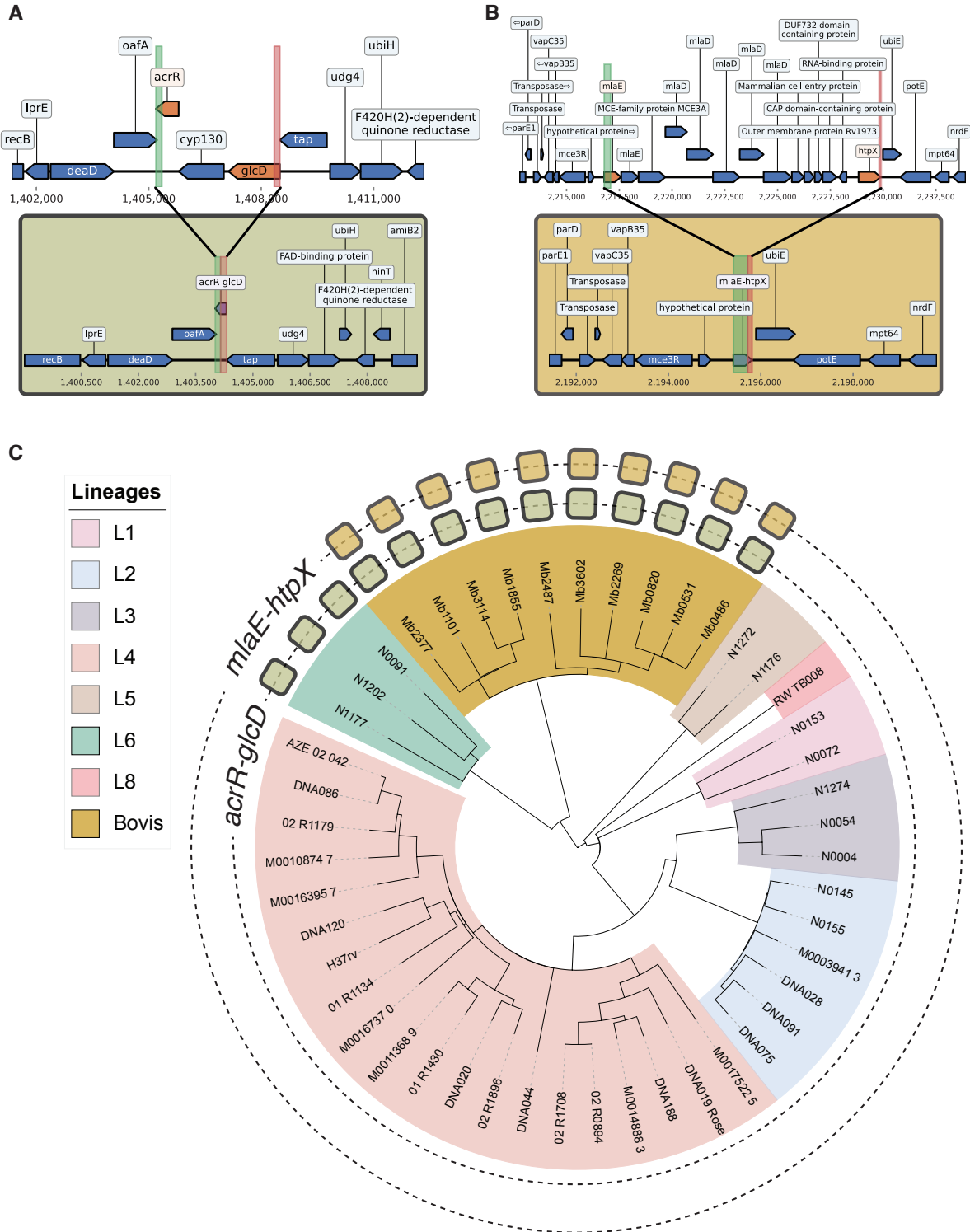
160 **Deletion-born fusions emerge during speciation in the *Mycobacterium Tuberculosis*** 161 **Complex**

162 We next turned to a well-characterized and deeply sequenced bacterial speciation event, the
163 *Mycobacterium tuberculosis-bovis* divergence, to investigate the potential for deletion-born
164 fusions to arise in complex natural systems.^{42,43}

165 We identified two putative deletion-born fusion genes: *acrR-glcD* (from a 2 kb deletion) and
166 *mlaE-htpX* (from a 12.5 kb deletion) arising during this divergence (**Figure 3**). Using a collection
167 of 47 long-read assemblies (37 *M. tb*, 10 *M. bovis*), we predicted all ORFs, clustered them into
168 orthologous groups, and searched for novel ORFs consistent with deletions (see METHODS).
169 Phylogenetic mapping showed that *mlaE-htpX* occurred exclusively in *M. bovis*, while *acrR-glcD*
170 was found in both *M. bovis* and *M. tb* Lineage 6, the closest relative of *M. bovis*.⁴⁴ These
171 monophyletic distributions suggest each fusion arose once, *mlaE-htpX* in the *M. bovis* split and
172 *acrR-glcD* in the *M. bovis*/L6 divergence.

173 Neither fusion appears to be immediately functional. Structurally, *acrR-glcD* maintains reading
174 frame continuity between the N-terminus of *glcD* and the C-terminus of *acrR*, whereas *mlaE-*
175 *htpX* fuses *mlaE* in-frame to an out-of-frame fragment of *htpX*. Pfam domain searches revealed
176 that both fusions lost the catalytic motifs of their ancestors and did not generate any new
177 conserved domains. While neither *acrR-glcD* nor *mlaE-htpX* exhibit nucleotide variation in the
178 genomes from Marin et al. and Charles et al.,^{42,43} a broader analysis (next section) revealed
179 both variation and signatures of selection.

180 The observation of deletion-born fusions in the twin contexts of laboratory evolution and
181 natural speciation in evolutionarily distant bacteria implies that this process may be
182 widespread across bacterial diversity.



183

184 **Figure 3. Two deletion-born fusions arose in the *M. tb*/*M. bovis* speciation event.**

185 (A) Top: representative pre-deletion *acrR-glcD* genomic neighborhood in *M. tb*. (genome: RW-TB008). Bottom:

186 representative post-deletion *acrR-glcD* genomic neighborhood in *M. bovis* + L6 *M. tb*. (genome: Mb1855).

187 (B) Top: pre-deletion *miaE-htpX* genomic neighborhood in *M. tb.* (genome: DNA019_Rose). Bottom: post-deletion *miaE-*
188 *htpX* genomic neighborhood in *M. bovis* (genome: Mb3602).

189 (C) Phylogenetic tree of the 47 analyzed samples. Clade colors denote the lineage each isolate falls within, and the filled
190 squares in each ring denote the genomes with *acrR-glcD* and *miaE-htpX* respectively.

191 **An alignment-free approach to characterize structural variation across a multi-million** 192 **bacterial genome collection**

193 To query structural variation at the scale of the bacterial tree of life, our previous alignment-
194 based techniques were computationally infeasible. We therefore developed an alignment-free
195 approach that queries short k-mers from the beginning (“prefix”) and the ending (“suffix”) of
196 candidate genes to infer structural rearrangements based upon the genomic distance between
197 these sequences (**Figure 4A**). Because this method relies only on the location of exact k-mer
198 matches, it can leverage FM-indices for sublinear query times,⁴⁵ making searches across
199 millions of genomes tractable. We applied this approach to AllTheBacteria (ATB), the largest
200 current bacterial genome collection, comprising 2.4 million uniformly assembled single isolate
201 genomes.

202 The distribution of observed prefix-suffix distances reflects structural variation within a given
203 gene. Nearly all insertion, deletion, or partial inversion events will change the distribution of
204 distances. For deletions, the ancestral pre-deletion sequence will have a larger prefix-suffix
205 distance than a fused post-deletion sequence. We used multi-modality in the prefix-suffix
206 distance distribution as an indicator of structural variation within a locus (see METHODS).

207 We first validated this prefix-suffix approach on the three previously identified fusion genes
208 and observed distinct peaks at the ancestral and fused distances in all cases (**Supplementary**
209 **Figure 4**). Notably, the scale of ATB revealed additional structural variation at these loci that
210 was invisible in our earlier, smaller-scale analyses. At the *yjcO-lysU* locus, we identified nine
211 distinct distance clusters spanning a range of structural variants. The 80 genomes in cluster 0
212 (prefix-suffix distance matching the fusion length) formed a monophyletic clade restricted to
213 LTEE-derived isolates. The remaining 116,892 genomes showed extended prefix-suffix
214 distances ranging from 29 to 96 kb, consistent with independent insertions and deletions at
215 this locus. Cluster 3, comprising 48,892 genomes with a prefix-suffix distance centered at 57
216 kb, matches the ancestral LTEE length; the remaining 68,000 genomes exhibit distinct
217 structural configurations, revealing additional diversity at this locus (**Supplementary Figure**
218 **4A**).

219 In *Mycobacteria*, the deeper sampling from ATB uncovered 4 distinct *acrR-glcD* alleles and 17
220 *miaE-htpX* alleles. Alignment-based Bayesian selection analysis showed strong evidence for
221 diversifying selection at codon 50 in *acrR-glcD* (posterior >0.9) and purifying selection at five
222 positions (87, 91, 92, 108, 125) in *miaE-htpX* (see METHODS). These patterns indicate that,
223 even absent recognizable catalytic domains, these fusions are experiencing selective pressure,
224 consistent with emerging or maintained function.

225 While we developed the prefix-suffix method to identify deletion-born fusions, the approach
226 detects any recurrent structural variation at a locus. In the next section we filter candidates to
227 deletion-born fusions, but many of the events we filter out are themselves biologically
228 interesting. We note the identification of internal deletions in the *mngB* gene in *E. coli*
229 (**Supplementary Figure 5A**), repeat prophage insertions into *rep13e12* gene in *M. tuberculosis*
230 (**Supplementary Figure 5B**), and variable gene cargo in an uncharacterized mobile element
231 disrupting the *pdp* gene in *S. pneumoniae* (**Supplementary Figure 5C**). Thus, the prefix-suffix
232 signal is generalizable to surveying structural variation beyond our specific application and,
233 because of its computational efficiency, scales readily to ever-expanding genome collections.

234 **Putative deletion-born fusions are found across the bacterial tree of life**

235 We next asked how pervasive deletion-born fusions are across diverse bacterial phyla. We
236 applied the prefix-suffix k-mer screen to all annotated protein-coding sequences from five type
237 strains spanning well-studied and diverse bacterial clades (*Escherichia coli* K12, *Mycobacterium*
238 *tuberculosis* H37Rv, *Neisseria gonorrhoeae* FA1090, *Campylobacter jejuni* NCTC1168, and
239 *Streptococcus pneumoniae* TIGR4). We selected these species since they span phylogenetically
240 distinct clades, and because their clinical importance has made them among the most deeply
241 sequenced bacteria, with each represented by at least 60,000 genomes in the ATB.

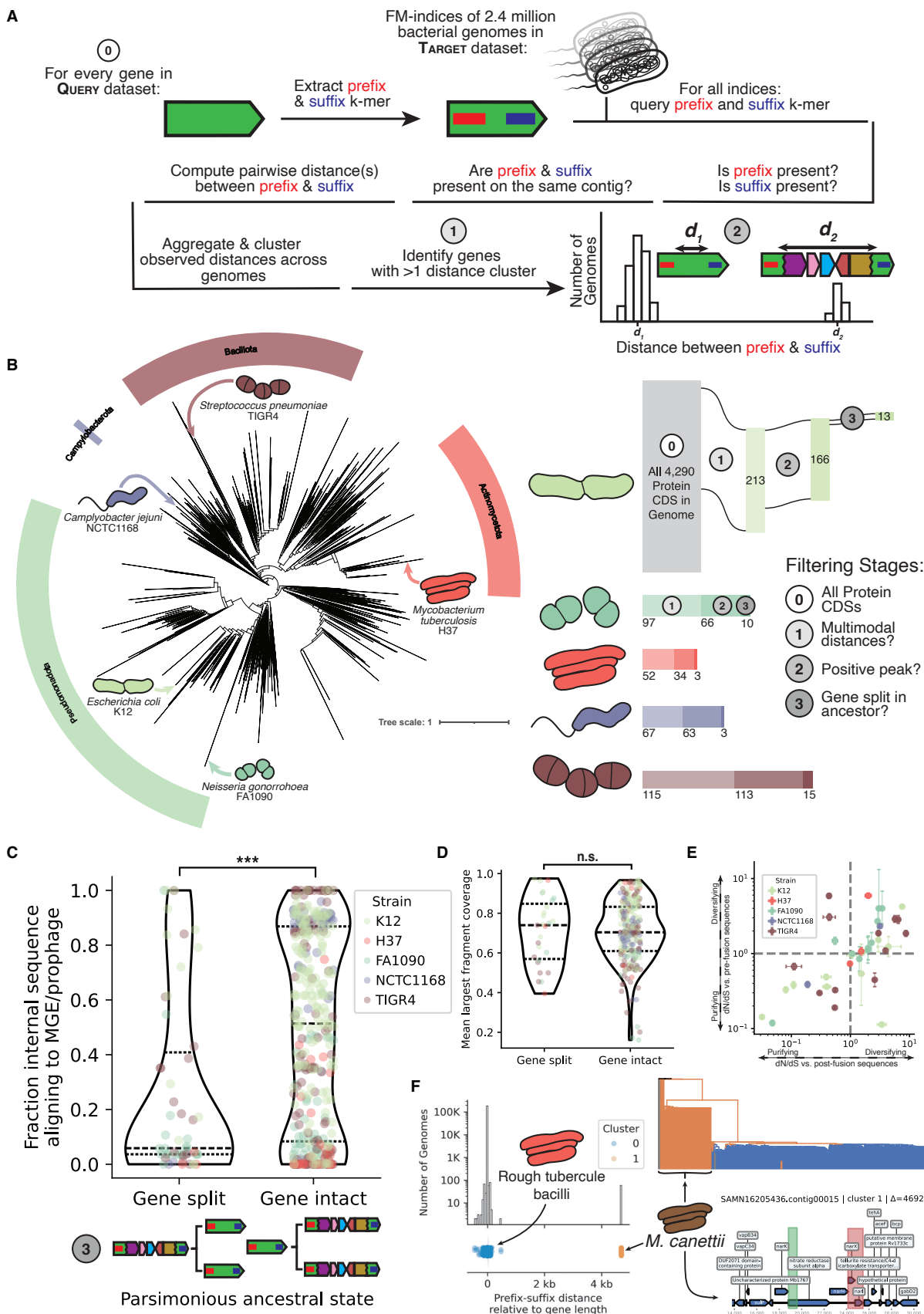
242 Across the strains, our pipeline resolved 44 putative deletion-born fusion candidates inferred
243 as “split” at the MRCA of the genomes sampled: 13 in *E. coli* K12, 3 in *M. tuberculosis* H37Rv,
244 10 in *N. gonorrhoeae* FA1090, 3 in *C. jejuni* NCTC1168, and 15 in *S. pneumoniae* TIGR4 (**Figure**
245 **4B, Supplementary Table 1**).

246 We identified these putative deletion-born fusions via a sequence of consecutive filters: We first
247 selected genes whose prefix-suffix distances were multimodal, indicative of structural variation.
248 We then retained candidates that showed one cluster centered near the length of the gene (a
249 difference of 0 implying an intact gene) and at least one more cluster at a larger, positive
250 distance (putative pre-deletion sequence) (see METHODS). By solely relying on the distance
251 between the prefix and suffix k-mers, we are unable to distinguish between insertion into a
252 gene and a deletion that led to the formation of that gene. To distinguish between these
253 possibilities, for each candidate gene, we extracted complete genomes with equal
254 representation from each cluster, added outgroups, built a k-mer-based distance tree, and
255 performed parsimony-based ancestral state reconstruction using cluster labels as discrete
256 states (see METHODS). Candidates were retained when the most parsimonious cluster
257 assignment for the most recent common ancestor (MRCA) of all sampled genomes was a cluster
258 with a positive distance length (split gene) implying that the ancestral state was an unfused
259 gene that later experienced a deletion and gene formation.

260 In genes whose parsimonious ancestral state is predicted to be intact, the observed positive
261 prefix-suffix distance peaks are significantly explained by the insertion of foreign elements,
262 particularly mobile genetic elements (MGEs) and prophages (**Figure 4C**). By contrast, at loci
263 whose MRCA is inferred not to have carried the intact gene, the intervening segments have
264 significantly fewer matches to MGEs or prophages, consistent with separation/fusion via
265 deletion or recombination rather than insertion.

266 To ensure the distance signals identified were not derived from spurious k-mer matches, we
267 sequence-aligned the putative fusion gene to genomes where it was predicted to be split. We
268 found that all putative deletion-born fusions have at least 80% gap-excluded identity to their
269 split ancestors, implying the prefix-suffix approach detected true sequence homology. Further,
270 the distribution of the relative contributions of the largest alignment fragment in putative
271 deletion-born fusions (MRCA = Gene split) matched that of disrupted genes (MRCA = Gene
272 intact, which we expect to be fully random), implying that spurious k-mer matches to either the
273 prefix or suffix alone are undetectably rare (**Figure 4D**).

274 The lack of spurious k-mer matches is likely attributed to the significant requirements we
275 enforce: at least 54 base pairs of exact nucleotides matching on the same contig separated by
276 roughly the same amount across numerous genomes. To test how robust this approach was to
277 varying lengths of k we also tested numerous values for three sets of 1,000 random bacterial
278 RefSeq proteins. We found that below k = 20, the number of genes with multimodal
279 distributions rises sharply (**Supplementary Figure 6**). Though some of these signals might be
280 real, we elected to continue with the more stringent value of k = 27 to ensure high confidence
281 matches.



283 **Figure 4. Deletion-born fusions are present across the bacterial tree of life.**

284 (A) Schematic of the prefix-suffix k-mer approach.

285 (B) Left: bacterial tree of life condensed to the genus level showing the five type strains queried. Right: filtering cascade
286 with counts of genes passing each stage: (0) all protein-coding ORFs, (1) multimodal prefix-suffix distances, (2) positive
287 distance peak present, (3) gene inferred as split in ancestor. Bacterial illustrations were traced from SEM images. *E.*
288 *coli*,⁴⁶ *N. gonorrhoeae*,⁴⁷ *M. tuberculosis*,⁴⁸ *C. jejuni*,⁴⁹ *S. pneumoniae*.⁵⁰

289 (C) Fraction of intervening sequence aligning to MGEs/prophages in genes with split versus intact ancestral states
290 (Mann-Whitney U test, $p < 10^{-5}$).

291 (D) Mean coverage of the largest alignment block between prefix and suffix k-mers; no significant difference between
292 distributions.

293 (E) dN/dS estimates for 44 candidate fusions. X-axis: dN/dS versus intact alleles; Y-axis: dN/dS versus reconstructed
294 pre-deletion sequences. Error bars show 95% confidence intervals.

295 (F) *narX* from *M. tuberculosis* H37Rv. Histogram and stripplot show prefix-suffix distances; Cluster 0 = rough tubercle
296 bacilli (including the fusion), Cluster 1 = *M. canettii* (putative ancestral split state). Phylogeny colored by cluster
297 assignment; bottom panel shows the 10 kb genomic neighborhood of the ancestral *narX* locus in *M. canettii*.

298 Signatures of positive selection are most common across the 44 putative deletion-born fusions
299 (**Figure 4E**). To distinguish between selection currently operating on the proteins and strong
300 selection immediately after the initial deletion event, we categorized synonymous and
301 nonsynonymous variation in two ways. First, we sampled genomes with the fusion gene and
302 measured their variation against the reference fusion gene (dN/dS vs. post-fusion sequences).
303 Second, we measured variation against “surrogate” genes created by merging the pre-fusion
304 fragments (dN/dS vs. pre-fusion sequences, see METHODS). We expect neofunctionalization to
305 be marked by initial diversifying selection followed by purifying selection; however, most
306 fusions identified exhibit signatures of positive selection under both regimes. We note that
307 uncorrected dN/dS is fairly coarse⁵¹ and the stipulation of an exact 54 base pair match throws
308 out any variation which might be occurring in the very beginning or end of given genes driving
309 down variation which might be occurring. Even with these caveats, the consistent identification
310 of signatures of positive selection seem to indicate many of these deletion-born fusions are still
311 in the process of exploring sequence space for function.

312 To illustrate one example in more depth, we turn to *narX*, a putative nitrate reductase gene in
313 MTBC genomes. In 182,154 Mycobacterial genomes, *narX* is 1,915 bp long, but in 59 other
314 genomes we identify the prefix and suffix 27-mers separated by 6.6 kb, a 4.6 kb increase from
315 the intact gene (**Figure 4F**). A rooted phylogeny built from sampled genomes shows that they
316 form a monophyletic clade, all belonging to *Mycobacterium canettii*, a smooth tubercle bacillus
317 regarded as one of the closest environmental relatives of transmissible MTBC members.⁵² When
318 we align the *M. tb.* H37Rv *narX* coding sequence to a representative *M. canettii* genome, we find
319 that the “split” *narX* corresponds to a composite of three genes in a nitrate reduction operon
320 and retains three recognizable domains (PF00384: molybdopterin oxidoreductase; PF02613
321 and PF02665: nitrate reductase delta/gamma subunits). This organization is consistent with a
322 scenario in which a deletion-driven fusion assembled *narX* from a previously distinct nitrate
323 reduction operon. Latent *M. tb.* infections are characterized by long-term survival in hypoxic
324 granulomas, where nitrate respiration supports energy generation in the absence of oxygen.⁵³⁻
325 ⁵⁵ In *M. tb.*, in-vitro experiments demonstrate nitrate reduction is dominated by *narGHJI*, with
326 a minimal role for *narX*.⁵⁴ Although no function has been identified for this fusion, the *narX*
327 example highlights how deletions can rewire core modules and may provide the raw material
328 for the development of novel coding sequences.

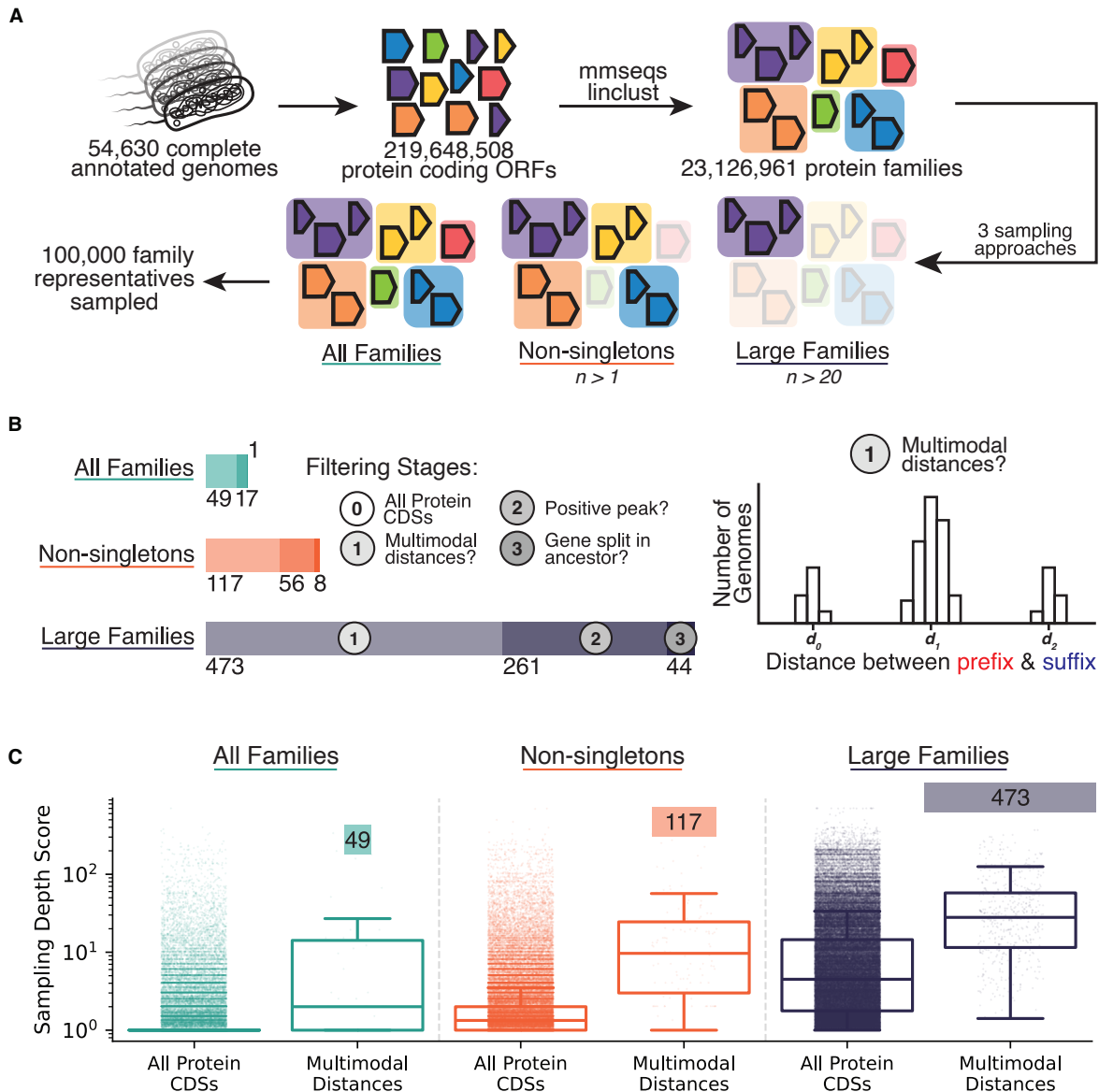
329 The identification of multiple putative deletion-born fusions across the bacterial tree of life
330 indicates the proposed mechanism is shared across bacteria; however, the type strains used
331 were specifically chosen because these species (*E. coli*, *M. tuberculosis*, *N. gonorrhoeae*, *C.*

332 *jejuni*, *S. pneumoniae*) are some of the most well-studied and consistently sequenced bacteria
333 on the planet. As a result, the genetic variation belonging to these clades is deeply sampled; a
334 reality sadly not shared by most bacteria.⁵⁶ To further clarify the role that sampling bias plays
335 in the ability of the prefix-suffix approach to capture deletion-born fusions across diverse
336 bacterial proteins we turned from a genome-first to a gene-first approach.

337 ***Detection of both deletion-born fusions and broader structural variation is dependent on***
338 ***sampling depth.***

339 We next applied the prefix-suffix approach to subsampled collections of 100,000 protein coding
340 ORFs extracted from all 54,630 “complete” RefSeq genomes available at the time of analysis
341 (**Figure 5A**). These genomes contained a total of 219,648,508 protein coding ORFs, these ORFs
342 were clustered into 23,126,961 protein families (see METHODS). Concordant with prior results,
343 we found that most (60.78%) of the protein families have only one sequence within them, so-
344 called “singletons” (**Supplementary Figure 7A**).⁵⁷⁻⁵⁹ This observation can partially be explained
345 by the bias of which bacterial genomes are selected for sequencing, assembled with high-
346 confidence, and then annotated. Indeed, we find that the top 20 represented species make up
347 about 35% of the 16,774 species sampled across the dataset (**Supplementary Figure 7B**).

348 We sampled 100,000 protein families and selected representatives from each in three distinct
349 ways. In the first, all protein families were given an equal chance of being drawn (“All
350 Families”). In the second, only those families with at least one other sampled member were
351 sampled (“Non-singletons”). In the third, only families with more than 20 sequences in their
352 family were considered (“Large Families”).



353

354 **Figure 5. Identification of deletion-born fusions scales with genomic sampling depth.**

355 (A) Sampling strategy schematic. Protein-coding ORFs from 54,630 complete RefSeq genomes were clustered into
 356 families, and 100,000 families were sampled under three strategies: "All Families" (uniform), "Non-singletons" (≥ 2
 357 members), and "Large Families" (>20 members).

358 (B) Filtering cascade showing genes passing each stage. "Large Families" yields 44 deletion-born fusions versus 8
 359 ("Non-singletons") and 1 ("All Families"). Cartoon schematic of multimodal distance filtering stage is depicted on the
 360 right.

361 (C) Sampling depth score (total genomes / total unique species) for all sampled proteins and multimodal proteins in
 362 each sampling strategy.

363 Applying the prefix-suffix approach to these 300,000 proteins, we identify increasing deletion-
 364 born fusions when moving to more represented protein families (see METHODS). Only 1 gene
 365 was identified in "All Families", 8 genes in "Non-singletons" and 44 genes in "Large Families"
 366 (Figure 5B, Supplementary Table 2). Importantly, these same trends are also observed for the

367 intermediate filtering stages: “Large Families” has not only the most putative deletion-born
368 fusions, but also the most proteins where some structural rearrangement was observed (473
369 multimodal distances vs. 49 in “All Families”; 217 putative gene insertions vs. 16 in “All
370 Families”).

371 We sought to examine the extent to which this disparity reflects sampling depth rather than
372 biological distribution by defining a sampling depth score as the total number of genomes
373 containing members of that family divided by the number of unique species represented. A
374 score of 1 indicates a protein family sampled either once or broadly across many species,
375 whereas higher scores indicate repeated sampling within the same species and thus deeper
376 lineage-specific coverage.

377 Larger protein families were enriched for more deeply sampled species. Protein families in the
378 “Large Families” category exhibited a mean sampling depth score of 15.64, indicating that a
379 typical protein in this set is represented, on average, by 15–16 genomes from the same species
380 (**Figure 5C**). This depth was significantly greater than that of the “Non-singletons” (mean =
381 3.31) and “All Families” (mean = 1.95) categories (Mann–Whitney U test, $p < 10^{-30}$ for all
382 comparisons), explaining the greater power to detect structural variation in more deeply
383 sampled protein families.

384 Proteins exhibiting multimodal prefix-suffix distance distributions showed significantly higher
385 sampling depth than the full set of proteins drawn with that strategy (“All Families” mean =
386 24.15, “Non-singletons” mean = 28.64, “Large Families” mean = 50.53; Mann–Whitney U test, p
387 $< 10^{-12}$ for all comparisons). This pattern indicates that, even controlling for sampling strategy,
388 the detection of structural variation is biased toward protein families with deeper within-
389 species sampling.

390 Together, these results demonstrate that the apparent enrichment of deletion-born fusions in
391 certain datasets is driven by sampling depth rather than biological prevalence. Detecting these
392 events requires observing both pre-deletion and post-deletion states within the same lineage,
393 which in turn demands not only broad phylogenetic sampling but also dense sampling within
394 species. As genomic databases continue to expand and sampling becomes deeper across the
395 bacterial tree of life, we anticipate that deletion-born fusions will be revealed as more
396 widespread than currently detectable.

397 **DISCUSSION**

398 In this work, we describe a distinct route to gene birth in which adaptive deletions generate
399 fusion ORFs as by-products. These “deletion-born fusion genes” inherit their initial frequency
400 from a beneficial structural change rather than drifting from rarity and are assembled from
401 pre-existing coding material rather than arising de novo.

402 This mechanism complements existing models while occupying a distinct niche. Duplication
403 and diversification require maintaining redundant copies in genomes under strong
404 streamlining pressure, overprinting must preserve the ancestral reading frame; and horizontal
405 gene transfer introduces novelty but defers origin to an external donor. Classical gene fusions
406 generally persist only when the fusion is beneficial.⁶⁰ By contrast, deletion-born fusions arise
407 as incidental consequences of selection acting at the level of genome architecture rather than
408 protein function. They require no external material, no long-term maintenance of redundancy,
409 and no immediate benefit of the fused product. Since these fusions are assembled from existing
410 sequences, they are more likely to produce biophysically viable proteins than sequences
411 emerging de novo from non-coding DNA, as proposed for novel genes arising in eukaryotes.^{61,62}
412 Our simulations formalize the hitchhiking advantage showing that elevated starting frequency
413 is the dominant determinant of whether functionalization occurs before loss.

414 We document this process across multiple evolutionary timescales. In the Lenski LTEE, we
415 observe a large deletion that rapidly sweeps to fixation, generating a novel fusion gene in the
416 process; a convergent deletion at the same locus suggests selection acted on the deletion rather
417 than the gene. At a longer timescale, we identify deletion-born fusions arising during the
418 *Mycobacterium tuberculosis*-*M. bovis* divergence, demonstrating that such fusions can persist
419 across speciation events. Finally, by screening millions of bacterial genomes, we identify
420 putative deletion-born fusions across diverse bacterial clades, indicating that this mechanism
421 reflects a general opportunity for bacterial genome innovation.

422 We have not, however, identified any clear novel beneficial function to an identified deletion-
423 born fusion. The *ujcO-lysU* fusion in the LTEE shows no nucleotide variation and the *acrR-glcD*
424 and *mIaE-htpX* fusions in the MTBC lack catalytic domains and show minimal variation. Most
425 candidates from our tree-of-life screen exhibit dN/dS ratios consistent with diversifying rather
426 than strong purifying selection, suggesting they remain in the early stages of sequence
427 exploration. The fusions identified here likely represent snapshots of this process, in which
428 proteins persist long enough to sample sequence space but have not yet undergone strong
429 functional refinement. Furthermore, while the genome streamlining literature provides strong
430 indirect evidence that deletions are often beneficial, we have not proven any specific deletion
431 was advantageous and its accompanying fusion was neutral or deleterious.

432 The absence of clearly functional fusions may reflect biology, but it may also reflect the
433 conservatism of our approach. Our survey provides a lower bound on the prevalence of
434 deletion-born fusions. A single nucleotide mutation in either the prefix or suffix k-mer is
435 sufficient to exclude a genome from the search, so fusions that have accumulated terminal
436 mutations are systematically missed. The problem is compounded by the biology of the events
437 we are seeking deletions often arise through recombination between repetitive sequences, yet
438 repeats are precisely where short-read assemblies tend to break, placing true deletion
439 junctions on contig edges and preventing us from measuring the distance between them.
440 Future studies employing more sensitive indices could relax the exact-match requirement,
441 potentially revealing functionalized fusions that have since diverged at their termini.

442 Sampling bias further constrains detection. We show that detection scales with sampling
443 depth: the “Large Families” category, drawn from deeply sampled protein families, yielded 44
444 deletion-born fusion candidates compared to just 1 in the uniformly sampled “All Families”
445 category. This disparity likely reflects the requirement to observe both pre- and post-deletion
446 states within the same dataset, rather than true biological enrichment in well-studied species.
447 Current genomic databases are dominated by culturable, human-associated pathogens;
448 environmental lineages remain sparsely sampled. This creates a double bind: we both lack the
449 breadth to detect fusions deep in the bacterial tree, and the depth to detect recent fusions in
450 poorly sampled lineages. We expect organisms whose lifestyles predispose them to deletions
451 (i.e. intracellular pathogens) to be enriched for deletion-born fusions, yet many remain too
452 sparsely sequenced to test this hypothesis.

453 This is not merely a limitation of our approach. Any broad computational analysis of existing
454 databases will inherit these gaps, a fact particularly pertinent given the recent proliferation of
455 protein language models. Our protein family analysis demonstrated that most bacterial
456 proteins in high-quality genome collections remain singletons, solely sampled once. We
457 attribute this to the same sampling bias: we simply have not conducted deep, high-quality
458 sampling across bacterial diversity. Targeted sequencing of underrepresented taxa will be
459 necessary to close these gaps, underscoring that large-scale computational analyses are only
460 as powerful as the underlying biological data relied upon.

461 The prefix-suffix k-mer approach developed here has utility beyond deletion-born fusions. It
462 enables detection of recurrent structural variation at specific loci across multi-million genome
463 collections neither relying on alignments nor prior knowledge of ancestral states. In this study
464 alone, we characterized internal deletions, repeat prophage insertions, and variable gene cargo
465 in an uncharacterized mobile element. More broadly, the method should be applicable to any
466 process that alters the genomic distance between conserved flanking sequences: novel bacterial
467 or phage introns,^{63,64} integron cassette expansion,⁶⁵ variable plasmid cargo,⁶⁶ or structural
468 variation in complex metagenomic communities.⁶⁷ We have not yet characterized which protein
469 domains or functional categories are enriched among structurally variable loci; such meta-
470 analyses are a natural extension. As sequence databases expand in scale and complexity,
471 alignment-free approaches may become the only tractable means of extracting biological signal
472 from the noisy chorus of bacterial diversity.

473 Our characterization of deletion-born fusions carries implications beyond evolutionary biology.
474 In clinical microbiology, large deletions are frequently observed during adaptation to host
475 environments;^{21,22} our results suggest that some of these events may generate novel proteins
476 with unpredictable properties. In synthetic biology, rational protein design often proceeds by
477 domain shuffling fusing functional modules from different proteins to create chimeras with new
478 activities.⁶⁸ Understanding how nature assembles and filters such chimeras may inform future
479 efforts to engineer functional novelty.

480 Our findings recast the deletional bias that pervades bacterial genome evolution as a possible
481 wellspring from which bacterial innovation may arise. By fusing distant sequences into new
482 combinations, deletions may contribute to the vast and still largely uncharacterized diversity of
483 bacterial proteins.

484

485 **RESOURCE AVAILABILITY**

486 **Lead contact**

487 Requests for further information and resources should be directed to and will be fulfilled by the
488 lead contact, Arya Kaul (arya_kaul@g.harvard.edu).

489 **Materials availability**

490 This study did not generate new unique reagents.

491 **Data and code availability**

492 Section 1: Data

- 493
- 494 • This paper analyzes existing, publicly available data.
 - 494 • Lenski LTEE
 - 495 ○ Metagenomic data downloaded from PRJNA380528
 - 496 ○ Clonal sequencing data downloaded from
 - 497 <https://github.com/barricklab/LTEE-Ecoli>
 - 498 • *Mycobacterium tuberculosis/bovis*
 - 499 ○ Genomes for *M. tb.* were downloaded from PRJNA719670, PRJNA480888,
 - 500 PRJNA436997 and PRJNA421446.
 - 501 ○ Genomes for *M. bovis* were downloaded from PRJNA832544.
 - 502 • Prefix-Suffix Screen
 - 503 ○ AllTheBacteria was downloaded from the OSF repository
 - 504 <https://osf.io/ZXFMY/overview>

505 Section 2: Code

- 506 • Code for analyses besides the prefix-suffix k-mer screen is available at:
507 <https://github.com/baymlab/deletion-born-fusion-manuscript>
- 508 • Code for the prefix-suffix k-mer screen is available at:
509 <https://github.com/aryakaul/prefixsuffix-kmer>

510 Section 3: Additional Information

- 511 • Any additional information required to reanalyze the data reported in this paper is
512 available from the lead contact upon request.

513

514 **ACKNOWLEDGMENTS**

515 The authors would like to warmly thank all past and present members of the Baym and
516 GenScale team for both illuminating conversations and additional scientific insight. This work
517 was supported by NIGMS of the National Institutes of Health (R35GM133700 and
518 R35GM156320), the David and Lucile Packard Foundation, the Pew Charitable Trusts, and the
519 Alfred P. Sloan Foundation. The prefix-suffix approach is based upon research performed in
520 France within the GenScale team at the Inria Center at Rennes University. The work was
521 supported by a Chateaubriand Fellowship of the Office for Science & Technology of the
522 Embassy of France in the United States and a mobility grant from the Collège doctoral de
523 Bretagne. This research was supported by the French National Research Agency (ANR) under
524 Grant ANR-24-CE45-1226 for the REALL project (KB).

525 **AUTHOR CONTRIBUTIONS**

526 Conceptualization, A.K., K.B., and M.B.; methodology, A.K., K.B., and M.B.; investigation, A.K.,
527 F.R., K.B., and M.B.; writing—original draft, A.K. writing—review & editing, A.K., F.R., K.B.,
528 and M.B.; funding acquisition, A.K., K.B., and M.B.; resources, K.B. and M.B.; supervision,
529 F.R., K.B., and M.B.

530 **DECLARATION OF INTERESTS**

531 None.

532 **DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES**

533 During the preparation of this work, the author(s) used ChatGPT and Claude to assist with
534 writing code for analysis. After using this tool or service, the author(s) reviewed and edited the
535 content as needed and take(s) full responsibility for the content of the publication.

536

539 **SUPPLEMENTARY NOTE(S)**

540 ***Estimation of *yjcO-lysU*/deletion selection coefficient***

541 For a haploid population in which a novel beneficial allele has fitness $1 + s$, and the wildtype
542 has a fitness of 1, the allele-frequency dynamics are given by:

543
$$p_{t+1} = \frac{p_t(1 + s)}{1 + sp_t}$$

544 Where p_t is the frequency of the novel allele at generation t . This can be approximated to the
545 differential equation, a classical result from Kimura:⁶⁹

546
$$\frac{dp}{dt} = sp(1 - p)$$

547 Integrating over this equation and solving for t gives:⁶⁹

548
$$t \approx \frac{1}{s} \ln \left(\frac{p_1(1 - p_0)}{p_0(1 - p_1)} \right)$$

549 Conditioning on the allele fixing, we can set $p_0 = \frac{1}{N_e}$ and $p_1 = 1 - \frac{1}{N_e}$. For large effective
550 populations, this yields the commonly used approximation:

551
$$t \approx \frac{2}{s} \ln(N_e)$$

552 Solving for s gives:

553
$$s \approx \frac{2 \ln(N_e)}{t}$$

554 From Good et al.,³⁹ we estimate $N_e = 10^7$, from our metagenomic analysis, we estimate $t = 500$:

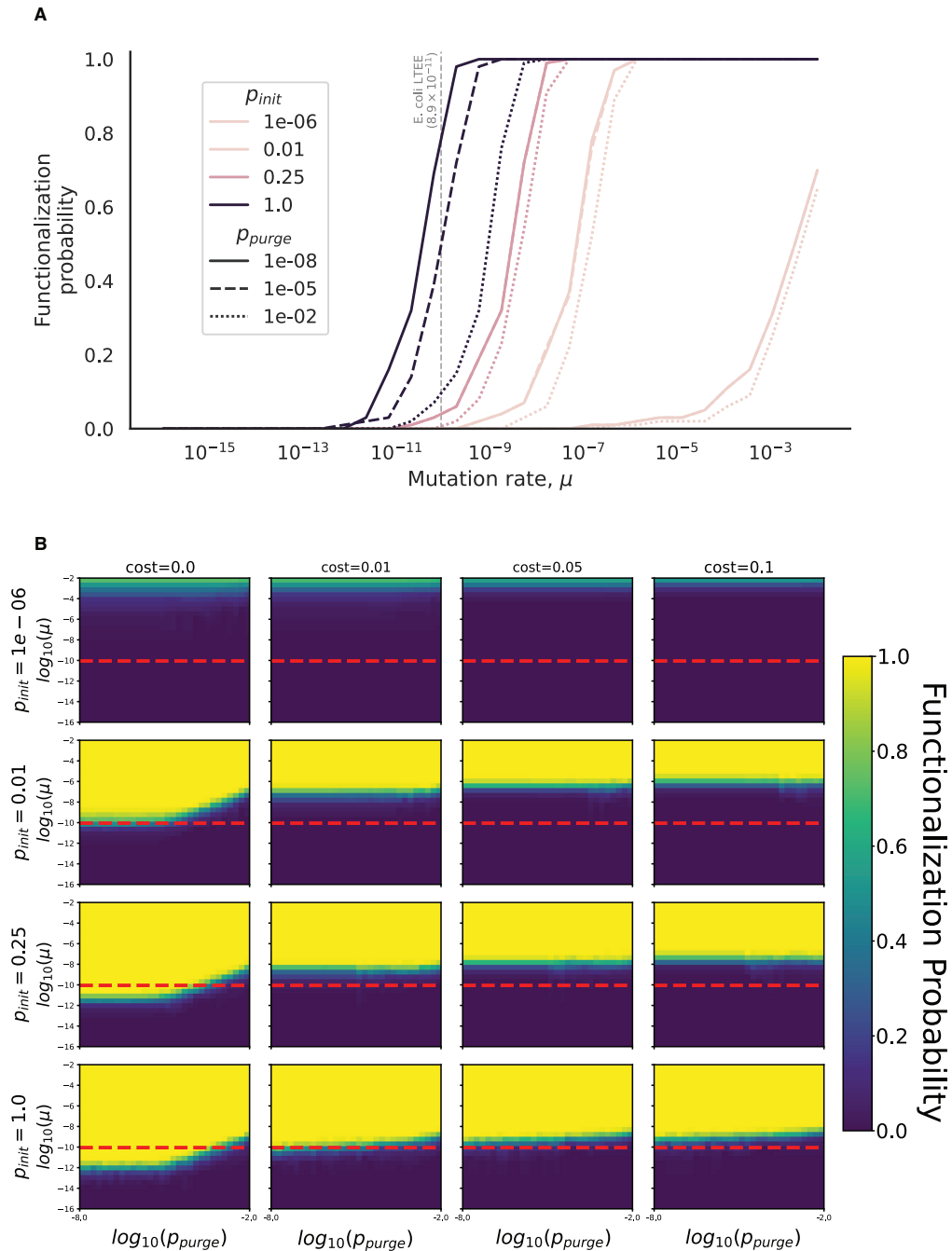
555
$$s \approx \frac{2 \ln(10^7)}{500} \approx 0.065$$

556 Yielding the estimated selection coefficient of 6.5% for the *yjcO-lysU* fusion/deletion.

557

558 **SUPPLEMENTAL INFORMATION**

559

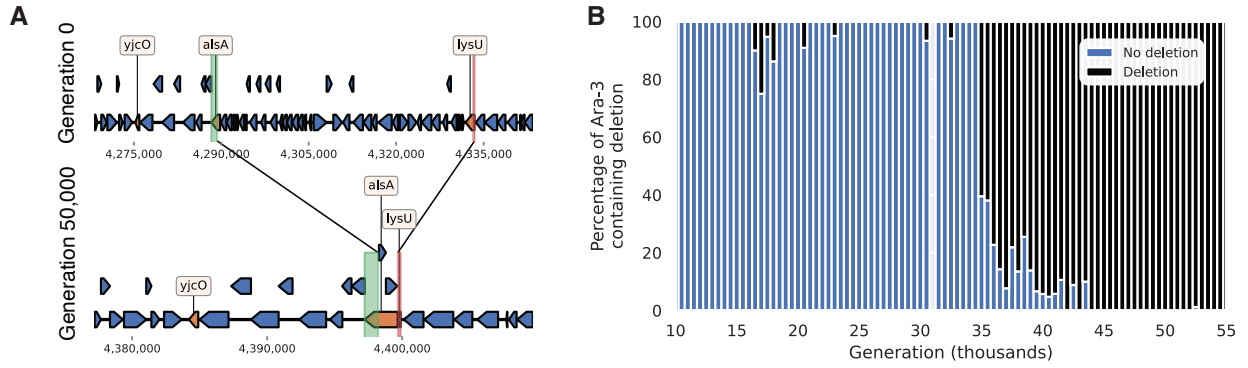


560

561 **Supplementary Figure 1. Forward simulations quantify the hitch-hiking advantage**

562 (A) Probability that at least one lineage functionalizes as a function of the mutation rate μ ; colors denote four starting
563 frequencies of the nascent fusion (p_{init}) and line styles three purge probabilities (p_{purge}). The dashed vertical line marks
564 the LTEE point-mutation rate. All curves are shown for a fusion gene with a pre-functionalized fitness cost of 0.01.

565 (B) Heat-maps show the same probability across grids of p_{purge} (x-axis, \log_{10} scale) and fitness cost of the unfixed fusion
566 (c, four columns) for the four p_{init} values (rows); the color of the cell indicates the proportion of times the fusion
567 functionalized before being purged. The dashed red line denotes the LTEE point-mutation rate. Each cell summarizes
568 100 Wright-Fisher runs of 10^6 haploids.



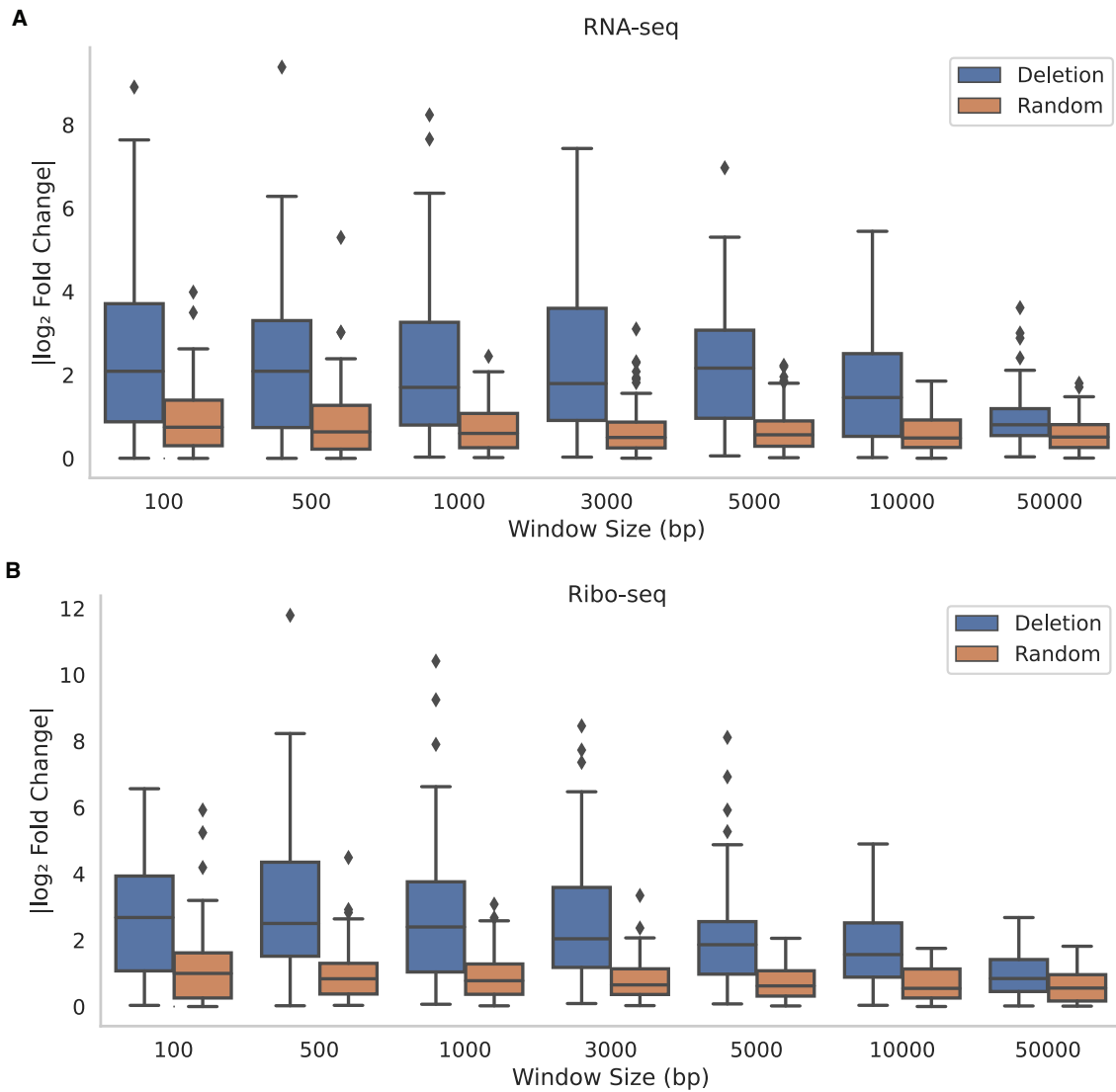
569

570 **Supplementary Figure 2. A convergent 43.4 kb deletion sweeps at the same locus in Ara-3.**
571 (A) Schematic of the deletion in Ara-3, orange genes highlight the resulting prior genes involved in the deletion. Green
572 and red bars correspond to BLAST alignments to this region.

573 (B) Metagenomic sequencing results showing the fraction of reads supporting the deletion or not.

574

575

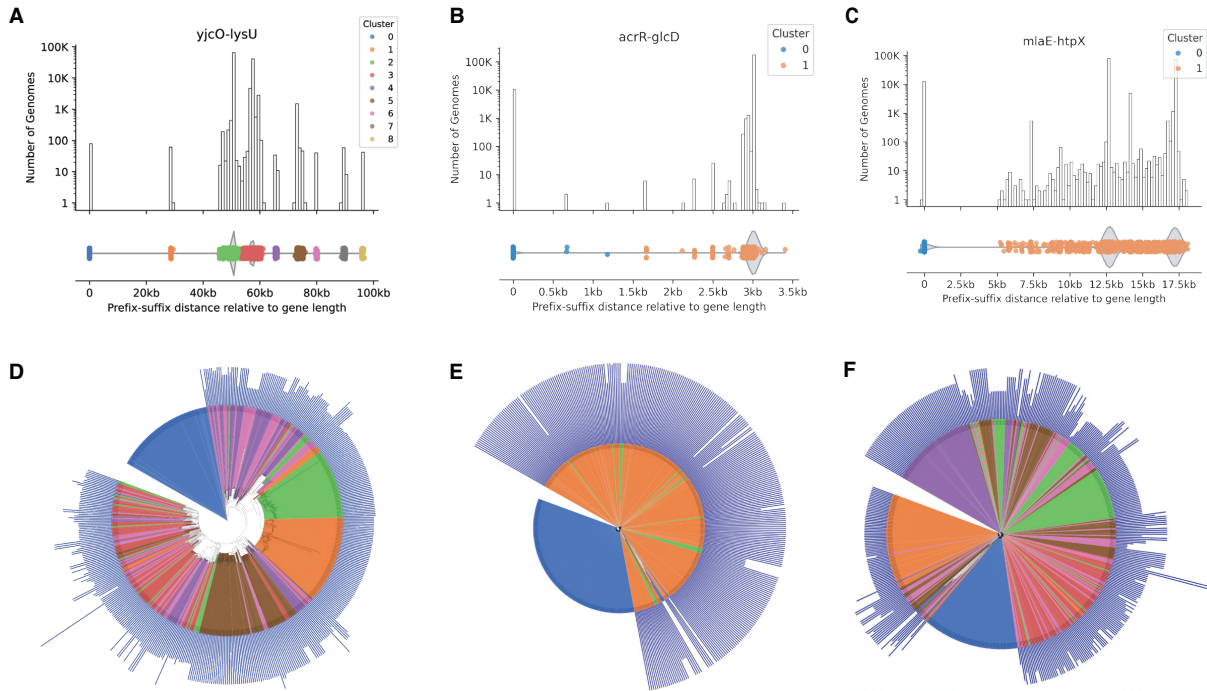


576

577 **Supplementary Figure 3. Large LTEE deletions are associated with significant changes to local**
578 **transcription and translation.**

579 (A) RNA-seq \log_2 fold changes for windows of varying size flanking ≥ 1 kb deletions (blue) and for randomly sampled
580 windows (orange). Fold changes are calculated between the ancestral strain and the evolved population at generation
581 50,000. Data downloaded from Favate et al. 2022.⁴⁰

582 (B) As in (A) but analyzing Ribo-seq data.



583

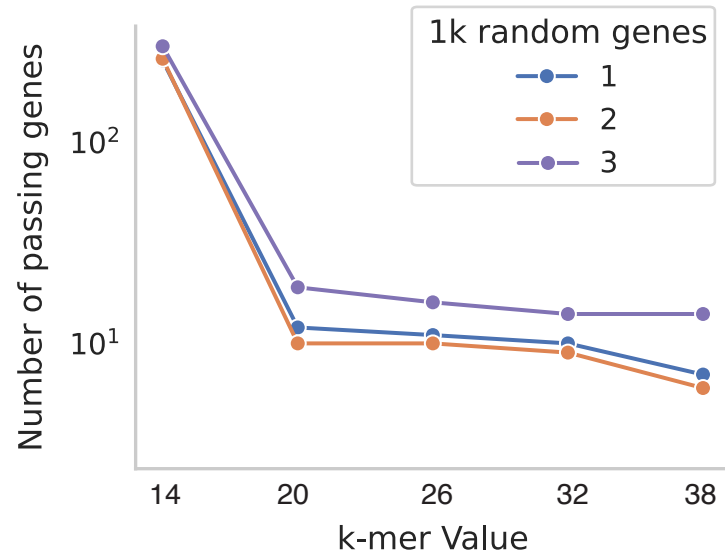
584 **Supplementary Figure 4. The prefix-suffix approach captures previously identified deletion-born**
585 **fusions and reveals additional structural variation.**

586 (A) Top: Log-scaled histogram of prefix-suffix distances relative to the length of the original *yjcO-lysU* fusion. Bottom:
587 underlying point distribution for histogram. Every point represents a single genome in the ATB dataset; colors
588 correspond to DBSCAN-derived clusters.

589 (B), (C) are the same as in (A) but with *acrR-glcD* and *mlaE-htpX* respectively.

590 (D) Distance-metric based phylogeny of 300 randomly sampled genomes with equal genomes sampled across the
591 number of clusters identified. Clades are colored by cluster membership, and blue vertical bars off leaves represent the
592 distance between the prefix-suffix found in that genome relative to the *yjcO-lysU* gene.

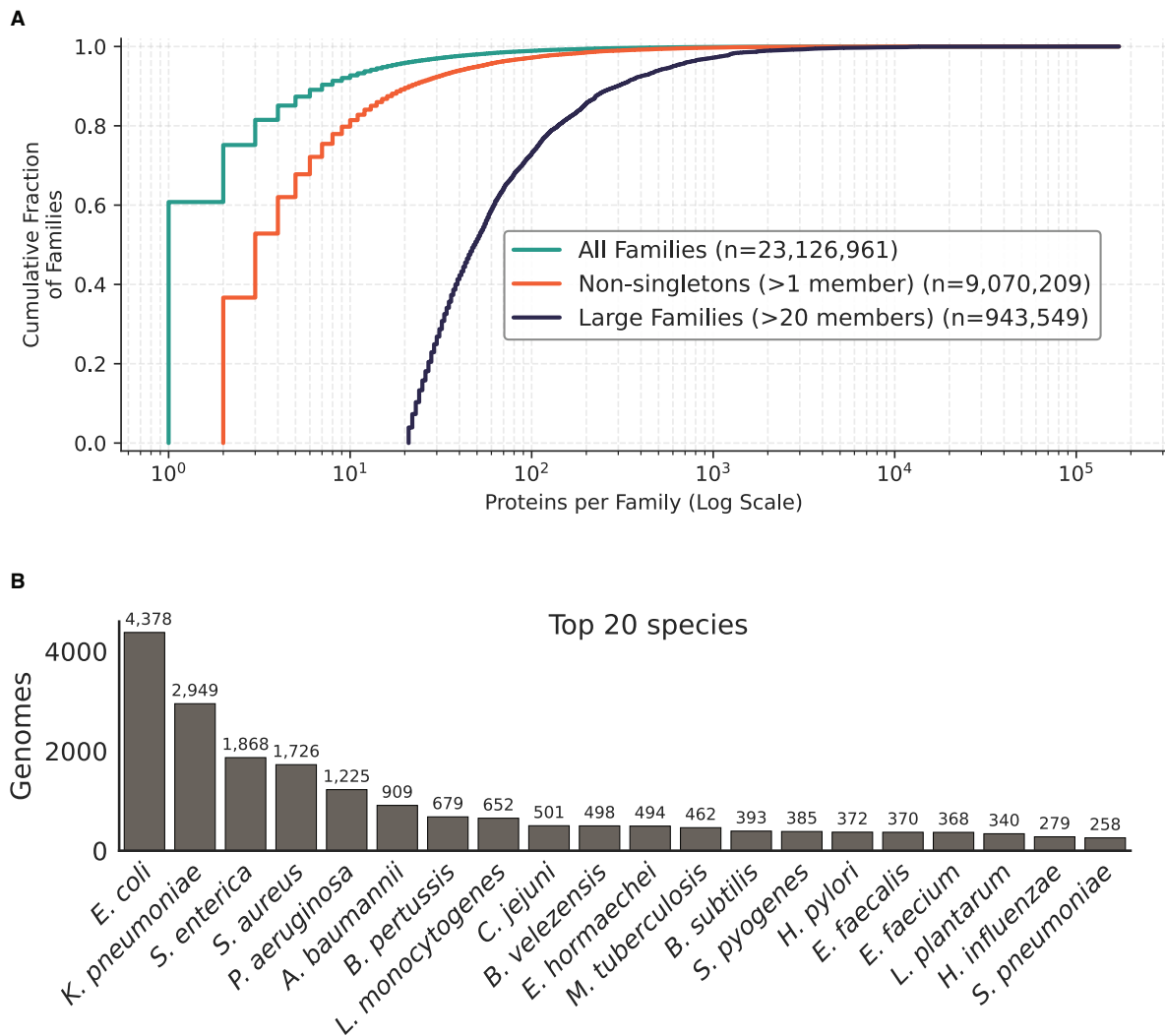
593 (E), (F) are the same as in (D) but with *acrR-glcD* and *mlaE-htpX* respectively.



603

604 **Supplementary Figure 6. Prefix-suffix approach is robust to values of $k \geq 20$.**

605 Number of genes with multimodal prefix-suffix distances detected across three random samples of 1,000 RefSeq genes
606 at varying k-mer lengths.



607

608 **Supplementary Figure 7. Representation of both protein families and bacterial species are highly**
609 **skewed in RefSeq complete genomes.**

610 (A) Cumulative distribution function of the number of protein members per family. Most represented proteins are
611 singletons

612 (B) Top 20 species represented in the 54,630 complete genomes proteins were pulled from. Count of each is displayed
613 above each bar.

614

615 **REFERENCES**

- 616 1. Hug, L.A., Baker, B.J., Anantharaman, K., Brown, C.T., Probst, A.J., Castelle, C.J.,
617 Butterfield, C.N., HERNSDORF, A.W., Amano, Y., Ise, K., et al. (2016). A new view of the tree of
618 life. *Nat Microbiol* 1, 1–6. <https://doi.org/10.1038/nmicrobiol.2016.48>.
- 619 2. Louca, S., Mazel, F., Doebeli, M., and Parfrey, L.W. (2019). A census-based estimate of
620 Earth's bacterial and archaeal diversity. *PLOS Biology* 17, e3000106.
621 <https://doi.org/10.1371/journal.pbio.3000106>.
- 622 3. Welch, R.A., Burland, V., Plunkett, G., Redford, P., Roesch, P., Rasko, D., Buckles, E.L.,
623 Liou, S.-R., Boutin, A., Hackett, J., et al. (2002). Extensive mosaic structure revealed by
624 the complete genome sequence of uropathogenic *Escherichia coli*. *Proceedings of the*
625 *National Academy of Sciences* 99, 17020–17024.
626 <https://doi.org/10.1073/pnas.252529799>.
- 627 4. Brockhurst, M.A., Harrison, E., Hall, J.P.J., Richards, T., McNally, A., and MacLean, C.
628 (2019). The Ecology and Evolution of Pangenomes. *Current Biology* 29, R1094–R1103.
629 <https://doi.org/10.1016/j.cub.2019.08.012>.
- 630 5. Koonin, E.V., Makarova, K.S., and Wolf, Y.I. (2021). Evolution of Microbial Genomics:
631 Conceptual Shifts over a Quarter Century. *Trends in Microbiology* 29, 582–592.
632 <https://doi.org/10.1016/j.tim.2021.01.005>.
- 633 6. Pasolli, E., Asnicar, F., Manara, S., Zolfo, M., Karcher, N., Armanini, F., Beghini, F.,
634 Manghi, P., Tett, A., Ghensi, P., et al. (2019). Extensive Unexplored Human Microbiome
635 Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age,
636 Geography, and Lifestyle. *Cell* 176, 649–662.e20.
637 <https://doi.org/10.1016/j.cell.2019.01.001>.
- 638 7. Baltoumas, F.A., Karatzas, E., Paez-Espino, D., Venetsianou, N.K., Aplakidou, E., Oulas,
639 A., Finn, R.D., Ovchinnikov, S., Pafilis, E., Kyrpides, N.C., et al. (2023). Exploring microbial
640 functional biodiversity at the protein family level—From metagenomic sequence reads to
641 annotated protein clusters. *Front. Bioinform.* 3.
642 <https://doi.org/10.3389/fbinf.2023.1157956>.
- 643 8. Álvarez-Lugo, A., and Becerra, A. (2021). The Role of Gene Duplication in the Divergence of
644 Enzyme Function: A Comparative Approach. *Front. Genet.* 12.
645 <https://doi.org/10.3389/fgene.2021.641817>.
- 646 9. Toll-Riera, M., Millan, A.S., Wagner, A., and MacLean, R.C. (2016). The Genomic Basis of
647 Evolutionary Innovation in *Pseudomonas aeruginosa*. *PLOS Genetics* 12, e1006005.
648 <https://doi.org/10.1371/journal.pgen.1006005>.
- 649 10. Sanchez-Herrero, J.F., Bernabeu, M., Prieto, A., Hüttner, M., and Juárez, A. (2020). Gene
650 Duplications in the Genomes of Staphylococci and Enterococci. *Front. Mol. Biosci.* 7.
651 <https://doi.org/10.3389/fmolb.2020.00160>.
- 652 11. Pavesi, A., Vianelli, A., Chirico, N., Bao, Y., Blinkova, O., Belshaw, R., Firth, A., and Karlin,
653 D. (2018). Overlapping genes and the proteins they encode differ significantly in their
654 sequence composition from non-overlapping genes. *PLOS ONE* 13, e0202513.
655 <https://doi.org/10.1371/journal.pone.0202513>.

- 656 12. Watson, A.K., Lopez, P., and Bapteste, E. (2022). Hundreds of Out-of-Frame Remodeled
657 Gene Families in the Escherichia coli Pangenome. *Molecular Biology and Evolution* 39,
658 msab329. <https://doi.org/10.1093/molbev/msab329>.
- 659 13. Irbäck, A., Peterson, C., and Potthast, F. (1996). Evidence for Non-Random Hydrophobicity
660 Structures in Protein Chains. Preprint at arXiv, [https://doi.org/10.48550/arXiv.chem-](https://doi.org/10.48550/arXiv.chem-ph/9512004)
661 [ph/9512004](https://doi.org/10.48550/arXiv.chem-ph/9512004) <https://doi.org/10.48550/arXiv.chem-ph/9512004>.
- 662 14. Foy, S.G., Wilson, B.A., Bertram, J., Cordes, M.H.J., and Masel, J. (2019). A Shift in
663 Aggregation Avoidance Strategy Marks a Long-Term Direction to Protein Evolution.
664 *Genetics* 211, 1345–1355. <https://doi.org/10.1534/genetics.118.301719>.
- 665 15. Pasek, S., Risler, J.-L., and Brézellec, P. (2006). Gene fusion/fission is a major contributor
666 to evolution of multi-domain bacterial proteins. *Bioinformatics* 22, 1418–1423.
667 <https://doi.org/10.1093/bioinformatics/btl135>.
- 668 16. Giovannoni, S.J., Cameron Thrash, J., and Temperton, B. (2014). Implications of
669 streamlining theory for microbial ecology. *ISME J* 8, 1553–1565.
670 <https://doi.org/10.1038/ismej.2014.60>.
- 671 17. Mira, A., Ochman, H., and Moran, N.A. (2001). Deletional bias and the evolution of
672 bacterial genomes. *Trends in Genetics* 17, 589–596. [https://doi.org/10.1016/S0168-](https://doi.org/10.1016/S0168-9525(01)02447-7)
673 [9525\(01\)02447-7](https://doi.org/10.1016/S0168-9525(01)02447-7).
- 674 18. Lynch, M. (2006). Streamlining and Simplification of Microbial Genome Architecture. *Annu.*
675 *Rev. Microbiol.* 60, 327–349. <https://doi.org/10.1146/annurev.micro.60.080805.142300>.
- 676 19. Wolf, Y.I., and Koonin, E.V. (2013). Genome reduction as the dominant mode of evolution.
677 *Bioessays* 35, 829–837. <https://doi.org/10.1002/bies.201300037>.
- 678 20. Giovannoni, S.J., Tripp, H.J., Givan, S., Podar, M., Vergin, K.L., Baptista, D., Bibbs, L.,
679 Eads, J., Richardson, T.H., Noordewier, M., et al. (2005). Genome Streamlining in a
680 Cosmopolitan Oceanic Bacterium. *Science* 309, 1242–1245.
681 <https://doi.org/10.1126/science.1114057>.
- 682 21. Ashish, A., Paterson, S., Mowat, E., Fothergill, J.L., Walshaw, M.J., and Winstanley, C.
683 (2013). Extensive diversification is a common feature of *Pseudomonas aeruginosa*
684 populations during respiratory infections in cystic fibrosis. *Journal of Cystic Fibrosis* 12,
685 790–793. <https://doi.org/10.1016/j.jcf.2013.04.003>.
- 686 22. Bottai, D., Frigui, W., Sayes, F., Di Luca, M., Spadoni, D., Pawlik, A., Zoppo, M., Orgeur,
687 M., Khanna, V., Hardy, D., et al. (2020). TbD1 deletion as a driver of the evolutionary
688 success of modern epidemic *Mycobacterium tuberculosis* lineages. *Nat Commun* 11, 684.
689 <https://doi.org/10.1038/s41467-020-14508-5>.
- 690 23. Lee, M.-C., and Marx, C.J. (2012). Repeated, Selection-Driven Genome Reduction of
691 Accessory Genes in Experimental Populations. *PLOS Genetics* 8, e1002651.
692 <https://doi.org/10.1371/journal.pgen.1002651>.
- 693 24. Raeside, C., Gaffé, J., Deatherage, D.E., Tenaillon, O., Briska, A.M., Ptashkin, R.N.,
694 Cruveiller, S., Médigue, C., Lenski, R.E., Barrick, J.E., et al. (2014). Large Chromosomal
695 Rearrangements during a Long-Term Evolution Experiment with *Escherichia coli*. *mBio* 5,
696 10.1128/mbio.01377-14. <https://doi.org/10.1128/mbio.01377-14>.

- 697 25. Rocha, E.P.C., Cornet, E., and Michel, B. (2005). Comparative and Evolutionary Analysis of
698 the Bacterial Homologous Recombination Systems. *PLOS Genetics* 1, e15.
699 <https://doi.org/10.1371/journal.pgen.0010015>.
- 700 26. Bichara, M., Wagner, J., and Lambert, I.B. (2006). Mechanisms of tandem repeat instability
701 in bacteria. *Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis*
702 598, 144–163. <https://doi.org/10.1016/j.mrfmmm.2006.01.020>.
- 703 27. Clayton, A.L., Jackson, D.G., Weiss, R.B., and Dale, C. (2016). Adaptation by Deletogenic
704 Replication Slippage in a Nascent Symbiont. *Mol Biol Evol* 33, 1957–1966.
705 <https://doi.org/10.1093/molbev/msw071>.
- 706 28. Hoff, G., Bertrand, C., Piotrowski, E., Thibessard, A., and Leblond, P. (2018). Genome
707 plasticity is governed by double strand break DNA repair in *Streptomyces*. *Sci Rep* 8, 5272.
708 <https://doi.org/10.1038/s41598-018-23622-w>.
- 709 29. Suzuki, N., Inui, M., and Yukawa, H. (2008). Random genome deletion methods applicable
710 to prokaryotes. *Appl Microbiol Biotechnol* 79, 519–526. [https://doi.org/10.1007/s00253-](https://doi.org/10.1007/s00253-008-1512-4)
711 [008-1512-4](https://doi.org/10.1007/s00253-008-1512-4).
- 712 30. Lynch, M., and Marinov, G.K. (2015). The bioenergetic costs of a gene. *Proceedings of the*
713 *National Academy of Sciences* 112, 15690–15695.
714 <https://doi.org/10.1073/pnas.1514974112>.
- 715 31. Lenski, R.E., Rose, M.R., Simpson, S.C., and Tadler, S.C. (1991). Long-Term Experimental
716 Evolution in *Escherichia coli*. I. Adaptation and Divergence During 2,000 Generations. *The*
717 *American Naturalist* 138, 1315–1341. <https://doi.org/10.1086/285289>.
- 718 32. uz-Zaman, M.H., D’Alton, S., Barrick, J.E., and Ochman, H. (2024). Promoter recruitment
719 drives the emergence of proto-genes in a long-term evolution experiment with *Escherichia*
720 *coli*. *PLOS Biology* 22, e3002418. <https://doi.org/10.1371/journal.pbio.3002418>.
- 721 33. Warsi, O., Knopp, M., Surkov, S., Jerlström Hultqvist, J., and Andersson, D.I. (2020).
722 Evolution of a New Function by Fusion between Phage DNA and a Bacterial Gene.
723 *Molecular Biology and Evolution* 37, 1329–1341.
724 <https://doi.org/10.1093/molbev/msaa007>.
- 725 34. Farr, A.D., Remigi, P., and Rainey, P.B. (2017). Adaptive evolution by spontaneous domain
726 fusion and protein relocalization. *Nat Ecol Evol* 1, 1562–1568.
727 <https://doi.org/10.1038/s41559-017-0283-7>.
- 728 35. Gallant, J., Mouton, J., Ummels, R., ten Hagen-Jongman, C., Kriel, N., Pain, A., Warren,
729 R.M., Bitter, W., Heunis, T., and Sampson, S.L. (2020). Identification of gene fusion events
730 in *Mycobacterium tuberculosis* that encode chimeric proteins. *NAR Genomics and*
731 *Bioinformatics* 2, lqaa033. <https://doi.org/10.1093/nargab/lqaa033>.
- 732 36. Bobay, L.-M., and Ochman, H. (2017). The Evolution of Bacterial Genome Architecture.
733 *Front. Genet.* 8. <https://doi.org/10.3389/fgene.2017.00072>.
- 734 37. Johnson, M.S., Martsul, A., Kryazhimskiy, S., and Desai, M.M. (2019). Higher-fitness yeast
735 genotypes are less robust to deleterious mutations. *Science* 366, 490–493.
736 <https://doi.org/10.1126/science.aay4199>.

- 737 38. Barrick, J.E., and Lenski, R.E. (2013). Genome dynamics during experimental evolution.
738 *Nat Rev Genet* 14, 827–839. <https://doi.org/10.1038/nrg3564>.
- 739 39. Good, B.H., McDonald, M.J., Barrick, J.E., Lenski, R.E., and Desai, M.M. (2017). The
740 dynamics of molecular evolution over 60,000 generations. *Nature* 551, 45–50.
741 <https://doi.org/10.1038/nature24287>.
- 742 40. Favate, J.S., Liang, S., Cope, A.L., Yadavalli, S.S., and Shah, P. (2022). The landscape of
743 transcriptional and translational changes over 22 years of bacterial adaptation. *eLife* 11.
744 <https://doi.org/10.7554/elife.81979>.
- 745 41. Couce, A., Limdi, A., Magnan, M., Owen, S.V., Herren, C.M., Lenski, R.E., Tenailon, O.,
746 and Baym, M. (2024). Changing fitness effects of mutations through long-term bacterial
747 evolution. *Science* 383, eadd1417. <https://doi.org/10.1126/science.add1417>.
- 748 42. Marin, M., Vargas, R., Harris, M., Jeffrey, B., Epperson, L.E., Durbin, D., Strong, M.,
749 Salfinger, M., Iqbal, Z., Akhundova, I., et al. (2022). Benchmarking the empirical accuracy
750 of short-read sequencing across the *M. tuberculosis* genome. *Bioinformatics* 38, 1781–
751 1787. <https://doi.org/10.1093/bioinformatics/btac023>.
- 752 43. Charles, C., Conde, C., Vorimore, F., Cochard, T., Michelet, L., Boschioli, M.L., and Biet,
753 F. (2023). Features of *Mycobacterium bovis* Complete Genomes Belonging to 5 Different
754 Lineages. *Microorganisms* 11, 177. <https://doi.org/10.3390/microorganisms11010177>.
- 755 44. Bespiatykh, D., Bespyatykh, J., Mokrousov, I., and Shitikov, E. (2021). A Comprehensive
756 Map of *Mycobacterium tuberculosis* Complex Regions of Difference. *mSphere* 6, e00535–21.
757 <https://doi.org/10.1128/mSphere.00535-21>.
- 758 45. Ferragina, P., and Manzini, G. (2000). Opportunistic data structures with applications. In
759 *Proceedings 41st Annual Symposium on Foundations of Computer Science*, pp. 390–398.
760 <https://doi.org/10.1109/SFCS.2000.892127>.
- 761 46. Nagarajan, D., Nagarajan, T., Roy, N., Kulkarni, O., Ravichandran, S., Mishra, M.,
762 Chakravorty, D., and Chandra, N. (2018). Computational antimicrobial peptide design and
763 evaluation against multidrug-resistant clinical isolates of bacteria. *J Biol Chem* 293, 3492–
764 3509. <https://doi.org/10.1074/jbc.M117.805499>.
- 765 47. Dekker, N.P., Lammel, C.J., and Brooks, Geo.F. (1991). Scanning electron microscopy of
766 piliated *Neisseria gonorrhoeae* processed with hexamethyldisilazane. *J. Elec. Microsc. Tech.*
767 19, 461–467. <https://doi.org/10.1002/jemt.1060190408>.
- 768 48. Dahl, J.L. (2004). Electron microscopy analysis of *Mycobacterium tuberculosis* cell division.
769 *FEMS Microbiology Letters* 240, 15–20. <https://doi.org/10.1016/j.femsle.2004.09.004>.
- 770 49. Hilbert, F., Scherwitzel, M., Paulsen, P., and Szostak, M.P. (2010). Survival of
771 *Campylobacter jejuni* under conditions of atmospheric oxygen tension with the support of
772 *Pseudomonas* spp. *Appl Environ Microbiol* 76, 5911–5917.
773 <https://doi.org/10.1128/AEM.01532-10>.
- 774 50. Straume, D., Piechowiak, K.W., Olsen, S., Stamsås, G.A., Berg, K.H., Kjos, M.,
775 Heggenhougen, M.V., Alcorlo, M., Hermoso, J.A., and Håvarstein, L.S. (2020). Class A PBPs
776 have a distinct and unique role in the construction of the pneumococcal cell wall.
777 *Proceedings of the National Academy of Sciences* 117, 6129–6138.
778 <https://doi.org/10.1073/pnas.1917820117>.

- 779 51. Lieberman, T.D. (2022). Detecting bacterial adaptation within individual microbiomes.
780 *Philosophical Transactions of the Royal Society B: Biological Sciences* 377, 20210243.
781 <https://doi.org/10.1098/rstb.2021.0243>.
- 782 52. Gagneux, S. (2018). Ecology and evolution of *Mycobacterium tuberculosis*. *Nat Rev*
783 *Microbiol* 16, 202–213. <https://doi.org/10.1038/nrmicro.2018.8>.
- 784 53. Gong, W., and Wu, X. (2021). Differential Diagnosis of Latent Tuberculosis Infection and
785 Active Tuberculosis: A Key to a Successful Tuberculosis Control Strategy. *Front. Microbiol.*
786 12. <https://doi.org/10.3389/fmicb.2021.745592>.
- 787 54. Sohaskey, C.D., and Wayne, L.G. (2003). Role of narK2X and narGHJI in Hypoxic
788 Upregulation of Nitrate Reduction by *Mycobacterium tuberculosis*. *Journal of Bacteriology*
789 185, 7247–7256. <https://doi.org/10.1128/jb.185.24.7247-7256.2003>.
- 790 55. Cunningham-Bussel, A., Zhang, T., and Nathan, C.F. (2013). Nitrite produced by
791 *Mycobacterium tuberculosis* in human macrophages in physiologic oxygen impacts
792 bacterial ATP consumption and gene expression. *Proc Natl Acad Sci U S A* 110, E4256–
793 E4265. <https://doi.org/10.1073/pnas.1316894110>.
- 794 56. Hunt, M., Lima, L., Anderson, D., Bouras, G., Hall, M., Hawkey, J., Schwengers, O., Shen,
795 W., Lees, J.A., and Iqbal, Z. (2025). AllTheBacteria – all bacterial genomes assembled,
796 available, and searchable. Preprint at bioRxiv,
797 <https://doi.org/10.1101/2024.03.08.584059>
798 <https://doi.org/10.1101/2024.03.08.584059>.
- 799 57. Ellrott, K., Jaroszewski, L., Li, W., Wooley, J.C., and Godzik, A. (2010). Expansion of the
800 Protein Repertoire in Newly Explored Environments: Human Gut Microbiome Specific
801 Protein Families. *PLOS Computational Biology* 6, e1000798.
802 <https://doi.org/10.1371/journal.pcbi.1000798>.
- 803 58. Yooseph, S., Sutton, G., Rusch, D.B., Halpern, A.L., Williamson, S.J., Remington, K.,
804 Eisen, J.A., Heidelberg, K.B., Manning, G., Li, W., et al. (2007). The Sorcerer II Global
805 Ocean Sampling Expedition: Expanding the Universe of Protein Families. *PLOS Biology* 5,
806 e16. <https://doi.org/10.1371/journal.pbio.0050016>.
- 807 59. Siew, N., and Fischer, D. (2003). Analysis of singleton ORFans in fully sequenced microbial
808 genomes. *Proteins* 53, 241–251. <https://doi.org/10.1002/prot.10423>.
- 809 60. Yanai, I., Wolf, Y.I., and Koonin, E.V. (2002). Evolution of gene fusions: horizontal transfer
810 versus independent events. *Genome Biol* 3, research0024.1. [https://doi.org/10.1186/gb-](https://doi.org/10.1186/gb-2002-3-5-research0024)
811 [2002-3-5-research0024](https://doi.org/10.1186/gb-2002-3-5-research0024).
- 812 61. Carvunis, A.-R., Rolland, T., Wapinski, I., Calderwood, M.A., Yildirim, M.A., Simonis, N.,
813 Charloteaux, B., Hidalgo, C.A., Barbette, J., Santhanam, B., et al. (2012). Proto-genes and
814 de novo gene birth. *Nature* 487, 370–374. <https://doi.org/10.1038/nature11184>.
- 815 62. Weisman, C.M., Murray, A.W., and Eddy, S.R. (2020). Many, but not all, lineage-specific
816 genes can be explained by homology detection failure. *PLOS Biology* 18, e3000862.
817 <https://doi.org/10.1371/journal.pbio.3000862>.
- 818 63. Lambowitz, A.M., and Zimmerly, S. (2004). Mobile Group II Introns. *Annual Review of*
819 *Genetics* 38, 1–35. <https://doi.org/10.1146/annurev.genet.38.072902.091600>.

- 820 64. Merk, L.N., Jones, T.A., and Eddy, S.R. (2025). Presence of group II introns in phage
821 genomes. *Nucleic Acids Res* 53, gkaf761. <https://doi.org/10.1093/nar/gkaf761>.
- 822 65. Labbate, M., Case, R.J., and Stokes, H.W. (2009). The Integron/Gene Cassette System: An
823 Active Player in Bacterial Adaptation. In *Horizontal Gene Transfer: Genomes in Flux*, M. B.
824 Gogarten, J. P. Gogarten, and L. C. Olendzenski, eds. (Humana Press), pp. 103–125.
825 https://doi.org/10.1007/978-1-60327-853-9_6.
- 826 66. Sereika, M., Mussig, A.J., Jiang, C., Knudsen, K.S., Jensen, T.B.N., Petriglieri, F., Yang, Y.,
827 Jørgensen, V.R., Delogu, F., Sørensen, E.A., et al. (2025). Genome-resolved long-read
828 sequencing expands known microbial diversity across terrestrial habitats. *Nat Microbiol* 10,
829 2018–2030. <https://doi.org/10.1038/s41564-025-02062-z>.
- 830 67. Chen, L., Zhao, N., Cao, J., Liu, X., Xu, J., Ma, Y., Yu, Y., Zhang, X., Zhang, W., Guan, X.,
831 et al. (2022). Short- and long-read metagenomics expand individualized structural
832 variations in gut microbiomes. *Nat Commun* 13, 3175. [https://doi.org/10.1038/s41467-](https://doi.org/10.1038/s41467-022-30857-9)
833 [022-30857-9](https://doi.org/10.1038/s41467-022-30857-9).
- 834 68. Kadelka, C., Wheeler, M., Veliz-Cuba, A., Murrugarra, D., and Laubenbacher, R. (2023).
835 Modularity of biological systems: a link between structure and function. *J R Soc Interface*
836 20, 20230505. <https://doi.org/10.1098/rsif.2023.0505>.
- 837 69. Kimura, M. (1983). *The Neutral Theory of Molecular Evolution* 1st ed. (Cambridge
838 University Press) <https://doi.org/10.1017/CBO9780511623486>.
- 839 70. Harris, C.R., Millman, K.J., van der Walt, S.J., Gommers, R., Virtanen, P., Cournapeau, D.,
840 Wieser, E., Taylor, J., Berg, S., Smith, N.J., et al. (2020). Array programming with NumPy.
841 *Nature* 585, 357–362. <https://doi.org/10.1038/s41586-020-2649-2>.
- 842 71. Waskom, M. (2021). seaborn: statistical data visualization. *JOSS* 6, 3021.
843 <https://doi.org/10.21105/joss.03021>.
- 844 72. Seemann, T. (2014). Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30,
845 2068–2069. <https://doi.org/10.1093/bioinformatics/btu153>.
- 846 73. Tonkin-Hill, G., MacAlasdair, N., Ruis, C., Weimann, A., Horesh, G., Lees, J.A., Gladstone,
847 R.A., Lo, S., Beaudoin, C., Floto, R.A., et al. (2020). Producing polished prokaryotic
848 pangenomes with the Panaroo pipeline. *Genome Biology* 21, 180.
849 <https://doi.org/10.1186/s13059-020-02090-4>.
- 850 74. Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and
851 Madden, T.L. (2009). BLAST+: architecture and applications. *BMC Bioinformatics* 10, 421.
852 <https://doi.org/10.1186/1471-2105-10-421>.
- 853 75. Zulkower, V., and Rosser, S. (2020). DNA Features Viewer: a sequence annotation
854 formatting and plotting library for Python. *Bioinformatics* 36, 4350–4352.
855 <https://doi.org/10.1093/bioinformatics/btaa213>.
- 856 76. BEDOPS: high-performance genomic feature operations | *Bioinformatics* | Oxford
857 Academic <https://academic.oup.com/bioinformatics/article/28/14/1919/218826>.
- 858 77. Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34,
859 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>.

- 860 78. Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A., Sonnhammer, E.L.L.,
861 Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J., et al. (2021). Pfam: The protein
862 families database in 2021. *Nucleic Acids Research* 49, D412–D419.
863 <https://doi.org/10.1093/nar/gkaa913>.
- 864 79. Eddy, Sean HMMER. <http://hmmer.org/>.
- 865 80. Kurtz, S., Phillippy, A., Delcher, A.L., Smoot, M., Shumway, M., Antonescu, C., and
866 Salzberg, S.L. (2004). Versatile and open software for comparing large genomes. *Genome*
867 *Biol* 5, R12. <https://doi.org/10.1186/gb-2004-5-2-r12>.
- 868 81. Larralde, M. (2022). Pyrodigal: Python bindings and interface to Prodigal, an efficient
869 method for gene prediction in prokaryotes. *Journal of Open Source Software* 7, 4296.
870 <https://doi.org/10.21105/joss.04296>.
- 871 82. Steinegger, M., and Söding, J. (2017). MMseqs2 enables sensitive protein sequence
872 searching for the analysis of massive data sets. *Nat Biotechnol* 35, 1026–1028.
873 <https://doi.org/10.1038/nbt.3988>.
- 874 83. Kille, B., Nute, M.G., Huang, V., Kim, E., Phillippy, A.M., and Treangen, T.J. (2024). Parsnp
875 2.0: scalable core-genome alignment for massive microbial datasets. *Bioinformatics* 40,
876 *btac311*. <https://doi.org/10.1093/bioinformatics/btac311>.
- 877 84. Letunic, I., and Bork, P. (2024). Interactive Tree of Life (iTOL) v6: recent updates to the
878 phylogenetic tree display and annotation tool. *Nucleic Acids Research* 52, W78–W82.
879 <https://doi.org/10.1093/nar/gkae268>.
- 880 85. Břinda, K., Lima, L., Pignotti, S., Quinones-Olvera, N., Salikhov, K., Chikhi, R., Kucherov,
881 G., Iqbal, Z., and Baym, M. (2025). Efficient and robust search of microbial genomes via
882 phylogenetic compression. *Nat Methods* 22, 692–697. [https://doi.org/10.1038/s41592-](https://doi.org/10.1038/s41592-025-02625-2)
883 [025-02625-2](https://doi.org/10.1038/s41592-025-02625-2).
- 884 86. Chambers, D. (2025). [d-chambers/dbscan1d](https://github.com/dchambers/dbscan1d).
- 885 87. Břinda, K. (2025). [karel-brinda/attotree](https://github.com/karel-brinda/attotree).
- 886 88. Katz, L., Griswold, T., Morrison, S., Caravas, J., Zhang, S., Bakker, H., Deng, X., and
887 Carleton, H. (2019). Mashtree: a rapid comparison of whole genome sequence files. *JOSS* 4,
888 1762. <https://doi.org/10.21105/joss.01762>.
- 889 89. Fitch, W.M. (1971). Toward Defining the Course of Evolution: Minimum Change for a
890 Specific Tree Topology. *Systematic Biology* 20, 406–416.
891 <https://doi.org/10.1093/sysbio/20.4.406>.
- 892 90. Johansson, M.H.K., Bortolaia, V., Tansirichaiya, S., Aarestrup, F.M., Roberts, A.P., and
893 Petersen, T.N. (2021). Detection of mobile genetic elements associated with antibiotic
894 resistance in *Salmonella enterica* using a newly developed web tool: MobileElementFinder.
895 *Journal of Antimicrobial Chemotherapy* 76, 101–109.
896 <https://doi.org/10.1093/jac/dkaa390>.
- 897 91. Dieppa-Colón, E., Martin, C., Kosmopoulos, J.C., and Anantharaman, K. (2025). Prophage-
898 DB: a comprehensive database to explore diversity, distribution, and ecology of prophages.
899 *Environmental Microbiome* 20, 5. <https://doi.org/10.1186/s40793-024-00659-1>.

- 900 92. Ranwez, V., Harispe, S., Delsuc, F., and Douzery, E.J.P. (2011). MACSE: Multiple
901 Alignment of Coding SEquences Accounting for Frameshifts and Stop Codons. PLOS ONE
902 6, e22594. <https://doi.org/10.1371/journal.pone.0022594>.
- 903 93. Murrell, B., Moola, S., Mabona, A., Weighill, T., Sheward, D., Kosakovsky Pond, S.L., and
904 Scheffler, K. (2013). FUBAR: A Fast, Unconstrained Bayesian AppRoximation for Inferring
905 Selection. *Molecular Biology and Evolution* 30, 1196–1205.
906 <https://doi.org/10.1093/molbev/mst030>.

907

908

909 **METHODS**

910 **Simulation Framework**

911 We simulated the fate of a novel gene using a minimal forward-time Wright–Fisher model with
912 a fixed haploid population size $N = 10^6$ and two initial states: 0 (no fusion), and 1 (non-
913 functional novel gene). State 1 carried a fitness cost c (relative fitness $1-c$), while states 0 had
914 fitness 1. State 2 is when the novel gene functionalizes, and no individual started at this point.

915 Simulations were initialized with the fusion in state 1 at frequency p_{init} and the rest at state 0.
916 Populations were updated each generation by fitness-proportionate multinomial sampling
917 (selection + drift). After “reproduction”, individuals in state 1 could functionalize (transition
918 from state 1 to state 2) with probability p_{func} per generation, and novel gene-bearing individuals
919 in state 1 can lose the gene (transition to state 0) with probability p_{purge} per generation. Each
920 replicate was run until the novel gene was either lost or functionalized (defined as the first
921 appearance of any state 2 individual, with a maximum of 100,000 generations.

922 For parameter sweeps (**Supplementary Figure 1**), we evaluated a log-spaced grid (20 points) of
923 p_{func} from 10^{-16} to 10^{-2} and p_{purge} from 10^{-8} to 10^{-2} for each combination of c and p_{init} running
924 100 replicate simulations per grid point with a fixed seed of 42. Simulations were implemented
925 in Python using numpy⁷⁰ and results were saved as matrices of functionalization probabilities.
926 Visualizations provided by seaborn.⁷¹ Code for simulations is available in the project repository.

927 **LTEE Data Analysis**

928 We obtained mutation call files for LTEE clonal isolates from the Barrick lab LTEE-Ecoli
929 repository (<https://github.com/barricklab/LTEE-Ecoli>). For each .gd file, we used gdttools
930 APPLY to generate an isolate-specific genome sequence by applying the curated variants to
931 REL606. We produced per-isolate FASTA outputs. These isolate-specific assemblies were used
932 for consistent gene calling and pangenome construction across all isolates.

933 To provide standardized gene models for pangenome inference, we annotated each isolate-
934 specific FASTA with Prokka,⁷² producing per-isolate GFF files. These GFFs were used as the
935 input to Panaroo.⁷³

936 For each of the twelve lineages, we ran Panaroo on the set of clonal isolate annotations for that
937 population together with the corresponding ancestral annotation (Anc+ or Anc-). Panaroo was
938 run with --merge_paralogs enabled and with stringent similarity thresholds (e.g., --
939 len_dif_percent 0.98, --threshold 0.98, --family_threshold 0.7), using the “moderate” clean
940 mode.

941 To identify candidate novel genes arising during evolution, we analyzed each population-
942 specific Panaroo run using a custom script. Briefly, isolate identifiers were mapped to LTEE
943 metadata (population and generation), and Panaroo’s gene_presence_absence.csv and
944 struct_presence_absence.Rtab were scanned to identify gene clusters whose first appearance
945 occurred after the ancestor and increased in presence among later isolates (“appearing genes”).

946 To distinguish gene families plausibly arising from structural rearrangement from those
947 reflecting annotation noise or near-identical ancestral sequence, we filtered candidate
948 appearing genes by sequence similarity to the REL606 ancestor. For each candidate gene
949 cluster, we extracted its representative nucleotide sequence from pan_genome_reference.fa and
950 aligned it to the ancestral REL606 genome with BLASTN.⁷⁴ We computed the fraction of the
951 candidate gene covered by its single largest BLAST match to REL606 and removed candidates
952 with >85% coverage by the best match, consistent with those being largely ancestral sequence
953 rather than novel junction-derived sequence.

954 *yjcO-lysU* was identified in this manner and its genomic context from the ancestor was
955 visualized using DNA Features Viewer.⁷⁵ RNA-sequencing and Ribosome profiling were
956 downloaded from Favate et al. and processed according with the same approach as
957 published.⁴⁰ DNA Features Viewer was again used to visualize the reads to the relevant clonal
958 genome.

959 For candidates that passed the REL606 similarity filter, we extracted a short sequence
960 spanning the putative novel junction. When BLAST produced two distinct hits to REL606, we
961 converted the two hit intervals to BED format and used bedops⁷⁶ with a 30 bp range to identify
962 the local overlap/adjacent junction region; when BLAST produced a single hit, we extracted a
963 30 bp window around the inferred boundary of the match (depending on hit
964 orientation/endpoint). The corresponding subsequence was extracted from the candidate gene
965 sequence and retained as a “junction” FASTA. To query the relative fraction of the population
966 with a specific variant, both the junction FASTA and the original sequence were aligned with
967 minimap2⁷⁷ to all metagenomic reads from Good et al.³⁹ The relative fraction of reads
968 supporting each variant was computed and results were visualized with seaborn.

969 Domain annotation was performed by using hmmsearch (version 3.4) on the translated protein
970 sequence against Pfam-A database (version 38).^{78,79}

971 To analyze the RNA-seq and Ribo-seq coverage changes around large structural variants, we
972 downloaded a LTEE structural variant table from <https://barricklab.org/shiny/LTEE-Ecoli/> on
973 2025.06.30. We first restricted to the 12 sequenced endpoint clones sequenced in the
974 RNA/Ribo dataset and selected events annotated as deletions or substitutions (DEL or SUB)
975 larger than 1 kb. Because the deletion coordinates were reported in REL606 reference
976 coordinates, we transferred breakpoints into each evolved clone’s coordinate system using
977 whole-genome alignment block coordinates (*.coords, generated by nucmer⁸⁰). For each clone,
978 we parsed alignment blocks describing reference-to-query interval mappings and mapped each
979 breakpoint by locating the corresponding block (allowing a 50 bp tolerance) and applying the
980 within-block offset; events with neither breakpoint mappable were excluded, while events with
981 only one breakpoint mappable were evaluated using a one-sided window around the mapped
982 breakpoint.

983 Per-base RNA-seq and Ribo-seq coverage was provided as strand-specific depth files for two
984 replicates per clone and for ancestral controls. For each deletion and each window size, we
985 extracted coverage across the mapped interval, averaged coverage across all available replicate
986 and strands to obtain a single mean coverage value and then computed the log₂ fold change of
987 evolved versus ancestral mean coverage in the same window. To assess spatial scale, we
988 repeated this analysis across multiple window sizes (100 bp to 50 kb). As a matched control,
989 for each clone and window size, we sampled the same number of random genomic windows,
990 excluding regions overlapping deletion windows (including flanks), mapped those windows into
991 clone coordinates using the same alignment-based procedure, and computed coverage log₂ fold
992 changes identically. Distributions of absolute log₂ fold change values were then compared
993 between deletion-flanking windows and matched random windows for RNA-seq and Ribo-seq.

994 ***Mycobacterium tuberculosis* Complex Analysis**

995 We first downloaded the 36 clinical *M. tb* genomes assembled in Marin et al. 2022 from NCBI
996 using BioProjects PRJNA719670, PRJNA480888, PRJNA436997 and PRJNA421446.⁴² The 10
997 *M. bovis* genomes used were also downloaded from NCBI using BioProjects PRJNA832544.⁴³
998 Genomes were retrieved in FASTA format and organized per isolate. The *M. tuberculosis* H37Rv
999 reference genome (accession AL123456) was also downloaded in GenBank format and used as
1000 a reference for downstream analyses.

1001 To enable consistent gene-content comparisons across isolates, we predicted protein-coding
1002 genes de novo for all genomes using pyrodigal⁸¹ and exported predicted proteins as amino acid
1003 FASTA files. All predicted proteins across genomes were then clustered into gene families using
1004 MMseqs2.⁸² We created a single MMseqs2 sequence database from all proteins, clustered
1005 sequences using identity-based clustering (minimum sequence identity 0.7), and exported
1006 cluster assignments as a tab-delimited table. These clusters were used to define gene families
1007 and identify accessory families whose presence varied across *M. tb* and *M. bovis* isolates.

1008 To prioritize candidate deletion-born fusions, we examined accessory gene families with
1009 lineage-restricted presence patterns and then performed nucleotide-level mapping in genomes
1010 lacking the gene family. For each candidate family, we extracted representative nucleotide
1011 sequences and aligned them to the corresponding genome FASTA files using BLASTN. BLAST
1012 hits were converted to genomic intervals, merged to identify discrete matching regions, and
1013 filtered to enrich for signatures consistent with deletion-born fusion rather than simple
1014 absence, fragmentation, or duplication. Specifically, retained candidates required at least 80%
1015 total aligned coverage across merged hits, a best single-hit length not exceeding 80% of the
1016 gene length, and multiple hits separated by at least 1 kb in the target genome. Candidates
1017 passing these criteria were treated as putative deletion-born fusion genes. Domain analysis
1018 was done as before for the *yjcO-lysU* fusion.

1019 To place candidate events in an evolutionary context, we constructed a core-genome alignment
1020 and phylogeny using Parsnp⁸³ with H37Rv as the reference, using the set of *M. tb* and *M. bovis*
1021 assemblies analyzed above. The resulting phylogeny was visualized in iTOL⁸⁴ and BLAST
1022 results were visualized using DNA Features Viewer.⁷⁵

1023 **Prefix-Suffix K-mer Screen**

1024 We extracted prefix and suffix k-mers from each query gene using coding nucleotide FASTA
1025 inputs. For each gene sequence, we took a k-mer of length k ($k = 27$) from near the 5' end and
1026 a second k-mer of length k from near the 3' end, excluding a small buffer region from each
1027 terminus to avoid start and stop codons and to shift the k-mers out of frame relative to the
1028 annotated coding sequence. Specifically, we used a fixed gap distance $g = 4$ bp: the prefix k-
1029 mer was taken from positions g to $g+k$ (4 to 31), and the suffix k-mer from positions $-(g+k)$ to
1030 $-g$.

1031 We queried these k-mers against the AllTheBacteria collection of 2.4 million bacterial isolate
1032 assemblies. Assemblies were processed in batches according to their Miniphy'd⁸⁵ output:
1033 genome FASTA files were stored in compressed tar.xz archives, and for each archive we
1034 concatenated subsets of 500 genomes into temporary multi-FASTA files and built BWA indices
1035 on these batches. Indexing was performed with `bwa index`, producing FM-indices for each
1036 genome batch.

1037 Exact k-mer placements were then obtained using `bwa fastmap`, run with the query k-mer
1038 FASTA and a matching k-mer length equal to k (the same value used in extraction). We used a
1039 large maximum hit window (`-w 99999`) to retain all exact match locations reported by fastmap.
1040 Fastmap outputs were gzip-compressed and parsed to recover, for each genome, all contig-level
1041 match positions for both prefix and suffix k-mers.

1042 For each genome and each query gene, we computed a prefix-suffix distance only when at least
1043 one prefix match and one suffix match occurred on the same contig. Distances were computed
1044 from the genomic coordinates of the matched k-mers on that contig. Matches split across
1045 contigs were ignored, and genomes lacking a same-contig prefix-suffix pair were treated as
1046 missing for that gene.

1047 **Multimodal Distance Detection and Clustering**

1048 We identified structurally variable genes by clustering the per-genome prefix-suffix distances
1049 for each query gene and testing whether the resulting distance distribution was multimodal.
1050 For each gene, we aggregated the set of observed “Difference” values across genomes (the
1051 prefix-suffix distance expressed relative to the expected gene length) along with their
1052 multiplicities, and clustered these one-dimensional values using a density-based algorithm
1053 (DBSCAN1D).⁸⁶ Clustering was run with epsilon = 800 (bp) and min_samples = 25, and we
1054 treated DBSCAN noise points as outliers (removed from downstream summaries). A gene was
1055 called “multimodal” only if DBSCAN identified at least two non-noise clusters, corresponding to
1056 two tight peaks in the distance distribution.

1057 To enrich specifically for candidates consistent with deletion-born fusion genes, we applied
1058 additional filters to multimodal loci. We required one cluster centered near 0 (the intact allele,
1059 where the observed prefix-suffix distance matches the reference gene length) and at least one
1060 additional cluster centered at a positive distance greater than 800 bp, consistent with a split
1061 ancestral state in which the prefix and suffix k-mers are separated by a substantial intervening
1062 segment. Genes with only small positive shifts (e.g., internal deletions) or without an intact-like
1063 cluster near 0 were excluded at this stage. The set of genes passing these clustering and
1064 distance-peak criteria was carried forward for phylogenetic filtering and downstream analyses.

1065 **Phylogenetic Reconstruction and Ancestral State Inference**

1066 To distinguish deletion-born fusions from gene disruption by insertion (**Figure 4B**, filter 3), we
1067 performed phylogenetic reconstruction and ancestral state inference using the DBSCAN cluster
1068 assignments from the prefix-suffix distance analysis. For each candidate gene passing the
1069 multimodality and peak-shape filters, we downsampled genomes to obtain a tractable but
1070 representative set for tree building by sampling equal numbers of genomes from each distance
1071 cluster (300 total genomes per gene, split evenly across clusters). Sampling was restricted to a
1072 high-quality genome set from the ATB to reduce artifacts from fragmented assemblies. For each
1073 sampled genome, we extracted the identified prefix and suffix k-mer, as well as the entire
1074 interleaving sequence to a new FASTA file and masked that same region in the whole genome
1075 sequence.

1076 Trees were constructed from masked sequences using attotree⁸⁷ (an optimized version of
1077 Mashtree⁸⁸) with default options on the masked whole genome sequences. To enable rooting,
1078 we also added two outgroup genomes per focal species, chosen as close relatives outside the
1079 focal clade: for *E. coli* K-12, *Klebsiella pneumoniae* MGH78578 (GCF_000016305.1) and
1080 *Salmonella enterica* serovar Typhimurium LT2 (GCF_000006945.2); for *M. tb.* H37Rv, *M.*
1081 *canettii* (NC_015848.1) and *M. caprae* (CP016401.1); for *N. gonorrhoeae* FA1090, *N. lactamica*
1082 (NC_014752.1) and *N. meningitidis* MC58 (NC_003112.2); for *C. jejuni* NCTC1168, *C. coli*
1083 (NC_022660.1) and *C. lari* (NC_012039.1); and for *S. pneumoniae* TIGR4, *S. mitis* (FN568063.1)
1084 and *S. oralis* (FR720602.1). For each of the two possible outgroups included, we identified the
1085 genome with the largest mean distance between it and the other leaves and chose that as the
1086 outgroup to root the resulting tree at.

1087 For each candidate gene passing the distance-based filters, we inferred whether the ancestral
1088 state was “split” or “intact” using a tree-based parsimony approach. Cluster labels were treated
1089 as discrete character states, and we reconstructed internal node states by multi-state Fitch
1090 parsimony.⁸⁹ The ancestral state for each gene was defined as the parsimony assignment at the
1091 most recent common ancestor (MRCA) of the ingroup genomes. Genes were classified as
1092 consistent with deletion-born fusions when the MRCA state corresponded to a “split” cluster (the
1093 cluster with a positive prefix-suffix difference) and at least one descendant clade carried an
1094 “intact” cluster (the cluster with mean difference near 0). Conversely, genes whose MRCA state
1095 was “intact” were interpreted as cases where the intact gene was ancestral, and the positive-

1096 distance cluster reflects disruption (often by insertion) and were excluded from the deletion-born
1097 fusion set.

1098

1099 **Mobile Genetic Element and Prophage Detection**

1100 To assess whether structurally variable loci were dominated by insertions of mobile genetic
1101 elements (MGEs) or prophages (**Figure 4C**), we aligned the sequence between the matched
1102 prefix and suffix k-mers to curated MGE and prophage databases. The MGE database was
1103 taken from MGEdb⁹⁰ and the prophage database from Prophage-DB⁹¹ (bacterial host
1104 prophages); the two FASTA sets were concatenated into a single nucleotide BLAST database.
1105 For each genome sampled, we extracted the intervening sequence between the matched prefix
1106 and suffix k-mers on the same contig (i.e., the sequence whose length drives the positive prefix-
1107 suffix distance signal) and queried it against the combined database using BLASTN. For each
1108 intervening sequence, we merged overlapping BLAST hit intervals along the query and
1109 computed the total number of query bases covered by any hit; the fraction of the intervening
1110 sequence explained by MGEs/prophages was then calculated as covered bases divided by
1111 query length. These per-genome fractions were summarized by locus and compared between
1112 loci whose inferred MRCA state was “split” versus “intact” using a two-sided Mann–Whitney U
1113 test.

1114 **Selection Analysis**

1115 For the selection analysis on *miaE-htpX* and *acrR-glcD*, genomes with the intact ORF were
1116 identified based on their membership to the cluster with a relative prefix-suffix “Distance” of 0.
1117 The prefix-suffix k-mer and intervening nucleotides were extracted, deduplicated, and codon-
1118 aligned using MACSE with default options.⁹² A phylogeny of the masked whole genome
1119 sequences was built using Parsnp and selection analysis was performed with FUBAR,
1120 implemented in the HyPhy package with default options.⁹³

1121 We assessed signatures of selection on candidate deletion-born fusions using two
1122 complementary dN/dS-style comparisons (**Figure 4E**). First, to estimate selection acting on the
1123 observed “intact” allele, we focused on genomes assigned to the cluster closest to zero
1124 difference (the intact-like cluster). For each genome we used the extracted locus sequence and
1125 deduplicated identical sequences by hashing, retaining a count of how many genomes shared
1126 each unique sequence. Each unique sequence was then codon-aligned to the canonical query
1127 CDS using MACSE,⁹² and we counted synonymous and nonsynonymous differences across
1128 aligned codons, ignoring codons overlapping gaps or ambiguous bases and tracking premature
1129 stop gains separately. Per-cluster estimates were computed by summing synonymous and
1130 nonsynonymous counts across unique sequences and weighting each sequence by the number
1131 of genomes in which it occurred; dN/dS was then calculated as weighted nonsynonymous
1132 divided by weighted synonymous changes.

1133 Second, to estimate the degree of divergence expected from the pre-deletion (“split”) state, we
1134 constructed “surrogate” genes for genomes assigned to the inferred ancestral cluster (the MRCA
1135 cluster from parsimony). For each genome we used the extracted intervening locus sequence
1136 and flanking sequence, deduplicated both sequence sets, and mapped each unique removed-
1137 region sequence into its parent flank sequence by exact substring matching (allowing reverse
1138 complement). We then aligned the canonical query CDS to each unique removed-region
1139 sequence using BLASTN and selected a set of high-scoring segment pairs that approximately
1140 tiled the query. Each BLAST segment was projected back onto the parent flank coordinates and
1141 intersected with ORFs predicted on the flank sequence using pyrodigal; for each segment we
1142 selected the ORF with the greatest overlap and stitched the resulting ORF nucleotide sequences
1143 together in query order, joining adjacent blocks with “NNN” to preserve codon-phase ambiguity.
1144 These stitched constructs were treated as surrogate pre-deletion sequences and were codon-

1145 aligned to the canonical query CDS again using MACSE, after which synonymous and
1146 nonsynonymous differences were counted as above and aggregated into a weighted dN/dS
1147 estimate.

1148 ***Protein Family Clustering and Sampling***

1149 We constructed a large-scale bacterial protein family catalogue from complete RefSeq genomes
1150 to quantify how sampling depth affects detection of deletion-born fusion genes (**Figure 5**). All
1151 complete bacterial genomes available in RefSeq at the time of analysis (54,630 assemblies) were
1152 downloaded using ncbi-genome-download in GenBank format. From each genome, we
1153 extracted all annotated protein-coding sequences by parsing the GenBank feature tables. For
1154 each CDS, we extracted both the nucleotide sequence and the corresponding protein sequence.
1155 When a translation was provided in the GenBank record it was used directly; otherwise, the
1156 nucleotide sequence was translated in-frame. Protein sequences were written to per-genome
1157 FASTA files and then concatenated into a single combined protein FASTA for clustering.

1158 Protein sequences were clustered into gene families using MMseqs2 Linclust.⁸² Clustering was
1159 performed with a minimum pairwise sequence identity of 80%, a minimum target coverage of
1160 80% (coverage mode 1), and 80 k-mers per sequence, using mmseqs easy-linclust. This
1161 procedure yielded 23,126,961 protein families, spanning both singleton and multi-member
1162 families. The resulting cluster table and representative sequences were used for downstream
1163 sampling and analysis.

1164 To evaluate how database sampling affects detection of structural variation, we defined three
1165 complementary protein-family sampling strategies. In the “All Families” strategy, we sampled
1166 protein families uniformly at random from the full set of MMseqs2 clusters, including
1167 singletons. In the “Non-singletons” strategy, we excluded singleton families and sampled
1168 uniformly from families containing at least two members. In the “Large Families” strategy, we
1169 restricted sampling to families with more than 20 members, enriching for deeply sampled
1170 lineages. For each strategy, we analyzed an equal number of protein families (100,000).

1171 For each protein family, we quantified sampling depth using a sampling depth score defined as
1172 the total number of genomes containing members of that family divided by the number of
1173 unique species represented among those genomes. A score near 1 indicates a family sampled
1174 broadly but shallowly across species, whereas higher values indicate repeated sampling of the
1175 same species (e.g., many isolates of a single lineage). Sampling depth scores were used to
1176 compare the effective detectability of deletion-born fusions across the three sampling
1177 strategies.