

A non-adaptive demographic mechanism for genome expansion in *Streptomyces*

Mallory J Choudoir^{a†}, Marko J Järvenpää^{b‡}, Pekka Marttinen^b, and *Daniel H Buckley^a

^aSchool of Integrative Plant Science, Cornell University, Ithaca, NY, USA

^bDepartment of Computer Science, Helsinki Institute for Information Technology HIIT, Aalto University, Espoo, Finland

[†]Present address: Department of Microbiology, University of Massachusetts Amherst, Amherst, MA, USA

[‡]Present address: Department of Biostatistics, University of Oslo, Oslo, Norway

*Corresponding author: Daniel H Buckley
School of Integrative Plant Science, Bradfield Hall
Cornell University, Ithaca, NY 14850
dhb28@cornell.edu, (607) 255 1716

Running title: Non-adaptive mechanism for genome expansion

1 **Abstract**

2 The evolution of microbial genome size is driven by gene acquisition and loss events that occur
3 at scales from individual genomes to entire pangenomes. The equilibrium between gene gain and
4 loss is shaped by evolutionary forces, including selection and drift, which are in turn influenced
5 by population demographics. There is a well-known bias towards deletion in microbial genomes,
6 which promotes genome streamlining. Less well described are mechanisms that promote genome
7 expansion, giving rise to the many microbes, such as *Streptomyces*, that have unusually large
8 genomes. We find evidence of genome expansion in *Streptomyces* sister-taxa, and we
9 hypothesize that a recent demographic range expansion drove increases in genome size through a
10 non-adaptive mechanism. These *Streptomyces* sister-taxa, NDR (northern-derived) and SDR
11 (southern-derived), represent recently diverged lineages that occupy distinct geographic ranges.
12 Relative to SDR genomes, NDR genomes are larger, have more genes, and their genomes are
13 enriched in intermediate frequency genes. We also find evidence of relaxed selection in NDR
14 genomes relative to SDR genomes. We hypothesize that geographic range expansion, coupled
15 with relaxed selection, facilitated the introgression of non-adaptive horizontally acquired genes,
16 which accumulated at intermediate frequencies through a mechanism known as genome surfing.
17 We show that similar patterns of pangenome structure and genome expansion occur in a
18 simulation that models the effects of population expansion on genome dynamics. We show that
19 non-adaptive evolutionary phenomena can explain expansion of microbial genome size, and
20 suggest that this mechanism might explain why some bacteria with large genomes can be found
21 in soil.

22

23 **Introduction**

24 Microbial genomes are extraordinarily dynamic. Genome size varies considerably, and gene
25 content in strains of the same species can differ dramatically, giving rise to the pangenome. The
26 pangenome concept has transformed our understanding of evolutionary processes in diverse taxa
27 [1–4]. The pangenome is the entire collection of genes in a microbial species, and is subdivided
28 into core genes present in all strains, dispensable or accessory genes present in some strains, and
29 strain-specific or unique genes [5, 6]. Rates of gene acquisition and gene loss determine the
30 individual genome size, and consequently, pangenome composition is shaped by evolutionary
31 mechanisms that alter gene frequencies in microbial populations [7–9].

32

33 Genome size varies by four orders of magnitude (10^4 – 10^7 kb) in eukaryotic organisms and two
34 orders of magnitude in prokaryotic organisms (from less than 150 kb in certain endosymbionts to
35 over 10 Mb for some free-living bacteria) [10]. Unlike eukaryotes, whose genomes contain large
36 portions of non-coding DNA, prokaryotic gene content is directly related to genome size because
37 bacterial and archaeal taxa have high coding density [11, 12]. While microbial genomes are
38 constantly in flux, deletion rates are approximately three-fold greater than rates of gene
39 acquisition [13]. Multiple factors contribute to the strong deletion bias in microbial genomes,
40 including selection for efficiency, “use it or lose it” purging of nonessential genes, and genetic
41 drift [14–17].

42

43 Because of this tendency towards deletion, microbial genome reduction has been examined in
44 greater detail than genome expansion. For example, the evolutionary mechanisms driving
45 genome reduction in obligate pathogens like *Rickettsia* and symbionts like *Buchnera* in aphids
46 are well described [17, 18]. The transition from a free-living to a host-associated lifestyle

47 involves substantial loss of superfluous genes, and generations of vertical transmission in small
48 asexual populations leads to gene inactivation and deletion accelerated by genetic drift (16).
49 Alternatively, genome streamlining leads to reduction of both genome and cell size through
50 selection for increased metabolic efficiency in free-living microbes with large populations [15].
51 Genome streamlining is historically associated with marine oligotrophic *Pelagibacter* [14, 19]
52 but has more recently been described for soil-dwelling *Verrucomicrobia* [20].

53

54 Large genomes are frequent among terrestrial free-living microbes, and must be the product of
55 evolutionary forces that drive genome expansion. A common, though relatively untested,
56 hypothesis to explain large genomes is that high environmental heterogeneity (a characteristic of
57 terrestrial habitats) selects for metabolic versatility afforded by gene gain, and thereby drives
58 genome expansion [21, 22]. For example, massive gene acquisition and adaptation to alkaline
59 conditions caused genome expansion in the myxobacterium *Sorangium cellulosum*, which at 15
60 Mb is one of the largest known bacterial genomes [23]. Mechanisms of gene gain include
61 duplication or horizontal gene transfer (HGT), and large genomes are enriched in functional
62 genes acquired from phylogenetically distant origins [24]. Much of the evolution of gene
63 families can be attributed to HGT rather than duplication events [25, 26], and HGT is a major
64 driver of genome expansion [27, 28]. While HGT-mediated gene acquisition occurs with great
65 frequency, microbial genomes remain relatively small, and genome size tends to be fairly
66 conserved within a species.

67

68 Gene frequencies at the population-level are governed by selection and drift, and these
69 evolutionary forces determine whether a newly acquired gene will be purged from the

70 pangenome or whether it will sweep to fixation. The strength of selection and drift varies
71 inversely, and their relative contributions are determined by a gene's selection coefficient and
72 effective population size (N_e) [29, 30]. Drift can exert large effects on populations with small N_e ,
73 but these effects decline as N_e increases and selection intensifies. Our ability to disentangle the
74 contributions of selection and drift to pangenome dynamics are complicated by the fact that it
75 remains difficult to estimate microbial N_e [31, 32] and to delimit microbial population and
76 species boundaries [33–35]. Another complication is that demographic models often include the
77 simplifying expectation that N_e is invariable over time.

78

79 Rapid changes in population size are typical in the evolutionary histories of many microbial
80 species, and fluctuations in N_e such as population bottlenecks or expansions can have profound
81 impacts on contemporary patterns of genomic diversity. For example, the population structure
82 for many pathogenic bacterial lineages is exemplified by episodes of rapid expansion of clonal
83 complexes repeated across space and time [36–38]. Microbial population expansions can also be
84 linked to ecological or geographical range expansions [39–42]. For instance, demographic
85 expansion in the oral bacteria *Streptococcus mutans* coincides with the origin of human
86 agricultural practices [41].

87

88 We find evidence for post-glacial range expansion in the genus *Streptomyces*, and these species
89 exhibit several of the genetic characteristics described in plant and animal species whose
90 biogeography was influenced by Pleistocene glaciation [43, 44]. By examining *Streptomyces*
91 isolated from sites across North America, we observed genetic evidence for dispersal limitation,
92 a latitudinal gradient in taxonomic richness, and a latitudinal gradient in genetic diversity [45,

93 46]. We also identified recently diverged sister-taxa comprising a more genetically diverse
94 southern-derived (SDR) clade and a more homogenous northern-derived (NDR) clade, which
95 occupied discrete geographic ranges spanning the boundary of glaciation [47]. We further
96 observed larger genomes in the northern clade compared to the southern clade.

97

98 We hypothesize that genome expansion in NDR is a consequence of demographic change driven
99 by post-Pleistocene range expansion. Here, we evaluate the effects of historical range expansion
100 on lineage divergence, genome size, and pangenome structure, and assess these data in the
101 context of the genome surfing hypothesis. Genome surfing is a non-adaptive mechanism which
102 describes the introgression of horizontally acquired genes facilitated by relaxed selection and
103 amplified by geographic expansion [48]. We hypothesize that range expansion, coupled with
104 relaxed selection, dampened gene loss thereby facilitating an increase in non-adaptive,
105 intermediate frequency genes in the NDR pangenome. We infer gene gain and loss dynamics by
106 evaluating patterns of shared gene content between strains. We predict that the contribution of
107 drift is greater in NDR compared to SDR, and determine the relative strength of selection by
108 comparing genome-wide rates of amino acid substitution between clades. Finally, we evaluate
109 our hypothesis by modeling population expansion under a regime of relaxed selection and ask
110 whether these demographic conditions increase retention of horizontally acquired genes at
111 intermediate frequencies, ultimately causing genome expansion.

112

113 **Material and Methods**

114 *Streptomyces* isolation and genomic DNA extraction

115 The strains in this study belong to a larger culture collection of *Streptomyces* isolated from
116 surface soils (0–5 cm) spanning sites across the United States (see [45, 46]) (Table S1). To
117 minimize the effects of environmental filtering in driving patterns of microbial diversity, we
118 selected sample locations with similar ecologies including meadow, pasture, or native grasslands
119 dominated by perennials and with moderately acidic to neutral soils (pH 6.0 ± 1.0 , mean \pm SD).

120

121 *Streptomyces* strains were isolated by plating air-dried soils on glycerol-arginine agar (pH 8.7)
122 plus cycloheximide and Rose Bengal [49, 50] as previously described [51]. Genomic DNA was
123 extracted with a standard phenol/chloroform/isoamyl alcohol protocol from 72 h liquid cultures
124 grown at 30°C with shaking in yeast extract-malt extract medium (YEME) + 0.5% glycine [52].

125

126 *Genome sequencing, assembly, and annotation*

127 Genome sequencing, assembly, and annotation is previously described (see [47]). Briefly, we
128 used the Nextera DNA Library Preparation Kit (Illumina, San Diego, CA, USA) to prepare
129 sequencing libraries. Genomes were sequenced on an Illumina HiSeq2500 instrument with
130 paired-end reads (2 x 100 bp). Genomes were assembled with the A5 pipeline [53] and annotated
131 with RAST [54]. This generated high quality draft genome assemblies with over 25X coverage
132 and estimated completeness > 99% as assessed with CheckM [55]. We used ITEP and MCL
133 clustering (inflation value = 2.0, cutoff = 0.04, maxbit score) [56] to identify orthologous
134 protein-coding gene clusters (i.e., genes). Genome sequences are available through NCBI under
135 BioProject PRJNA401484 accession numbers SAMN07606143–SAMN07606166.

136

137 *Phylogeny*

138 Phylogenetic relationships were reconstructed from whole genome alignments. We used Mugsy
139 [57] to generate multiple genome nucleotide alignments and trimAl v1.2 [58] for automatic
140 trimming of poorly aligned regions. Maximum likelihood (ML) trees were built using the
141 generalized time reversible nucleotide substitution model [59] with gamma distributed rate
142 heterogeneity among sites (GTRGAMMA) in RAxML v7.3.0[60], and bootstrap support was
143 determined following 20 ML searches with 100 inferences using the RAxML rapid bootstrapping
144 algorithm [61]. Average nucleotide identity (ANI) was calculated from whole genome nucleotide
145 alignments using mothur [62].

146

147 *Pangenome and population genetics analyses*

148 The pangenome was determined from the gene content of 24 *Streptomyces* genomes (Table S2).
149 Strains in this collection were initially chosen for whole genome sequencing based on their
150 genetic similarity at house-keeping loci (see [46]). Subsequent analyses focused on recently
151 diverged sister-taxa clades of 10 genomes each, the northern-derived (NDR) and southern-
152 derived (SDR) lineages. Gene content patterns between strains and pangenome gene frequency
153 distributions were determined from gene presence/absence data.

154

155 Gene-level attributes across gene pools were determined from the average of all nucleotide
156 sequences within an orthologous protein-coding gene cluster (see above). GC content was
157 calculated for each gene using the R package Biostrings [63]. Codon usage bias was calculated
158 for each gene using the R package cordon [64]. Clade-level population genetic traits were
159 evaluated using 2,778 single-copy genes conserved across all 24 genomes. For each core gene,
160 nucleotide sequences were aligned using MAFFT v.7 [65], and Gblocks [66] removed poorly

161 aligned positions. PAL2NAL [67] generated codon alignments, and SNAP [68] calculated intra-
162 clade non-synonymous (K_A) and synonymous (K_S) substitution rates (values > 2 were filtered
163 prior to plotting and statistical analysis).

164

165 *Demographic simulation*

166 We assumed that the SDR pangenome approximates the gene frequency distribution of the last
167 common ancestor of NDR and SDR. For the starting generation 0, we used the model from
168 Marttinen *et al.* [69] to simulate a population of sequences and learn parameter values for rates
169 of gene acquisition and deletion that produced the frequency distribution for SDR. To model
170 range expansion demographics (i.e., severe bottleneck followed by exponential growth), we
171 sampled 5 strains from generation 0 as the founding population for the subsequent generation,
172 and simulated this for 100 generations. The simulated population had a growth rate of 5% per
173 generation until a maximum of 100 individuals was reached. We varied the initial sizes of the
174 founding population as well as the growth rate, and observed qualitatively similar results.

175

176 The model included gene acquisition events and deletion events similar to Marttinen *et al.* [69]
177 but modified to allow for multiple changes. Instead of acquisitions/deletions happening
178 independently, there were $k=20$ simultaneous acquisitions/deletions per strain per generation.
179 The previous model [69] included a multiplicative fitness penalty of 0.99 for each gene
180 exceeding a pre-specified genome size threshold. During the expansion, we relaxed the penalty
181 for excess genes to $0.99^{(\text{current size}/\text{max size})}$ allowing for genome size variation.

182

183 **Results**

184 *Streptomyces sister-taxa*

185 We sequenced the genomes of 20 *Streptomyces* strains isolated from ecologically similar
186 grasslands sites across the United States (Table S1, Table S2). These genomes derive from sister-
187 taxa comprising a northern-derived (NDR) and southern-derived clade (SDR), which originate
188 from sites spanning the historical extent of glaciation (Figure S1, see [45]). These sister-taxa
189 represent closely related but genetically distinct microbial species. Genomes within NDR share
190 $97.8 \pm 1.3\%$ (mean \pm SD) ANI and those within SDR share $97.6 \pm 0.1\%$ (mean \pm SD) ANI,
191 while inter-clade genomic ANI is $93.0 \pm 0.14\%$ (mean \pm SD). An ANI of 93–96% is typically
192 indicative of taxonomic species boundaries [70, 71]. For comparative purposes, we also
193 sequenced the genomes of four strains that co-localized with the sister-taxa. The closest
194 taxonomic neighbor to these 24 strains is *Streptomyces griseus* subsp. *griseus* NBRC 13350,
195 although all strains share $< 95\%$ ANI with this type strain (Figure S1).

196

197 *Genomic attributes and gene content*

198 NDR genomes are larger (8.70 ± 0.23 Mb, mean \pm SD) than SDR genomes (7.87 ± 0.19 Mb,
199 mean \pm SD), and this difference is significant (Mann Whitney U test; $P < 0.0001$) (Figure 1a).
200 NDR genomes also have also have more orthologous protein-coding gene clusters (hereby
201 referred to as genes) ($7,775 \pm 196$ genes, mean \pm SD) than SDR genomes ($7,093 \pm 205$ genes,
202 mean \pm SD), and this difference is also significant (Mann Whitney U test; $P < 0.0001$) (Figure
203 1b). As expected, there is a strong positive correlation between genome size and gene content
204 ($R^2 = 0.95$, $P < 0.0001$), but coding density did not differ between clades (Figure S2). SDR
205 genomes are more genetically diverse than NDR. Nucleotide diversity (π) across conserved,
206 single-copy core genes is greater in SDR than NDR, and this difference is significant (Mann

207 Whitney U test; $P < 0.0001$) (see [47]). Finally, NDR genomes have slightly lower genome-wide
208 GC content ($71.50 \pm 0.087\%$, mean \pm SD) than SDR genomes ($71.62 \pm 0.11\%$, mean \pm SD), and
209 this difference is significant (Mann Whitney U test; $P = 0.017$) (Figure 1c). Shared gene content
210 between strains correlates with genomic similarity as measured by ANI (NDR: $R^2 = 0.82$, $P <$
211 0.0001 ; SDR: $R^2 = 0.64$, $P < 0.0001$) (Figure 2). However, gene content varies more in NDR
212 than in SDR, and there is a significant interaction between genomic similarity and clade with
213 respect to gene content shared between strains (Table S3). This interaction comes from shared
214 gene content between strains increasing more rapidly over recent phylogenetic timescales in
215 NDR compared to SDR (Figure 3).

216

217 *Pangenome structure and dynamics*

218 The 24 *Streptomyces* genomes (Table S2) contain 22,055 total orthologous protein-coding gene
219 clusters (i.e., genes), and 42% (9,285 genes) are strain-specific. All 24 genomes share 3,234
220 (2,778 single-copy) genes, which represent 40–48% of the total gene content per strain. While
221 NDR has a smaller core genome than SDR (4,234 and 4,400 genes, respectively), its pangenome
222 is larger (13,681 genes in NDR versus 12,259 genes in SDR) and contains a greater number of
223 clade-specific genes (5,647 genes unique to NDR versus 4,308 genes unique to SDR) (Figure 3,
224 Figure 4, Table S4).

225

226 For most microbial species, pangenome frequency distributions are U-shaped, reflecting high
227 proportions of both strain-specific genes and core genes [72]. While the pangenome structures of
228 our *Streptomyces* sister-taxa generally conform to this shape, the NDR pangenome is enriched in
229 intermediate frequency accessory genes relative to SDR (Figure 4). The proportion of

230 intermediate-low frequency (i.e., present in 3–5 strains) accessory genes is higher in NDR than
231 in SDR (19% of total genes for NDR versus 9.2% of total genes for SDR) (Table S4), and this
232 difference is statistically significant (two proportion z-test; $P < 0.0001$). Conversely, the
233 proportion of intermediate-high frequency (i.e., present in 6–8 strains) accessory genes is
234 equivalent (6.9% of total genes for NDR versus 7.2% of total genes for SDR; two proportion z-
235 test; $P = 0.26$) (Table S4).

236

237 Next, we determined if genes across different gene pools, binned according to their pangenome
238 frequencies, differed in genetic attributes including per-gene GC content and codon usage bias.
239 GC content differs between gene pools for both NDR and SDR pangenomes (ANOVA; $F_{3, 25932}$
240 $= 267.5$, P -value < 0.0001) (Figure S3). In general, GC content is greater in high frequency and
241 core genes compared to rare and intermediate frequency genes for both sister-taxa. Codon usage
242 bias as measured by the effective number of codons (ENC) [73] also differs between gene pools
243 for both NDR and SDR pangenomes (ANOVA; $F_{3, 21624} = 1862.7$, P -value < 0.0001) (Figure
244 S4). Rare and intermediate frequency genes exhibit less overall codon bias compared to high
245 frequency and core genes, which tend to use codons more preferentially.

246

247 *Historical population demography*

248 Due to founder effects occurring at the edge of an expanding population, N_e is dramatically
249 reduced during geographic range expansion [74]. Consequently, relaxed selection will
250 accompany range expansion since the contribution of selection scales directly with N_e . Based on
251 the theory of neutral molecular evolution, which states that selection on synonymous sites is
252 negligible [75], the ratio of non-synonymous to synonymous amino acid substitutions (K_A/K_S)

253 reflects the relative strength of selection acting on a sequence. When assessed at the level of
254 single-copy genes conserved between the sister taxa (2,444 genes), we observe that genome-wide
255 K_A/K_S tends to be higher in NDR than in SDR (Figure 5), and this difference is significant
256 (Mann-Whitney U test; $P < 0.0001$). This result indicates that selection is weaker and genetic
257 drift stronger in NDR relative to SDR.

258

259 We used a population model (modified from [69]) to determine whether demographic expansion
260 could produce increased intermediate gene frequencies and result in genome expansion. We
261 simulated gene gain and loss events in a population undergoing exponential growth over 100
262 generations, and determined changes in pangenome structure and genome size. To approximate
263 relaxed selection during the population expansion, we imposed a fitness penalty for newly
264 acquired genes that scaled inversely with population size. At the beginning of expansion, most
265 genes were present at high frequencies due to strong founder effects (Figure S5, top and middle).
266 During the expansion, we observed a transient enrichment of intermediate frequency genes
267 within the pangenome (Figure S5, top and middle). Total gene content also increased during
268 population expansion due to relaxed selection pressure when N_e was small, which allowed for
269 the persistence of newly HGT-acquired genes. Genome size stabilized when N_e reached
270 maximum size, and selection pressure balanced HGT-mediated gene gain with simultaneous
271 gene loss (Figure S5, bottom).

272

273 **Discussion**

274 We have hypothesized that the biogeography of our *Streptomyces* sister-taxa is explained by
275 historical demographic change driven by geologic and climatic events that occurred in the late

276 Pleistocene [46, 47]. Following the last glacial maxima, North American plant and animal
277 species rapidly colonized glacial retreat zones, and the genetic consequences of post-glacial
278 expansion are well documented and include northern-ranged populations with low diversity that
279 established vast geographic extent [43, 44]. We hypothesize that the recent common ancestor of
280 NDR and SDR inhabited southern glacial refugia prior to the last glacial maxima (LGM). Post
281 glaciation, NDR dispersed northward and colonized the latitudinal range it occupies today (see
282 [46]). We previously described patterns of gene flow, genomic diversity, and ecological
283 adaptation in these sister-taxa, with both adaptive and non-adaptive processes likely reinforcing
284 lineage divergence [47]. Here, we evaluate the outcomes of historical range expansion on sister-
285 taxa pangenome structure and genome size.

286

287 Expanding populations experience repeated founder effects as individuals along the leading edge
288 disperse and colonize new landscapes, creating spatial patterns of genetic diversity akin to
289 genetic drift [74]. Allele surfing, or gene surfing, is a non-adaptive mechanism that propagates
290 rare alleles along an expanding edge such that neutral, or even deleterious, variants ‘surf’ to
291 higher frequencies than would be expected under population equilibrium [76–78]. When applied
292 to expanding microbial populations, gene surfing can facilitate genome surfing, a neutral
293 mechanism acting at the pangenome level that causes rare genes to surf to higher frequencies
294 independent of natural selection [48]. Below, we outline how historical range expansion and
295 genome surfing could give rise to genome expansion in *Streptomyces*.

296

297 Genome surfing is most likely to occur in microbial populations with intermediate levels of
298 dispersal and in taxa capable of HGT. Bacteria in the genus *Streptomyces* are ubiquitous in soil

299 and produce desiccation and starvation resistant spores which are easily disseminated [52],
300 making them ideal for studying patterns of biogeography dependent on dispersal limitation.
301 Rates of HGT in *Streptomyces* are among the highest estimated across a range of bacterial
302 species [51, 79, 80]. In many instances, HGT events occurred in ancestral lineages creating
303 patterns of shared genetic ancestry and reticulate evolution in many extant *Streptomyces* species
304 [81]. We previously observed a distance decay relationship between sites up to 6,000 km apart,
305 indicative of dispersal limitation at intermediate spatial scales that allows detection of geographic
306 patterns of diversity across the sampled range [45, 46]. We also found evidence of restricted
307 gene flow between the core genomes of NDR and SDR [47].

308

309 Since NDR and SDR sister-taxa share a recent common ancestor (Figure S1), they must also
310 share a common ancestral genome size. Hence, differences in genome size accompanying
311 lineage divergence resulted from either genome expansion in NDR or genome reduction in SDR.
312 Given that changes in genome size are ultimately the result of gene gain and loss, we first
313 evaluated differences in shared gene content between NDR and SDR strains. We find greater
314 variability in shared gene content in NDR compared to SDR (Figure 2, Figure 3). This result
315 suggests relative gene content stability for SDR and gene content instability for NDR, most
316 notably in recent phylogenetic history (Figure 3). Likewise, the pangenome of NDR exceeds that
317 of SDR by over 1,000 genes. Evidence suggests that during range expansion, founders at the
318 expansion edge disperse into new habitats and acquire genes from local gene pools
319 asymmetrically at unequal rates, and gene flow is almost exclusively from local to invading
320 genomes [82]. These data are consistent with the observation that that NDR has a larger, more
321 diverse, and more dynamic pangenome than SDR due to introgression from local gene pools.

322 Regardless of their origin, most novel horizontally-acquired genes are neutral or nearly neutral
323 [83]. In most situations, selection will balance gene gain with gene deletion, and genome size
324 will remain relatively constant.

325

326 Genetic diversity in individuals at the leading edge of an expanding population is dramatically
327 reduced, and their genomes experience relaxed selection pressure due to consecutive population
328 bottlenecks and low N_e [84]. We find that NDR has lower genetic diversity [47] and higher rates
329 of K_A/K_S across its core genome relative to SDR (Figure 5), which is consistent with the
330 prediction that NDR has experienced a period of relaxed selection relative to SDR. A positive
331 correlation is observed between GC content and selection pressure on microbial genomes [85,
332 86], and genome expansion in *Chlamydia* has been linked to relaxed selection resulting in a
333 decrease in genome-wide GC content [87]. We likewise observe a decrease in genome-wide GC
334 content in NDR relative to SDR (Figure 1). Relaxed selection pressure in NDR would mitigate
335 the natural bias towards deletion and permit genes acquired by HGT to persist in the genome,
336 regardless of their adaptive coefficient. Microbial sectoring that accompanies geographic range
337 expansion [88] would then allow these newly acquired genes to accumulate at intermediate
338 frequencies in the pangenome. The fact that NDR has larger overall genome size and that relative
339 selection pressure is lower in NDR than SDR, is contrary to the predictions of the metabolic
340 versatility hypothesis of large genomes.

341

342 We hypothesize that relaxed selection and drift caused genome expansion in NDR. While these
343 same mechanisms are known to promote genome reduction in endosymbionts and obligate
344 pathogens [17, 18], it is important to recognize that these outcomes are not contradictory (Figure

345 6). Genome size is regulated by rates of gene gain and loss, the selective coefficient for each
346 gene in the genome, and the strength of selection. Endosymbionts and obligate intracellular
347 pathogens have small population sizes and accordingly, relaxed selection and stronger drift.
348 Relaxed selection pressure should lessen deletion bias. But under these conditions, host
349 compensation for microbial gene function radically alters selective coefficients of core genes,
350 thereby favoring genome reduction, and slightly deleterious mutations accumulate over time via
351 Muller's ratchet [89, 90]. In addition, rates of HGT from non-host sources are essentially zero,
352 since there is little opportunity for endosymbionts to interact with other microbial cells, resulting
353 in a one way track to genome erosion. In contrast, for free-living microbes relaxed selection
354 pressure should bring about genome expansion by shifting the selective coefficients of accessory
355 genes towards neutral. For example, genome expansion in *Chlamydia* was driven by relaxed
356 selection, recombination, and introgression [87]. In this way small population size can favor
357 genome erosion in endosymbionts, while also favoring genome expansion in free-living
358 organisms (Figure 6). Meanwhile, free-living organisms that have large population sizes and
359 high selection pressure will experience high rates of deletion that purge unnecessary genes in
360 order to promote genome streamlining [14, 15].

361
362 Newly acquired genes tend to occur at low frequency in a population unless they provide an
363 adaptive benefit [91], while adaptive genes will increase rapidly in frequency to join the core
364 genome. These dynamics are believed to explain the characteristic U-shape of pangenome gene
365 frequency distributions [72, 92]. Deviations from U-shape expectations, including increased
366 intermediate frequency genes, can result from changes in selection coefficients of genes or under
367 conditions where HGT exceeds deletion rates [93]. Alternatively, negative frequency dependent

368 selection can cause highly beneficial genes to occur at low and intermediate frequencies [94, 95].

369 A large portion of rare genes in microbial pangenomes are hypothetical proteins or genes of
370 unknown function acquired through HGT [96, 97]. For both NDR and SDR, approximately 60%
371 of unique-rare genes (i.e., present in 1–2 strains) are annotated as hypothetical proteins. Nearly
372 half of the 2,596 genes in NDR’s intermediate-low frequency gene pool (i.e., present in 3–5
373 strains) are also hypothetical genes. Furthermore, intermediate-low frequency genes are similar
374 to rare frequency genes in regards to GC content (Figure S3) and codon usage (Figure S4). These
375 data are consistent with our hypotheses that NDR intermediate frequency genes represent
376 evolutionarily recent HGT-gene acquisitions, which increased in frequency as a result of genome
377 surfing.

378
379 HGT-mediated genome expansion supplies a reservoir of novel genetic material for the evolution
380 of gene families [25, 26], biosynthetic pathways [98], and formation of new metabolic networks
381 [99]. Hence, the metabolic versatility of large genomes might be a classic example of an
382 evolutionary spandrel [100], an adaptive trait associated with large genomes that originated not
383 because of selection for versatility, but rather because the acquisition of diverse metabolic
384 pathways is a byproduct of non-adaptive evolutionary process that cause genome expansion.

385
386 We show that pangenome analysis of *Streptomyces* sister-taxa verifies several predictions of the
387 hypothesis that genome expansion within this clade was enabled by non-adaptive evolutionary
388 processes, most likely driven by late Pleistocene demography. We hypothesize that small
389 effective population size and relaxed selection, a consequence of geographic range expansion,
390 allowed for genes newly acquired by HGT to increase in frequency within the NDR pangenome

391 as a result of genome surfing. Further amplifying this effect is introgression of genes from local
392 gene pools encountered following dispersal into new environments. Non-adaptive genome
393 expansion is inherently a non-equilibrium process driven by a transient period of relaxed
394 selection, and population stabilization will re-impose selection pressures that favor deletion. At
395 this point, intermediate frequency genes will either be lost to deletion or fixed if they provide
396 adaptive benefits, and these processes will shift the pangenome structure back to U-shaped
397 expectations. These insights highlight the importance of considering population demography and
398 the profound influence of historical contingency on contemporary patterns of microbial genome
399 diversity.

400

401 **Data Availability**

402 *Streptomyces* genome sequences are available through NCBI under BioProject PRJNA401484
403 accession numbers SAMN07606143–SAMN07606166.

404

405 **Acknowledgements**

406 This work was supported by the National Science Foundation under Grant No. DEB-1456821
407 awarded to Daniel H Buckley.

408

409 **Competing Interests**

410 The authors claim no conflicts of interest nor have competing interests.

411

412 **References**

413 1. Vernikos G, Medini D, Riley DR, Tettelin H. Ten years of pan-genome analyses. *Curr Opin*

- 414 *Microbiol* 2015; **23**: 148–154.
- 415 2. Rasko DA, Rosovitz MJ, Myers GSA, Mongodin EF, Fricke WF, Gajer P, et al. The
416 pangenome structure of *Escherichia coli*: comparative genomic analysis of *E. coli*
417 commensal and pathogenic isolates. *J Bacteriol* 2008; **190**: 6881–6893.
- 418 3. Reno ML, Held NL, Fields CJ, Burke PV, Whitaker RJ. Biogeography of the *Sulfolobus*
419 *islandicus* pan-genome. *Proc Natl Acad Sci U S A* 2009; **106**: 8605–8610.
- 420 4. Lefébure T, Bitar PDP, Suzuki H, Stanhope MJ. Evolutionary dynamics of complete
421 *Campylobacter* pan-genomes and the bacterial species concept. *Genome Biol Evol* 2010; **2**:
422 646–655.
- 423 5. Medini D, Donati C, Tettelin H, Massignani V, Rappuoli R. The microbial pan-genome.
424 *Curr Opin Genet Dev* 2005; **15**: 589–594.
- 425 6. Tettelin H, Massignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al. Genome
426 analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the
427 microbial ‘pan-genome’. *Proc Natl Acad Sci U S A* 2005; **102**: 13950–13955.
- 428 7. McInerney JO, McNally A, O’Connell MJ. Why prokaryotes have pangenomes. *Nature*
429 *Microbiology* 2017; **2**: 17040.
- 430 8. Brockhurst MA, Harrison E, Hall JPJ, Richards T, McNally A, MacLean C. The ecology
431 and evolution of pangenomes. *Curr Biol* 2019; **29**: R1094–R1103.
- 432 9. Azarian T, Huang I-T, Hanage WP. Structure and Dynamics of Bacterial Populations:
433 Pangenome Ecology. In: Tettelin H, Medini D (eds). *The Pangenome: Diversity, Dynamics*
434 *and Evolution of Genomes*. 2020. Springer, Cham (CH).
- 435 10. Lynch M. Streamlining and simplification of microbial genome architecture. *Annu Rev*
436 *Microbiol* 2006; **60**: 327–349.

- 437 11. Fraser CM, Eisen J, Fleischmann RD, Ketchum KA, Peterson S. Comparative genomics and
438 understanding of microbial biology. *Emerg Infect Dis* 2000; **6**: 505–512.
- 439 12. Kuo C-H, Moran NA, Ochman H. The consequences of genetic drift for bacterial genome
440 complexity. *Genome Res* 2009; **19**: 1450–1454.
- 441 13. Puigbò P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin EV. Genomes in turmoil:
442 quantification of genome dynamics in prokaryote supergenomes. *BMC Biol* 2014; **12**: 66.
- 443 14. Viklund J, Ettema TJG, Andersson SGE. Independent genome reduction and phylogenetic
444 reclassification of the oceanic SAR11 clade. *Mol Biol Evol* 2012; **29**: 599–615.
- 445 15. Giovannoni SJ, Cameron Thrash J, Temperton B. Implications of streamlining theory for
446 microbial ecology. *ISME J* 2014; **8**: 1553–1565.
- 447 16. Mira A, Ochman H, Moran NA. Deletional bias and the evolution of bacterial genomes.
448 *Trends Genet* 2001; **17**: 589–596.
- 449 17. Moran NA. Microbial minimalism: Genome reduction in bacterial pathogens. *Cell* 2002;
450 **108**: 583–586.
- 451 18. McCutcheon JP, Moran NA. Extreme genome reduction in symbiotic bacteria. *Nat Rev*
452 *Microbiol* 2011; **10**: 13–26.
- 453 19. Giovannoni SJ, Tripp HJ, Givan S, Podar M, Vergin KL, Baptista D, et al. Genome
454 streamlining in a cosmopolitan oceanic bacterium. *Science* 2005; **309**: 1242–1245.
- 455 20. Brewer TE, Handley KM, Carini P, Gilbert JA, Fierer N. Genome reduction in an abundant
456 and ubiquitous soil bacterium ‘*Candidatus Udaeobacter copiosus*’. *Nat Microbiol* 2016; **2**:
457 16198.
- 458 21. Konstantinidis KT, Tiedje JM. Trends between gene content and genome size in prokaryotic
459 species with larger genomes. *Proc Natl Acad Sci U S A* 2004; **101**: 3160–3165.

- 460 22. Dini-Andreote F, Andreote FD, Araújo WL, Trevors JT, van Elsas JD. Bacterial genomes:
461 habitat specificity and uncharted organisms. *Microb Ecol* 2012; **64**: 1–7.
- 462 23. Han K, Li Z-F, Peng R, Zhu L-P, Zhou T, Wang L-G, et al. Extraordinary expansion of a
463 *Sorangium cellulosum* genome from an alkaline milieu. *Sci Rep* 2013; **3**: 2101.
- 464 24. Cordero OX, Hogeweg P. The impact of long-distance horizontal gene transfer on
465 prokaryotic genome size. *Proc Natl Acad Sci U S A* 2009; **106**: 21748–21753.
- 466 25. Lerat E, Daubin V, Ochman H, Moran NA. Evolutionary origins of genomic repertoires in
467 bacteria. *PLoS Biol* 2005; **3**: e130.
- 468 26. Treangen TJ, Rocha EPC. Horizontal transfer, not duplication, drives the expansion of
469 protein families in prokaryotes. *PLoS Genet* 2011; **7**: e1001284.
- 470 27. Bohlin J, Brynildsrud OB, Sekse C, Snipen L. An evolutionary analysis of genome
471 expansion and pathogenicity in *Escherichia coli*. *BMC Genomics* 2014; **15**: 882.
- 472 28. Tsai Y-M, Chang A, Kuo C-H. Horizontal gene acquisitions contributed to genome
473 expansion in insect-Symbiotic *Spiroplasma clarkii*. *Genome Biol Evol* 2018; **10**: 1526–
474 1532.
- 475 29. Wright S. Evolution in Mendelian populations. *Genetics* 1931; **16**: 97–159.
- 476 30. Kimura M. Evolutionary rate at the molecular level. *Nature* 1968; **217**: 624–626.
- 477 31. Bobay L-M, Ochman H. The Evolution of Bacterial Genome Architecture. *Front Genet*
478 2017; **8**: 72.
- 479 32. Rocha EPC. Neutral theory, microbial practice: Challenges in bacterial population genetics.
480 *Mol Biol Evol* 2018; **35**: 1338–1347.
- 481 33. Roselló-Mora R, Amann R. The species concept for prokaryotes. *FEMS Microbiol Rev*
482 2001; **25**: 39–67.

- 483 34. Achtman M, Wagner M. Microbial diversity and the genetic nature of microbial species.
484 *Nat Rev Microbiol* 2008; **6**: 431–440.
- 485 35. Shapiro BJ. What microbial population genomics has taught us about speciation. In: Polz
486 MF, Rajora OP (eds). *Population Genomics: Microorganisms*. 2019. Springer International
487 Publishing, Cham, pp 31–47.
- 488 36. Smith NH, Dale J, Inwald J, Palmer S, Gordon SV, Hewinson RG, et al. The population
489 structure of *Mycobacterium bovis* in Great Britain: clonal expansion. *Proc Natl Acad Sci U*
490 *S A* 2003; **100**: 15271–15275.
- 491 37. Achtman M. Population structure of pathogenic bacteria revisited. *Int J Med Microbiol*
492 2004; **294**: 67–73.
- 493 38. Nübel U, Dordel J, Kurt K, Strommenger B, Westh H, Shukla SK, et al. A timescale for
494 evolution, population expansion, and spatial spread of an emerging clone of methicillin-
495 resistant *Staphylococcus aureus*. *PLoS Pathog* 2010; **6**: e1000855.
- 496 39. Wirth T, Hildebrand F, Allix-Béguet C, Wölbeling F, Kubica T, Kremer K, et al. Origin,
497 spread and demography of the *Mycobacterium tuberculosis* complex. *PLoS Pathog* 2008; **4**:
498 e1000160.
- 499 40. Takuno S, Kado T, Sugino RP, Nakhleh L, Innan H. Population genomics in bacteria: a case
500 study of *Staphylococcus aureus*. *Mol Biol Evol* 2012; **29**: 797–809.
- 501 41. Cornejo OE, Lefebure T, Pavinski 2. Paulina, Lang P, Richards 2. Vincent P., Eilertson K,
502 et al. Evolutionary and population genomics of the cavity causing bacteria *Streptococcus*
503 *mutans*. *Evolution* ; **30**: 881–893.
- 504 42. Montano V, Didelot X, Foll M, Linz B, Reinhardt R, Suerbaum S, et al. Worldwide
505 population structure, long-term demography, and local adaptation of *Helicobacter pylori*.

- 506 *Genetics* 2015; **200**: 947–963.
- 507 43. Hewitt G. Some genetic consequences of ice ages, and their role in divergence and
508 speciation. *Biological Journal of the Linnean Society* 1996; **58**: 247–276.
- 509 44. Hewitt GM. Genetic consequences of climatic oscillations in the Quaternary. *Philos Trans*
510 *R Soc Lond B Biol Sci* 2004; **359**: 183–195.
- 511 45. Andam CP, Doroghazi JR, Campbell AN, Kelly PJ, Choudoir MJ, Buckley DH. A
512 latitudinal diversity gradient in terrestrial bacteria of the genus *Streptomyces*. *MBio* 2016; **7**:
513 e02200–15.
- 514 46. Choudoir MJ, Doroghazi JR, Buckley DH. Latitude delineates patterns of biogeography in
515 terrestrial *Streptomyces*. *Environ Microbiol* 2016; **18**: 4931–4945.
- 516 47. Choudoir MJ, Buckley DH. Phylogenetic conservatism of thermal traits explains dispersal
517 limitation and genomic differentiation of *Streptomyces* sister-taxa. *ISME J* 2018; **12**: 2176–
518 2186.
- 519 48. Choudoir MJ, Panke-Buisse K, Andam CP, Buckley DH. Genome surfing as driver of
520 microbial genomic diversity. *Trends Microbiol* 2017; **25**: 624–636.
- 521 49. El-Nakeeb MA, Lechevalier HA. Selective isolation of aerobic Actinomycetes. *Appl*
522 *Microbiol* 1963; **11**: 75–77.
- 523 50. Ottow JCG. Rose Bengal as a selective aid in the isolation of fungi and actinomycetes from
524 natural sources. *Mycologia* 1972; **64**: 304.
- 525 51. Doroghazi JR, Buckley DH. Widespread homologous recombination within and between
526 *Streptomyces* species. *ISME J* 2010; **4**: 1136.
- 527 52. Kieser, T, Bibb, MJ, Buttner, MJ, Charter, KF, Hopwood, DA. *Practical Streptomyces*
528 Genetics. 2000. John Innes Foundation, Norwich, UK.

- 529 53. Tritt A, Eisen JA, Facciotti MT, Darling AE. An integrated pipeline for *de novo* assembly
530 of microbial genomes. *PLoS One* 2012; **7**: e42304.
- 531 54. Aziz RK, Bartels D, Best AA, DeJongh M, Disz T, Edwards RA, et al. The RAST server:
532 Rapid annotations using subsystems technology. *BMC Genomics* 2008; **9**: 75.
- 533 55. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the
534 quality of microbial genomes recovered from isolates, single cells, and metagenomes.
535 *Genome Res* 2015; **25**: 1043–1055.
- 536 56. Benedict MN, Henriksen JR, Metcalf WW, Whitaker RJ, Price ND. ITEP: an integrated
537 toolkit for exploration of microbial pan-genomes. *BMC Genomics* 2014; **15**: 8.
- 538 57. Angiuoli SV, Salzberg SL. Mugsy: fast multiple alignment of closely related whole
539 genomes. *Bioinformatics* 2011; **27**: 334–342.
- 540 58. Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. trimAl: a tool for automated
541 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 2009; **25**: 1972–
542 1973.
- 543 59. Tavaré S. Some probabilistic and statistical problems in the analysis of DNA sequences.
544 Lectures Math Life Sci 17: 57-86 58. Warscheid T, Braams J (2000) Biodeterioration of
545 stone: a review. *Int Biodeterior Biodegradation* 1986; **46**: 343–368.
- 546 60. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with
547 thousands of taxa and mixed models. *Bioinformatics* 2006; **22**: 2688–2690.
- 548 61. Stamatakis A, Hoover P, Rougemont J, Renner S. A rapid bootstrap algorithm for the
549 RAxML web servers. *Syst Biol* 2008; **57**: 758–771.
- 550 62. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing
551 mothur: Open-source, platform-independent, community-supported software for describing

- 552 and comparing microbial communities. *Appl Environ Microbiol* 2009; **75**: 7537–7541.
- 553 63. Pagès HA, Gentleman P, DebRoy R. Biostrings: Efficient manipulation of biological
554 strings. *R package version 2.59*. 2020.
- 555 64. Elek A, Kuzman M, Vlahoviček K. coRdon: codon usage analysis and prediction of gene
556 expressivity. *R package version 1.8.0*. 2019.
- 557 65. Katoh K, Standley DM. MAFFT Multiple sequence alignment software version 7:
558 Improvements in performance and usability. *Mol Biol Evol* 2013; **30**: 772–780.
- 559 66. Talavera, Gerard, Castresana, Jose. Improvement of phylogenies after removing divergent
560 and ambiguously aligned blocks from protein sequence alignments. *Syst Biol* 2007; **56**:
561 564–577.
- 562 67. Suyama M, Torrents D, Bork P. PAL2NAL: robust conversion of protein sequence
563 alignments into the corresponding codon alignments. *Nucleic Acids Res* 2006; **34**: W609–
564 12.
- 565 68. Korber BT. HIV Signature and Sequence Variation Analysis. In: Allen G. Rodrigo and
566 Gerald H. Learn (ed). *Computational Analysis of HIV Molecular Sequences*. 2000.
567 Dordrecht, Netherlands: Kluwer Academic Publishers, p Chapter 4, pages 55–72.
- 568 69. Marttinen P, Croucher NJ, Gutmann MU, Corander J, Hanage WP. Recombination
569 produces coherent bacterial species clusters in both core and accessory genomes. *Microbial*
570 *Genomics* 2015; **1**.
- 571 70. Kim M, Oh H-S, Park S-C, Chun J. Towards a taxonomic coherence between average
572 nucleotide identity and 16S rRNA gene sequence similarity for species demarcation of
573 prokaryotes. *Int J Syst Evol Microbiol* 2014; **64**: 346–351.
- 574 71. Ciufu S, Kannan S, Sharma S, Badretdin A, Clark K, Turner S, et al. Using average

- 575 nucleotide identity to improve taxonomic assignments in prokaryotic genomes at the NCBI.
576 *Int J Syst Evol Microbiol* 2018; **68**: 2386–2392.
- 577 72. Haegeman B, Weitz JS. A neutral theory of genome evolution and the frequency
578 distribution of genes. *BMC Genomics* 2012; **13**: 1.
- 579 73. Wright F. The ‘effective number of codons’ used in a gene. *Gene* 1990; **87**: 23–29.
- 580 74. Slatkin M, Excoffier L. Serial founder effects during range expansion: a spatial analog of
581 genetic drift. *Genetics* 2012; **191**: 171–181.
- 582 75. Kimura M. The Neutral Theory of Molecular Evolution. 1983. Cambridge University Press.
- 583 76. Edmonds CA, Lillie AS, Cavalli-Sforza LL. Mutations arising in the wave front of an
584 expanding population. *Proc Natl Acad Sci U S A* 2004; **101**: 975–979.
- 585 77. Travis JMJ, Münkemüller T, Burton OJ, Best A, Dytham C, Johst K. Deleterious mutations
586 can surf to high densities on the wave front of an expanding population. *Mol Biol Evol*
587 2007; **24**: 2334–2343.
- 588 78. Chuang A, Peterson CR. Expanding population edges: theories, traits, and trade-offs. *Glob*
589 *Chang Biol* 2016; **2**: 494–512.
- 590 79. Doroghazi JR, Buckley DH. Intraspecies comparison of *Streptomyces pratensis* genomes
591 reveals high levels of recombination and gene conservation between strains of disparate
592 geographic origin. *BMC Genomics* 2014; **15**: 970.
- 593 80. Cheng K, Rong X, Huang Y. Widespread interspecies homologous recombination reveals
594 reticulate evolution within the genus *Streptomyces*. *Mol Phylogenet Evol* 2016; **102**: 246–
595 254.
- 596 81. Andam CP, Choudoir MJ, Vinh Nguyen A, Sol Park H, Buckley DH. Contributions of
597 ancestral inter-species recombination to the genetic diversity of extant *Streptomyces*

- 598 lineages. *ISME J* 2016; **10**: 1731–1741.
- 599 82. Currat M, Ruedi M, Petit RJ, Excoffier L. The hidden side of invasions: massive
600 introgression by local genes. *Evolution* 2008; **62**: 1908–1920.
- 601 83. Gogarten JP, Townsend JP. Horizontal gene transfer, genome innovation and evolution. *Nat*
602 *Rev Microbiol* 2005; **3**: 679–687.
- 603 84. Excoffier L, Foll M, Petit RJ. Genetic consequences of range expansions. *Annu Rev Ecol*
604 *Evol Syst* 2009; **40**: 481–501.
- 605 85. Hildebrand F, Meyer A, Eyre-Walker A. Evidence of selection upon genomic GC-content
606 in bacteria. *PLoS Genet* 2010; **6**: e1001107.
- 607 86. Raghavan R, Kelkar YD, Ochman H. A selective force favoring increased G+C content in
608 bacterial genes. *Proc Natl Acad Sci U S A* 2012; **109**: 14504–14507.
- 609 87. Bohlin J. Genome expansion in bacteria: the curious case of *Chlamydia trachomatis*. *BMC*
610 *Res Notes* 2015; **8**: 512.
- 611 88. Hallatschek O, Hersen P, Ramanathan S, Nelson DR. Genetic drift at expanding frontiers
612 promotes gene segregation. *Proc Natl Acad Sci U S A* 2007; **104**: 19926–19930.
- 613 89. Moran NA. Accelerated evolution and Muller’s ratchet in endosymbiotic bacteria. *Proc Natl*
614 *Acad Sci U S A* 1996; **93**: 2873–2878.
- 615 90. Rispe C, Moran NA. Accumulation of deleterious mutations in endosymbionts: Muller’s
616 ratchet with two levels of selection. *Am Nat* 2000; **156**: 425–441.
- 617 91. Kuo C-H, Ochman H. The fate of new bacterial genes. *FEMS Microbiol Rev* 2009; **33**: 38–
618 43.
- 619 92. Lobkovsky AE, Wolf YI, Koonin EV. Gene frequency distributions reject a neutral model
620 of genome evolution. *Genome Biol Evol* 2013; **5**: 233–242.

- 621 93. Domingo-Sananes MR, McInerney JO. Selection-based model of prokaryote pangenomes.
622 *bioRxiv*. 2019; doi:10.1101/782573.
- 623 94. Cordero OX, Ventouras L-A, DeLong EF, Polz MF. Public good dynamics drive evolution
624 of iron acquisition strategies in natural bacterioplankton populations. *Proc Natl Acad Sci U*
625 *SA* 2012; **109**: 20059–20064.
- 626 95. McNally A, Kallonen T, Connor C, Abudahab K, Aanensen DM, Horner C, et al.
627 Diversification of colonization factors in a multidrug-resistant *Escherichia coli* lineage
628 evolving under negative frequency-dependent selection. *MBio* 2019; **10**: e00644–19.
- 629 96. Mira A, Klasson L, Andersson SGE. Microbial genome evolution: sources of variability.
630 *Curr Opin Microbiol* 2002; **5**: 506–512.
- 631 97. Mira A, Martín-Cuadrado AB, D’Auria G, Rodríguez-Valera F. The bacterial pan-genome:
632 a new paradigm in microbiology. *Int Microbiol* 2010; **13**: 45–57.
- 633 98. Boucher Y, Doolittle WF. The role of lateral gene transfer in the evolution of isoprenoid
634 biosynthesis pathways. *Mol Microbiol* 2000; **37**: 703–716.
- 635 99. Pál C, Papp B, Lercher MJ. Adaptive evolution of bacterial metabolic networks by
636 horizontal gene transfer. *Nat Genet* 2005; **37**: 1372–1375.
- 637 100. Gould SJ, Lewontin RC. The spandrels of San Marco and the Panglossian paradigm: a
638 critique of the adaptationist programme. *Proc R Soc Lond B Biol Sci* 1979; **205**: 581–598.

639

640

641

642

643

644 **Figure Legends**

645 **Figure 1.** Genomic attributes of NDR and SDR sister-taxa. NDR genomes are larger, have more
646 genes, and have lower GC content compared to SDR genomes. Plots show the distributions of
647 genome size in Mb (a), number of genes (b), and genome-wide GC content (%) (c) for
648 *Streptomyces* sister-taxa. Boxplots show the clade-level medians, interquartile ranges, and 1.5
649 times interquartile ranges. Colored circles illustrate the values for individual genomes belonging
650 to the NDR clade (blue) or the SDR clade (green).

651
652 **Figure 2.** Genomic similarity versus shared gene content for NDR and SDR. Differences in
653 shared gene content across increasing average nucleotide identity (ANI) are greater within the
654 NDR clade compared to the SDR clade (Table S3). Circles show pairwise comparisons of the
655 number of shared genes between two strains versus ANI and are colored by clade according to
656 the legend. Dashed lines show linear regressions, and the shaded area is the 95% confidence
657 interval.

658
659 **Figure 3.** Presence/absence of genes across phylogeny. Gene content changes more rapidly
660 across ancestral phylogenetic nodes for NDR genomes compared to SDR genomes. Tree is made
661 from whole genome nucleotide alignments, and the scale bar shows nucleotide substitutions per
662 site (see Figure S1). Branch colors reflect clade membership. Phylogenetic nodes are labeled
663 with the number of genes conserved in all members of descendent nodes. Gray pie charts at tree
664 tips show the portion of total genes per genome that are strain-specific (black slice). Right panel
665 plots the differences in gene content across the phylogeny beginning at the shared ancestral node
666 and ending with extant taxa at the terminal tips for NDR (blue-solid) and SDR (green-dashed)

667 lineages. Multiple lines represent monophyletic lineages.

668

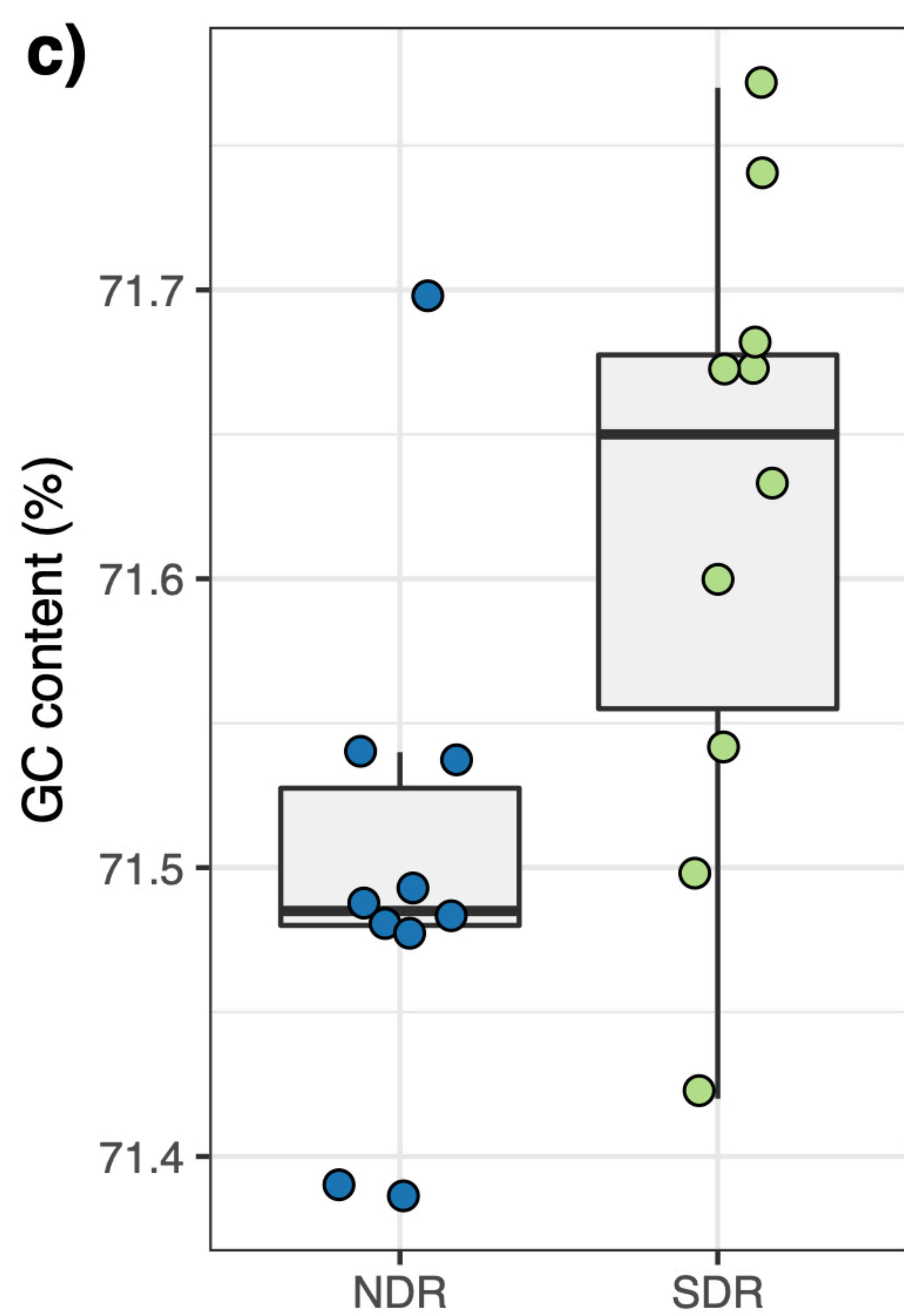
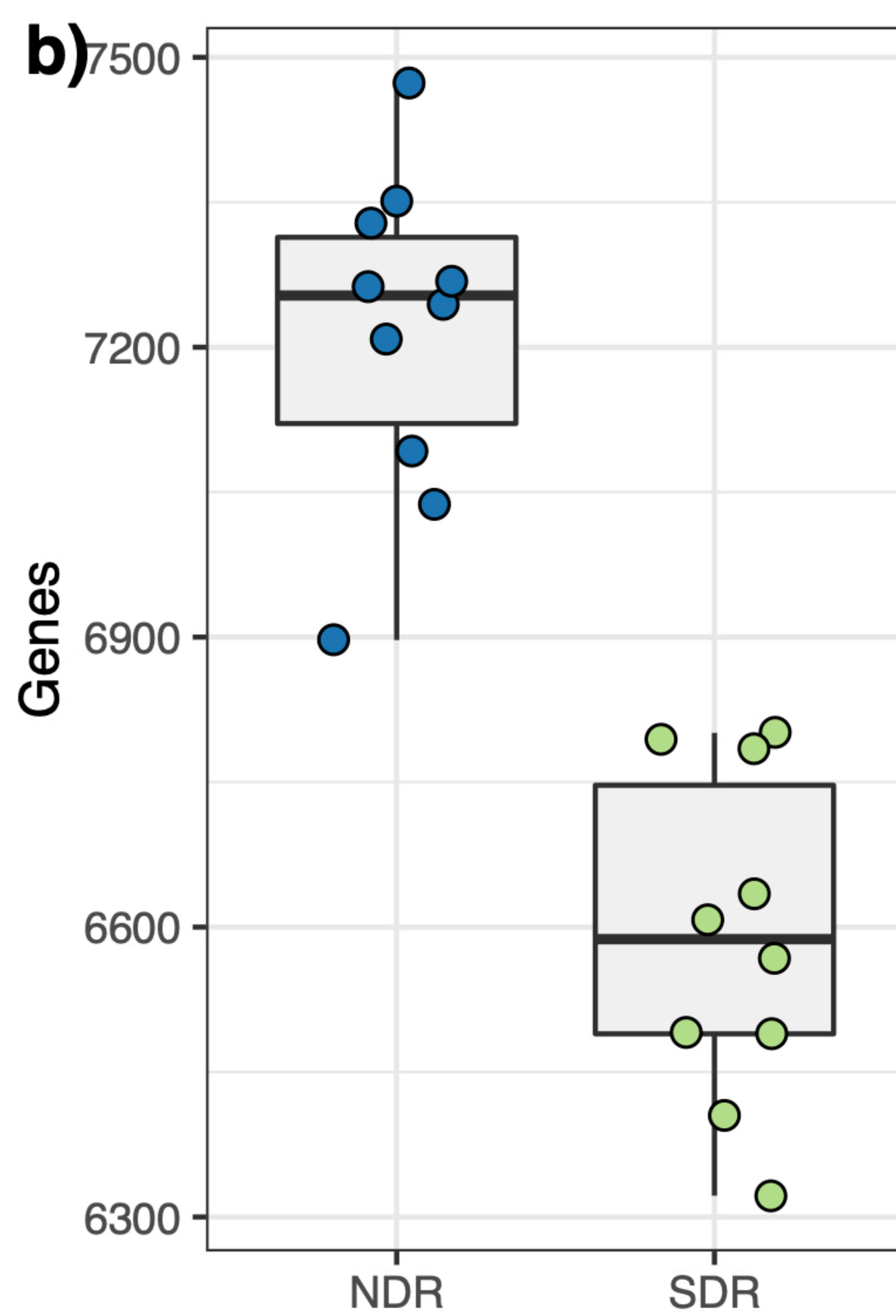
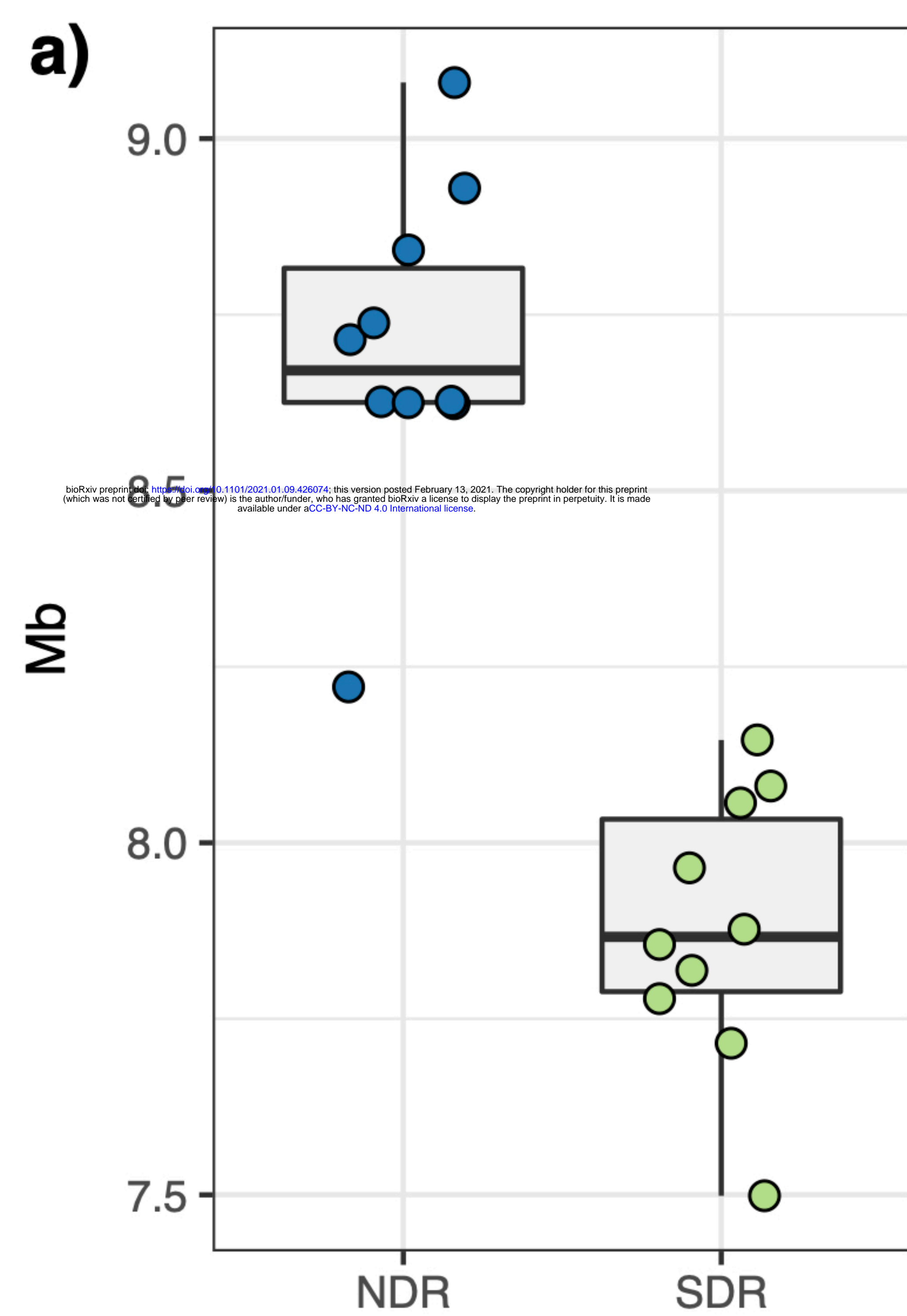
669 **Figure 4.** Pangenome gene frequency distributions. NDR genomes are enriched in intermediate
670 frequency genes. Plots show the pangenome gene frequency distributions for NDR (left) and
671 SDR (right). Bars show the population-level sums of genes present in 1–10 genomes. See Table
672 S3 for raw values and proportions.

673

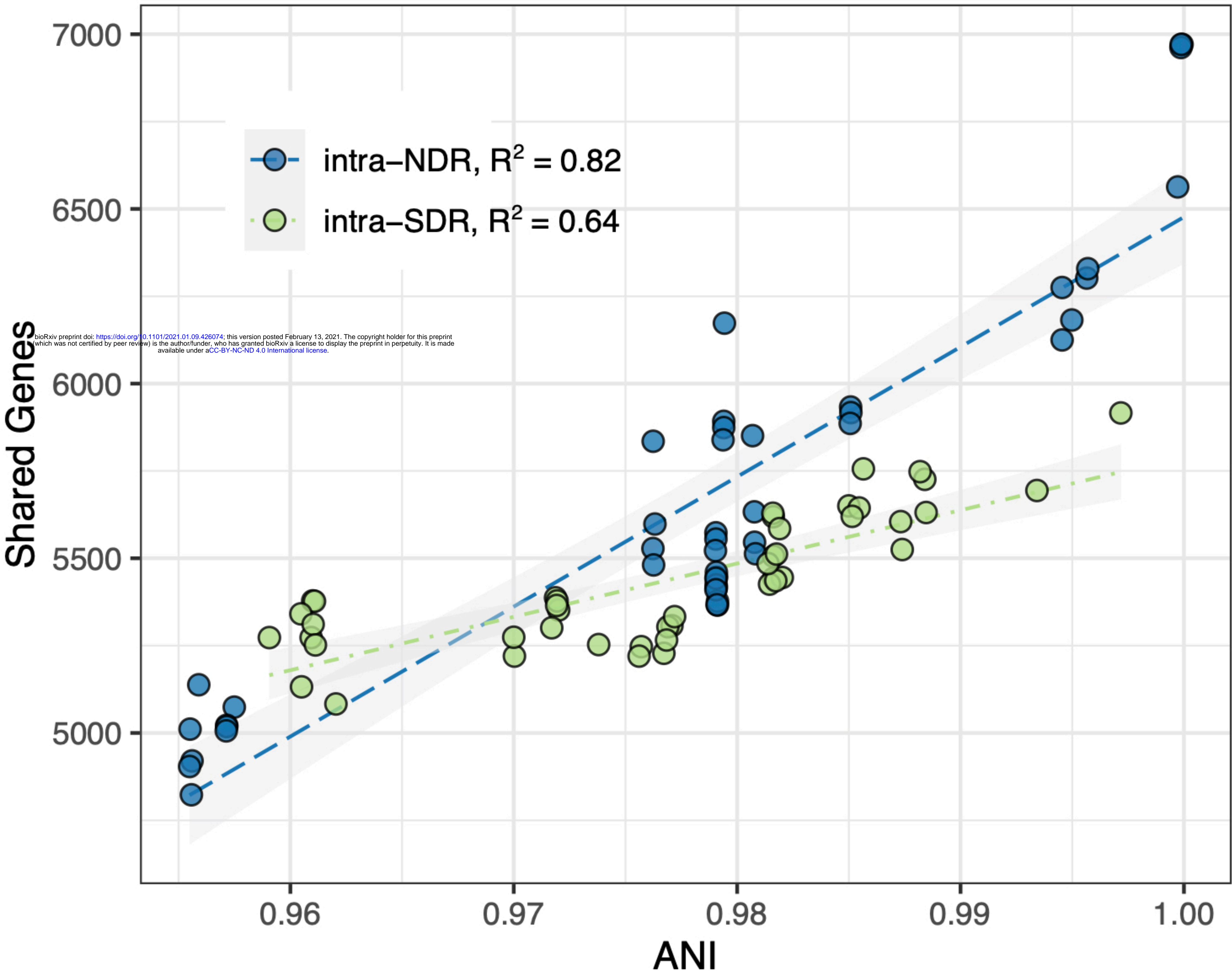
674 **Figure 5.** K_A/K_S values between the NDR and SDR sister-taxa core genome. NDR core genes
675 have, on average, greater rates of non-synonymous to synonymous amino acid substitutions
676 compared to SDR core genes. Circles plot clade-level rates of non-synonymous to synonymous
677 amino acid substitutions (K_A/K_S) for each of 2,444 single-copy core genes for NDR (y-axis) and
678 SDR (x-axis). Axes are logarithmic scale. The black dashed line is a slope of 1, and points along
679 this line are genes with equal K_A/K_S mean values in both clades. K_A/K_S is proportional to the
680 relative strength of genetic drift and inversely proportional to the relative strength of selection.

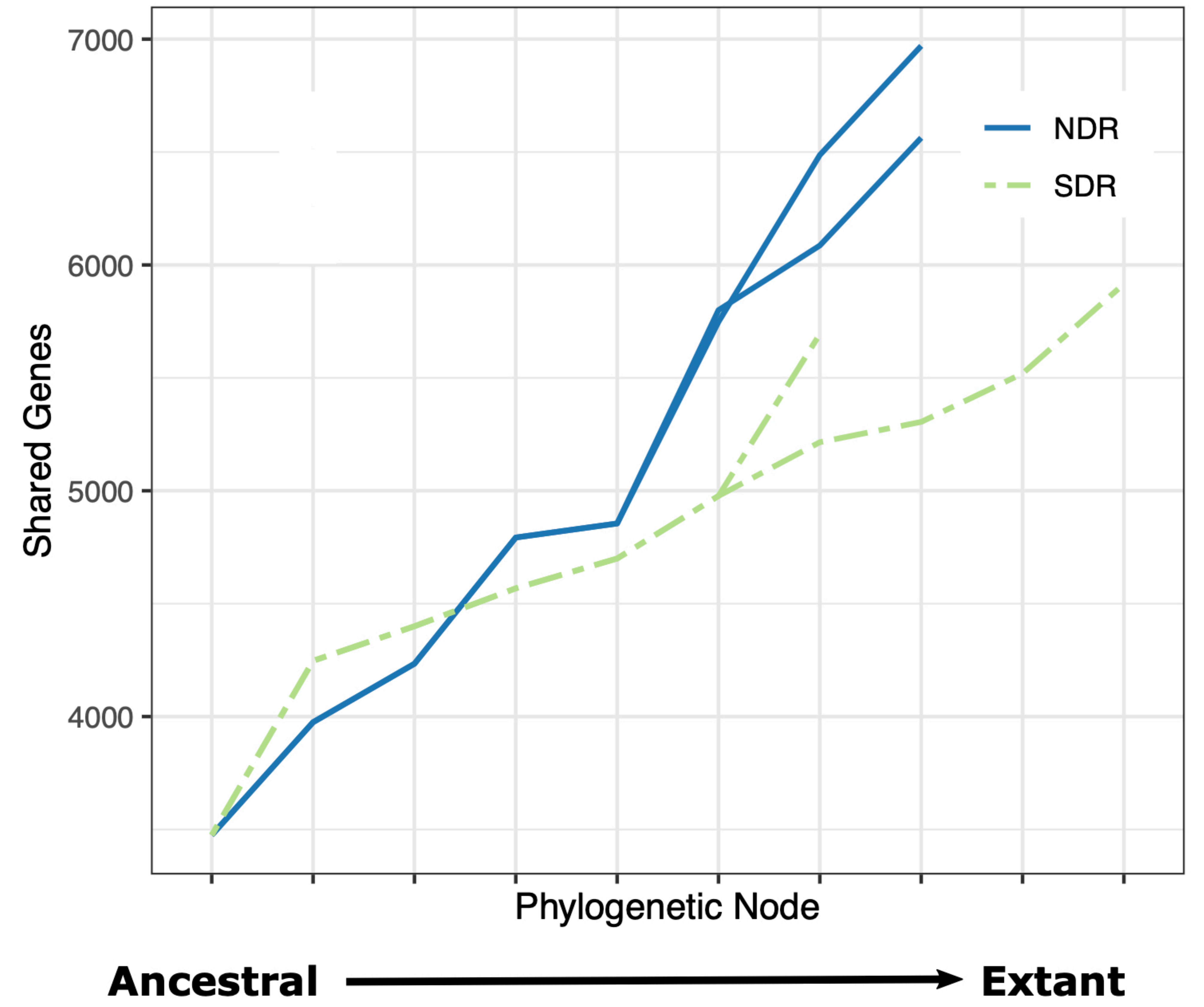
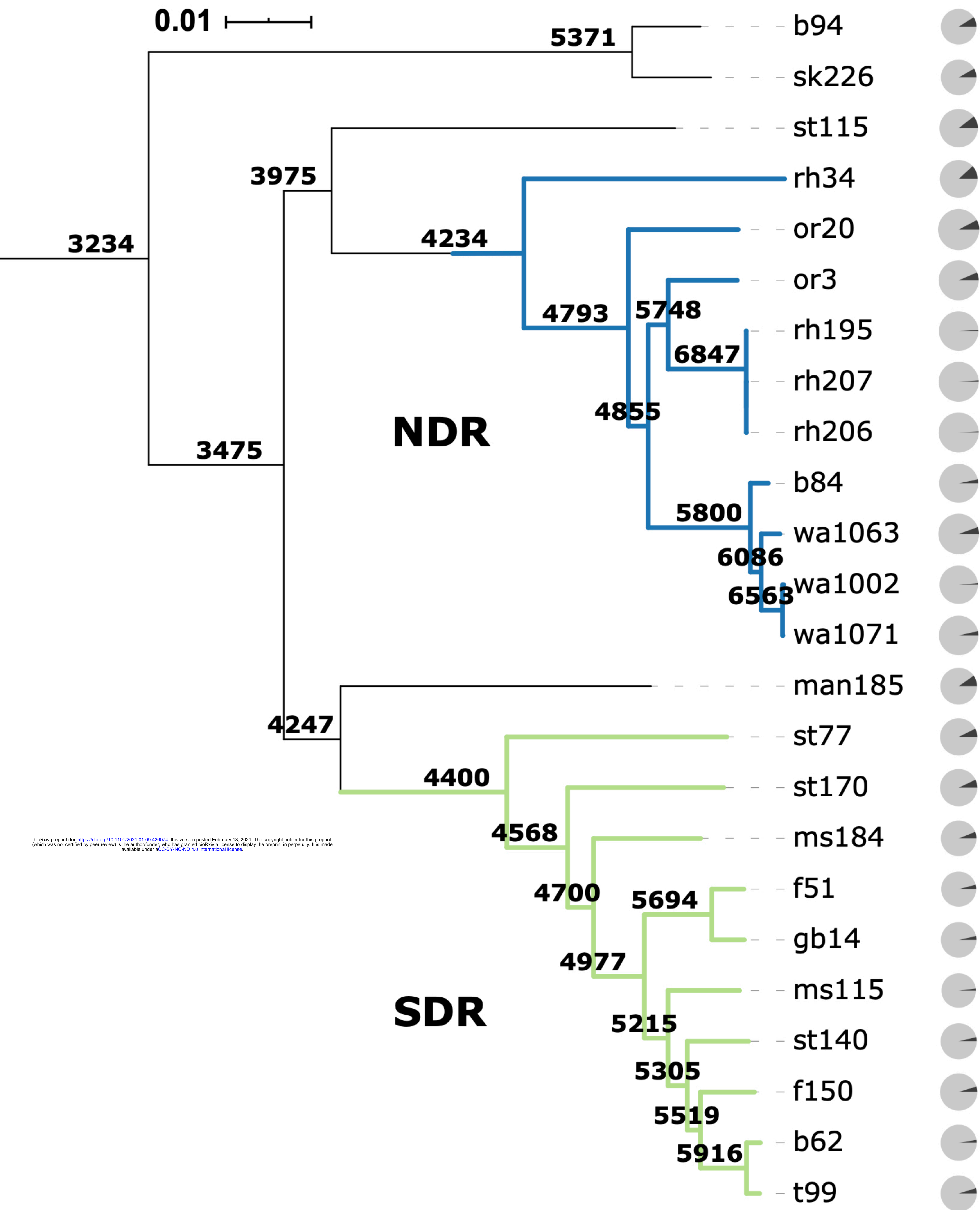
681

682 **Figure 6.** Conceptual overview of the evolutionary processes and demographic conditions that
683 support changes in genome size. Genome erosion (left) in endosymbionts is the result of small
684 N_e and strong genetic drift, with host compensation lowering costs of deletion while restricting
685 gene flow (HGT). Genome streamlining (middle) in free-living microbes with large populations
686 like *Pelagibacter* involves strong selection and elimination of non-adaptive genes. Genome
687 expansion (right) in *Streptomyces* is facilitated by high rates of HGT and relaxed selection,
688 allowing for the accumulation of non-adaptive genes and ultimately larger genomes.

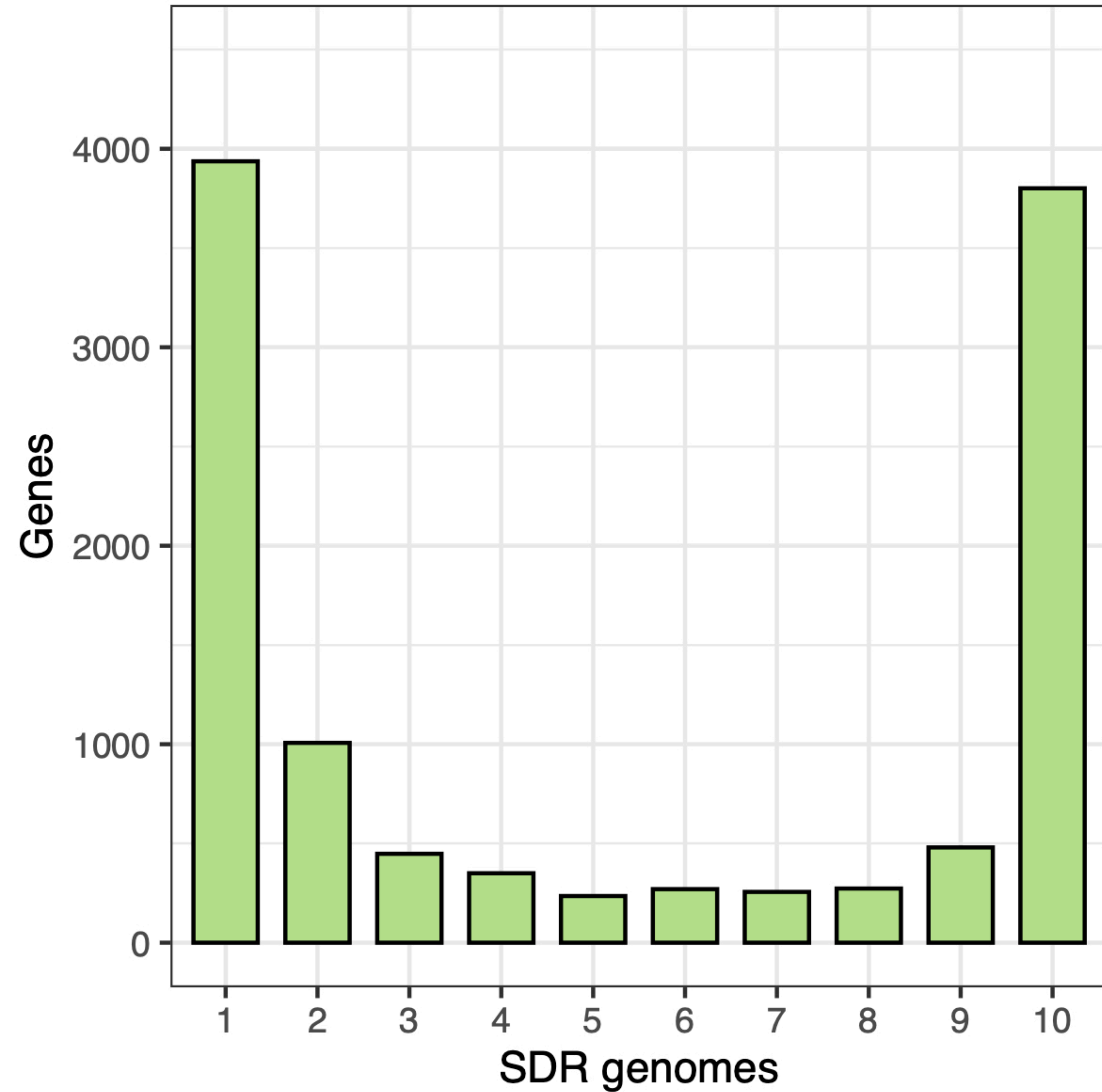
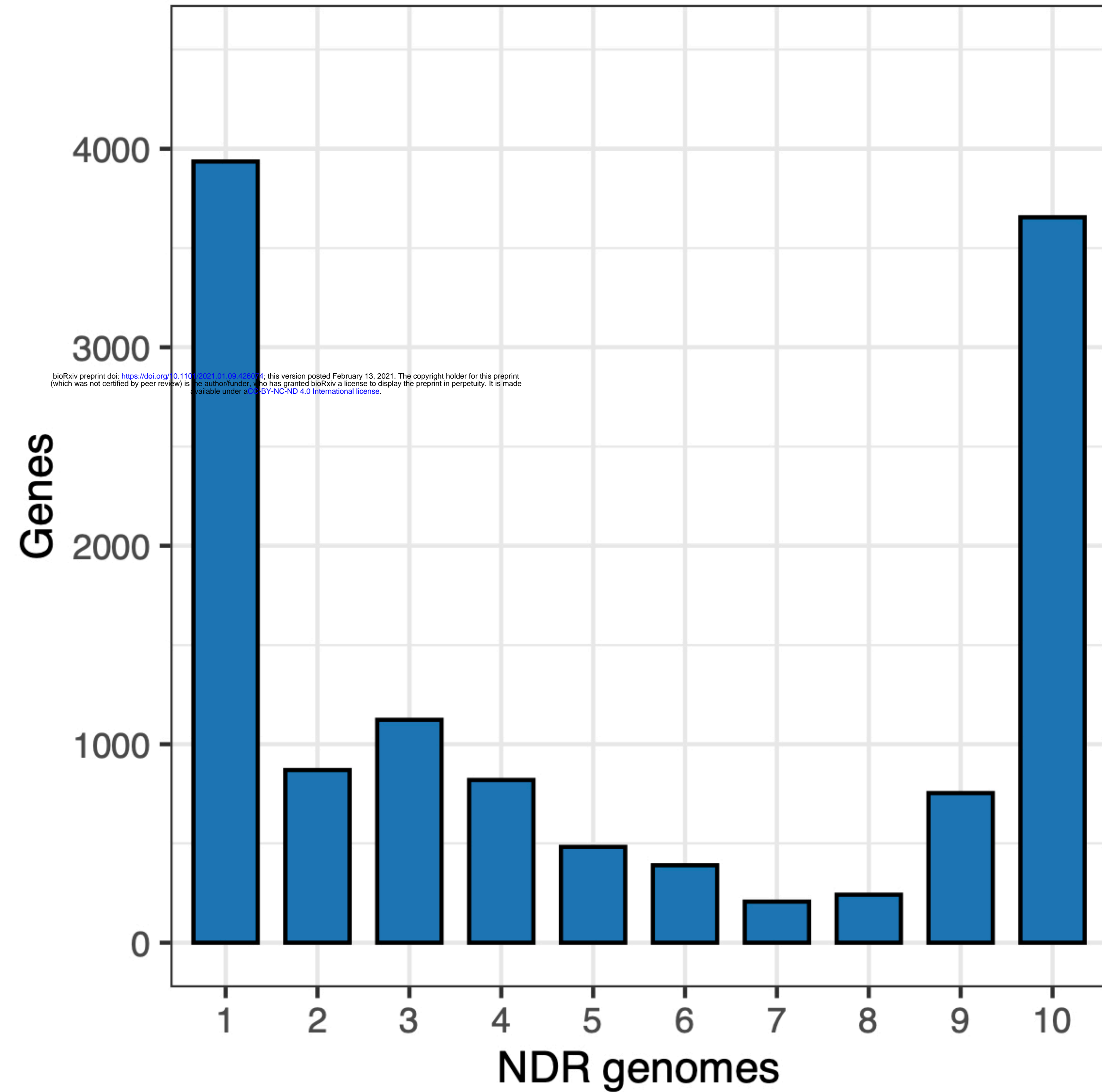


bioRxiv preprint doi: <https://doi.org/10.1101/2021.01.09.426074>; this version posted February 13, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.





bioRxiv preprint doi: <https://doi.org/10.1101/2021.01.09.429074>; this version posted February 13, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



NDR $\log K_A/K_S$

1.00

0.10

0.01

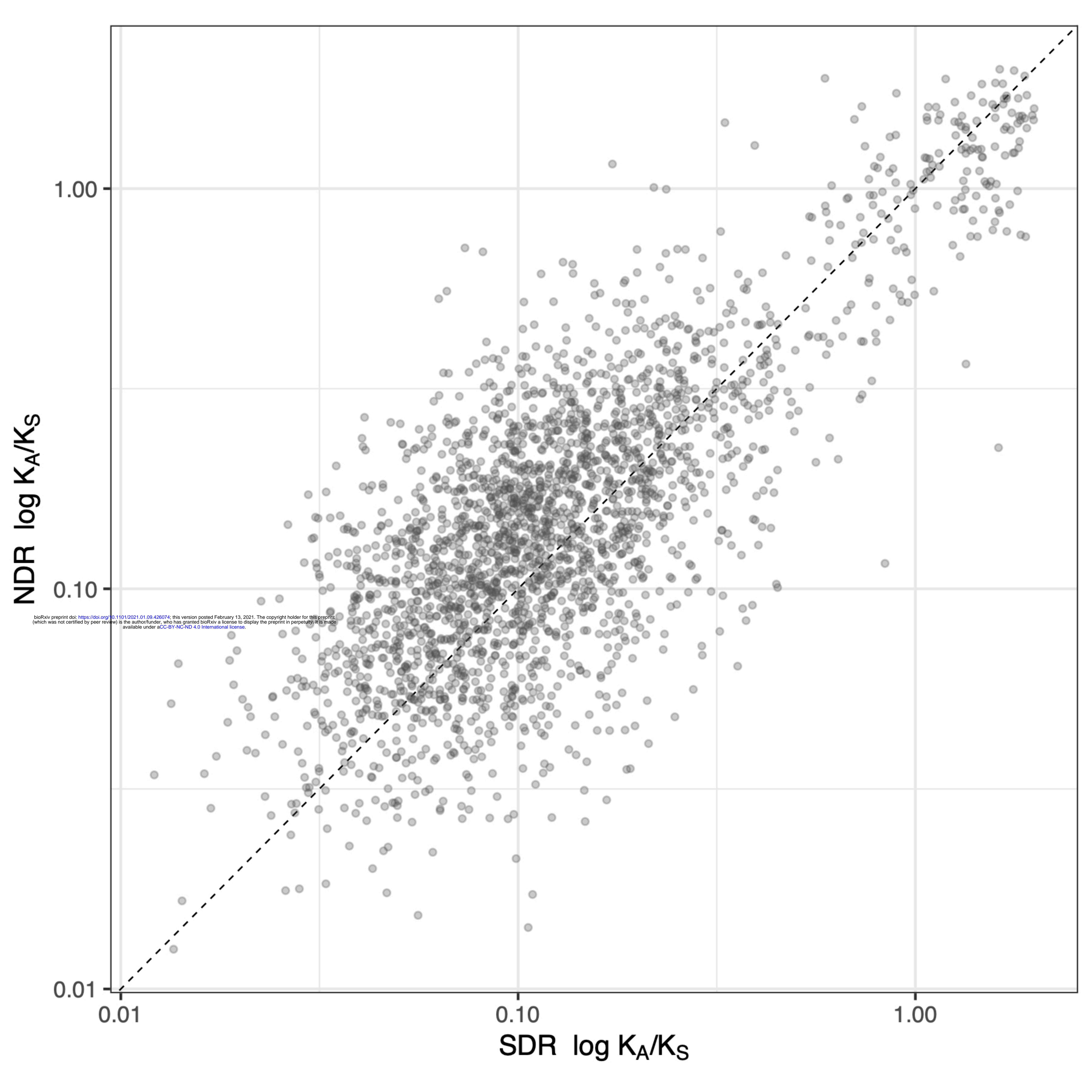
0.01

0.10

1.00

SDR $\log K_A/K_S$

bioRxiv preprint doi: <https://doi.org/10.1101/2021.01.09.426074>; this version posted February 13, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



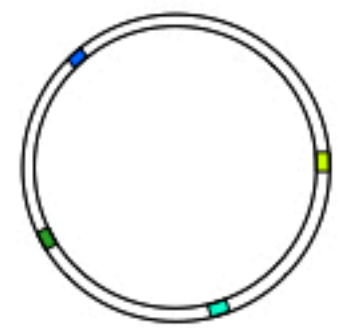
Genome Erosion (e.g. Endosymbionts)

bioRxiv preprint doi: <https://doi.org/10.1101/2021.01.09.426074>; this version posted February 13, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

Host environment restricts HGT, lowers cost of deletion

deletion \gg acquisition

Elimination of redundant genes



Negligible HGT, small N_e , strong genetic drift, favors **genome reduction**

Genome Streamlining (e.g. *Pelagibacter*)

HGT

deletion \geq acquisition

Elimination of non-adaptive genes

HGT, large N_e , strong selection, favors **genome reduction**

Genome Expansion (e.g. *Streptomyces*)

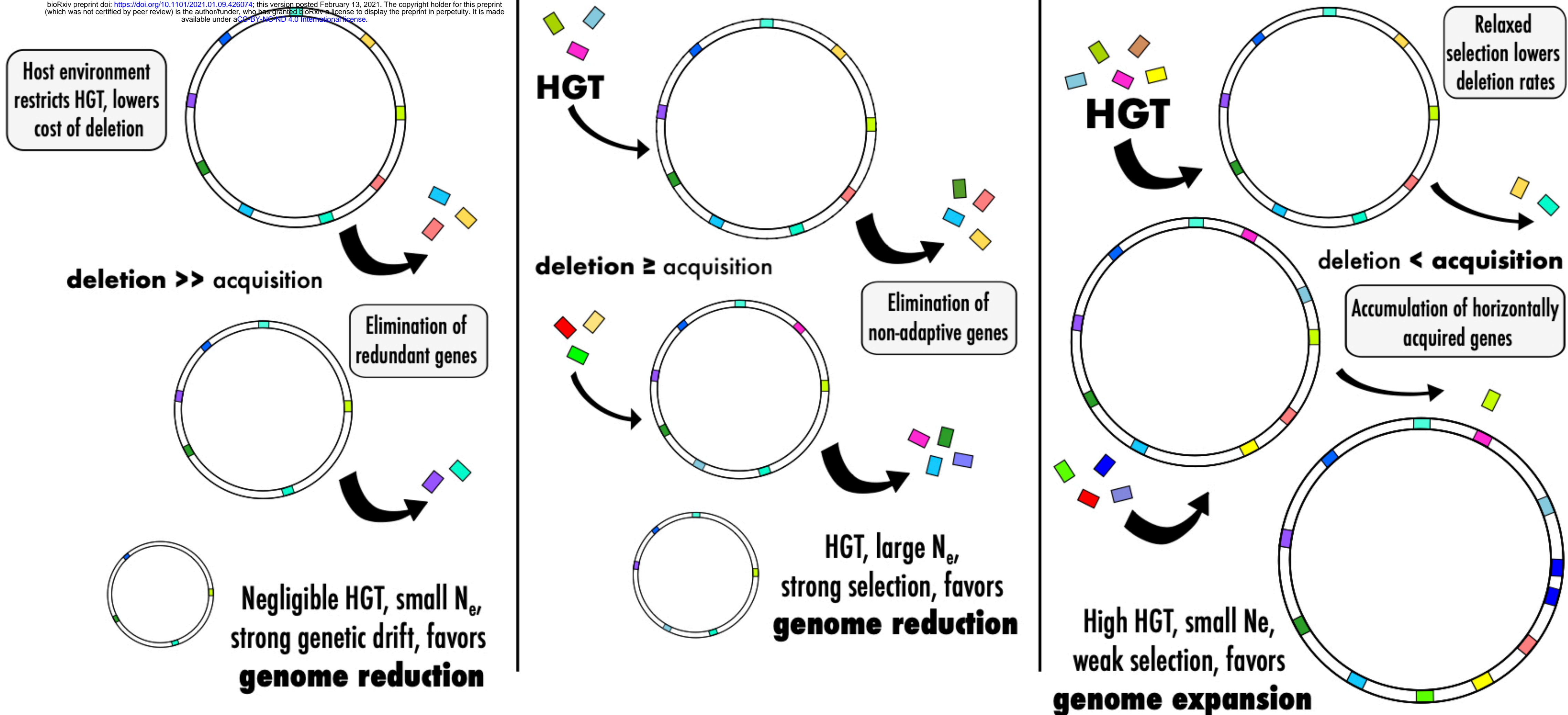
HGT

deletion \lt acquisition

Accumulation of horizontally acquired genes

High HGT, small N_e , weak selection, favors **genome expansion**

Relaxed selection lowers deletion rates



Supplementary Information for

A non-adaptive demographic mechanism for genome expansion in *Streptomyces*

MJ Choudoir, MJ Järvenpää, P Marttinen, and DH Buckley

Table S1. *Streptomyces* strains were isolated from 11 grassland sites across the United States.

Strain names begin with the site code referencing their isolation location. Mean annual temperature (MAT) reflects the 30-year average reported by NOAA.

| Sample Site Location | Code | Latitude (°N) | Longitude (°W) | MAT (°C) |
|------------------------|------------|---------------|----------------|----------|
| Manley Hot Springs, AK | man | 63.9 | -149.0 | -3.9 |
| Bothell, WA | wa | 47.7 | -122.2 | 11.4 |
| Astoria, OR | or | 46.2 | -123.9 | 10.8 |
| Rhinelander, WI | rh | 45.6 | -89.3 | 5.4 |
| Bear Creek, WI | sk | 43.4 | -90.1 | 7.6 |
| Brookfield, WI | b | 43.1 | -88.1 | 7.3 |
| Palo Alto, CA | st | 37.4 | -122.2 | 14.5 |
| Greensboro, NC | gb | 36.1 | -79.9 | 15.1 |
| Starkville, MS | ms | 33.5 | -88.8 | 17.2 |
| Austin, TX | t | 30.2 | -97.8 | 21.0 |
| Fort Pierce, FL | f | 27.5 | -80.4 | 22.9 |

Table S2. Descriptive attributes for 24 *Streptomyces* genomes (previously described in (1)).

NCBI accession numbers are associated with BioProject PRJNA401484. Strain names reflect sample sites (see Table S1). Clade membership includes the northern-derived (NDR) and southern-derive (SDR) clades, and the remaining four strains belong to independent (IND) lineages. Genome size in Mb. Genome-wide GC content (%). Number of open reading frames (ORFs). Protein-coding orthologous gene clusters (Genes).

| Strain | NCBI Accession | Clade | Size (Mb) | GC (%) | ORFs | Genes |
|--------|----------------|-------|-----------|--------|------|-------|
| b62 | SAMN07606143 | SDR | 7.82 | 71.63 | 7073 | 6568 |
| b84 | SAMN07606144 | NDR | 8.71 | 71.49 | 7657 | 7092 |
| b94 | SAMN07606145 | IND | 8.03 | 72.56 | 6939 | 6502 |
| f150 | SAMN07606147 | SDR | 8.06 | 71.60 | 7320 | 6794 |
| f51 | SAMN07606146 | SDR | 7.72 | 71.77 | 6851 | 6405 |
| gb14 | SAMN07606148 | SDR | 7.88 | 71.67 | 6998 | 6489 |
| man185 | SAMN07606149 | IND | 8.07 | 71.60 | 7244 | 6701 |
| ms115 | SAMN07606150 | SDR | 7.50 | 71.74 | 6776 | 6322 |
| ms184 | SAMN07606151 | SDR | 7.86 | 71.68 | 6958 | 6491 |
| or20 | SAMN07606153 | NDR | 9.08 | 71.54 | 8087 | 7474 |
| or3 | SAMN07606152 | NDR | 8.93 | 71.70 | 7956 | 7329 |
| rh195 | SAMN07606155 | NDR | 8.63 | 71.49 | 7804 | 7245 |
| rh206 | SAMN07606156 | NDR | 8.62 | 71.48 | 7817 | 7262 |
| rh207 | SAMN07606175 | NDR | 8.62 | 71.48 | 7825 | 7268 |
| rh34 | SAMN07606154 | NDR | 8.22 | 71.54 | 7392 | 6897 |
| sk226 | SAMN07606158 | IND | 7.96 | 72.51 | 6930 | 6505 |
| st115 | SAMN07606160 | IND | 8.35 | 71.75 | 7498 | 6966 |
| st140 | SAMN07606161 | SDR | 7.96 | 71.50 | 7184 | 6635 |

| | | | | | | |
|--------|--------------|-----|------|-------|------|------|
| st170 | SAMN07606162 | SDR | 8.08 | 71.42 | 7351 | 6801 |
| st77 | SAMN07606159 | SDR | 8.15 | 71.54 | 7346 | 6784 |
| t99 | SAMN07606163 | SDR | 7.78 | 71.67 | 7076 | 6607 |
| wa1002 | SAMN07606164 | NDR | 8.63 | 71.48 | 7577 | 7038 |
| wa1063 | SAMN07606165 | NDR | 8.84 | 71.39 | 7879 | 7351 |
| wa1071 | SAMN07606166 | NDR | 8.74 | 71.39 | 7755 | 7208 |

Table S3. Linear model summary. Table reports coefficient, standard error, t statistic, and *P*-value for explanatory variables. See Figure 2.

| | B | S.E. | t | <i>P</i>-value |
|-------------|----------|-------------|----------|-----------------------|
| (Constant) | -30696 | 2022 | -15.18 | < 0.001 |
| ANI | 37173 | 2066 | 17.99 | < 0.001 |
| Clade | 21220 | 3370 | 6.30 | < 0.001 |
| ANI x Clade | -21906 | 3448 | -6.35 | < 0.001 |

$R^2 = 0.82$, Adjusted $R^2 = 0.81$

Table S4. Pangenome frequency distributions for NDR and SDR clades. Table reports the number of genes (*n*) and the proportion (*prop*) of the total pangenome across frequencies for NDR and SDR. Frequency refers to the number of genomes a gene is present in, ranging from 1–10. Gene pools are categorized by gene frequencies. For example, intermediate-low genes are present in 3–4 strains, and intermediate-high genes are present in 6–8 strains. *P*-values are from a two proportion z-test, with Bonferroni adjustment for multiple comparisons, evaluating the null hypothesis that the proportion of genes at each frequency is the same for NDR and SDR. Significant *P*-values (< 0.0001) are in bold italics.

| Gene Pool | Frequency | NDR | | SDR | | <i>P</i> -value |
|-------------------------------|-----------|----------|-------------|----------|-------------|-----------------------|
| | | <i>n</i> | <i>prop</i> | <i>n</i> | <i>prop</i> | |
| Unique | 1 | 4132 | 0.302 | 4188 | 0.342 | <i>9.8E-11</i> |
| Rare | 2 | 926 | 0.068 | 1097 | 0.089 | <i>7.3E-10</i> |
| Intermediate - <i>Low</i> | 3 | 1182 | 0.086 | 494 | 0.040 | <i>3.3E-50</i> |
| | 4 | 881 | 0.064 | 377 | 0.031 | <i>3.3E-35</i> |
| | 5 | 533 | 0.039 | 260 | 0.021 | <i>1.5E-15</i> |
| Intermediate - <i>High</i> | 6 | 430 | 0.031 | 294 | 0.024 | 0.0032 |
| | 7 | 230 | 0.017 | 289 | 0.024 | 0.0012 |
| | 8 | 279 | 0.020 | 303 | 0.025 | 0.21 |
| Near Core | 9 | 854 | 0.062 | 557 | 0.045 | <i>2.0E-08</i> |
| Core | 10 | 4234 | 0.309 | 4400 | 0.359 | <i>3.7E-16</i> |

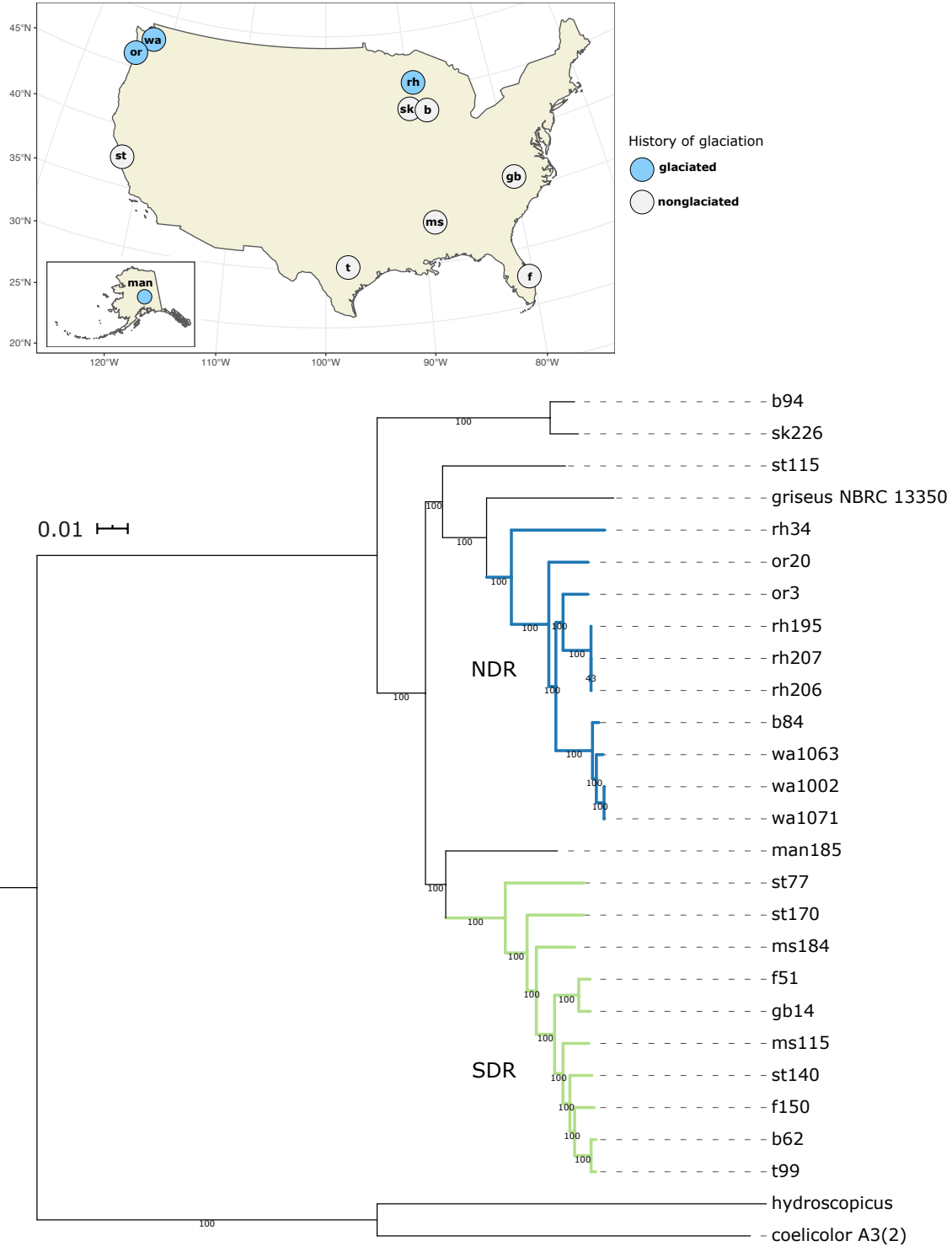


Figure S1. Whole genome phylogeny. *Streptomyces* northern-derived (NDR) and southern-derived (SDR) clades are recently diverged and originate from discrete latitudinal geographic ranges. Phylogenetic relationships were reconstructed from whole genome alignments using maximum likelihood and a GTRGAMMA model of evolution. The scale bar indicates nucleotide substitution per site, and nodes are labeled with bootstrap values. Strains are named according to their sample location (see Table S1). NDR clade branches are blue, and SDR clade branches are green. *Streptomyces griseus* NBRC 13350 is included as the closest taxonomic neighbor. The tree was rooted with *Streptomyces hydroscopicus* and *Streptomyces coelicolor* A3(2). Map shows the sample locations which are labeled with site codes (see Table S1) and colored according to the extent of historical glaciation (see (2)).

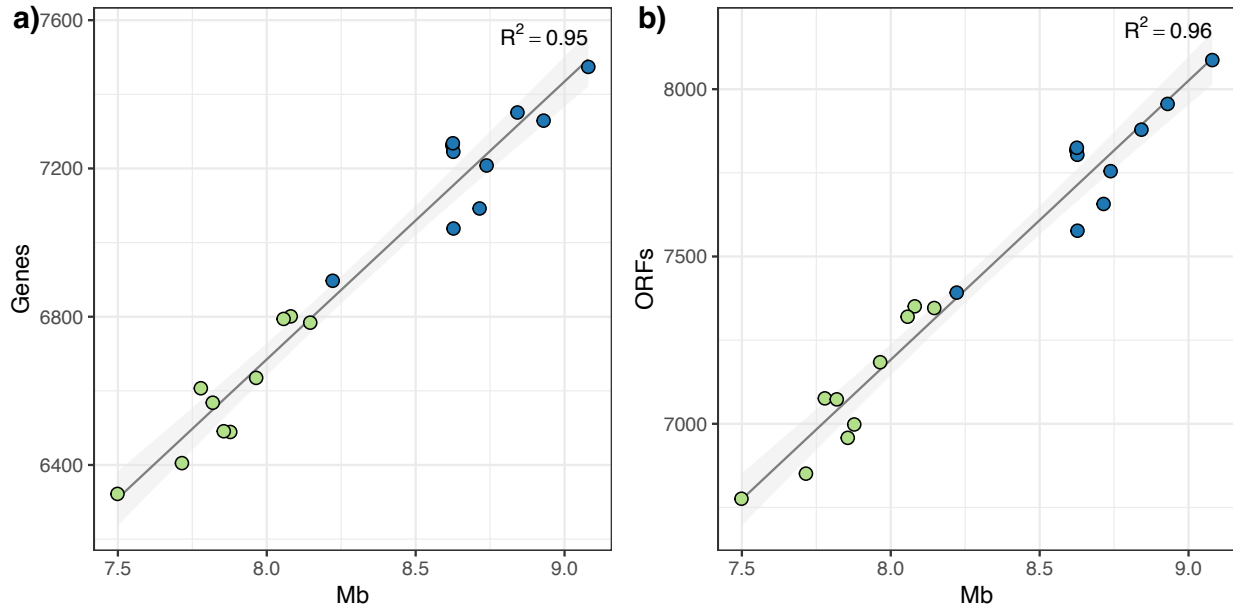


Figure S2. Genome size correlates positively with gene content. Plots show the relationship between genome size in Mb and number of genes (i.e., protein-coding orthologous gene clusters) (a) and open reading frames (ORFs) (b). Circles show values for individual *Streptomyces* genomes (NDR in blue and SDR in green). Lines show the linear regression, and the shaded area is the 95% confidence interval.

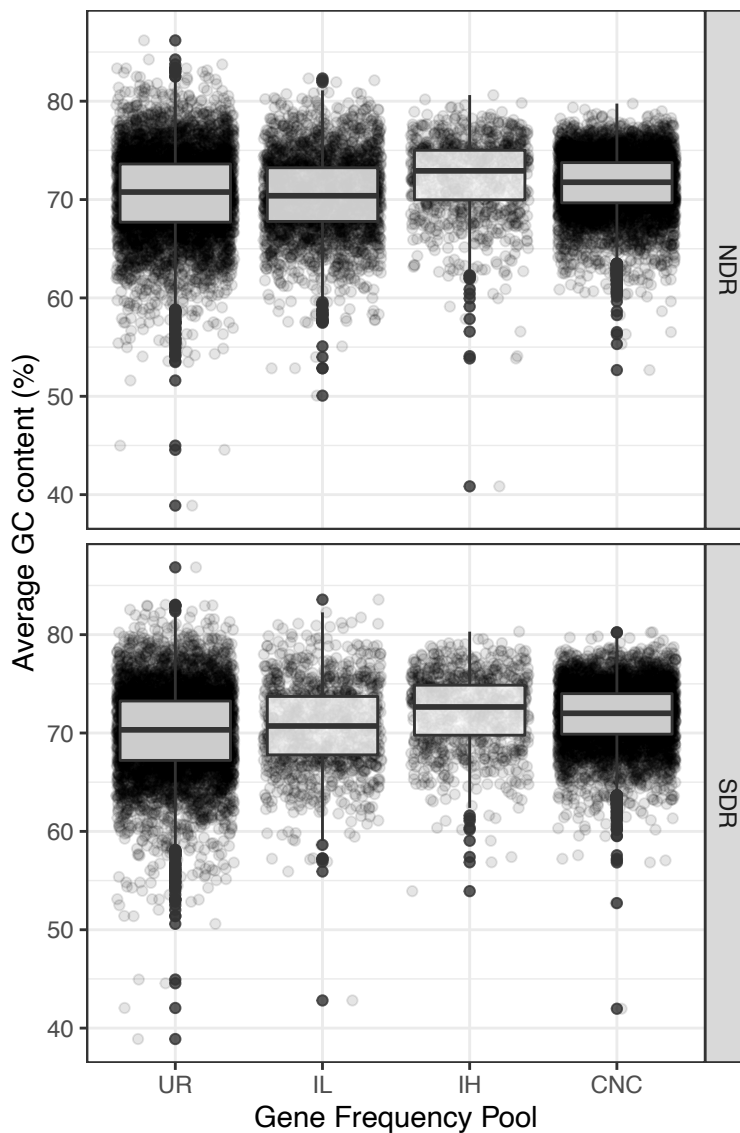


Figure S3. Per-gene GC content. For each gene frequency pool, plots illustrate the distributions of GC skew for NDR (top) and SDR (bottom). Circles show the average GC content (%) for each gene, and boxplots show the medians, interquartile ranges, and 1.5 times interquartile ranges. Gene pools are defined by frequency and include unique-rare (UR) (present in 1–2 strains), intermediate-low (IL) (present in 3–5 strains), intermediate-high (IH) (present in 6–8), and core-near-core (CNC) (present in 9–10 strains) (Table S4).

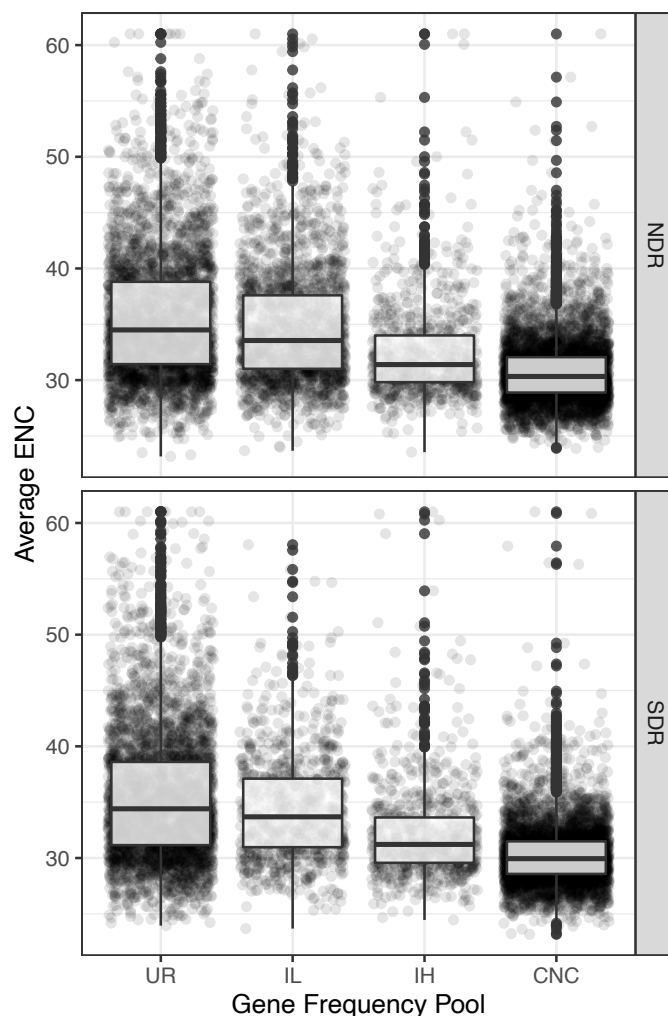


Figure S4. Per-gene codon usage bias. For each gene frequency pool, plots illustrate the distributions of codon usage bias as measured by the effective number of codons (ENC) (3) for NDR (top) and SDR (bottom). Circles show the average ENC for each gene, and boxplots show the medians, interquartile ranges, and 1.5 times interquartile ranges. ENC values range from 61 (indicating no codon usage bias, or all synonymous codons are used in equal frequency), to 2 (indicating extreme codon bias, or extreme favoring of certain codons). Gene pools are defined by frequency and include unique-rare (UR) (present in 1–2 strains), intermediate-low (IL) (present in 3–5 strains), intermediate-high (IH) (present in 6–8), and core-near-core (CNC) (present in 9–10 strains) (Table S4).

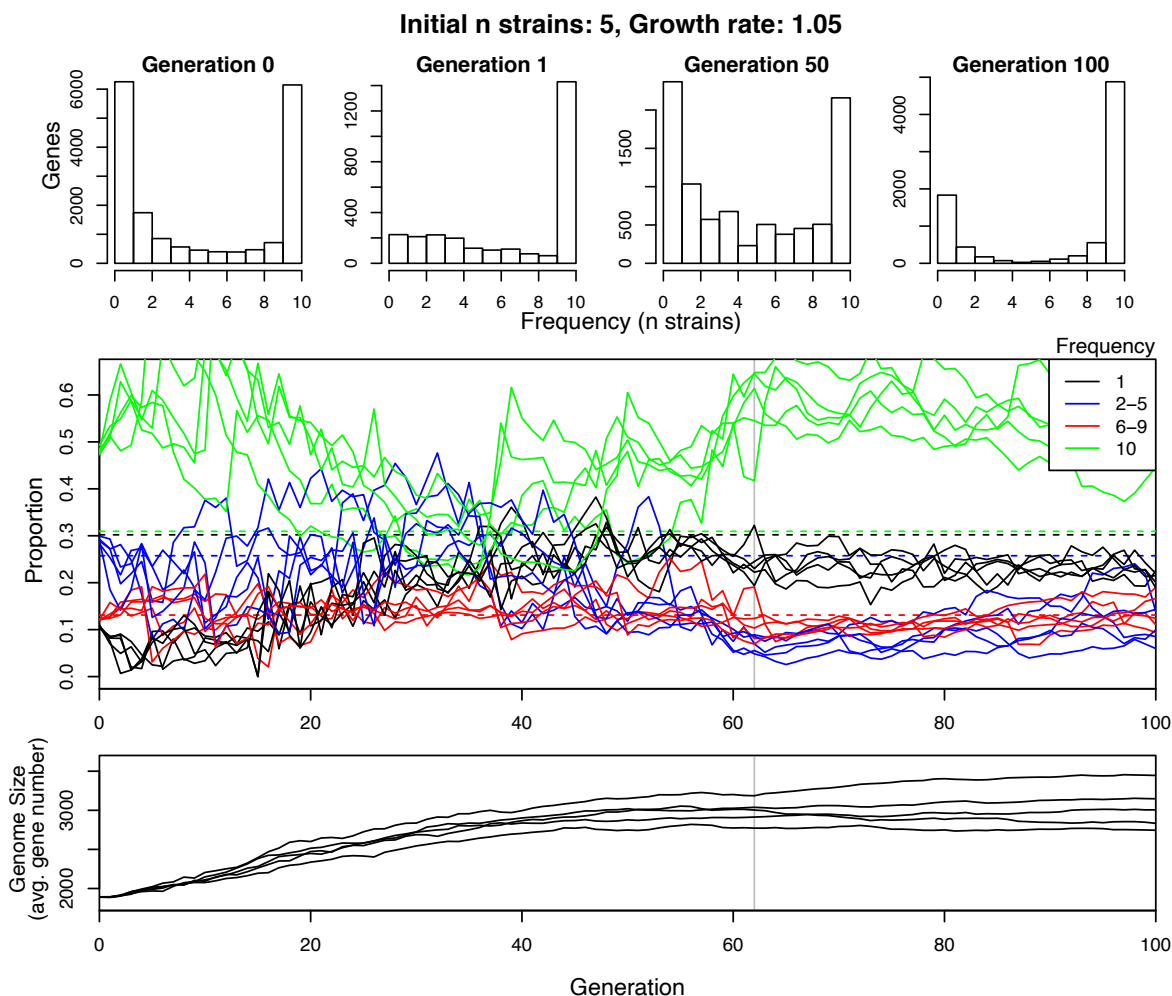


Figure S5. Demographic simulation. We simulated a population range expansion and modeled pangenome dynamics and genome size for 100 generations. *Top:* Pangenome gene frequency distributions from a single randomly selected simulation. Generation 0 shows the frequency distribution of the initial population before the bottleneck. Generation 1 shows the frequency distribution right after the bottleneck. Generations 50 and 100 show the frequency distributions at the subsequent generations. *Middle:* Trajectories of gene frequencies across 100 generations for 5 independent simulations. Colored lines show the proportion of genes within the pangenome present at different frequencies as according to legend. For example, blue lines are the proportions

of total genes present in 2–5 strains. Dashed horizontal lines are the observed proportions of the SDR clade (see Table S3). The vertical gray line indicates the generation when the population reached the maximum size. *Bottom*: Genome size. Lines show the average genome size for 5 independent simulations.

Supplementary References

1. Choudoir MJ, Buckley DH. 2018. Phylogenetic conservatism of thermal traits explains dispersal limitation and genomic differentiation of *Streptomyces* sister-taxa. *ISME J* 12:2176–2186.
2. Andam CP, Doroghazi JR, Campbell AN, Kelly PJ, Choudoir MJ, Buckley DH. 2016. A Latitudinal Diversity Gradient in Terrestrial Bacteria of the Genus *Streptomyces*. *MBio* 7:e02200–15.
3. Wright F. 1990. The “effective number of codons” used in a gene. *Gene* 87:23–29.