**A natural definition for a bacterial strain and clonal complex**.

Luis M. Rodriguez-R[1*], Roth E. Conrad[2,*], Dorian J. Feistel[2], Tomeu Viver[3], Ramon Rosselló-Móra[3], and Konstantinos T. Konstantinidis[2]

[1]Department of Microbiology, and Digital Science Center (DiSC), University of Innsbruck, Innsbruck, Tyrol, Austria.

[2]School of Civil and Environmental Engineering, and School of Biological Sciences, Georgia Institute of Technology, Atlanta, GA, USA.

[3]Marine Microbiology Group, Department of Animal and Microbial Biodiversity, Mediterranean Institutes for Advanced Studies (IMEDEA, CSIC-UIB), Esporles, Spain

Correspondence should be addressed to K.T.K. (email: kostas.konstantinidis@gatech.edu)

*equal first authors.

**Classification:** Microbiology
**Keywords:** ANI, strain definition, micro-diversity, epidemiology, clonal complex

**1 Abstract**

2 Large scale surveys of prokaryotic communities (metagenomics) as well as isolate

3 genomes have recently revealed that prokaryotic diversity is predominantly organized in

4 sequence-discrete units that may be equated to species. Specifically, genomes of the same unit

5 (or species) commonly show >95% genome-aggregate average nucleotide identity (ANI) to each

6 other and <90% ANI to members of other species, while genomes showing 90-95% ANI are

7 comparatively rare. However, it remains unclear if such "discontinuities" or gaps in ANI values

8 can be observed within species and thus, be used to define strains or clonal complexes; two

9 cornerstone concepts for microbiology that remain ill-defined. By analyzing 18,123 complete

10 isolate genomes from 330 bacterial species with at least ten genome representatives each, we

11 show that such a natural discontinuity exists at around 99.5% ANI. Further, we show that the

12 99.5% ANI threshold is largely consistent with how clonal complexes have been defined in

13 previous epidemiological studies but provides clusters with ~20% higher accuracy in terms of

14 evolutionary relatedness of the grouped genomes and greater homogeneity in gene content.

15 Collectively, our results should facilitate future micro-diversity studies across clinical or

16 environmental settings because they provide a more natural definition of a clonal complex and

17 strain.

18

**19 Main**

20 The strain represents a fundamental unit of microbial diversity that is commonly used

21 across medical or environment studies to designate the smallest distinguishable unit. This term

22 is based on the concept of a pure culture (1), and has been defined by the Bacteriological Code

23 as the group of cultures made "*of the descendants of a single isolation in pure culture*" (2).

24 Accordingly, a strain is expected to represent a genome or a collection of genomes that have no

25 single-nucleotide or gene content differences or, if such differences exist, they are expected to

26 not encode for important phenotypic differences (3). This strain definition, while being

27 commonly used across microbiological fields, remains problematic, however, and often causes

28 confusion in communication. Most notably, it is not clear when two distinct genomes or cells

29 should be considered the same or separate strains since cell ancestry information is often

30 missing, and the existence of phenotypic similarities (or differences) may vary, depending on the

2

31 growth conditions. Additionally, the isolation of an organism (wild-type) in the laboratory is

32 frequently accompanied by phenotypic changes due to adaptation to the laboratory conditions;

33 yet, the wild-type and the lab-adapted cells or types are typically considered the same strain (1).

34

35 A related term that is also commonly used to catalogue intra-species diversity is the clonal

36 complex (CC). The CC has been used, especially in medical microbiology and epidemiological

37 studies, to identify -for instance- an outbreak caused by a specific pathogenic strain or groups

38 (complexes) of highly related strains. Unlike the more relaxed, context-dependent definition of a

39 "strain", a CC is typically defined as a collection of genomes with no nucleotide sequence diversity

40 (zero single-nucleotide polymorphisms) in 6-7 genetic loci (4). These loci are typically distributed

41 across the genome, to avoid co-selection evolutionary events, and represent fragments

42 (amplicons) of genes shared by most members of a species; that is, core genes. While this

43 definition is pragmatic and operational, it is also problematic because it is arbitrary (e.g.,

44 sometimes CC definitions allow 0- or 1-point mutations), and core genes tend to be more

45 conserved than the genome average. Thus, it remains somewhat speculative how similar (or not)

46 isolates of the same CC may be in the rest of their genome, and this may also depend, at least

47 partly, on the exact loci used in the analysis (5, 6). If the intra-species diversity was organized in

48 discrete units that can be equated to CC or strains, this would have provided for more natural

49 and meaningful definitions for CCs and strains compared to the existing definitions and thus,

50 improved communication about intra-species diversity. While it has been recently recognized

51 that prokaryotic organisms may form such discrete units at the species level (7), it remains

52 unclear whether or not such units also exist within species.

53

54 Specifically, culture-independent (metagenomic) studies of natural microbial populations

55 during the past decade revealed that bacteria and archaea predominantly form sequence-

56 discrete populations with intra-population genomic sequence relatedness typically ranging from

57 ~95% to ~100% genome-aggregate average nucleotide identity (or ANI) depending on the

58 population considered (e.g., younger populations since the last population diversity sweep event

59 show lower levels of intra-population diversity). In contrast, ANI values between distinct

60    populations are typically lower than 90% (8). Intermediate identity genotypes, for example,

61    sharing 85–95% ANI, when present, are generally ecologically differentiated and scarcer in

62    abundance, and thus should probably be considered distinct species (7, 9, 10) rather than

63    representing cultivation or other sampling biases (11). Such sequence-discrete populations have

64    been recovered from many different habitats, including marine, freshwater, soils, human gut,

65    and biofilms, and were usually persistent over time and space [e.g., (12-16)] indicating that they

66    are not ephemeral but long-lived entities. Further, these sequence-discrete populations

67    commonly harbor substantial intra-population gene content diversity (i.e., they are rarely clonal)

68    (12, 15). Therefore, these populations appear to be "species-like" and may constitute important

69    units of microbial communities. Moreover, the 95% ANI threshold appears to be largely

70    consistent with how genomes have been classified into (named) species in the last couple

71    decades; that is, ~97% of named species include only organisms with genomes sharing >95% ANI

72    (17). In summary, it appears that a natural gap in ANI values can be used to define prokaryotic

73    species and has been largely consistent with how species are recognized (17) [Discontinuity or

74    gap here refers to the small number of genome pairs showing 85–95% ANI relative to counts of

75    pairs showing >95% and <85% ANI]. Whether or not a similar ANI gap exists and can be used to

76    define strains or CC has not been evaluated yet.

77

78    **Results/Discussion**

79    **An ANI gap within species around 99.5%**

80          In the process of assessing cultivation biases as a possible explanation for the ANI-based

81    sequence discrete populations (18), we observed another discontinuity (or gap) in ANI values that

82    may be used to define the smallest unit within a species, that of the strain or CC. Specifically, the

83    analysis of 18,123 complete genomes from 330 species available in NCBI's Assembly database

84    with at least 10 such genome representatives each revealed a clear bimodal distribution in ANI

85    values within named species or 95% ANI-defined groups of genomes (genomospecies). That is,

86    there is a scarcity of genomes pairs showing 99.2-99.8% ANI (average around 99.5% ANI) in

87    contrast to genome pairs showing >99.8% or <99.2% ANI. Specifically, among the 18,123

88    complete genomes in our dataset, there are 4,280,133 genome pairs showing >96% ANI, which

89    would translate to about 107,000 pairs per every 0.1 percent unit of ANI if there was no bimodal

90    distribution but the ANI values among these genome pairs were evenly distributed between 96%

91    and 100% ANI. Our analysis revealed only 235,527 genome pairs between 99.2% and 99.8% ANI,

92    which is three-fold fewer data points than expected by chance alone in a uniform ANI value

93    distribution (642,000 pairs expected). No other ANI range within 96-100% had such a strong bias

94    based on our dataset. That is, a pronounced gap in ANI values is observed among very closely

95    related members of a species around 99.5% ANI (Fig. 1). Importantly, this ANI gap appears to be

96    consistent across phylogenetically diverse species from a dozen of distinct bacterial phyla

97    evaluated, including gram-negative and gram-positive, and does not seem to be driven by a

98    couple or a few species based on a sub-sampling of all species to the same number of genomes

99    (n=10) (Fig. S1). Instead, it represents a universal property of the 330 species evaluated (see also

100    Fig. S2 for specific species examples). Therefore, it appears that another important level of

101    genomic differentiation may exist within species that can be used to define strains and CC.

102

103        It is unlikely that this 99.5% intra-species ANI gap is due to cultivation or classification

104    biases due the reasons mentioned previously, such as that cultivation media should not

105    distinguish between members of the same species or closely related groups of organisms (18),

106    and that random subsampling provided similar patterns (Fig. S1). In fact, it is highly likely that the

107    intra-species ANI gap may be even more pronounced in nature because very closely related

108    genomes (e.g., showing >99.8% ANI to each other) are often selected against for genome

109    sequencing (and thus, are likely underrepresented in our collection) based on pre-screening using

110    fingerprinting techniques (e.g., RARP-PCR, MLST) in order to avoid sequencing of redundant

111    genomes. Despite this bias against very closely related genomes, which in all probability exists

112    but we are currently unable to estimate its magnitude, several species have enough very closely

113    related genomes sequenced in our dataset for robust evaluation of patterns in their ANI value

114    distributions (Fig. 1 and S2). We were not able to identify clear exceptions to this 99.5% intra-

115    species ANI gap when examining individual species with enough sequenced representatives,

116    although such exceptions likely exist. For instance, several species in our collection did not have

117    enough very highly related genome representatives (showing >99% ANI to each other) to assess

118    the critical area of ANI value distribution (i.e., the 99-100% range), and this could be due to the

119    pre-screening biases mentioned above or reflect their actual natural diversity patterns (see also

120    Fig. 2, Panel D for an example within the *E. coli*). Further, for a few species (n < 10) such as *Listeria*

121    *monocytogenes* and *Bordetella bronchiseptica*, the intra-species ANI gap appears to exist but is

122    shifted compared to the 99.5% ANI that characterized most well-sampled species (Fig. S2).

123    Therefore, for future studies, we suggest evaluating the ANI value distribution for the species of

124    interest, and if the data indicate so, to adjust the ANI threshold to match the gap in the observed

125    ANI value distribution. The 99.5% ANI should work for most species based on the dataset

126    evaluated here.

127

128    **Gene content diversity within 99.5% ANI clusters**

129    Another notable observation from the data from all species comparisons is that shared

130    gene content decreases, on average, as ANI distance (or genomic divergence) increases within

131    the 95% ANI clusters, but the decrease is biphasic. That is, shared gene content decreases quickly

132    among genome pairs sharing 99.0-100% ANI but then, the decrease is less dramatic in genome

133    pairs sharing between 96.0%-99.0% ANI. In other words, genome pairs sharing between 99.0%-

134    100% ANI (i.e., one ANI unit range) may differ in their total gene content by up to 10% (average

135    values of genome pairs showing ~99.0% ANI) and more divergent genomes of the same species

136    (i.e., showing 96.0%>ANI>99.0%) may differ by up to 20% (average values of genome pairs

137    showing ~96.0% ANI), adding another ~10% of gene content differences for 3 additional units of

138    ANI (vs. 1 unit in the 99-100% range). Collectively, these results show that genome pairs showing

139    >99.5% ANI are also expected to be much more similar in gene-content compared to more

140    divergent genomes of the same species. Not only do these results quantify the amount of gene

141    content diversity expected in comparisons of genomes within the same species, but they are also

142    consistent with the notion that members of the same CC should be highly similar in shared

143    functional gene content and thus, phenotype. While the high gene content diversity within

144    species has been observed previously, even based on the very first few bacterial genomes

145    sequenced (19), the analysis presented here provides robust quantification of shared gene

146    content as a function of ANI and CC or strain designation, which should be useful in future studies

147    that aim to quantify the size and value of bacterial species pangenomes.

148

149    **Comparison to clonal complexes (or Sequence Types, STs)**

150    We also assessed how consistent the 99.5% ANI threshold is with the assignment of

151    genomes to the same CC, the latter defined as identical sequences for 6-7 genetic loci. We focus

152    this analysis on the *E. coli* species because it is a good representative of the ANI patterns observed

153    within other species in our dataset, the large number of *E. coli* genomes available (n=975), and

154    the availability of a robust Multi-Locus Sequence Typing/Analysis (MLST/MLSA) scheme (20) that

155    has been used for at least two decades to provide below-species resolution and identify

156    outbreaks of *E. coli* pathogens. Under the *E. coli* MLST scheme, genomes are assigned to the same

157    CC, also commonly called Sequence Type (ST), based on identical sequences in seven *E. coli* core

158    genes (namely, *adk*, *fumC*, *gyrB*, *icd*, *mdh*, *purA*, and *recA*) (20). Our evaluation showed that the

159    99.5% ANI threshold is largely consistent with how genomes are assigned to STs; that is, ~80% of

160    ST assignments, for the top four most abundant STs (n=615 genomes), were supported by the

161    99.5% ANI threshold (Fig. 2). In other words, only about 20% of the genomes that were assigned

162    to the same ST showed <99.5% ANI among themselves (high recall). Further, ~2% of the genomes

163    with >99.5% were assigned to different ST, revealing even higher precision (Table S1; Suppl. Fig.

164    S3). The high recall (and precision) of the 99.5% ANI threshold is due, at least in part, to the fact

165    that the core genes used in current MLST schemes tend to be more highly conserved, at the

166    sequence level, than the genome average (represented by ANI), or have been exchanged

167    horizontally very recently among the evaluated genomes. It is also important to note that

168    hundreds, if not thousands (e.g., at least 3,000 for the *E. coli* pairs), of genes are used in each ANI

169    calculation at this level of high relatedness vs. 6-7 genes in MLST schemes and that's another

170    reason for the higher robustness of the former metric against horizontal gene transfer or other

171    evolutionary events that could affect sequence identity. These results are also consistent with

172    our previous conclusion that the 95% ANI threshold should be adjusted upwards if the ANI is

173    based on a few universal or core genes (17). In any case, the 99.5% ANI threshold is a property

174    that emerges from the data themselves as opposed to a manmade, arbitrary identity threshold

175   (i.e., identical sequences in seven loci used in MLST applications) and thus, it should capture the

176   natural diversity patterns better. Consistent with this interpretation, our preliminary results from

177   applying the 99.5% ANI threshold to a collection of *E. coli* isolate genomes collected over a period

178   of 18 months in Northern coastal Ecuador as part of the EcoZUR study (for "*E. coli* en Zonas

179   Urbanas y Rurales") (21) shows that the 99.5%-ANI-defined CCs map well to local outbreaks of

180   pathogenic *E. coli* (Feistel et al., in preparation). Therefore, using the 99.5% ANI to define CCs will

181   prevent inaccurate calls and provide more biologically informed means to define CCs. Another

182   important advantage of the 99.5% ANI is that it can be automatically implemented and thus, does

183   not require manual curation, which is the case when establishing new ST numbers when novel

184   (meaning, not seen previously) sequences become available (4).

185

186   **What are the underlying mechanisms for the 99.5% ANI gap?**

187   The mechanism(s) that underlie the 99.5% ANI gap (or the earlier 95% ANI gap for the

188   species level) remain essentially speculative and should be the subject of future research in order

189   to further advance the mechanistic understanding of the microbial diversity patterns observed

190   in nature. Most notable is the idea that members of a population cohere together via means of

191   unbiased (random) genetic exchange which is more frequent within vs. between populations or

192   CCs (i.e., *the biological or sexual species concept*). A competing hypothesis is that several

193   members of the species are functionally differentiated from each other either due to

194   specialization for different growth conditions or different affinities for the same energy substrate

195   and thus, selection over time for these functions purge diversity (i.e., *the ecological species*

196   *concept*) (22-24). It is intriguing to note that the ecological explanation is also consistent with the

197   notion that CCs or different strains of the same species are somewhat ecologically and/or

198   functionally distinguishable from each other. Notably, given an estimated mutation rate of ~$4 \times 10^{-10}$

199   per nucleotide per generation (25) and between 100 to 300 generations per year (26), it would

200   take two distinct *E. coli* lineages or CCs at least forty thousand years since their last common

201   ancestor to accumulate 1% difference (i.e., fixed mutations) in their core genes or 99% ANI.

202   Therefore, there is enough time, at least theoretically, for the ecological purging of diversity to

203   take place at around the 99.5% ANI level and thus, account for the ANI patterns observed herein.

8

204    Intriguingly, it has been shown that the explicit inclusion of extinction events in a neutral model

205    of evolution can also result in punctuated distributions of genetic differentiation, opening up a

206    third possibility of historical contingency from stochastic events (27). However, we note that

207    while stochasticity can explain bimodal (or multimodal) distance distributions, a scarcity of ANI

208    values in the exact same range (i.e., around 99.5% ANI) would be unlikely to repeatedly emerge

209    by chance alone across many different species with distinct lifestyles and evolutionary tempo, as

210    opposed to this range varying between species. In any case, the data available in support of one

211    of these (or another) hypotheses remain sparse and/or anecdotal to date, to the best of our

212    knowledge, and the analysis presented in this study did not aim to advance this issue further.

213

214    **Conclusion**

215    Regardless of what the underlying mechanisms are for the 99.5% ANI gap, the results

216    presented here show that the patterns of natural diversity among thousands of sequenced

217    genomes are consistent with a 99.5% ANI threshold that can be used to identify CCs and strains

218    more reliably and precisely compared to the current practice. Regarding the use of this threshold

219    to define CCs vs. strains, we believe that the threshold is highly appropriate, as well as it matches

220    well the intended meaning and use of CCs, and thus its application to CC definition is

221    straightforward. For the strain level, 0.5% or 0.2% difference in ANI (correspondingly, 99.5% and

222    99.8% ANI) represents substantial, non-trivial, genomic divergence that, in most cases, would

223    likely encompass several genomes with at least some phenotypic differences (due to substantial

224    sequence or gene content differences among the genomes; see Fig. 1). Thus, multiple strains will

225    be likely grouped together under the same 99.5% ANI cluster in such cases, and strain, in general,

226    represents a more fine-grained level of resolution than the 99.5% ANI level. That said, we also

227    expect that the latter would somewhat depend on the context of the study and the existence or

228    not of phenotypic differences. For instance, it is possible that all sequence and gene-content

229    differences within some 99.5% ANI genome clusters to be neutral (or not functional), for at least

230    some growth conditions and habitats. Therefore, in such cases, a strain could be defined at the

231    99.5% ANI level and thus, include high(er) intra-strain sequence and/or gene-content diversity.

232    Hence, the 99.5% ANI level is also a good starting point in trying to define strains and

233    subsequently assess their gene and phenotypic differences (or lack of). Collectively, we expect

234    that the findings reported here will advance the molecular toolbox for accurately delineating and

235    following the important units of diversity withing prokaryotic species and thus, would greatly

236    facilitate future epidemiological and micro-diversity studies.

237

238    **Material and Methods**

239        Step by step methods, including how average trendlines were fit to the data, custom

240    Python code, NCBI Assembly accession numbers for selected genomes, and plots for each

241    selected                species                are                available                from:

242    https://github.com/rotheconrad/bacterial_strain_definition. Briefly, all genome sequences were

243    obtained from NCBI's RefSeq Assembly database on April 20th, 2022 and were labeled as

244    "complete" and "latest". ANI values and the shared genome fractions were directly obtained

245    from the output of FastANI version 1.32 "One to Many" mode with default settings (Jain et al

246    2018). Results were concatenated to create within species all vs. all output. Self matches were

247    removed, and genome pairs were filtered by minimum ANI values according to the axes of each

248    figure. Selected individual species plots (i.e., all vs. all output) of shared genome fraction vs. ANI

249    are shown in the Supplementary Material. *E. coli* genomes were assigned to sequence types (ST)

250    using the using the command-line tool mlst (https://github.com/tseemann/mlst) version 2.19.0

251    (20) with default settings.

252

253    **Acknowledgments**

254    This work has been supported by the US National Science Foundation (Award No 1759831 and

255    2129823) to KTK.

256

257    **Code and data availability**

258    All            code            and            data            details            are            available            from

259    https://github.com/rotheconrad/bacterial_strain_definition.

260

261    **Competing interests**

262 The authors declare no competing interests.

263

264 **Author contact list**

265 Luis M. Rodriguez-R: lmrodriguezr@gmail.com

266 Roth E. Conrad: rotheconrad@gatech.edu

267 Dorian J. Feistel: dfeistel3@gatech.edu

268 Tomeu Viver: tviver@imedea.uib-csic.es

269 Ramon Rosselló-Móra: ramon@imedea.uib-csic.es

270 Konstantinos T. Konstantinidis: kostas.konstantinidis@gatech.edu

271

272 **Figure 1. ANI vs. shared gene content for the 17,283 complete genomes used in this study.** Each

273 datapoint represents a comparison between a pair of genomes. FastANI (17) was used to

274 generate ANI values between the genomes of a pair (x-axis) and their shared genome fraction (y-

275 axis). The shared genome fraction was calculated by dividing the number of bidirectional

276 fragment mappings over the total query fragments determined by FastANI. Only a single set of

277 values is reported per pair, the one that used the longer genome as the reference (and the

278 reverse comparison was omitted). Note that only datapoints representing genome pairs sharing

279 ANI >95% are shown, and that panel B is a zoomed-in version of panel A. The main scatter plot is

280 shaded by density of the points using the Datashader package in Python with Matplotlib. The

281 trendline was calculated using linearGAM from pyGAM and includes the 95% confidence interval.

282 The marginal plots outside the two axes show histograms for the density of datapoints of each

283 axis. Note the low-density region in the ANI value distribution around 99.2-99.8% (Panel A), which

284 becomes more obvious when zooming in to the 98-100% ANI range (Panel B).
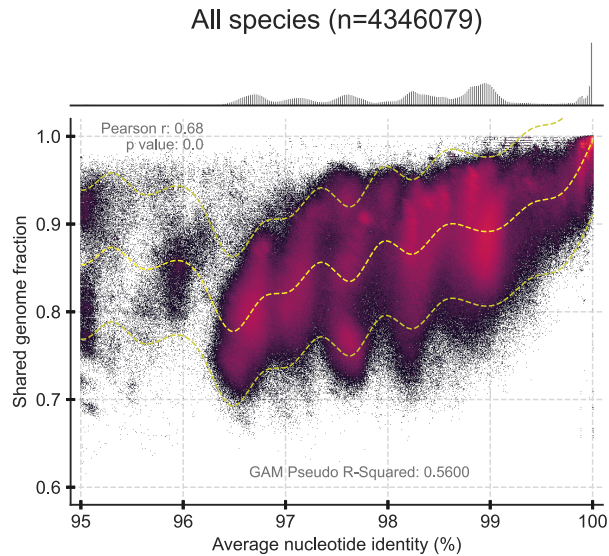
285 **Figure 1, panel A**.
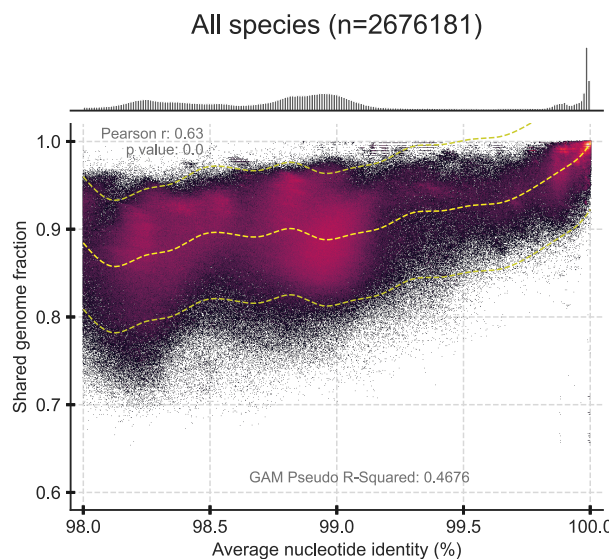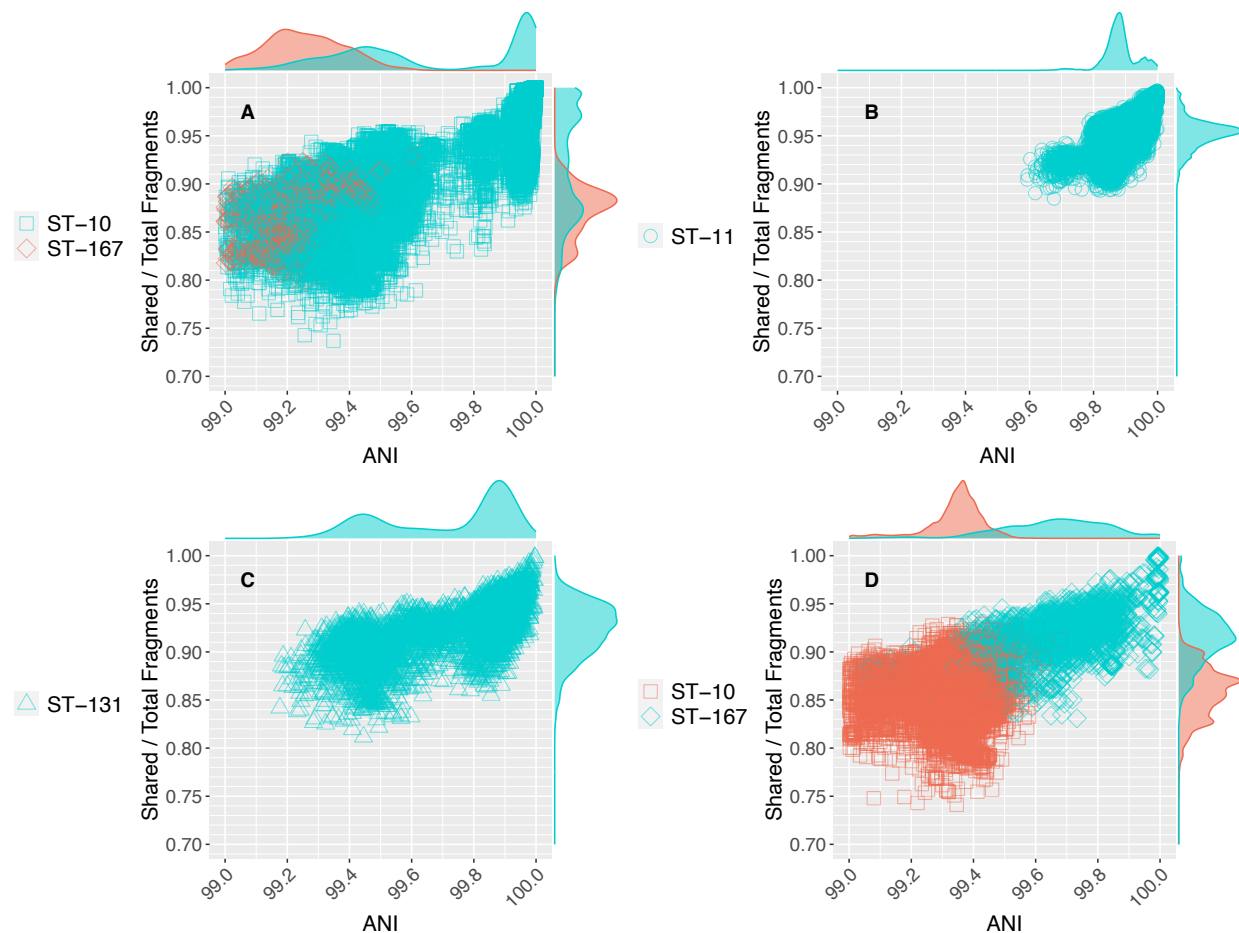
286

**Figure 1, panel B**.



288

**Figure 2. Comparison of the 99.5% ANI threshold to available clonal complexes of *E. coli*.** All *E. coli* genomes were assigned to a Sequence Type, ST (or clonal complex) using the command-line tool mlst (https://github.com/tseemann/mlst) version 2.19.0 (20). The four panels show the four most abundant STs based on the number of genomes assigned to them (see Table S1 and Figure S1 for underlying data and F1 statistic, respectively). Each datapoint is a comparison between two genomes, similar to those shown in Figure 1; the marginal plots show the kernel density estimate of datapoints for each axis. Datapoints are cyan if both genomes in the pair were assigned to the same, reference ST; red datapoints represent pairs for which one of the genomes

297 in the pair is assigned to a closely related, yet distinct ST than the reference ST. Note that for ST-

298 11, recall and precision of 99.5% ANI vs. CC is perfect because there are no genomes, and thus

299 STs, that are closely related to ST-11, which is also consistent with a pronounced ANI gap at 99.5%

300 for ST-11, and the substantial overlap in terms of ANI values between the closely related ST-10

301 and ST167 (low recall). ST-131 and ST-10 appear to harbor too much genomic diversity and could

302 be split in more than one CC based on the 99.5% ANI criterion (and the bimodal ANI value

303 distribution around 99.6% ANI).

304



305

## References

1. L. Dijkshoorn, B. M. Ursing, J. B. Ursing, Strain, clone and species: comments on three basic concepts of bacteriology. *J Med Microbiol* **49**, 397-401 (2000).

2. C. T. Parker, B. J. Tindall, G. M. Garrity, International Code of Nomenclature of Prokaryotes. *Int J Syst Evol Microbiol* 10.1099/ijsem.0.000778 (2015).

3. F. C. Tenover *et al.*, Interpreting chromosomal DNA restriction patterns produced by pulsed-field gel electrophoresis: criteria for bacterial strain typing. *J Clin Microbiol* **33**, 2233-2239 (1995).

4. M. C. Maiden *et al.*, Multilocus sequence typing: a portable approach to the identification of clones within populations of pathogenic microorganisms. *Proc Natl Acad Sci U S A* **95**, 3140-3145 (1998).

5. D. A. Baltrus, K. Dougherty, S. M. Beckstrom-Sternberg, J. S. Beckstrom-Sternberg, J. T. Foster, Incongruence between multi-locus sequence analysis (MLSA) and whole-genome-based phylogenies: Pseudomonas syringae pathovar pisi as a cautionary tale. *Mol Plant Pathol* **15**, 461-465 (2014).

6. K. T. Konstantinidis, A. Ramette, J. M. Tiedje, Toward a more robust assessment of intraspecies diversity, using fewer genetic markers. *Appl Environ Microbiol* **72**, 7286-7293 (2006).

7. R. L. Rodriguez, C. Jain, R. E. Conrad, S. Aluru, K. T. Konstantinidis, Reply to: "Re-evaluating the evidence for a universal genetic boundary among microbial species". *Nat Commun* **12**, 4060 (2021).

8. A. Caro-Quintero, K. T. Konstantinidis, Bacterial species may exist, metagenomics reveal. *Environ Microbiol* **14**, 347-355 (2012).

9. T. Viver *et al.*, Distinct ecotypes within a natural haloarchaeal population enable adaptation to changing environmental conditions without causing population sweeps. *ISME J* 10.1038/s41396-020-00842-5 (2020).

10. R. Conrad *et al.*, Toward quantifying the adaptive role of bacterial pangenomes during environmental perturbations. *The ISME Journal* https://doi.org/10.1101/2021.03.15.435471, Accepted (2021).

11. C. S. Murray, Y. Gao, M. Wu, Re-evaluating the evidence for a universal genetic boundary among microbial species. *Nat Commun* **12**, 4059 (2021).

12. K. T. Konstantinidis, E. F. DeLong, Genomic patterns of recombination, clonal divergence and environment in marine microbial populations. *ISME J* **2**, 1052-1065 (2008).

13. M. L. Bendall *et al.*, Genome-wide selective sweeps and gene-specific sweeps in natural bacterial populations. *ISME J* **10**, 1589-1601 (2016).

14. E. R. Johnston *et al.*, Metagenomics Reveals Pervasive Bacterial Populations and Reduced Community Diversity across the Alaska Tundra Ecosystem. *Front Microbiol* **7**, 579 (2016).

15. A. Meziti *et al.*, Quantifying the changes in genetic diversity within sequence-discrete bacterial populations across a spatial and temporal riverine gradient. *ISME J* **13**, 767-779 (2019).

16. M. R. Olm *et al.*, Consistent Metagenome-Derived Metrics Verify and Delineate Bacterial Species Boundaries. *mSystems* **5** (2020).

17.    C. Jain, R. L. Rodriguez, A. M. Phillippy, K. T. Konstantinidis, S. Aluru, High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* **9**, 5114 (2018).

18.    L. M. Rodriguez-R, C. Jain, R. E. Conrad, S. Aluru, K. T. Konstantinidis, Rebuttal to Murray and colleagues Nature Communication Matters Arising article. . *Narure Communications*, In press (2021).

19.    R. A. Welch *et al.*, Extensive mosaic structure revealed by the complete genome sequence of uropathogenic Escherichia coli. *Proc Natl Acad Sci U S A* **99**, 17020-17024 (2002).

20.    T. Wirth *et al.*, Sex and virulence in Escherichia coli: an evolutionary perspective. *Mol Microbiol* **60**, 1136-1151 (2006).

21.    A. Pena-Gonzalez *et al.*, Metagenomic Signatures of Gut Infections Caused by Different Escherichia coli Pathotypes. *Appl Environ Microbiol* **85** (2019).

22.    B. J. Shapiro *et al.*, Population genomics of early events in the ecological differentiation of bacteria. *Science* **336**, 48-51 (2012).

23.    C. Fraser, W. P. Hanage, B. G. Spratt, Recombination and the nature of bacterial speciation. *Science* **315**, 476-480 (2007).

24.    G. W. Tyson *et al.*, Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**, 37-43 (2004).

25.    J. W. Drake, B. Charlesworth, D. Charlesworth, J. F. Crow, Rates of spontaneous mutation. *Genetics* **148**, 1667-1686 (1998).

26.    Gibbons RJ, K. B, Estimates of the overall rate of growth of the intestinal microflora of hamsters, guinea pigs, and mice. *Journal of Bacteriology* **93**, 510–512 (1967).

27.    T. J. Straub, O. Zhaxybayeva, A null model for microbial diversification. *Proc Natl Acad Sci U S A* **114**, E5414-E5423 (2017).