# A genomic perspective on genome size distribution across Earth's microbiomes reveals a tendency to gene loss — Source link ⧉

Alejandro Rodríguez-Gijón, Julia K. Nuy, Maliheh Mehrshad, Moritz Buck ...+3 more authors

**Institutions:** Stockholm University, Swedish University of Agricultural Sciences, Joint Genome Institute

Related papers:

- A genomic perspective across Earth's microbiomes reveals that genome size in Archaea and Bacteria is linked to ecosystem type and trophic strategy

- Ecological selection for small microbial genomes along a temperate-to-thermal soil gradient

- Diversity and Genomic Characterization of a Novel Parvarchaeota Family in Acid Mine Drainage Sediments.

- Metagenome-assembled genomes uncover a global brackish microbiome

- Correlations between bacterial ecology and mobile DNA.

1 **Title: A genomic perspective on genome size distribution across Earth's microbiomes**
2 **reveals a tendency to gene loss**
3
4 **Running title: Archaea and Bacteria genome size distribution**
5

6 Authors: Alejandro Rodríguez-Gijón #[a] (alejandro.gijon@su.se), Julia K. Nuy #[a]
7 (julia.nuy@su.se), Maliheh Mehrshad [b] (maliheh.mehrshad@slu.se), Moritz Buck [b]
8 (moritz.buck@slu.se), Frederik Schulz [c] (fschulz@lbl.gov), Tanja Woyke [c] (twoyke@lbl.gov),
9 Sarahi L. Garcia [a*] (sarahi.garcia@su.se).
10
11 [a] Department of Ecology, Environment, and Plant Sciences, Science for Life Laboratory,
12 Stockholm University, 10691 Stockholm, Sweden
13 [b] Department of Aquatic Sciences and Assessment, Swedish University of Agricultural Sciences,
14 Uppsala, Lennart Hjelms väg 9, 75651 Uppsala, Sweden
15 [c] DOE Joint Genome Institute, Berkeley, CA 94720, USA
16
17 #Equal contributions
18 *Corresponding author
19 Article type: Review article
20
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  **Abstract**
48  Our view of genome size distribution in Bacteria and Archaea has remained skewed as the data
49  used to paint its picture has been dominated by genomes of microorganisms that can be
50  cultivated under laboratory settings. However, the continuous effort to catalogue the genetic
51  make-up of Earth's microbiomes, specifically propelled by recent extensive work on
52  uncultivated microorganisms, provides a unique opportunity to revise our perspective on genome
53  size distribution. Genome size is largely a function of the expansion and contraction, by gain or
54  loss of DNA elements. While genome expansion provides microorganisms the capability to
55  acquire a wide repertoire of ecological functions, genome reduction increases the fitness of the
56  microorganisms to very specific niches. Capitalizing on a recently released large catalog of tens
57  of thousands of metagenome-assembled genomes, we here provide a comprehensive overview of
58  genome size distributions, suggesting that the known phylogenetic diversity of environmental
59  microorganisms possess significantly smaller genomes (aquatic bacteria average 3.1 Mb, host-
60  associated bacterial genomes average 3.0 Mb, and terrestrial bacteria average 3.8 Mb) than the
61  collection of laboratory isolated microorganisms (average 4.4 Mb). Moreover, the variation in
62  genome sizes across different types of environments reflects the different ecological and
63  evolutionary strategies used by microorganisms to thrive in their native environment. Finally, the
64  fact that genome sizes in Bacteria and Archaea remain relatively small might be a reflection of
65  the constraints imposed by selection and an overall dominance of gene loss as a survival strategy.
66
67  **Introduction**
68  Genomes are dynamic databases that encode the machinery behind evolution and adaptation of
69  living organisms to environmental settings. In brief, a genome encompasses all genetic material
70  present in one organism and includes both its genes and its non-coding DNA. Genome size is
71  largely a function of expansion and contraction by gain or loss of DNA fragments. The genomes
72  of extant organisms are the result of a long evolutionary history. In eukaryotes, an organism's
73  complexity is not directly proportional to its genome size which can have variations over 64,000-
74  fold (1, 2). However, the genome size ranges in Bacteria and Archaea are smaller and the
75  genomes are information-rich (3), and known to range from 100 kb to 16 Mb (4, 5). While
76  subject to genetic drift bacterial and eukaryotic genomes evolve in opposite directions. Bacteria
77  exhibit a mutational bias that deletes superfluous sequences, whereas Eukaryotes are biased
78  toward large insertions (6). In Bacteria and Archaea, evolutionary studies have revealed
79  extremely rapid and highly variable flux of genes (7) with evolutionary forces acting on
80  individual genes (8). On one hand, mechanisms for genome expansion are promoting the gain of
81  new functions through horizontal gene transfer, *de novo* gene birth, and gene duplications (9,
82  10). On the other hand, it is known that the primary driving forces for genome reduction are
83  metabolic and spatial economy and cell multiplication speed (11, 12).
84
85  Our overall view of the diversity, distribution, and genome characteristics of Bacteria and
86  Archaea have remained biased for most of the microbial ecology history. These biases stem
87  chiefly from the "great plate count anomaly" because for more than a decade, the genomes that
88  were sequenced were primarily from laboratory isolates (13, 14). More than two decades have
89  passed since the first bacterial genomes were completely sequenced (15, 16). In the first decade
90  of genome sequencing about 300 bacterial genomes and two metagenomic projects with
91  assembled genomes were published (17). Since then, rapid advances in metagenome sequencing
92  and data analyses have enabled large-scale cataloging of bacterial and archaeal genomes from a

93    wide range of environments (18-20). With all sequencing efforts the representation of number of
94    genomes from phylogenetically diverse groups of Bacteria and Archaea has greatly increased.
95    The Genome Taxonomy Database (GTDB) include 194,600 genomes, with 31,910 of those
96    being species representatives and 8,792 of those species representatives are based on published
97    named species (21). Genome catalogs such as Genomes from Earth's Microbiomes (GEMs)
98    contain ~52,500 genomes all of them being metagenome-assembled genomes. Using these novel
99    resources, it is now possible to obtain an updated view of microbial genome characteristics,
100   diversity, and distribution of microbes in the environment.

101

102   Genome size and its evolution has been studied by many researchers who each focused on
103   different taxonomic lineages or different ecological or evolutionary backgrounds (8, 12, 22-25).
104   As microbial researchers, how do we define what is a small genome or a big genome? Perhaps,
105   researchers working on model organisms such as *Escherichia coli* with a genome size of ~5 Mb
106   (26), would define 'big' or 'small' very differently to researchers working on *Prochlorococcus*
107   with a genome size of ~2 Mb (27), soil-dwelling *Minicystis rosea* with a genome size of 16 Mb
108   (5) or bacterial endosymbionts of insects that may have genomes merely larger than 100 kb (4).
109   The recently published expanded database of environmental bacterial and archaeal genomes (18)
110   allows us to revisit and acquire a more complete understanding of genome size distribution
111   across different environments in higher resolution. In this review, we provide an overview of the
112   evolutionary and ecological drivers behind the different genome sizes of Bacteria and Archaea.
113   Moreover, we offer an overview of the distribution of genome sizes of all known bacterial and
114   archaeal phyla across different environments. We found that while there are phyla with
115   consistently smaller genome sizes (< 2 Mb), such as Caldisericota, Aenigmarchaeota,
116   Micrarchaeota, Nanohaloarchaeota, and Ianarchaeota, 78.4% of bacterial and archaeal genomes
117   recovered through genome-resolved metagenomics represent estimated genome sizes below 4
118   Mb.

119

120   **Extant genome size distribution in the environment**
121   The current state of environmental sequencing, assembly, and binning technologies allows us to
122   review and renew our view of bacterial and archaeal genome size distribution on Earth (18). To
123   minimize representation biases (28), from the ~52,500 genomes we included one representative
124   per mOTU, defined by 95% average nucleotide identity (ANI), from the GEMs environmental
125   MAGs resulting in ~15,000 MAGs (Figure 1A). We complemented these data by adding ~8,000
126   species cluster representatives from >90% complete genomes of isolates from GTDB (Figure 1).
127   GEMs reported that MAGs in the same species than isolate genomes were consistent in size
128   (average estimated genome length per OTU MAGs = -0.17 + 1.01 average estimated genome
129   length per OTU isolates, r=0.95) (18).While this suggests that there is not a big bias in
130   metagenome assembly and binning it is important to keep in mind that MAG assembly might
131   discriminate against ribosomal RNAs, transfer RNAs, mobile element functions and genes of
132   unknown function (29).

133

134   Furthermore, we compared the genome size distribution of all environmental MAGs versus that
135   of their taxonomic relatives from cultivated isolates, as derived from the GTDB. The genomes
136   from bacterial isolates have the average genome size of 4.4 Mb. When comparing this genome
137   size distribution from isolates with that of the environmental MAGs, the first striking observation
138   is that environmental MAGs have significantly lower genome sizes (t-test Bacteria p>2e-16

139 Archaea p>2e-16). Environmental aquatic bacteria average 3.1 Mb, host-associated
140 metagenome-assembled bacterial genomes average 3.0 Mb, and terrestrial bacteria average 3.8
141 Mb (Figure 1A). A reason for the difference in genome size between isolates and environmental
142 microorganisms might be the tendency to sample different types of microorganisms with culture
143 dependent and independent methods (30). For example, it is known that current cultivation
144 techniques with rich media bias cultivation towards copiotroph microorganisms (31). Moreover,
145 microorganisms in nature do not live in isolation but instead have coevolved with other
146 microorganisms and might have specific requirements that are hard to meet in batch-culture
147 standard-media isolation techniques (32). Other reasons for biases in cultivation include slow
148 growth of microorganisms (33), host dependence (34), dormancy (35), and microorganisms with
149 very limited metabolic capacity (36) among others. Innovations to culturing the uncultured
150 microbial majority might help breach this genome size gap in the future (37).
151
152 Placing bacterial and archaeal genome sizes in the context of a phylogenetic tree (Figure 2 and 3)
153 shows that the distribution of representative genomes and its sizes vary widely not only between
154 different phyla but widely within different phyla. To this view, we want to bring into the
155 discussion the biphasic model of evolution (25). In this model it is discussed that genome
156 evolution occurs in two phases. One phase involves gene gains that occur in bursts and are
157 associated with the emergence of novel microbial groups. The other phase involves gene loss
158 that occurs gradually. In the extant phylogenetic tree of Bacteria and Archaea it is noticeable
159 how closely related species are shaped by the different genetic processes that influence genome
160 size (Figure 2 and 3).
161
162 **Genetic processes that shape genome size**
163 The variability in genome size that we observe across different microbial taxa is the result of the
164 reached equilibrium between gains and losses of genetic information (Figure 4). The
165 evolutionary events that drive these changes are diverse. Some lineages follow a highly
166 mutational mode of evolution (38) while other lineages have recombination as a stronger
167 evolutionary force (39-42). The acquisition of new genetic information and metabolic capacities
168 is often accompanied by the expansion of gene families. *In silico* studies indicate that the
169 acquisition of new genes could have a vital role in adaptation (43). Moreover, a strong
170 correlation has been observed between genome size with gene family expansions and length of
171 non-coding sequences in complex cyanobacteria (44). The most important evolutionary events
172 involved in genome expansion processes are *de novo* gene birth, the duplication of genes, and
173 Lateral/Horizontal Gene Transfer (LGT/HGT). *De novo* gene birth is the process by which new
174 genes emerge from non-genic DNA sequences (45). However, most of the known examples of
175 this process are found in eukaryotes. Furthermore, comparative genomics of some bacterial
176 taxonomic lineages has suggested that HGT is more relevant on the expansion of bacterial
177 metabolic networks than gene duplications (46). HGT can foster the acquisition of new functions
178 while duplications relate to a higher gene dosage (47). However, phylogenomic analysis of other
179 lineages such as Nitrososphaerales (Thermoproteota) indicate a predominant role of gene
180 duplications over HGT (48). These examples highlight how HGT and duplications aid
181 microorganisms into adaptation to their niches. For example, different Archaea have shown
182 modifications in their metabolic potential through these genome expansion processes (49-51). In
183 a nutshell, genome expansion processes can influence gene dosage, acquisition of new ecological
184 capacities and adaptation in both, Archaea and Bacteria.

185     Conversely, genome reduction fosters the development of more compact genomes (Figure 4).
186     There are three main processes involved in gene loss: genetic drift, pseudogenization and
187     streamlining. Genetic drift describes stochastic changes on the gene repertoire variants.
188     Mutations which are biased towards deletions over time promote genome reduction (8). Genetic
189     drift is more pronounced in species that have a small effective population size such as host-
190     associated endosymbiotic microorganisms. As an example, endosymbiotic lineages of
191     Gammaproteobacteria such as in *Buchnera aphidicola* have lost ample genes that have already
192     reached stasis (52). When a gene loses its original function, it is often turned into a pseudogene
193     (53). A pseudogene is a derived form of regular genes that might present a different function or
194     turn obsolete. Comparative analyses of archaeal genomes show that up to 8.6% of their genomes
195     are constituted by pseudogenes which usually present at least one inactivating mutation (54).
196     Moreover, pseudogenization has been suggested to be a special type of gene loss when
197     adaptation to new ecological niches is needed. In the *Roseobacter* lineage (class
198     Alphaproteobacteria), this process was correlated to switches in resource recovery, energy
199     conservation, stress tolerance and different metabolic pathways (55). Finally, streamlining is the
200     process of gene loss through selection and it is mainly observed in free-living microorganisms
201     with high effective population sizes. Streamlining creates a series of distinct patterns, such as
202     increase in nutritional connectivity between individuals, reduction of genome size, lower GC
203     content and higher coding gene density (12). Aquatic microorganisms have been used as
204     exemplary cases of streamlining in which many have gone through community adaptive
205     selections and gene loss (56). In fact, their gene loss goes so far  that these free-living aquatic
206     microorganisms depend on community associations and thus thrive in functional cohorts (57).
207     The renewed view of genome sizes and characteristics confirms that genomes from aquatic
208     microorganisms have a higher coding density compared to those from other ecosystems (Figure
209     1).

210     Both genome expansion and reduction have a vital role in the evolution of microorganisms.
211     However, these two processes are not in perfect equilibrium. While genome expansion might
212     allow cells to become highly flexible in terms of developmental capacities and physiological
213     performance, gene loss allows cells to become highly successful in particular niches (44).
214     Moreover, gene loss dominates in the evolutionary history of Bacteria and Archaea (25). For
215     example, *in silico* studies of 34 bacterial genera and one archaeal genus show that the rate of
216     gene gains is three times lower than that of gene loss (7). In the same study, highly dynamic
217     genomes were found presenting these evolutionary events 25 times more often than the most
218     stable genomes. This tendency that gene loss is more prevalent than gene gain has also been
219     described in short term in host-associated *Pseudomonas aeruginosa* (Gammaproteobacteria)
220     (58). There, clinical isolates of *P. aeruginosa* indicate gene loss rates six times greater than gene
221     acquisition during the first year of a chronical infection as an adaptive strategy to avoid the
222     host's immune response. Even evolutionary reconstructions of the Last Common Ancestor of
223     Archaea show that genomes of early Archaea were more complex and thus gene loss played
224     likely a critical role in their evolution (59, 60). In summary, these examples illustrate the synergy
225     of both evolutionary processes, with genome expansion providing microorganisms the capability
226     to acquire a wide repertoire of ecological adaptations and genome reduction increasing the
227     fitness of the microorganisms to very specific niches (Figure 4).

228
229

230 **Environmental impact on genome size in different taxonomic lineages**
231 The most up-to-date view of genome sizes on Earth provided here shows that the distribution of
232 genomes from terrestrial environments average at size of 3.8 Mb (Figure 1). The sub-ecosystems
233 considered in this view are soil, and deep subsurface among others (Figure 5). Terrestrial
234 microorganism's genome size is larger than what is commonly found in aquatic and host-
235 associated ecosystems. However, it is smaller than expected based on previous metagenomic
236 predictions which placed the genome size of soil bacteria at 4.74 Mb (61). Trends of larger
237 genome sizes in soil have been hypothesized to be related to scarcity and diversity of nutrients,
238 fluctuating environment combined with little penalty for the slow growth rate (23, 62, 63). In
239 fact, although terrestrial or soil environments are physically structured, they are generally
240 characterized by two to three orders of magnitude greater variations (in temperature and
241 currents) than marine environments (64). *In silico* studies predict that large genome sizes could
242 be the result of higher environmental variability (65). A recent example showed that isolates of
243 terrestrial Cyanobacteria have larger genomes (6.0-8.0. Mb), as compared to their freshwater
244 counterparts (4.0-6.0 Mb) and their relatives originating from the marine environment (1.5-2.5
245 Mb) (62). The general theory of genome expansion states that the genetic repertoire increases to
246 allow microorganisms to gain adaptive capacities to face perturbations and survive in variable
247 environments. Despite these general trends showing larger genome sizes in terrestrial
248 environments, it is worth noting that streamlined microorganisms such Patescibacteria (Fig 1B)
249 as '*Candidatus* Udaeobacter copiosus' (Verrucomicrobiota) are abundant in soils (66).
250
251 Some of the most numerically abundant and streamlined microorganisms known to date, such as
252 Pelagibacter (class Alphaproteobacteria) (12), marine methylotrophs (class
253 Gammaproteobacteria) (67), *Prochlorococcus* (phylum Cyanobacteria) (27) Thermoproteota (68)
254 and Patescibacteria (69) are commonly found in aquatic niches. This is well reflected in the
255 MAG data, illustrating that genomes from aquatic sites are among the smallest (Figure 1).
256 Aquatic environments are less physically structured than soils. However, there is some vertical
257 structure in physicochemical parameters connected to depth variables such as light penetration,
258 temperature, oxygen, and nutrient gradients, as well as microscale spatial structure due to the
259 presence of heterogeneous particles. Aquatic structures are drivers of the genetic repertoire of
260 aquatic microorganisms. For instance, metagenomic sequencing reported the increase of genome
261 sizes for Bacteria and Archaea with increasing depths (70). While several hypotheses have been
262 proposed as drivers of such evolutionary trends, nutrient limitation might be one of the central
263 factors determining genomic properties (71) (Table 1). Temperature might be as important, for
264 example, a study based on twenty-one Thermoproteota and Euryarchaeota fosmids
265 (Euryarchaetoa is now reclassified into Methanobacteriota, Halobacteriota and
266 Nanohaloarchaeota) showed high rates of gene gains through HGT to adapt to cold and nutrient-
267 depleted marine environments (72). Moreover, aquatic hyperthermophilic microbes show
268 reduced genomes compared to those of microorganisms adapted to very cold environments (73-
269 75) supporting this negative relation between temperature and genome size. One last driver we
270 want to point out in aquatic environments is light which decreases with depth. Photosynthetic
271 bacteria such as *Prochlorococcus* spp. are well differentiated into a high-light adapted ecotype
272 with smaller genome sizes (average 1.6 Mb), and a low-light-adapted ecotype with slightly
273 bigger genome size (average 1.9 Mb) (76).
274
275 In host-associated microbiomes, microorganisms are shaped in their ecological and evolutionary

6

276     history by the differing levels of intimacy they might have with their host. For example, within
277     the Chlamydiaceae family there are lineages that have evolved intracellular associations with
278     eukaryotes (77, 78). Recent metagenomic studies uncovered that extensive metabolic capabilities
279     were present in the common ancestor of environmental Chlamydiia (class) and subsequently lost
280     in Chlamydiaceae (79). Moreover, host-associated bacterial genomes show a variation in size
281     depending on the type of host (plant, animal, etc.) and the type of association they have with the
282     host (endosymbiotic, ectobiotic or epibiotic). Generally, microorganisms associated with
283     Arthropoda (52), humans (80) and other mammals show smaller genomes whereas protist- and
284     plant-associated bacteria present bigger genomes (81) (Figure 5). *In silico* studies of
285     Alphaproteobacteria show massive genome expansions diversifying plant-associated Rhizobiales
286     and extreme gene losses in the ancestor of the intracellular lineages Rickettsia, Wolbachia,
287     Bartonella and Brucella that are animal- and human-associated (82). Within the Chloroflexota,
288     genomes associated with plants or algae range between 4.75 and 7.5 Mb, and genomes
289     associated with Arthropoda range between 0.75 and 1.75 Mb (Figure S1). Although
290     microorganisms that are host-associated are widely known for their reduced genomes, the
291     characteristics of metagenomic host-associated bacterial genomes show lower coding density
292     than streamlined genomes in aquatic environments in the genome sizes ranged 1 – 4 Mb (Figure
293     1F). However, at size range below 1 Mb the MAGs and available genomes of endosymbionts are
294     often reduced and at same time have high coding density of ~91% (Figure 1F) (83).
295
296     **Table 1**
297

| Chemical, physical or biological variable influencing genome size | | Taxa | References |
|---|---|---|---|
| Temperature | *Literature review indicates a negative correlation between genome size and temperature.* | | |
| | Comparative genomic of genomes of hyperthermophilic microorganisms shows average genome sizes of about 2.3 Mb with very active horizontal gene transfer (HGT) mechanisms | *Thermus thermophilus* (phylum Deinococcota) *Thermus spp.* | (73, 74) |
| | Metagenomics suggest that gene gains would have played an important role in adaptation to low temperature and oligotrophic deep marine environments | Thermoproteota and Euryarchaeota (phyla) | (72) |
| | Comparative genomics of isolates in one genus indicate larger genomes in colder environments. | *Janthinobacterium* spp. (class Gammaproteobacteria) | (75) |
| | Environmental samples indicate that hypersaline environments could increase gene gain via HGT, whereas thermal environments decrease it. | Halobacteria and Thermoproteia (class) | (84) |
| Nutrients | *When talking about nutrients, diversity and quantity of nutrients are two factors that drive ecology and evolution. Some literature present conflicting results on how these two dimensions of nutrient influence genome sizes.* | | |
| | Metagenomics indicate dominance of reduced genomes in the Baikal Lake. Small genomes are thought to reflect the extremely oligotrophic conditions. | Actinobacteria, Bacteroidetes, Cyanobacteria Verrucomicrobia and Thermoproteota | (85) |
| | Online databases indicate that larger | 70 closely related bacterial | (23) |

| | | | |
|---|---|---|---|
| | genome-sized species may dominate environments where resources are scarce but diverse. | genomes | |
| | Phylogenomics of isolates show gene loss in functions like resource scavenging and energy acquisition when adapting to nutrient-rich environments in algae and corals. | Roseobacter spp. (class Alphaproteobacteria) | (55) |
| | Oceanic metagenomic data show positive correlation between nutrient concentration and genome size. | Different bacteria phyla | (86) |
| | Metagenomics indicates small genomes in mesopelagic environments are the result of adaptation to energy scarcity. | Some Thermoproteota (phylum) | (68) |
| | Whole-genome shotgun sequencing indicated that deep oligotrophic marine environments are dominated by large genomes with high GC content. | Lactobacillales (phylum Firmicutes) | (87) |
| | Oceanic metagenomic samples suggest that deeper areas with more nitrate and phosphate as nutrients are dominated by large genomes and high GC content. | Bacteria (SAR11, *Prochlorococcus* spp., *Roseobacter* spp., etc.,) and Archaea (Thermoproteota and Euryarchaeota) | (70) |
| Light | *In oxygenic phototrophs there is negative correlation between light irradiance and the genome size.* | | |
| | Genomes of cultures and single cells show high-light-adapted ecotypes with smaller genome sizes and low-light-adapted ecotypes with bigger genomes. | *Prochlorococcus* spp. (phylum Cyanobacteria) | (27, 76, 88) |
| Particles | *Microorganisms with particle associated lifestyle tend to have larger genome sizes.* | | |
| | Comparison of metagenomes in coastal ecosystems show larger genome sizes for particle associated microorganisms than free-living. | Metagenomic data | (89) |
| | Particle associated microbes have larger genome sizes than free-living bacteria. | Cyanobacteria and Bacteroidetes | (86) |
| Host-association | *Host-associated bacterial genomes show a variation in size depending on the type of host (plant, animal, etc.) and the type of association they have with the host (endosymbiotic, ectobiotic or epibiotic)* | | |
| | In silico studies indicate massive genome expansions in plant-associated bacteria. | Alphaproteobacteria (class) | (82) |
| | Isolates from sugarcane (*Saccharum* sp.) rhizosphere and endophytic roots and stalks show 26 individual genomes of associated bacteria whose genomes ranged from 3.9 to 7.5 Mbp. | Diverse bacterial taxa (Burkholderiaceae, Rhizobiaceae, Caulobacteraceae, Xanthomonadaceae, etc.) | (90) |
| | Genomic comparison of 3837 bacterial genomes identified thousands of plant-associated gene clusters and found genomes of plant associated microorganisms tended to be larger | Diverse bacterial taxa | (81) |
| | Intense genome reduction in isolates of microbes associated with aphids (Arthropoda). | *Buchnera aphidicola* (class Gammaproteobacteria) | (52) |

8

| | *In vitro* cultures and metagenomic datasets indicate reduced genome sizes in microbes associated with humans and other mammmals | *Salmonella enterica* (class Gammaproteobacteria) Patescibacteria (phylum) | (80, 91) |
|---|---|---|---|
| | Environmental samples indicate that symbionts and epibionts of other microbes present highly reduced genomes. | Bacteria of the CPR clade (such as *Vampirococcus lugosii*) and Archaea of the DPANN | (92, 93) |
| Viruses | Marine isolates support the "Cryptic Escape Theory". In here small cell size is a strategy to minimize viral predation. This article also finds a correlation between genome size and cell size. | Different bacteria lineages (Cyanobacteria, Proteobacteria, Actinobacteria, among others) | (94) |

298
299 **Conclusion**
300 Since the sequencing of the first isolate bacterial genomes in 1995, profound improvements in
301 both sequencing technologies and bioinformatic analysis tools have accelerated our access to the
302 genetic make-up of the uncultivated majority. This allowed us for the first time to provide a more
303 global view of the distribution of bacterial and archaeal genomes from a wide array of
304 microbiomes on Earth. In this review, we offer an overview where genomes obtained from
305 environmental samples show to be smaller than those obtained from laboratory isolates. This is
306 not because isolates and MAGs from the same species differed in size but because cultivation
307 methods bias the sampling of nature towards obtaining copiotrophs, fast growers, and more
308 metabolically independent microorganisms. Moreover, we find the distribution of genome sizes
309 across the phylogenetic tree of Bacteria and Archaea reflects that genome evolution occurs in a
310 gene gain phase and gene loss phase, as the biphasic model theory suggests. Finally, we review
311 the ecological and evolutionary effectors causing the varying sizes of genomes in different
312 environments. Soils might have the microorganisms with the bigger genome sizes due to higher
313 environmental variability. Genomes in aquatic environments might be shaped by vertical
314 stratification in nutrients, particles, and light penetration. Host-associations might shape genomes
315 differentially based on the kind of relationship between the microorganisms and the host. We
316 expect that as the microbial ecology field keeps moving forward, we get a deeper resolution on
317 physicochemical, spatial, and biological drivers of bacterial and archaeal genome sizes.

318
319
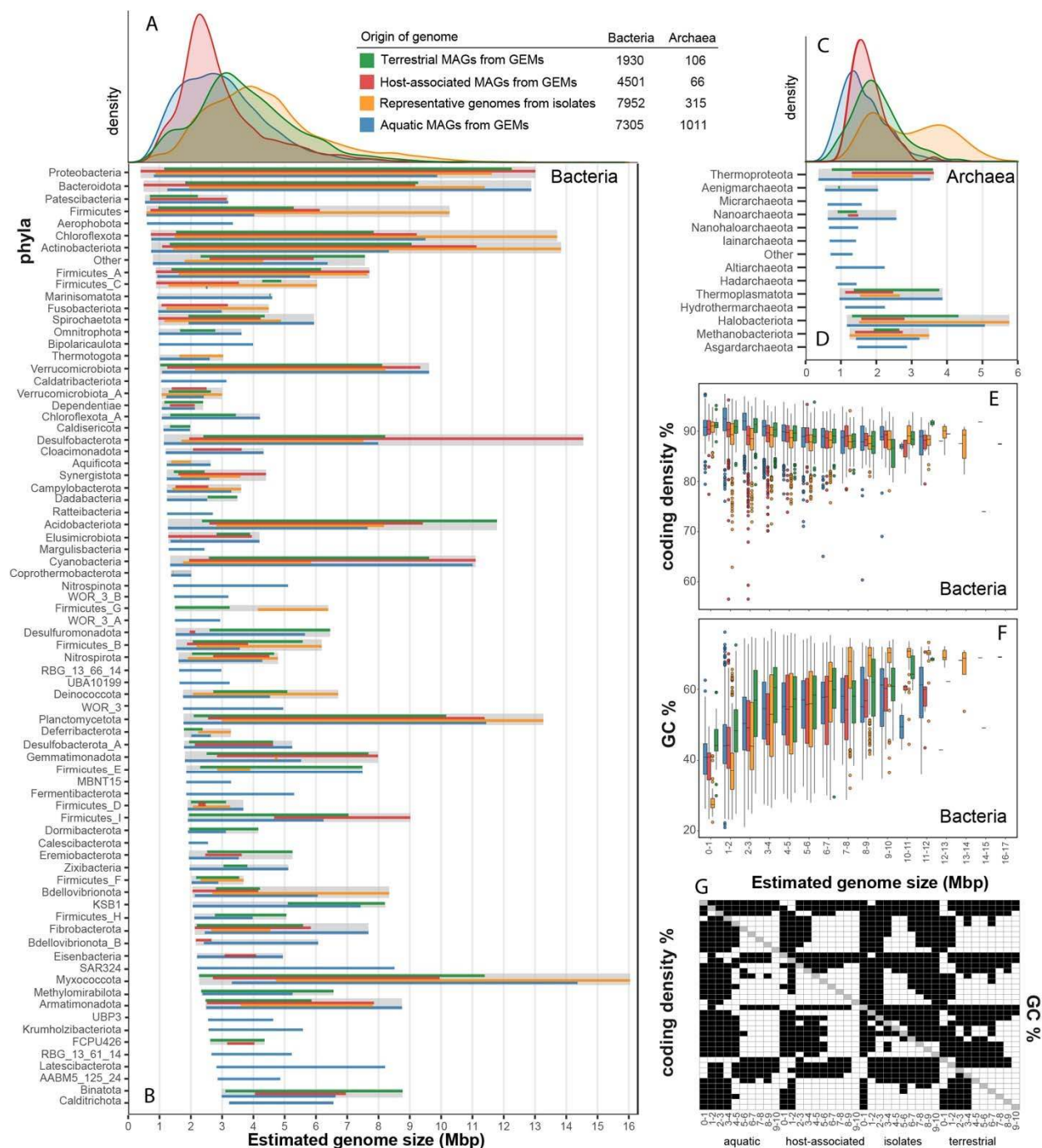320
321
322
323
324
325
326
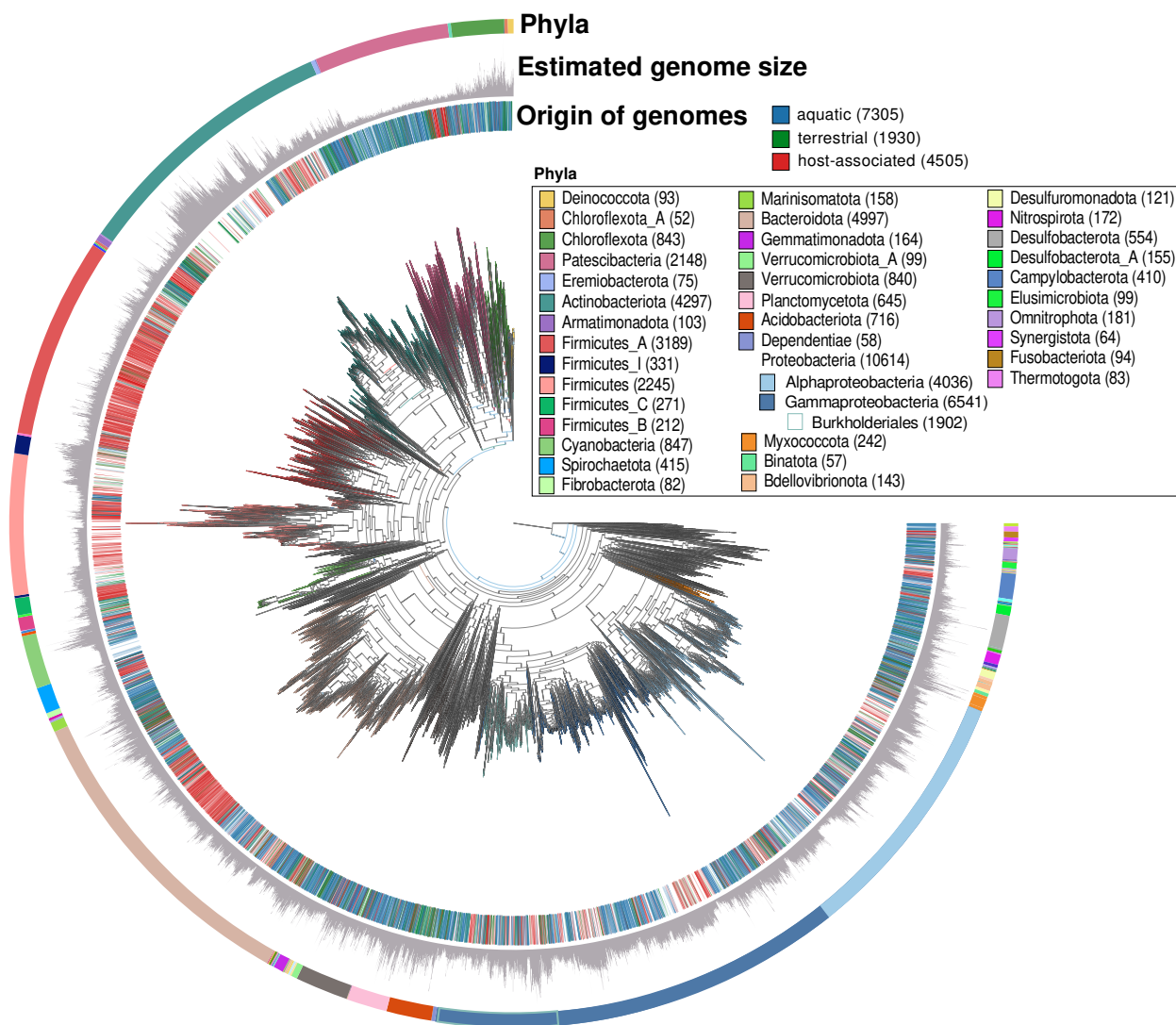327
328
329
330
331 **Figures**
332

9

333
334
335 **Figure 1**. Overview of the genome size distribution across Earth's microbiomes. Genome size
336 distribution of Bacteria [A] and Archaea [C] from different environmental sources and across
337 different bacterial [B] and archaeal [D] phyla is shown for a total of 23,186 genomes. The coding
338 density [E] and GC content (%) [F] is shown for the bacterial MAGs across different
339 environments and isolates. Pair-wise t-test was performed in all variables of panel E and F and if
340 the pair-wise comparison was significant (p < 0.05) it is shown in panel G in black. The figure
341 was constructed in R (95) using representative isolate genomes from GTDB database as well as
342 MAGs (metagenome assembled genomes) from GEMs catalog. The GEMs genomes available

10

343     were clustered into mOTUs (metagenomic operational taxonomic unit) at the threshold of the
344     operational definition of species (95% ANI). To eliminate over-representation biases for some
345     mOTUs, we used only one representative genome per mOTU from the GEMs catalog in the
346     plots. We addressed the same bias for the GTDB database by selecting the representative isolate
347     genome per species cluster that were circumscribed based on the ANI (>=95%) and alignment
348     fraction ((AF) >65%) between genomes (21). To construct the figures, we plotted the estimated
349     genome sizes which was calculated based on the genome assembly size and completeness
350     estimation provided. In panel B, 'other' includes 45 phyla all with less than 5 genomes. For a
351     complete list of bacterial phyla please see Figure S2. In panel D, 'other' includes 2 phyla all with
352     less than 5 genomes. For a complete list of archaeal phyla please see Figure 3.
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
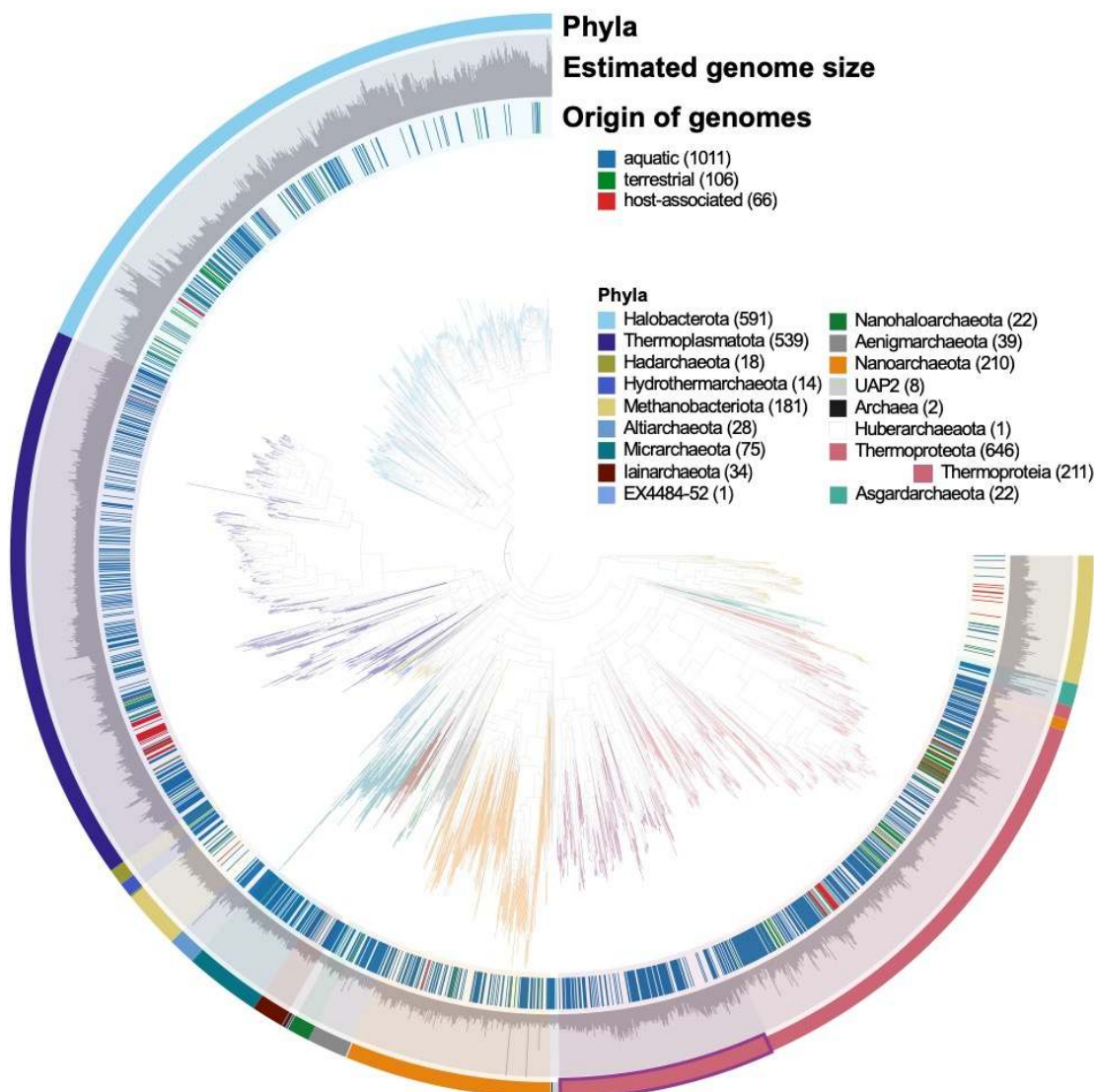375
376
377
378
379
380
381
382

11

**Figure 2**. Phylogenetic tree of bacterial representative genomes shows variation in genome size between and within phyla. Tree was constructed using GTDB-tk and aligned concatenated set of 120 single copy marker proteins for Bacteria (96). Estimated genome size shows distribution of larger and smaller genomes sizes are non-monophyletic. The tree shows origin of the genomes: aquatic, terrestrial and host-associated genomes are MAGs from GEMs database. The backbone genomes were added by GTDB-tk and it consists of their representative genomes. Estimated genome size scale is from 0 Mb to 14 Mb. Phyla are color-coded and legend includes the phyla with most representatives. Phyla with less than 50 genomes are not included in the legend. For full legend please refer to Figure S2. Burkholderiales is the Order with most genomes.
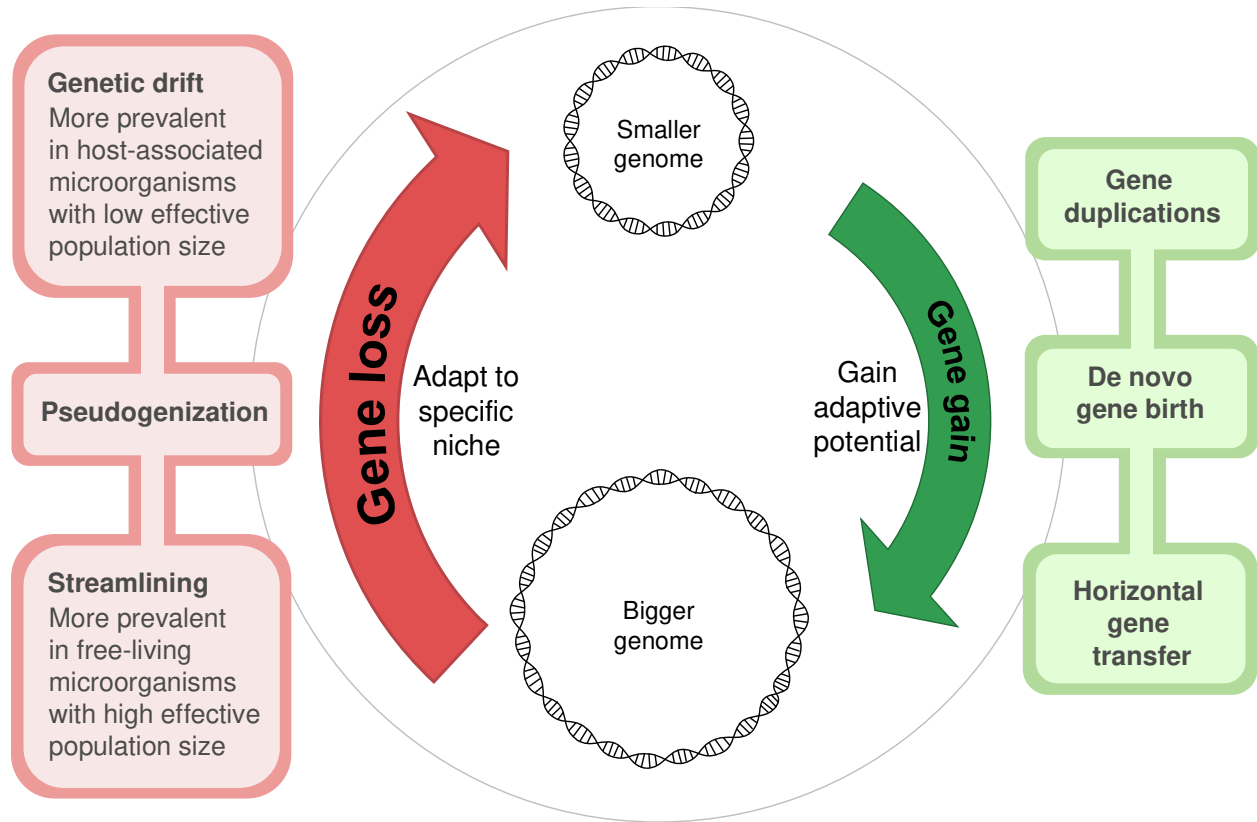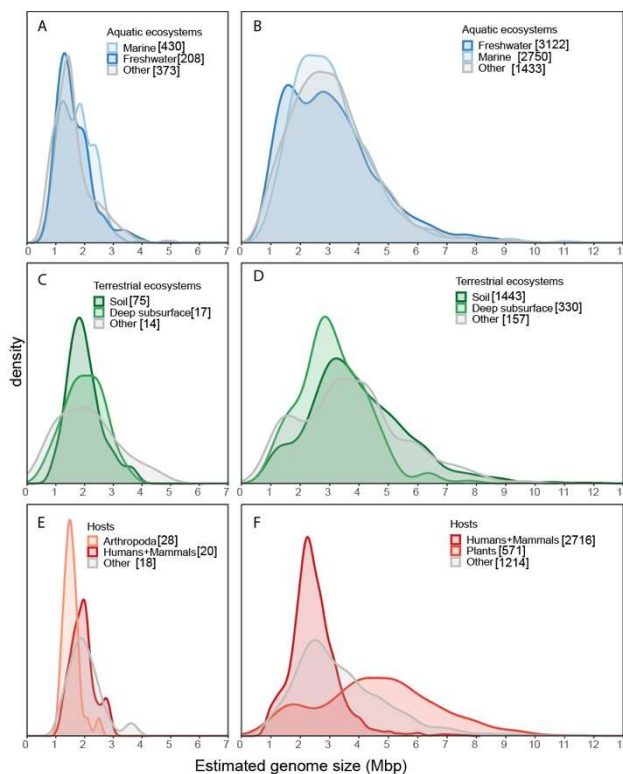
12

401
402



403
404
**Figure 3**. Phylogenetic tree of archaeal representative genomes shows variation in genome size between and within phyla. Tree was constructed using GTDB-tk and aligned concatenated set of 122 single copy marker proteins for Archaea (96). Estimated genome size shows distribution of larger and smaller genomes sizes are non-monophyletic. The tree shows origin of the genomes: aquatic, terrestrial and host-associated genomes are MAGs from GEMs database. The backbone genomes were added by GTDB-tk and it consists of their representative genomes. Estimated genome size scale is from 0 Mb to 6 Mb. Phyla are color-coded.
412

**Figure 4.** Conceptual figure of the evolutionary forces driving the expansion and reduction of genome sizes. Gene loss is represented with a bigger arrow because it dominates the evolutionary history we know based on extant microorganisms.

**Figure 5**. Genome size distribution in different sub-categories of environments. [A] Aquatic archaeal genomes, [B] aquatic bacterial genomes, [C] terrestrial archaeal genomes, [D] terrestrial bacterial genomes, [E] hos-associated archaeal genomes and [F] host-associated bacterial genomes. Inside the parenthesis is stated the number of MAGs per sub-environment.

**References**

1.      Gregory TR. Genome Size Evolution in Animals.  The Evolution of the Genome2005. p. 3-87.

2.      Pellicer J, Hidalgo O, Dodsworth S, Leitch IJ. Genome Size Diversity and Its Impact on the Evolution of Land Plants. Genes (Basel). 2018;9(2).

3.      Kirchberger PC, Schmidt ML, Ochman H. The Ingenuity of Bacterial Genomes. Annu Rev Microbiol. 2020;74:815-34.

4.      Moran NA, Bennett GM. The tiniest tiny genomes. Annu Rev Microbiol. 2014;68:195-215.

444    5.    Garcia R, Gemperlein K, Muller R. Minicystis rosea gen. nov., sp. nov., a

445    polyunsaturated fatty acid-rich and steroid-producing soil myxobacterium. Int J Syst Evol

446    Microbiol. 2014;64(Pt 11):3733-42.

447    6.    Bobay LM, Ochman H. The Evolution of Bacterial Genome Architecture. Front Genet.

448    2017;8:72.

449    7.    Puigbo P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin EV. Genomes in turmoil:

450    quantification of genome dynamics in prokaryote supergenomes. BMC biology. 2014;12:66.

451    8.    Kuo CH, Moran NA, Ochman H. The consequences of genetic drift for bacterial genome

452    complexity. Genome Res. 2009;19(8):1450-4.

453    9.    Gillings MR. Lateral gene transfer, bacterial genome evolution, and the Anthropocene.

454    Annals of the New York Academy of Sciences. 2017;1389(1):20-36.

455    10.    Han K, Li ZF, Peng R, Zhu LP, Zhou T, Wang LG, et al. Extraordinary expansion of a

456    Sorangium cellulosum genome from an alkaline milieu. Scientific reports. 2013;3:2101.

457    11.    Cavalier-Smith T. Economy, speed and size matter: evolutionary forces driving nuclear

458    genome miniaturization and expansion. Ann Bot. 2005;95(1):147-75.

459    12.    Giovannoni SJ, Cameron Thrash J, Temperton B. Implications of streamlining theory for

460    microbial ecology. ISME J. 2014;8(8):1553-65.

461    13.    Land M, Hauser L, Jun SR, Nookaew I, Leuze MR, Ahn TH, et al. Insights from 20 years

462    of bacterial genome sequencing. Funct Integr Genomics. 2015;15(2):141-61.

463    14.    Staley JT, Konopka A. Measurement of in Situ Activities of Nonphotosynthetic

464    Microorganisms in Aquatic and Terrestrial Habitats. Annual Review of Microbiology.

465    1985;39(1):321-46.

466   15.   Fleischmann RD, Adams MD, White O, Clayton RA, Kirkness EF, Kerlavage AR, et al.

467   Whole-genome random sequencing and assembly of Haemophilus influenzae Rd. Science.

468   1995;269(5223):496-512.

469   16.   Fraser CM, Gocayne JD, White O, Adams MD, Clayton RA, Fleischmann RD, et al. The

470   minimal gene complement of Mycoplasma genitalium. Science. 1995;270(5235):397-403.

471   17.   Binnewies TT, Motro Y, Hallin PF, Lund O, Dunn D, La T, et al. Ten years of bacterial

472   genome sequencing: comparative-genomics-based discoveries. Funct Integr Genomics.

473   2006;6(3):165-85.

474   18.   Nayfach S, Roux S, Seshadri R, Udwary D, Varghese N, Schulz F, et al. A genomic

475   catalog of Earth's microbiomes. Nat Biotechnol. 2020.

476   19.   Parks DH, Rinke C, Chuvochina M, Chaumeil PA, Woodcroft BJ, Evans PN, et al.

477   Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life.

478   Nat Microbiol. 2017;2(11):1533-42.

479   20.   Anantharaman K, Brown CT, Hug LA, Sharon I, Castelle CJ, Probst AJ, et al. Thousands

480   of microbial genomes shed light on interconnected biogeochemical processes in an aquifer

481   system. Nat Commun. 2016;7:13219.

482   21.   Parks DH, Chuvochina M, Chaumeil PA, Rinke C, Mussig AJ, Hugenholtz P. A

483   complete domain-to-species taxonomy for Bacteria and Archaea. Nat Biotechnol.

484   2020;38(9):1079-86.

485   22.   Batut B, Knibbe C, Marais G, Daubin V. Reductive genome evolution at both ends of the

486   bacterial population size spectrum. Nature reviews Microbiology. 2014;12(12):841-50.

487   23.   Konstantinidis KT, Tiedje JM. Trends between gene content and genome size in

488   prokaryotic species with larger genomes. Proc Natl Acad Sci U S A. 2004;101(9):3160-5.

489    24.    Lynch M. Streamlining and simplification of microbial genome architecture. Annu Rev

490    Microbiol. 2006;60:327-49.

491    25.    Wolf YI, Koonin EV. Genome reduction as the dominant mode of evolution. Bioessays.

492    2013;35(9):829-37.

493    26.    Abram K, Udaondo Z, Bleker C, Wanchai V, Wassenaar TM, Robeson MS, et al. What

494    can we learn from over 100,000 Escherichia coli genomes? bioRxiv. 2020.

495    27.    Rocap G, Larimer FW, Lamerdin J, Malfatti S, Chain P, Ahlgren NA, et al. Genome

496    divergence in two Prochlorococcus ecotypes reflects oceanic niche differentiation. Nature.

497    2003;424(6952):1042-7.

498    28.    Gweon HS, Bailey MJ, Read DS. Assessment of the bimodality in the distribution of

499    bacterial genome sizes. ISME J. 2017;11(3):821-4.

500    29.    Nelson WC, Tully BJ, Mobberley JM. Biases in genome reconstruction from

501    metagenomic data. PeerJ. 2020;8:e10119.

502    30.    Shade A, Hogan CS, Klimowicz AK, Linske M, McManus PS, Handelsman J. Culturing

503    captures members of the soil rare biosphere. Environmental Microbiology. 2012;14(9):2247-52.

504    31.    Swan BK, Tupper B, Sczyrba A, Lauro FM, Martinez-Garcia M, Gonzalez JM, et al.

505    Prevalent genome streamlining and latitudinal divergence of planktonic bacteria in the surface

506    ocean. Proc Natl Acad Sci U S A. 2013;110(28):11463-8.

507    32.    Garcia SL. Mixed cultures as model communities: hunting for ubiquitous

508    microorganisms, their partners, and interactions. Aquatic Microbial Ecology. 2016;77(2):79-85.

509    33.    Imachi H, Nobu MK, Nakahara N, Morono Y, Ogawara M, Takaki Y, et al. Isolation of

510    an archaeon at the prokaryote-eukaryote interface. Nature. 2020;577(7791):519-25.

511   34.    Cross KL, Campbell JH, Balachandran M, Campbell AG, Cooper SJ, Griffen A, et al.

512   Targeted isolation and cultivation of uncultivated bacteria by reverse genomics. Nat Biotechnol.

513   2019;37(11):1314-21.

514   35.    Hoehler TM, Jorgensen BB. Microbial life under extreme energy limitation. Nature

515   reviews Microbiology. 2013;11(2):83-94.

516   36.    Figueroa-Gonzalez PA, Bornemann TLV, Adam PS, Plewka J, Révész F, von Hagen CA,

517   et al. Saccharibacteria as Organic Carbon Sinks in Hydrocarbon-Fueled Communities. Frontiers

518   in microbiology. 2020;11.

519   37.    Lewis WH, Tahon G, Geesink P, Sousa DZ, Ettema TJG. Innovations to culturing the

520   uncultured microbial majority. Nature reviews Microbiology. 2020.

521   38.    Marais GA, Calteau A, Tenaillon O. Mutation rate and genome reduction in

522   endosymbiotic and free-living bacteria. Genetica. 2008;134(2):205-10.

523   39.    Croucher NJ, Harris SR, Barquist L, Parkhill J, Bentley SD. A high-resolution view of

524   genome-wide pneumococcal transformation. PLoS pathogens. 2012;8(6):e1002745.

525   40.    de Vries J, Wackernagel W. Integration of foreign DNA during natural transformation of

526   Acinetobacter sp. by homology-facilitated illegitimate recombination. Proc Natl Acad Sci U S A.

527   2002;99(4):2094-9.

528   41.    Shapiro BJ, Polz MF. Microbial Speciation. Cold Spring Harb Perspect Biol.

529   2015;7(10):a018143.

530   42.    Zaremba-Niedzwiedzka K, Viklund J, Zhao WZ, Ast J, Sczyrba A, Woyke T, et al.

531   Single-cell genomics reveal low recombination frequencies in freshwater bacteria of the SAR11

532   clade. Genome Biology. 2013;14(11).

533    43.    Sela I, Wolf YI, Koonin EV. Theory of prokaryotic genome evolution. Proc Natl Acad

534    Sci U S A. 2016;113(41):11399-407.

535    44.    Larsson J, Nylander JA, Bergman B. Genome fluctuations in cyanobacteria reflect

536    evolutionary, developmental and adaptive traits. BMC evolutionary biology. 2011;11:187.

537    45.    Van Oss SB, Carvunis AR. De novo gene birth. PLoS Genet. 2019;15(5):e1008160.

538    46.    Pal C, Papp B, Lercher MJ. Adaptive evolution of bacterial metabolic networks by

539    horizontal gene transfer. Nat Genet. 2005;37(12):1372-5.

540    47.    Treangen TJ, Rocha EP. Horizontal transfer, not duplication, drives the expansion of

541    protein families in prokaryotes. PLoS Genet. 2011;7(1):e1001284.

542    48.    Sheridan PO, Raguideau S, Quince C, Holden J, Zhang L, Thames C, et al. Gene

543    duplication drives genome expansion in a major lineage of Thaumarchaeota. Nat Commun.

544    2020;11(1):5494.

545    49.    Deppenmeier U, Johann A, Hartsch T, Merkl R, Schmitz RA, Martinez-Arias R, et al.

546    The genome of Methanosarcina mazei: evidence for lateral gene transfer between bacteria and

547    archaea. J Mol Microbiol. 2002;4(4):453-61.

548    50.    Lurie-Weinberger MN, Peeri M, Tuller T, Gophna U. Extensive Inter-Domain Lateral

549    Gene Transfer in the Evolution of the Human Commensal Methanosphaera stadtmanae. Front

550    Genet. 2012;3:182.

551    51.    Nelson-Sathi S, Dagan T, Landan G, Janssen A, Steel M, McInerney JO, et al.

552    Acquisition of 1,000 eubacterial genes physiologically transformed a methanogen at the origin of

553    Haloarchaea. Proc Natl Acad Sci U S A. 2012;109(50):20537-42.

554    52.    Tamas I, Klasson L, Canback B, Naslund AK, Eriksson AS, Wernegreen JJ, et al. 50

555    million years of genomic stasis in endosymbiotic bacteria. Science. 2002;296(5577):2376-9.

556    53.    Andersson JO, Andersson SG. Pseudogenes, junk DNA, and the dynamics of Rickettsia

557    genomes. Molecular biology and evolution. 2001;18(5):829-39.

558    54.    van Passel MW, Smillie CS, Ochman H. Gene decay in archaea. Archaea. 2007;2(2):137-

559    43.

560    55.    Chu X, Li S, Wang S, Luo D, Luo H. Gene loss through pseudogenization contributes to

561    the ecological diversification of a generalist Roseobacter lineage. ISME J. 2020.

562    56.    Morris JJ, Lenski RE, Zinser ER. The Black Queen Hypothesis: evolution of

563    dependencies through adaptive gene loss. Mbio. 2012;3(2).

564    57.    Mondav R, Bertilsson S, Buck M, Langenheder S, Lindstrom ES, Garcia SL. Streamlined

565    and Abundant Bacterioplankton Thrive in Functional Cohorts. 2020.

566    58.    Gabrielaite M, Johansen HK, Molin S, Nielsen FC, Marvig RL. Gene Loss and

567    Acquisition in Lineages of Pseudomonas aeruginosa Evolving in Cystic Fibrosis Patient

568    Airways. Mbio. 2020;11(5).

569    59.    Csuros M, Miklos I. Streamlining and large ancestral genomes in Archaea inferred with a

570    phylogenetic birth-and-death model. Molecular biology and evolution. 2009;26(9):2087-95.

571    60.    Wolf YI, Makarova KS, Yutin N, Koonin EV. Updated clusters of orthologous genes for

572    Archaea: a complex ancestor of the Archaea and the byways of horizontal gene transfer. Biology

573    direct. 2012;7:46.

574    61.    Raes J, Korbel J, Lercher M, von Mering C, Bork P. Prediction of effective genome size

575    in metagenomic samples. Genome Biology. 2007;8(1):R10.

576    62.    Chen MY, Teng WK, Zhao L, Hu CX, Zhou YK, Han BP, et al. Comparative genomics

577    reveals insights into cyanobacterial evolution and habitat adaptation. ISME J. 2020.

578    63.    Cobo-Simon M, Tamames J. Relating genomic characteristics to environmental

579    preferences and ubiquity in different microbial taxa. BMC genomics. 2017;18(1):499.

580    64.    Steele JH, Brink KH, Scott BE. Comparison of marine and terrestrial ecosystems:

581    suggestions of an evolutionary perspective influenced by environmental variation. ICES Journal

582    of Marine Science. 2019;76(1):50-9.

583    65.    Bentkowski P, Van Oosterhout C, Mock T. A Model of Genome Size Evolution for

584    Prokaryotes in Stable and Fluctuating Environments. Genome biology and evolution.

585    2015;7(8):2344-51.

586    66.    Brewer TE, Handley KM, Carini P, Gilbert JA, Fierer N. Genome reduction in an

587    abundant and ubiquitous soil bacterium 'Candidatus Udaeobacter copiosus'. Nature

588    Microbiology. 2016;2:16198.

589    67.    Giovannoni SJ, Hayakawa DH, Tripp HJ, Stingl U, Givan SA, Cho JC, et al. The small

590    genome of an abundant coastal ocean methylotroph. Environ Microbiol. 2008;10(7):1771-82.

591    68.    Aylward FO, Santoro AE. Heterotrophic Thaumarchaea with Small Genomes Are

592    Widespread in the Dark Ocean. mSystems. 2020;5(3).

593    69.    Tian R, Ning D, He Z, Zhang P, Spencer SJ, Gao S, et al. Small and mighty: adaptation

594    of superphylum Patescibacteria to groundwater environment drives their genome simplicity.

595    Microbiome. 2020;8(1):51.

596    70.    Mende DR, Bryant JA, Aylward FO, Eppley JM, Nielsen T, Karl DM, et al.

597    Environmental drivers of a microbial genomic transition zone in the ocean's interior. Nat

598    Microbiol. 2017;2(10):1367-73.

599    71.    Grzymski JJ, Dussaq AM. The significance of nitrogen cost minimization in proteomes

600    of marine microorganisms. ISME J. 2012;6(1):71-80.

601    72.    Brochier-Armanet C, Deschamps P, Lopez-Garcia P, Zivanovic Y, Rodriguez-Valera F,

602    Moreira D. Complete-fosmid and fosmid-end sequences reveal frequent horizontal gene transfers

603    in marine uncultured planktonic archaea. ISME J. 2011;5(8):1291-302.

604    73.    Blesa A, Averhoff B, Berenguer J. Horizontal Gene Transfer in Thermus spp. Curr Issues

605    Mol Biol. 2018;29:23-36.

606    74.    Borges KM, Bergquist PL. Genomic restriction map of the extremely thermophilic

607    bacterium Thermus thermophilus HB8. J Bacteriol. 1993;175(1):103-10.

608    75.    Dieser M, Smith HJ, Ramaraj T, Foreman CM. Janthinobacterium CG23_2: Comparative

609    Genome Analysis Reveals Enhanced Environmental Sensing and Transcriptional Regulation for

610    Adaptation to Life in an Antarctic Supraglacial Stream. Microorganisms. 2019;7(10).

611    76.    Berube PM, Biller SJ, Hackl T, Hogle SL, Satinsky BM, Becker JW, et al. Single cell

612    genomes of Prochlorococcus, Synechococcus, and sympatric microbes from diverse marine

613    environments. Sci Data. 2018;5:180154.

614    77.    Toft C, Andersson SG. Evolutionary microbial genomics: insights into bacterial host

615    adaptation. Nature reviews Genetics. 2010;11(7):465-75.

616    78.    Collingro A, Kostlbacher S, Horn M. Chlamydiae in the Environment. Trends Microbiol.

617    2020;28(11):877-88.

618    79.    Dharamshi JE, Tamarit D, Eme L, Stairs CW, Martijn J, Homa F, et al. Marine Sediments

619    Illuminate Chlamydiae Diversity and Evolution. Current biology : CB. 2020;30(6):1032-48 e7.

620    80.    McLean JS, Bor B, Kerns KA, Liu Q, To TT, Solden L, et al. Acquisition and Adaptation

621    of Ultra-small Parasitic Reduced Genome Bacteria to Mammalian Hosts. Cell reports.

622    2020;32(3):107939.

623     81.     Levy A, Salas Gonzalez I, Mittelviefhaus M, Clingenpeel S, Herrera Paredes S, Miao J,

624     et al. Genomic features of bacterial adaptation to plants. Nat Genet. 2017;50(1):138-50.

625     82.     Boussau B, Karlberg EO, Frank AC, Legault BA, Andersson SG. Computational

626     inference of scenarios for alpha-proteobacterial genome evolution. Proc Natl Acad Sci U S A.

627     2004;101(26):9722-7.

628     83.     Serra V, Gammuto L, Nitla V, Castelli M, Lanzoni O, Sassera D, et al. Morphology,

629     ultrastructure, genomics, and phylogeny of Euplotes vanleeuwenhoeki sp. nov. and its ultra-

630     reduced endosymbiont "Candidatus Pinguicoccus supinus" sp. nov. Scientific reports.

631     2020;10(1):20311.

632     84.     Rhodes ME, Spear JR, Oren A, House CH. Differences in lateral gene transfer in

633     hypersaline versus thermal environments. BMC evolutionary biology. 2011;11:199.

634     85.     Cabello-Yeves PJ, Zemskaya TI, Rosselli R, Coutinho FH, Zakharenko AS, Blinov VV,

635     et al. Genomes of Novel Microbial Lineages Assembled from the Sub-Ice Waters of Lake

636     Baikal. Appl Environ Microbiol. 2018;84(1).

637     86.     Allen LZ, Allen EE, Badger JH, McCrow JP, Paulsen IT, Elbourne LD, et al. Influence of

638     nutrients and currents on the genomic composition of microbes across an upwelling mosaic.

639     ISME J. 2012;6(7):1403-14.

640     87.     Makarova K, Slesarev A, Wolf Y, Sorokin A, Mirkin B, Koonin E, et al. Comparative

641     genomics of the lactic acid bacteria. Proc Natl Acad Sci U S A. 2006;103(42):15611-6.

642     88.     Dufresne A, Salanoubat M, Partensky F, Artiguenave F, Axmann IM, Barbe V, et al.

643     Genome sequence of the cyanobacterium Prochlorococcus marinus SS120, a nearly minimal

644     oxyphototrophic genome. Proc Natl Acad Sci U S A. 2003;100(17):10020-5.

645  89.  Smith MW, Zeigler Allen L, Allen AE, Herfort L, Simon HM. Contrasting genomic

646  properties of free-living and particle-attached microbial assemblages within a coastal ecosystem.

647  Frontiers in microbiology. 2013;4:120.

648  90.  de Souza RSC, Armanhi JSL, Damasceno NB, Imperial J, Arruda P. Genome Sequences

649  of a Plant Beneficial Synthetic Bacterial Community Reveal Genetic Features for Successful

650  Plant Colonization. Frontiers in microbiology. 2019;10:1779.

651  91.  Nilsson AI, Koskiniemi S, Eriksson S, Kugelberg E, Hinton JC, Andersson DI. Bacterial

652  genome size reduction by experimental evolution. Proc Natl Acad Sci U S A.

653  2005;102(34):12112-6.

654  92.  Moreira D, Le Guyader H, Philippe H. The origin of red algae and the evolution of

655  chloroplasts. Nature. 2000;405(6782):69-72.

656  93.  Castelle CJ, Banfield JF. Major New Microbial Groups Expand Diversity and Alter our

657  Understanding of the Tree of Life. Cell. 2018;172(6):1181-97.

658  94.  Yooseph S, Nealson KH, Rusch DB, McCrow JP, Dupont CL, Kim M, et al. Genomic

659  and functional adaptation in surface ocean planktonic prokaryotes. Nature. 2010;468(7320):60-6.

660  95.  Team RC. R: A language and environment for statistical computing. 2020 [R Foundation

661  for Statistical Computing , Vienna, Austria.:[Available from: http://www.R-project.org/.

662  96.  Chaumeil PA, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: a toolkit to classify

663  genomes with the Genome Taxonomy Database. Bioinformatics. 2020;36(6):1925-7.

664