

1 **Title: Mutational spectra analysis reveals bacterial niche and transmission routes**

2
3 **Authors:** Christopher Ruis^{1,2,3}, Aaron Weimann^{1,2,3}, Gerry Tonkin-Hill⁴, Arun Prasad
4 Pandurangan⁵, Marta Matuszewska^{2,6}, Gemma G. R. Murray⁷, Roger C. Lévesque⁸, Tom L.
5 Blundell⁵, R. Andres Floto^{1,3,9*†}, Julian Parkhill^{2*†}

6 **Affiliations:**

7 ¹Molecular Immunity Unit, University of Cambridge Department of Medicine, MRC-
8 Laboratory of Molecular Biology; Cambridge, UK.

9 ²Department of Veterinary Medicine, University of Cambridge; Cambridge, UK.

10 ³Cambridge Centre for AI in Medicine, University of Cambridge; Cambridge, UK.

11 ⁴Department of Biostatistics, University of Oslo; Blindern, Norway.

12 ⁵Department of Biochemistry, Sanger Building, University of Cambridge; Cambridge, UK.

13 ⁶Department of Medicine, University of Cambridge; Cambridge, UK.

14 ⁷Parasites and Microbes Programme, Wellcome Sanger Institute; Wellcome Genome
15 Campus, Cambridge, UK.

16 ⁸Institut de biologie intégrative et des systèmes (IBIS), Université Laval; Québec City,
17 Québec, Canada.

18 ⁹Cambridge Centre for Lung Infection, Papworth Hospital; Cambridge, UK.

19 † These Authors contributed equally to this work

20 * Corresponding authors. Email: arf27@cam.ac.uk, jp369@cam.ac.uk

21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

49 **Abstract:**

50 As observed in cancers, individual mutagens and defects in DNA repair create distinctive
51 mutational signatures that combine to form context-specific spectra within cells. We reasoned
52 that similar processes must occur in bacterial lineages, potentially allowing decomposition
53 analysis to identify disrupted DNA repair processes and niche-specific mutagen exposure.
54 Here we reconstructed mutational spectra for 84 clades from 31 diverse bacterial species,
55 assigned signatures to specific DNA repair pathways using hypermutator lineages, and, by
56 comparing mutational spectra of clades from different environmental and biological locations,
57 extracted reproducible niche-associated mutational signatures. We show that mutational
58 spectra can predict general and specific bacterial niches and therefore reveal the site of
59 infection and types of transmission routes for established and emergent human bacterial
60 pathogens.

61
62 **One sentence summary:** Variable mutagen exposure and DNA repair drive differential
63 mutational spectra between bacteria and enable niche inference

64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96

97 **Main text:**

98 Work using human cells and tissues has demonstrated that mutagens induce highly specific
99 context-dependent patterns of base substitutions termed mutational signatures, which
100 combine to form a mutational spectrum (1-6). However, these patterns are compounded with
101 the signatures of endogenous mutations and DNA repair, which also exhibit specific mutational
102 signatures (7-9).

103
104 Reconstructing the set of mutations and signatures within cancers has enabled inference of
105 the drivers of tumourigenesis (1, 2, 7). We therefore reasoned that reconstructing mutational
106 spectra in bacteria, differentiating them into different signatures, and correlating these with
107 known DNA repair defects and environmental exposures, should allow the association of
108 specific DNA signatures with bacterial niches. These signatures could then be used to predict
109 niche or infection sites and to identify defects in DNA repair when niche is known. To test this,
110 we undertook the first large-scale comparison of mutational spectra and their underlying
111 signatures across bacteria, correlating the results with DNA repair pathways and niche.

112
113 We used whole genome sequence alignments and phylogenetic trees to reconstruct single
114 base substitution (SBS) mutational spectra of 84 phylogenetic clades from 31 diverse bacterial
115 species, implemented in a specifically-developed open-source bioinformatic tool, MutTui (**fig.**
116 **S1; fig. S2; table S1; table S2; Supplementary Methods**). SBS spectra were rescaled by
117 genomic nucleotide composition to enable direct comparison between bacteria. We find that
118 such spectra are highly diverse, both in the nucleotide mutations themselves and their
119 surrounding context (**Fig. 1; fig. S2**). However, several generalisable properties could be
120 identified. We found that transition mutations are more common than transversion mutations
121 (10) in all cases (ranging from 52-55% in *Klebsiella pneumoniae* to >90% in *Campylobacter*
122 *jejuni*; **fig. S2**). Cytosine to thymine (C>T) was typically the most common mutation type
123 identified (in 69 of 84 SBS spectra examined), potentially due to cytosine deamination (11).
124 Genomic G+C content exhibits a negative correlation with proportion of C>A/T mutations but
125 a positive correlation with proportion of C>G mutations (**fig. S3**). Finally, transition mutations
126 exhibit enriched context specificity compared to transversion mutations while several
127 contextual mutations are significantly elevated across datasets (**fig. S4**). UMAP clustering
128 revealed groups of similar SBS spectra across bacterial clades (**Fig. 1**). We observe a strong
129 correlation between phylogenetic relatedness and spectrum similarity (Tukey HSD corrected
130 ANOVA $P < 0.001$; **fig. S5**) and spectra are typically conserved across highly-related clades
131 where there has likely been no change of niche or DNA repair capacity (**fig. S6**).

132
133 We reasoned that bacterial SBS spectra can be decomposed into combinations of mutational
134 signatures, each driven by distinct defects in DNA repair or by endogenous processes or
135 specific mutagens, as has previously been achieved for cancer-associated mutations (1, 3, 7).
136 Therefore, we first extracted mutational signatures associated with distinct DNA repair
137 pathways by calculating SBS spectra of 50 naturally-occurring hypermutator lineages across
138 four bacterial species (**Fig. 2A**). By identifying the genes most likely responsible for
139 hypermutation, we were able to attribute mutational signatures to defects in 11 DNA repair
140 genes that function in mismatch repair (MMR), base excision repair (BER), or homologous
141 recombination (HR) (**Fig. 2B-D; Supplementary Methods**).

142

143 Mutations of MMR genes result in high levels of context-specific C>T and T>C mutations (**Fig.**
144 **2B, fig. S7**) (1, 8, 12) which likely represent the error profile of DNA Polymerase III that is
145 usually repaired by functional MMR. While context specificity is highly similar between species,
146 the relative rates of C>T and T>C differ between *Pseudomonas aeruginosa* and *Burkholderia*
147 *cenocepacia* (**Fig. 2B; fig. S8**), likely reflecting distinct polymerase error profiles (a possibility
148 supported by structural modelling analysis; **fig. S9**).

149
150 Mutations in distinct base excision repair (BER) components results in characteristic gene-
151 specific patterns (**Fig. 2C**), as expected from the diverse repair functions of proteins within this
152 pathway (11). We identified *P. aeruginosa* hypermutators for each component of the GO repair
153 pathway (*mutT*, *mutY* and *mutM*) that prevents 8-oxoguanine (8-oxo-G)-induced mutations
154 (13). Mutation of *mutT*, whose product degrades 8-oxo-G monomers to prevent their
155 incorporation into DNA (13), results in non-specific T>G mutations (**Fig. 2C**), suggesting
156 incorporation of 8-oxo-G opposite adenine is context-independent. Conversely, mutation of
157 *mutY* which excises adenine opposite 8-oxo-G (13), results in C>A mutations predominantly
158 in CpCpN and TpCpN contexts (**Fig. 2C**), indicating context-specific mutation of incorporated
159 guanine to 8-oxo-G. This likely represents the pattern of reactive oxygen species (ROS)
160 damage, of which 8-oxo-G is a major mutagenic lesion (9). The C>A contexts differ between
161 the *P. aeruginosa mutY* signature and human cell signatures of ROS exposure (5) and
162 knockout of either the *mutY* homologue or OGG1 (7, 9) (**fig. S10**), suggesting differential repair
163 of these lesions by other proteins. Mutation of *mutM* results in C>G mutations in ApCpN
164 contexts (**Fig. 2C**). While the mechanism of C>G mutations is unclear, the lack of C>A
165 mutations in *mutM* knockouts is potentially due to functional MutY being sufficient to repair
166 mutagenic 8-oxo-G lesions (14). We additionally identify PA4172 in *P. aeruginosa* whose
167 knockout exhibits C>A mutations in CpCpN and TpCpN contexts similar to *mutY* (Pearson's *r*
168 $P < 0.001$; **Fig. 2C; fig. S10**), suggesting that its product may similarly repair mutagenic 8-
169 oxo-G lesions.

170
171 Disruption of *ung*, whose product removes uracil from DNA (11), results in similar patterns of
172 context-specific C>T mutations in *P. aeruginosa* and *Mycobacterium abscessus* (Pearson's *r*
173 $P < 0.001$; **Fig. 2C; fig. S11**). This bacterial signature exhibits subtle contextual differences
174 compared with *ung* knockout in human cells (9), particularly through enriched mutations in
175 NpCpG contexts (**fig. S11**), suggesting differential patterns of uracil incorporation in humans
176 and bacteria.

177
178 Mutation of *nth*, whose product Endonuclease III removes damaged pyrimidines, results in
179 C>T mutations in multiple *Mycobacteria* species and human cells (8) but with different context
180 specificity (Pearson's *r* $P > 0.05$; **Fig. 2C; fig. S12**). Disruption of the apurinic-apyrimidinic
181 (AP) endonuclease *xthA* results in mutations in multiple specific contexts (**Fig. 2C**), particularly
182 transversions in [C,G,T]p[C,T]pG contexts, indicating repair of a broad range of specific
183 lesions. Finally, hypermutators resulting from mutation of the homologous recombination
184 pathway components *recF* and *recN* exhibit context-specific transition mutations (**Fig. 2D**).
185 Recombination is known to drive GC-biased gene conversion (15) and this may contribute to
186 this signature.

187
188 We subsequently tested whether we could detect differences in SBS spectra from bacteria
189 with different repair capabilities occupying a similar niche (and therefore exposed to similar
190 sets of niche-specific mutagens). Decomposition analysis showed that almost all mutations

191 elevated in *C. jejuni* compared with the gastrointestinal *Escherichia coli* lineage 34 (16) can
192 be explained by a failure to repair deaminated cytosines and a lack of MMR (**Fig. 2E**; **fig.**
193 **S13**); pathways which are known to be absent in *C. jejuni* (17). Our results indicate that
194 differences in DNA repair can be inferred by comparing bacteria from a similar niche.

195
196 We then proceeded to extract further bacterial signatures through a decomposition analysis
197 employing nonnegative matrix factorisation (NMF) (18, 19) on SBS spectra datasets from a
198 range of species and genera (**table S3**). We extracted 33 SBS signatures and collapsed these
199 into a final set of 24 (named with the prefix Bacteria_SBS) by combining highly similar
200 signatures (with cosine similarity of 0.95 or greater) (**fig. S14**; **table S4**). The extracted
201 signatures exhibit divergent base mutations and contexts and are differentially present across
202 bacteria (**fig. S14**), supporting differential activity of mutagens and repair between clades. An
203 exception to this pattern was signature Bacteria_SBS15 that was extracted from the
204 *Staphylococcus* genus dataset and the *Enterococcus faecalis*, *Streptococcus pneumoniae*
205 and *Streptococcus agalactiae* species datasets (**fig. S14**), indicating broad distribution across
206 Bacillota. As these bacteria inhabit different niches, this signature likely represents phylum-
207 specific endogenous mutations and/or DNA repair profiles.

208
209 We next explored the influence of pathogen niche on mutational spectrum, focussing on
210 *Mycobacteria* and *Burkholderia*, genera that contain both clades that are transmitted from
211 person-to-person and clades that are acquired from environmental sources (20, 21). We find
212 that known lung and environmental clades cluster separately based on SBS spectrum
213 composition (**Fig. 3A**). Spectrum subtractions consistently revealed elevated C>A and C>T
214 mutations in lung bacteria and higher levels of T>C in environmental bacteria (**Fig. 3B, C**).
215 Lung and environmental bacteria additionally exhibit different contextual patterns within C>A
216 and T>C mutations (**fig. S15**). Decomposition of niche-specific mutations from subtracted
217 spectra using known human mutagen signatures suggests that higher C>A in lung bacteria is
218 likely driven by tobacco smoke (2) (found in human but not animal infecting lung clades) and
219 exposure to reactive oxygen species (ROS), while higher T>C within the environment is
220 probably caused by exposure to alkylating agents and nitro-polycyclic aromatic hydrocarbons
221 (5) (**Fig. 3D**), mutagens known to be present in the environment (5). It is also possible that the
222 long-term evolutionary selection towards GC richness seen in some bacterial genomes (22)
223 may contribute to the observed environmental signature.

224
225 We further examined niche signatures through a targeted NMF decomposition of the
226 *Mycobacteria* and *Burkholderia* spectra and were able to extract a lung-associated mutational
227 signature consisting of multiple mutation types that we term Bacteria_Lung1 (**Fig. 3E, F**).
228 Building on these different patterns, we developed a set of leave-one-out classifiers and found
229 that they were able to robustly predict known lung or environmental niche based on either SBS
230 spectrum, proportion of the six mutation types or cosine similarity between SBS spectra (**table**
231 **S5**; **fig. S16**).

232
233 Due to their success in predicting known niches, we next used SBS spectra to infer niche for
234 several *Mycobacteria* clades where this was previously unknown. We find strong evidence
235 that *Mycobacterium leprae* and the dominant circulating clones (DCCs) of *M. abscessus* (23)
236 replicate within the lung as they: cluster with known human lung bacteria based on their SBS
237 spectrum (**Fig. 3A**); exhibit lung-like contextual patterns of C>A and T>C mutations (**fig. S15**);
238 exhibit high levels of C>A and low levels of T>C (**Fig. 3B, C**); and exhibit signature

239 Bacteria_Lung1 at similar levels to known lung bacteria (**Fig. 3F**). These observations suggest
240 human-to-human lung transmission and are supported by reports that *M. leprae* can replicate
241 within human alveolar epithelial cells *in vitro* and can infect mouse lung macrophages and
242 epithelial cells during *in vivo* challenge (24), and that the *M. abscessus* DCCs have spread
243 through global transmission chains involving Cystic Fibrosis (CF) and non-CF individuals (20,
244 23).

245
246 The *Mycobacterium kansasii* main cluster (MKMC) causes the majority of *M. kansasii*
247 infections (25) and exhibits characteristics of both lung and environmental spectra.
248 Specifically, the MKMC exhibits lung-like C>A patterns but environmental-like T>C patterns
249 (**fig. S15**) and is therefore intermediate between known human lung and environmental
250 spectra in the SBS clustering (**Fig. 3A**) and C>A vs T>C comparison (**Fig. 3B**). Together,
251 these results suggest that the MKMC is exposed to both lung and environmental mutagens
252 and therefore replicates within (and is potentially acquired from) both niches.

253
254 We next examined multiple niches within the same host by comparing human *Salmonella*
255 lineages that cause enteric infection with those that have adapted to cause invasive disease
256 (**table S1 & S6**). The SBS spectra cluster by niche (**Fig. 4A**; association index $P < 0.001$),
257 rather than by phylogeny. While we could not identify clear and conserved differences between
258 SBS spectra in the different niches by eye (**fig. S17**), we were again able to develop classifiers
259 that could robustly predict enteric or invasive niche (**table S5; fig. S16**), suggesting that even
260 subtle spectrum differences can be sufficient to reliably distinguish niche.

261
262 Finally, we tested whether mutational spectra could distinguish sub-niches within the same
263 host niche. We find a high level of CC>TT double mutations characteristic of UV-light damage
264 (5) in the pan-skin bacterium *Cutibacterium acnes* that is not present in *Staphylococcus*
265 *epidermidis* which preferentially inhabits moist (26), and therefore less sun-exposed, skin sites
266 (**Fig. 4B**). Together, these results demonstrate that mutational spectra can predict bacterial
267 niche with very high levels of spatial resolution.

268
269 In conclusion, we show that we can reconstruct mutational spectra from bacterial phylogenies
270 and decompose these into specific signatures. We can ascribe some of these signatures to
271 defects in DNA repair pathways, and others to exposure to location-dependent mutagens
272 which can be used to predict the niche in which bacteria replicate and infer transmission
273 routes. We anticipate that identification of signatures at different levels in bacterial phylogenies
274 will identify ancestral niches and therefore sources of emergent human pathogens, reveal
275 routes of acquisition of infection permitting targeted interventions, and provide a mechanism
276 to monitor pathogenic evolution and host adaptation. We envisage that mutational spectra
277 analysis could be applied to viruses and parasites, enabling similar predictions.

278
279
280
281
282
283
284
285
286

287 **References and Notes:**

- 288 1. L. B. Alexandrov, S. Nik-Zainal, D. C. Wedge, S. A. J. R. Aparicio, S. Behjati, A. V.
289 Biankin, G. R. Bignell, N. Bolli, A. Borg, A.-L. Børresen-Dale, S. Boyault, B. Burkhardt, A. P.
290 Butler, C. Caldas, H. R. Davies, C. Desmedt, R. Eils, J. E. Eyfjörd, J. A. Foekens, M. Greaves, F.
291 Hosoda, B. Hutter, T. Ilcic, S. Imbeaud, M. Imielinski, N. Jäger, D. T. W. Jones, D. Jones, S.
292 Knappskog, M. Kool, S. R. Lakhani, C. López-Otín, S. Martin, N. C. Munshi, H. Nakamura, P. A.
293 Northcott, M. Pajic, E. Papaemmanuil, A. Paradiso, J. V. Pearson, X. S. Puente, K. Raine, M.
294 Ramakrishna, A. L. Richardson, J. Richter, P. Rosenstiel, M. Schlesner, T. N. Schumacher, P.
295 N. Span, J. W. Teague, Y. Totoki, A. N. J. Tutt, R. Valdés-Mas, M. M. van Buuren, L. van 't
296 Veer, A. Vincent-Salomon, N. Waddell, L. R. Yates, J. Zucman-Rossi, P. A. Futreal, U.
297 McDermott, P. Lichter, M. Meyerson, S. M. Grimmond, R. Siebert, E. Campo, T. Shibata, S.
298 M. Pfister, P. J. Campbell, M. R. Stratton, Signatures of mutational processes in human
299 cancer. *Nature*. **500**, 415–421 (2013).
- 300 2. L. B. Alexandrov, Y. S. Ju, K. Haase, P. V. Loo, I. Martincorena, S. Nik-Zainal, Y. Totoki,
301 A. Fujimoto, H. Nakagawa, T. Shibata, P. J. Campbell, P. Vineis, D. H. Phillips, M. R. Stratton,
302 Mutational signatures associated with tobacco smoking in human cancer. *Science*. **354**, 618–
303 622 (2016).
- 304 3. L. B. Alexandrov, J. Kim, N. J. Haradhvala, M. N. Huang, A. W. Tian Ng, Y. Wu, A. Boot,
305 K. R. Covington, D. A. Gordenin, E. N. Bergstrom, S. M. A. Islam, N. Lopez-Bigas, L. J.
306 Klimczak, J. R. McPherson, S. Morganella, R. Sabarinathan, D. A. Wheeler, V. Mustonen, G.
307 Getz, S. G. Rozen, M. R. Stratton, The repertoire of mutational signatures in human cancer.
308 *Nature*. **578**, 94–101 (2020).
- 309 4. A. Degasperi, T. D. Amarante, J. Czarnecki, S. Shooter, X. Zou, D. Glodzik, S.
310 Morganella, A. S. Nanda, C. Badja, G. Koh, S. E. Momen, I. Georgakopoulos-Soares, J. M. L.
311 Dias, J. Young, Y. Memari, H. Davies, S. Nik-Zainal, A practical framework and online tool for
312 mutational signature analyses show intertissue variation and driver dependencies. *Nat.*
313 *Cancer*. **1**, 249–263 (2020).
- 314 5. J. E. Kucab, X. Zou, S. Morganella, M. Joel, A. S. Nanda, E. Nagy, C. Gomez, A.
315 Degasperi, R. Harris, S. P. Jackson, V. M. Arlt, D. H. Phillips, S. Nik-Zainal, A Compendium of
316 Mutational Signatures of Environmental Agents. *Cell*. **177**, 821-836.e16 (2019).
- 317 6. S. Nik-Zainal, J. E. Kucab, S. Morganella, D. Glodzik, L. B. Alexandrov, V. M. Arlt, A.
318 Weninger, M. Hollstein, M. R. Stratton, D. H. Phillips, The genome as a record of
319 environmental exposure. *Mutagenesis*. **30**, 763–770 (2015).
- 320 7. A. Degasperi, X. Zou, T. Dias Amarante, A. Martinez-Martinez, G. C. C. Koh, J. M. L.
321 Dias, L. Heskin, L. Chmelova, G. Rinaldi, V. Y. W. Wang, A. S. Nanda, A. Bernstein, S. E.
322 Momen, J. Young, D. Perez-Gil, Y. Memari, C. Badja, S. Shooter, J. Czarnecki, M. A. Brown, H.
323 R. Davies, Genomics England Research Consortium, S. Nik-Zainal, Substitution mutational
324 signatures in whole-genome–sequenced cancers in the UK population. *Science*. **376**,
325 abl9283.
- 326 8. J. Drost, R. van Boxtel, F. Blokzijl, T. Mizutani, N. Sasaki, V. Sasselli, J. de Ligt, S.
327 Behjati, J. E. Grolleman, T. van Wezel, S. Nik-Zainal, R. P. Kuiper, E. Cuppen, H. Clevers, Use

- 328 of CRISPR-modified human stem cell organoids to study the origin of mutational signatures
329 in cancer. *Science*. **358**, 234–238 (2017).
- 330 9. X. Zou, G. C. C. Koh, A. S. Nanda, A. Degasperi, K. Urgo, T. I. Roumeliotis, C. A. Agu, C.
331 Badja, S. Momen, J. Young, T. D. Amarante, L. Side, G. Brice, V. Perez-Alonso, D. Rueda, C.
332 Gomez, W. Bushell, R. Harris, J. S. Choudhary, J. Jiricny, W. C. Skarnes, S. Nik-Zainal, A
333 systematic CRISPR screen defines mutational mechanisms underpinning signatures caused
334 by replication errors and endogenous DNA damage. *Nat. Cancer*. **2**, 643–657 (2021).
- 335 10. H. Lee, E. Popodi, H. Tang, P. L. Foster, Rate and molecular spectrum of spontaneous
336 mutations in the bacterium *Escherichia coli* as determined by whole-genome sequencing.
337 *Proc. Natl. Acad. Sci.* **109**, E2774–E2783 (2012).
- 338 11. K. J. Wozniak, L. A. Simmons, Bacterial DNA excision repair pathways. *Nat. Rev.*
339 *Microbiol.*, 1–13 (2022).
- 340 12. B. Meier, N. V. Volkova, Y. Hong, P. Schofield, P. J. Campbell, M. Gerstung, A.
341 Gartner, Mutational signatures of DNA mismatch repair deficiency in *C. elegans* and human
342 cancers. *Genome Res.* **28**, 666–675 (2018).
- 343 13. A. V. Endutkin, D. O. Zharkov, GO System, a DNA Repair Pathway to Cope with
344 Oxidative Damage. *Mol. Biol.* **55**, 193–210 (2021).
- 345 14. L. H. Sanders, J. Sudhakaran, M. D. Sutton, The GO System Prevents ROS-Induced
346 Mutagenesis and Killing in *Pseudomonas aeruginosa*. *FEMS Microbiol. Lett.* **294**, 89–96
347 (2009).
- 348 15. F. Lassalle, S. Périan, T. Bataillon, X. Nesme, L. Duret, V. Daubin, GC-Content
349 Evolution in Bacterial Genomes: The Biased Gene Conversion Hypothesis Expands. *PLOS*
350 *Genet.* **11**, e1004941 (2015).
- 351 16. G. Horesh, G. A. Blackwell, G. Tonkin-Hill, J. Corander, E. Heinz, N. R. Y. 2021
352 Thomson, A comprehensive and high-quality collection of *Escherichia coli* genomes and
353 their genes. *Microb. Genomics.* **7**, 000499.
- 354 17. J. Parkhill, B. W. Wren, K. Mungall, J. M. Ketley, C. Churcher, D. Basham, T.
355 Chillingworth, R. M. Davies, T. Feltwell, S. Holroyd, K. Jagels, A. V. Karlyshev, S. Moule, M. J.
356 Pallen, C. W. Penn, M. A. Quail, M.-A. Rajandream, K. M. Rutherford, A. H. M. van Vliet, S.
357 Whitehead, B. G. Barrell, The genome sequence of the food-borne pathogen *Campylobacter*
358 *jejuni* reveals hypervariable sequences. *Nature*. **403**, 665–668 (2000).
- 359 18. L. B. Alexandrov, S. Nik-Zainal, D. C. Wedge, P. J. Campbell, M. R. Stratton,
360 Deciphering Signatures of Mutational Processes Operative in Human Cancer. *Cell Rep.* **3**,
361 246–259 (2013).
- 362 19. S. M. A. Islam, Y. Wu, M. Díaz-Gay, E. N. Bergstrom, Y. He, M. Barnes, M. Vella, J.
363 Wang, J. W. Teague, P. Clapham, S. Moody, S. Senkin, Y. R. Li, L. Riva, T. Zhang, A. J. Gruber,
364 R. Vangara, C. D. Steele, B. Otlu, A. Khandekar, A. Abbasi, L. Humphreys, N. Syulyukina, S. W.
365 Brady, B. S. Alexandrov, N. Pillay, J. Zhang, D. J. Adams, I. Martincorena, D. C. Wedge, M. T.
366 Landi, P. Brennan, M. R. Stratton, S. G. Rozen, L. B. Alexandrov, *bioRxiv*, in press,
367 doi:10.1101/2020.12.13.422570.

- 368 20. J. M. Bryant, D. M. Grogono, D. Rodriguez-Rincon, I. Everall, K. P. Brown, P. Moreno,
369 D. Verma, E. Hill, J. Drijkoningen, P. Gilligan, C. R. Esther, P. G. Noone, O. Giddings, S. C. Bell,
370 R. Thomson, C. E. Wainwright, C. Coulter, S. Pandey, M. E. Wood, R. E. Stockwell, K. A.
371 Ramsay, L. J. Sherrard, T. J. Kidd, N. Jabbour, G. R. Johnson, L. D. Knibbs, L. Morawska, P. D.
372 Sly, A. Jones, D. Bilton, I. Laurenson, M. Ruddy, S. Bourke, I. C. J. W. Bowler, S. J. Chapman,
373 A. Clayton, M. Cullen, O. Dempsey, M. Denton, M. Desai, R. J. Drew, F. Edenborough, J.
374 Evans, J. Folb, T. Daniels, H. Humphrey, B. Isalska, S. Jensen-Fangel, B. Jönsson, A. M. Jones,
375 T. L. Katzenstein, T. Lillebaek, G. MacGregor, S. Mayell, M. Millar, D. Modha, E. F. Nash, C.
376 O'Brien, D. O'Brien, C. Ohri, C. S. Pao, D. Peckham, F. Perrin, A. Perry, T. Pressler, L. Prtak, T.
377 Qvist, A. Robb, H. Rodgers, K. Schaffer, N. Shafi, J. van Ingen, M. Walshaw, D. Watson, N.
378 West, J. Whitehouse, C. S. Haworth, S. R. Harris, D. Ordway, J. Parkhill, R. A. Floto,
379 Emergence and spread of a human-transmissible multidrug-resistant nontuberculous
380 mycobacterium. *Science*. **354**, 751–757 (2016).
- 381 21. C. Chewapreecha, A. E. Mather, S. R. Harris, M. Hunt, M. T. G. Holden, C. Chaichana,
382 V. Wuthiekanun, G. Dougan, N. P. J. Day, D. Limmathurotsakul, J. Parkhill, S. J. Peacock,
383 Genetic variation associated with infection and the environment in the accidental pathogen
384 *Burkholderia pseudomallei*. *Commun. Biol.* **2**, 1–11 (2019).
- 385 22. F. Hildebrand, A. Meyer, A. Eyre-Walker, Evidence of Selection upon Genomic GC-
386 Content in Bacteria. *PLOS Genet.* **6**, e1001107 (2010).
- 387 23. C. Ruis, J. M. Bryant, S. C. Bell, R. Thomson, R. M. Davidson, N. A. Hasan, J. van Ingen,
388 M. Strong, R. A. Floto, J. Parkhill, Dissemination of *Mycobacterium abscessus* via global
389 transmission networks. *Nat. Microbiol.*, 1–10 (2021).
- 390 24. C. A. M. Silva, L. Danelishvili, M. McNamara, M. Berredo-Pinho, R. Bildfell, F. Biet, L.
391 S. Rodrigues, A. V. Oliveira, L. E. Bermudez, M. C. V. Pessolani, Interaction of *Mycobacterium*
392 *leprae* with Human Airway Epithelial Cells: Adherence, Entry, Survival, and Identification of
393 Potential Adhesins by Surface Proteome Analysis. *Infect. Immun.* **81**, 2645–2659 (2013).
- 394 25. T. Luo, P. Xu, Y. Zhang, J. L. Porter, M. Ghanem, Q. Liu, Y. Jiang, J. Li, Q. Miao, B. Hu,
395 B. P. Howden, J. A. M. Fyfe, M. Globan, W. He, P. He, Y. Wang, H. Liu, H. E. Takiff, Y. Zhao, X.
396 Chen, Q. Pan, M. A. Behr, T. P. Stinear, Q. Gao, Population genomics provides insights into
397 the evolution and adaptation to humans of the waterborne pathogen *Mycobacterium*
398 *kansasii*. *Nat. Commun.* **12**, 2491 (2021).
- 399 26. J. Oh, A. L. Byrd, C. Deming, S. Conlan, H. H. Kong, J. A. Segre, Biogeography and
400 individuality shape function in the human skin metagenome. *Nature*. **514**, 59–64 (2014).
- 401 27. N. J. Croucher, A. J. Page, T. R. Connor, A. J. Delaney, J. A. Keane, S. D. Bentley, J.
402 Parkhill, S. R. Harris, Rapid phylogenetic analysis of large samples of recombinant bacterial
403 whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15–e15 (2015).
- 404 28. A. Stamatakis, RAxML version 8: a tool for phylogenetic analysis and post-analysis of
405 large phylogenies. *Bioinformatics*. **30**, 1312–1313 (2014).
- 406 29. P. Sagulenko, V. Puller, R. A. Neher, TreeTime: Maximum-likelihood phylodynamic
407 analysis. *Virus Evol.* **4**, vex042 (2018).

- 408 30. L. McInnes, J. Healy, N. Saul, L. Großberger, UMAP: Uniform Manifold Approximation
409 and Projection. *J. Open Source Softw.* **3**, 861 (2018).
- 410 31. J. M. Bryant, K. P. Brown, S. Burbau, I. Everall, J. M. Belardinelli, D. Rodriguez-
411 Rincon, D. M. Grogono, C. M. Peterson, D. Verma, I. E. Evans, C. Ruis, A. Weimann, D. Arora,
412 S. Malhotra, B. Bannerman, C. Passemar, K. Templeton, G. MacGregor, K. Jiwa, A. J. Fisher, T.
413 L. Blundell, D. J. Ordway, M. Jackson, J. Parkhill, R. A. Floto, Stepwise pathogenic evolution
414 of *Mycobacterium abscessus*. *Science*. **372** (2021), doi:10.1126/science.abb8699.
- 415 32. R. Nikam, A. Kulandaisamy, K. Harini, D. Sharma, M. M. Gromiha, ProThermDB:
416 thermodynamic database for proteins and mutants revisited after 15 years. *Nucleic Acids*
417 *Res.* **49**, D420–D424 (2021).
- 418 33. A. Waterhouse, M. Bertoni, S. Bienert, G. Studer, G. Tauriello, R. Gumienny, F. T.
419 Heer, T. A. P. de Beer, C. Rempfer, L. Bordoli, R. Lepore, T. Schwede, SWISS-MODEL:
420 homology modelling of protein structures and complexes. *Nucleic Acids Res.* **46**, W296–
421 W303 (2018).
- 422 34. J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K.
423 Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, A. Bridgland, C. Meyer, S. A. A. Kohl, A. J.
424 Ballard, A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D.
425 Reiman, E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein,
426 D. Silver, O. Vinyals, A. W. Senior, K. Kavukcuoglu, P. Kohli, D. Hassabis, Highly accurate
427 protein structure prediction with AlphaFold. *Nature*. **596**, 583–589 (2021).
- 428 35. E. F. Pettersen, T. D. Goddard, C. C. Huang, E. C. Meng, G. S. Couch, T. I. Croll, J. H.
429 Morris, T. E. Ferrin, UCSF ChimeraX: Structure visualization for researchers, educators, and
430 developers. *Protein Sci.* **30**, 70–82 (2021).
- 431 36. G. Tonkin-Hill, J. A. Lees, S. D. Bentley, S. D. W. Frost, J. Corander, Fast hierarchical
432 Bayesian analysis of population structure. *Nucleic Acids Res.* **47**, 5539–5549 (2019).

433

434 **Acknowledgements:**

435 The Authors would like to thank all researchers who helped to obtain published datasets
436 used in this study, including Uzma Basit Khan, Christopher Beaudoin, Sophie Belman,
437 Stephen Bentley, Sebastian Bruchmann, Jessica Calland, Claire Chewapreecha, Jukka
438 Corander, Dorota Jamrozy, Anna Kaarina Pöntinen, Noémie Lefrancq, Stephanie Lo, Neil
439 MacAlasdair, Samuel Sheppard, Andries van Tonder and Lucy Weinert.

440

441 **Funding:**

442 Funding for this work was provided by The Wellcome Trust through Investigator awards
443 107032/Z/15/Z (RAF, CR, AW) and 200814/Z/16/Z (TLB, APP), Fondation Botnar
444 (Programme grant 6063; RAF, JP, TLB, CR, AW) and the UK CF Trust (Innovation Hub
445 Award 001; Strategic Research Centre SRC010; CR, AW, TLB, RAF, JP).

446

447 **Author contributions:**

448 Conceptualization: CR, RAF, JP

449 Methodology: CR, GTH

450 Investigation: CR, AW, APP, MM, GGRM, RCL

451 Visualization: CR, APP
452 Funding acquisition: TB, RAF, JP
453 Project administration: RAF, JP
454 Supervision: TB, RAF, JP
455 Writing - original draft: CR, RAF, JP
456 Writing - review & editing: CR, AW, GTH, APP, MM, GGRM, RCL, TB, RAF, JP

457

458 **Competing interests:**

459 Authors declare that they have no competing interests.

460

461 **Data and materials availability:**

462 All data and code used for data analysis is available at

463 https://github.com/chrisruis/Mutational_spectra_data. The MutTui pipeline used to
464 reconstruct pathogen mutational spectra is available at <https://github.com/chrisruis/MutTui>.

465

466 **Supplementary Materials:**

467 Materials and Methods

468 Figs. S1 to S17

469 Tables S1 to S8

470 References (27-36)

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

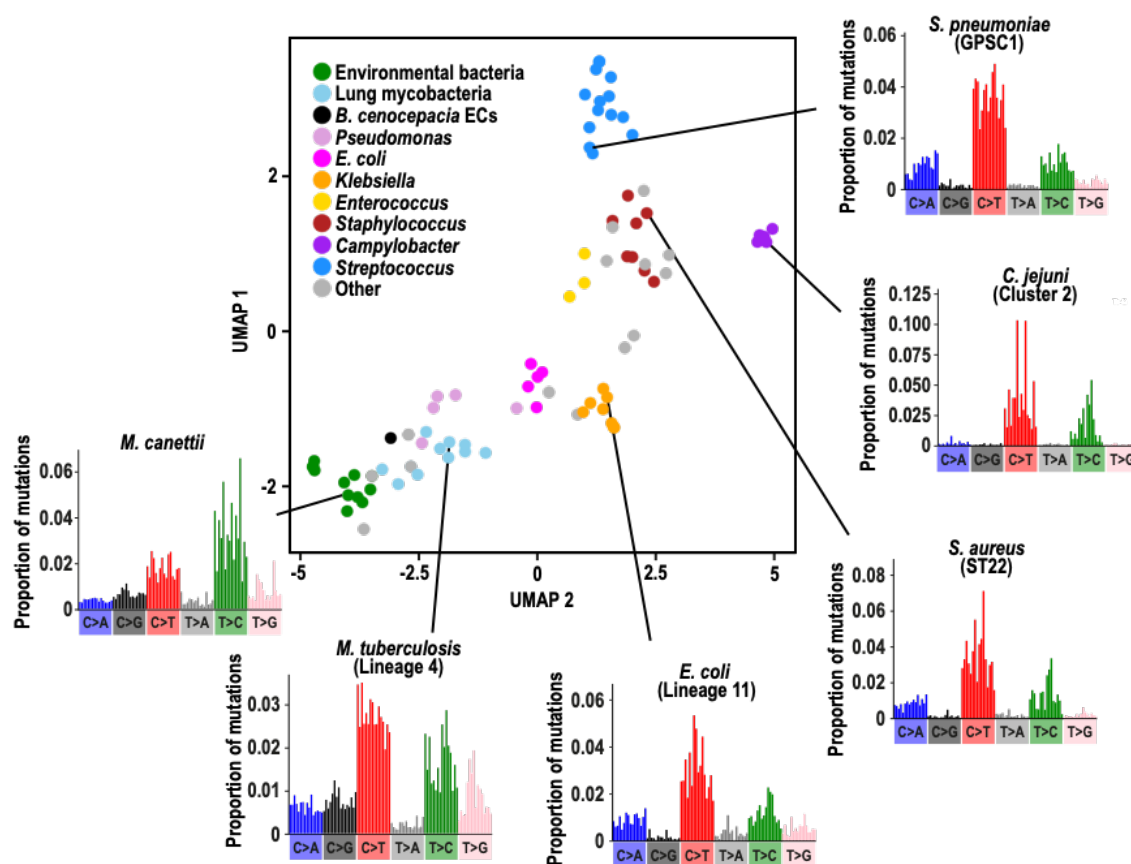
495

496

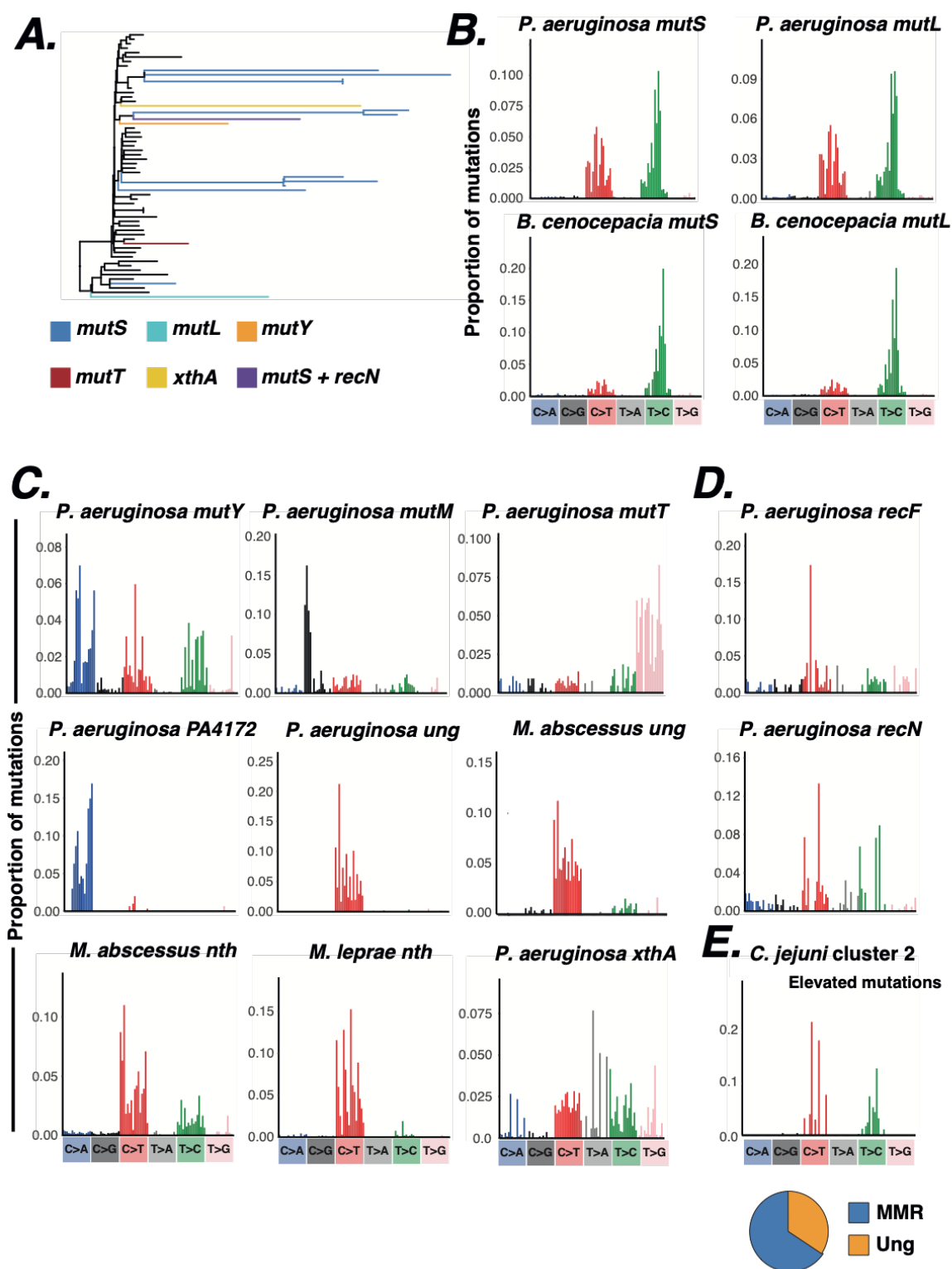
497

498

499
500

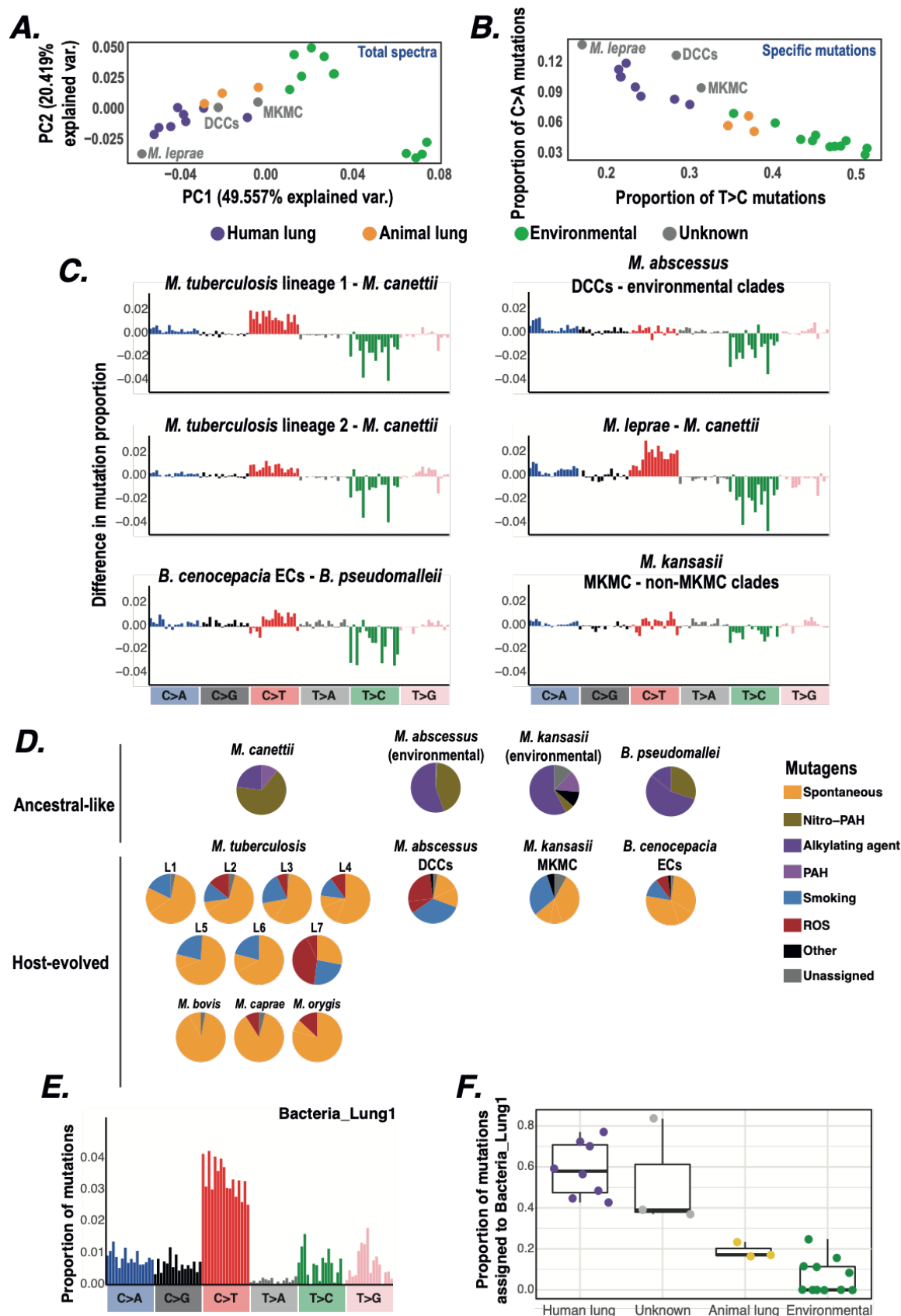


501
502 **Fig. 1. Clustering of bacterial SBS spectra.** UMAP clustering based on contextual mutation
503 proportions within the 84 SBS spectra across 31 bacterial species. Selected groups are
504 coloured. The environmental bacteria label includes *Burkholderia pseudomallei* and known
505 environmental *Mycobacteria*. Example SBS spectra are shown for selected groups.
506



507
 508 **Fig. 2. Mutational signatures associated with DNA repair genes.** (A) Example *P. aeruginosa*
 509 phylogenetic tree (ST274) showing hypermutator branches and the inferred responsible
 510 genes. Hypermutator branches were identified based on branch length and the ratio of
 511 transition and transversion mutations. Responsible genes were identified as DNA repair
 512 genes exhibiting a mutation on the long phylogenetic branch or ancestral branch. Black
 513 branches are background non-hypermutator branches that did not contribute to

514 hypermutator spectra. **(B)** Mutational signatures associated with MMR genes. **(C)**
515 Mutational signatures associated with BER genes. **(D)** Mutational signatures associated with
516 genes involved in homologous recombination. **(E)** Top panel shows the mutations elevated
517 in *C. jejuni* cluster 2 compared with *E. coli* lineage 34, calculated by subtracting each
518 respective mutation proportion in the SBS spectra. The pie chart shows the proportion of
519 mutations elevated in *C. jejuni* cluster 2 that are assigned to each bacterial DNA repair gene
520 signature in a decomposition analysis.
521
522



523

524

525

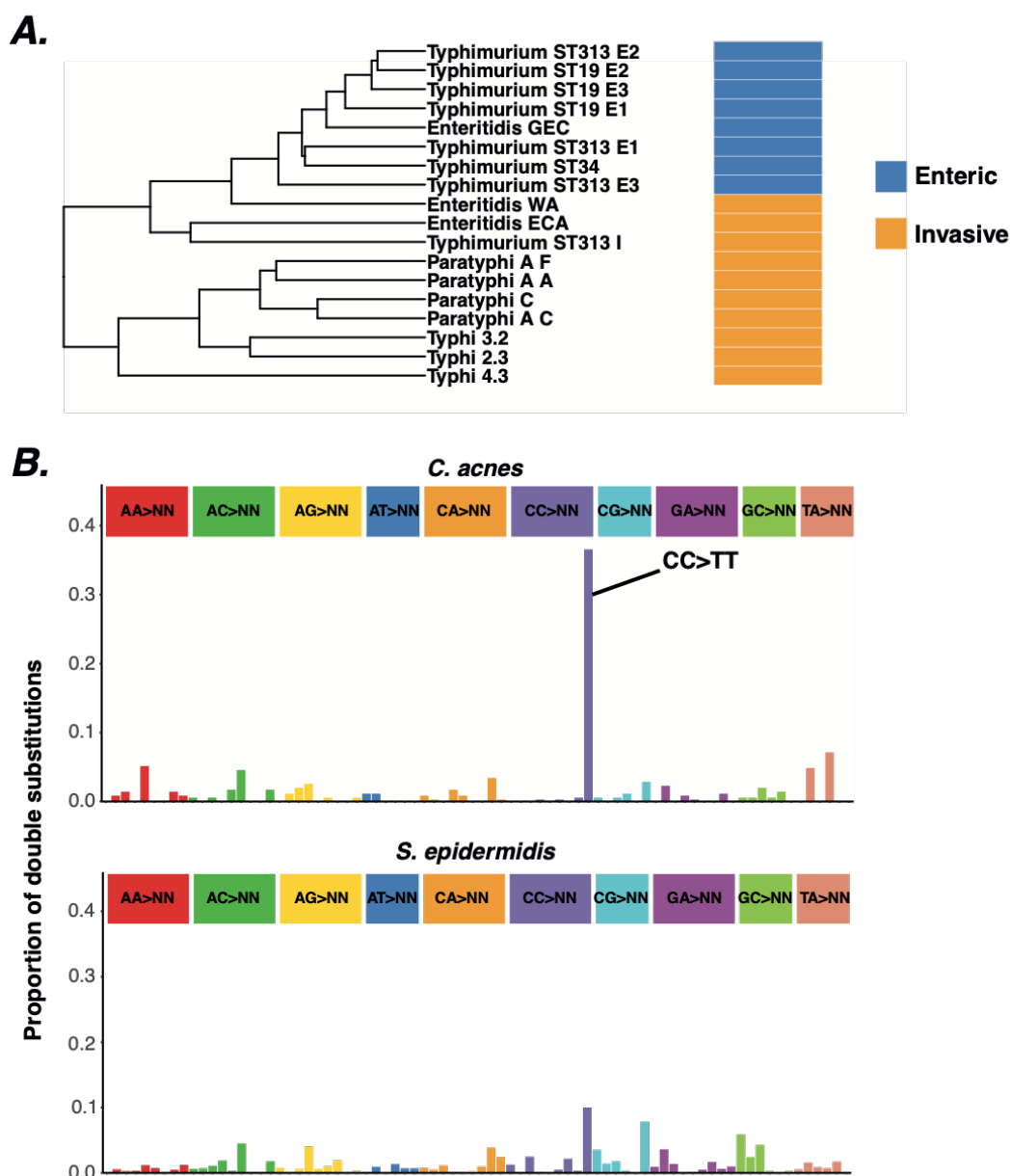
526

Fig. 3. Comparison of mutational spectra between lung and environmental niches. (A)

Principal component analysis on mutation proportions in the SBS spectra across

Mycobacteria and *Burkholderia*. Axes labels include the inferred proportion of variance each

527 principal component describes. Points are coloured by niche; clades with a previously
528 unknown niche are labelled. Environmental includes *B. pseudomallei* and known
529 environmental clades of *Mycobacteria*. **(B)** Comparison of the proportion of T>C and
530 proportion of C>A mutations in *Mycobacteria* and *Burkholderia* SBS spectra, coloured as in
531 **A**. **(C)** Subtraction of mutation proportions in SBS spectra between closely related bacterial
532 clades. Each comparison subtracts the SBS spectrum of a known environmental clade from
533 the SBS spectrum of a clade either known to reside within the lung or with an unknown
534 niche. **(D)** Decomposition of mutational spectra into their underlying components. Only
535 mutations elevated within the respective clade compared to a closely related clade in a
536 different niche were included. Known environmental clades were decomposed into the set
537 of previously extracted environmental mutagen signatures (5) while known lung clades and
538 clades with unknown niche were decomposed into the set of previously extracted lung
539 signatures from human data. *B. cenocepacia* ECs: *B. cenocepacia* epidemic clones. Nitro-
540 PAH: nitro-polycyclic aromatic hydrocarbons; PAH: polycyclic aromatic hydrocarbons; ROS:
541 reactive oxygen species. **(E)** Composition of signature Bacteria_Lung1 extracted from NMF
542 decomposition of *Mycobacteria* and *Burkholderia* SBS spectra. **(F)** The proportion of
543 mutations within each *Mycobacteria* and *Burkholderia* SBS spectrum assigned to signature
544 Bacteria_Lung1.
545



546
 547 **Fig. 4. Comparison of mutational spectra between niches.** (A) Hierarchical clustering of
 548 mutation proportions in salmonella SBS spectra labelled by niche. The enteric clades within
 549 Typhimurium ST19 and ST313 were split into three subclades for this analysis, labelled E1-
 550 E3. Typhimurium ST313 also contains an invasive clade labelled I. GEC - Global Enteric clade;
 551 WA - West Africa clade; ECA - East and Central Africa clade. Three genotypes of Paratyphi A
 552 were included - A, C and F. (B) Comparison of double substitution spectra between *C. acnes*
 553 and *S. epidermidis* (phylogenetic groups A-C combined). CC>TT is indicated in the *C. acnes*
 554 spectrum and is a classic signature of exposure to UV light (5).