



HAL
open science

A rapid and simple method for assessing and representing genome sequences relatedness

Martial Briand, Mariam Bouzid, Gilles Hunault, Marc Legeay, Marion Le Saux, Matthieu Barret

► **To cite this version:**

Martial Briand, Mariam Bouzid, Gilles Hunault, Marc Legeay, Marion Le Saux, et al.. A rapid and simple method for assessing and representing genome sequences relatedness. *BioRxiv*, 2019, 10.1101/569640 . hal-02171444

HAL Id: hal-02171444

<https://hal.science/hal-02171444>

Submitted on 2 Jul 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

1**Title: A rapid and simple method for assessing and representing genome sequences**

2**relatedness.**

3

4**Authors and affiliations :** Briand M^{1*}, Bouzid M^{1†}, Hunault G², Legeay M³, Fischer-Le Saux

5M¹, Barret M¹

6

7¹IRHS, Agrocampus-Ouest, INRA, Université d'Angers, SFR4207 QuaSaV, 49071

8Beaucouzé, France.

9² Université d'Angers, Laboratoire d'Hémodynamique, Interaction Fibrose et Invasivité

10tumorale hépatique, UPRES 3859, IFR 132, F-49045 Angers, France

11³Novo Nordisk Foundation Center for Protein Research, University of Copenhagen, Denmark

12

13**Authors email addresses:**

14Martial Briand: martial.briand@inra.fr

15Mariam Bouzid : mariam.bouzid@icloud.com

16Gilles Hunault : gilles.hunault@univ-angers.fr

17Marc Legeay : legeay.marc@free.fr

18Marion Fischer-Le Saux : marion.le-saux@inra.fr

19Matthieu Barret: matthieu.barret@inra.fr

20

21*Corresponding author

22Phone: +33 (0)2 41 22 57 29

23Fax: +33 (0)2 41 22 57 55

24 Abstract

25 Coherent genomic groups are frequently used as a proxy for bacterial species delineation
26 through computation of overall genome relatedness indices (OGRI). Average nucleotidic
27 identity (ANI) is the method of choice for estimating relatedness between genomic
28 sequences. However, pairwise comparisons of genome sequences based on ANI is relatively
29 computationally intensive and therefore precludes analyses of large datasets composed of
30 thousand genomes sequences.

31 In this work we evaluated an alternative OGRI based on *k*-mers counts to study prokaryotic
32 species delimitation. A dataset containing more than 3,500 *Pseudomonas* genome
33 sequences was successfully classified in few hours with the same precision than ANI. A new
34 visualization method based on zoomable circle packing was employed for assessing
35 relationships between the 350 cliques generated. Amendment of databases with these
36 *Pseudomonas* cliques greatly improve classification of metagenomic read sets with *k*-mers-
37 based classifier.

38 The developed workflow was integrated in the user-friendly KI-S tool that is available at the
39 following address: <https://iris.angers.inra.fr/galaxypub-cfbp>.

40

41 **Keywords : ANI, *k*-mers, circle packing, *Pseudomonas*, metagenome**

42

43

44 **Background**

45 Species is the unit of biological diversity. Species delineation of *Bacteria* and *Archaea*
46 historically relies on a polyphasic approach based on a range of genotypic, phenotypic and
47 chemo-taxonomic (e.g. fatty acid profiles) data of cultured specimen. According to the List of
48 Prokaryotic Names with Standing in Nomenclature (LPSN), approximately 15,500 bacterial
49 species names have been currently validated within this theoretical framework [1]. Since the
50 number of bacterial species inhabiting planet Earth is predicted to range between 10^7 to 10^{12}
51 species according to different estimates [2,3], the genomics revolution provides an
52 opportunity to accelerate the pace of species description.

53 Prokaryotic species are primarily described as cohesive genomic groups and
54 approaches based on similarity of whole genome sequence, also known as overall genome
55 relatedness indices (OGRI), have been proposed for delineating species. Average nucleotidic
56 identity (ANI) is nowadays the mostly acknowledged OGRI for assessing relatedness
57 between genomic sequences. Distinct ANI algorithms such as ANI based on BLAST (ANIb
58 [4]), ANI based on MUMmer (ANIm [5]) or ANI based on orthologous gene (OrthoANIb [6];
59 OrthoANLu [7]; gANI,AF [8]), which differ in their precisions but more importantly on their
60 calculation times [7], have been developed. Indeed, improvement of calculation time for
61 whole genomic comparison of large datasets is an essential parameter. As of November
62 2018, the total number of prokaryotic genome sequences publicly available in the NCBI
63 database is 170,728. Considering an average time of 1 second for calculating ANI values of
64 one pair of genome sequence, it would take approximately 1,000 years for obtaining ANI
65 values for all pairwise comparisons.

66 The number of words of length k (k -mers) shared between read sets [9] or genomic
67 sequences [10] is an alignment-free alternative for assessing the dis(similarities) between
68 entities. Methods based on k -mers counts, such as SIMKA [9], can quickly compute pairwise
69 comparison of multiple metagenome read sets with high accuracy. In addition, specific k -
70 mers profiles are now routinely employed by multiple read classifiers for estimating the

71 taxonomic structure of metagenome read sets [11–13]. While these k -mers based classifiers
72 differ in term of sensitivity and specificity [14], they rely on accurate genome databases for
73 affiliating read to a taxonomic rank.

74 The objective of the current work was to evaluate an alternative method based on k -
75 mers counts to study species delimitation on extensive genome datasets. We therefore
76 decided to employ k -mers counts for assessing similarity between genome sequences
77 belonging to the *Pseudomonas* genus. Indeed, this genus contains an important diversity of
78 species ($n = 207$), whose taxonomic affiliation is under constant evolution [15–21], and
79 numerous genome sequences are available in public databases. We also proposed an
80 original visualization tool based on D3 Zoomable Circle Packing
81 (<https://gist.github.com/mbostock/7607535>) for assessing relatedness of thousands of
82 genomes sequences. Finally, the benefit of taxonomic curation of reference database on the
83 taxonomic affiliation of metagenomics read sets was assessed. The developed workflow was
84 integrated in the user-friendly KI-S tool which is available in the galaxy toolbox of CIRM-
85 CFBP (<https://iris.angers.inra.fr/galaxypub-cfbp>).

86

87 **Methods**

88

89 **Genomic dataset**

90 All genome sequences ($n=3,623$ as of April 2017) from *Pseudomonas* genus were
91 downloaded from the NCBI database
92 (<https://www.ncbi.nlm.nih.gov/genome/browse#!/overview/>).

93

94 **Calculation of Overall Genome Relatedness Indices**

95 The percentage of shared k -mers between genome sequences was calculated with Simka
96 version 1.4 [9] with the following parameters (abundance-min 1 and k -mers length ranging
97 from 10 to 20). The percentage of shared k -mers was compared to ANIb values calculated
98 with PYANI version 0.2.3 (<https://github.com/widowquinn/pyani>). Due to the computing time
99 required for ANIb calculation, only a subset of *Pseudomonas* genomic sequences ($n=934$)
100 was selected for this comparison. This subset was composed of genome sequences
101 containing less than 150 scaffolds.

102

103 **Development of KI-S tool**

104 An integrative tool named KI-S was developed. The number of shared k -mers between
105 genome sequences is first calculated with Simka [9]. A custom R script is then employed to
106 cluster the genome sequences according to their connected components at different selected
107 threshold (e.g. 50% of shared 15-mers). The clustering result is visualized with Zoomable
108 Circle Packing representation with the D3.js JavaScript library
109 (<https://gist.github.com/mbostock/7607535>). The source code of the KI-S tool is available at
110 the following address: <https://sourcesup.renater.fr/projects/ki-s/>. A wrapper for accessing KI-S
111 in a user-friendly Galaxy tool is also available at the following address:
112 <https://iris.angers.inra.fr/galaxypub-cfbp>.

113

114 **Taxonomic inference of metagenomic read sets**

115 The taxonomic profiles of 9 metagenome read sets derived from seed, germinating seeds
116 and seedlings of common bean (*Phaseolus vulgaris* var. Flavert) were estimated with Clark
117 version 1.2.4 [13]. These metagenome datasets were selected because of the high relative
118 abundance of reads affiliated to *Pseudomonas* [22]. The following Clark default parameters –
119 `k 31 -t <minFreqTarget> 0` and `-o <minFreqObject> 0` were used for the taxonomic profiling.
120 Three distinct Clark databases were employed: (i) the original Clark database from
121 NCBI/RefSeq at the species level (ii) the original Clark database supplemented with the
122 3,623 *Pseudomonas* genome sequences and their original NCBI taxonomic affiliation (iii) the
123 original Clark database supplemented with the 3,623 *Pseudomonas* genome sequences
124 whose taxonomic affiliation was corrected according to the reclassification based on the
125 number of shared *k*-mers. For this third database, genome sequences were clustered at
126 >50% of 15-mers.

127

128 **Results**

129 **Selection of optimal *k*-mers size and percentage of shared *k*-mers**

130 Using the percentage of shared *k*-mers as an OGRI for species delineation first required to
131 determine the optimal *k*-mers size. This was performed by comparing the percentage of
132 shared *k*-mers to a widely acknowledged OGRI, ANIb [4], between 934 *Pseudomonas*
133 genome sequences. Since species delineation threshold was initially proposed following the
134 observation of a gap in the distribution of pairwise comparison values [23], the distribution
135 profiles obtained with *k*-mers lengths ranging from 10 to 20 were compared to ANIb values.
136 Short *k*-mers ($k < 12$) were evenly shared by most strains and then not discriminative (**Fig.**
1371). As the size of the *k*-mers increased, a multimodal distribution based on four peaks were
138 observed (**Fig. 1**). The first peak is related to genomes sequences that do not belong to the
139 same species. Then, depending on *k* length, the second and third peaks (e.g. 50% and 80%
140 for $k = 15$) corresponded to genome sequences associated to the same species and
141 subspecies, respectively. The fourth peak at 100% of shared *k*-mers was related to identical
142 genome sequences.

143 Fifty percent of 15-mers is closed to ANIb value of 0.95 (**Fig. 2**), a threshold
144 commonly employed for delineating bacterial species level [4]. More precisely the median
145 percentage of shared 15-mers is 49% [34%-66%] for ANIb value ranging from 0.94 to 0.96. In
146 addition, 15-mers allows the investigation of inter-and infra-specific relationship at lower and
147 higher percentage of shared 15-mers, respectively.

148 Computing time of 15-mers for 934 genome sequences was 4 hours on a DELL
149 Power Edge R510 server, while it took approximately 3 months for obtaining all ANIb pairwise
150 comparisons (500-fold decrease of computing time).

151

152 **Classification of *Pseudomonas* genomes**

153 The percentage of shared 15-mers was then used to investigate relatedness between 3,623
154 *Pseudomonas* genome sequences publicly available. At a threshold of 50% of 15-mers, we

155identified 350 cliques. The clique containing the most abundant number of genome
156sequences was by far related to *P. aeruginosa* species ($n = 2,341$), followed by the
157phylogroups PG1 ($n = 111$), PG3 ($n = 92$) and PG2 ($n = 74$) of *P. syringae* species complex
158([16]; **Table S1**). At the clustering threshold employed, 185 cliques were composed of a
159single genome sequence, therefore highlighting the high *Pseudomonas* strain diversity.
160Moreover, according to Chao1 index, *Pseudomonas* species richness is estimated at 629
161cliques [± 57], which indicates that additional strain isolations and sequencing effort are
162needed to cover the whole diversity of this bacterial genus. Graphical representation of
163hierarchical clustering by dendrogram for a large dataset is generally not optimal. Here we
164employed Zoomable circle packing as an alternative to dendrogram for representing similarity
165between genome sequences (**Fig. 3** and **FigS1.html**). The different clustering thresholds that
166can be superimposed on the same graphical representation allow the investigation of inter-
167and intra- groups relationships (**Fig. 3** and **FigS1.html**). This is useful for affiliating specific
168clique to a group or subgroup of *Pseudomonas* species.

169

170**Improvement of taxonomic affiliation of metagenomic read sets.**

171The taxonomic composition of metagenome read sets is frequently estimated with k -mers
172based classifiers. While these k -mers based classifiers differ in term of sensitivity and
173specificity, they all rely on accurate genome databases for affiliating read to a taxonomic
174rank. Here, we investigated the impact of database content and curation on taxonomic
175affiliation. Using Clark [13] as a taxonomic profiler with the original Clark database, we
176classified metagenome read sets derived from bean seeds, germinating seeds and seedlings
177[22]. Adding the 3,623 *Pseudomonas* genomes with their original taxonomic affiliation from
178NCBI to the original Clark database did not increase the percentage of classified reads (**Fig.**
179**4**). However, adding the same genome sequences reclassified in cliques according to their
180percentage of shared k -mers ($k=15$; threshold= 50%) increased 1.4-fold on average the
181number of classified reads (**Fig. 4**).

182

183

184 Discussion

185 Classification of bacterial strains on the basis on their genome sequences similarities has
186 emerged since a decade as an alternative to the cumbersome DNA-DNA hybridizations [24].
187 Although ANIb is the current gold-standard method for investigating these genomic
188 relatedness, its intensive computational time prohibited its used for comparing large genome
189 datasets [7]. In contrast, investigating the percentage of shared k -mers is scalable for
190 comparing thousands of genome sequences.

191 In a method based on k -mers counts, choosing the length of k is a compromise
192 between accuracy and speed. The distribution of shared k -mers values between genome
193 sequences is impacted by k length. For $k = 15$, four peaks were observed at 15%, 50%, 80%
194 and 100% of shared k -mers. The second peak is closed to ANIb value of 0.95 and falls in the
195 so called grey or fuzzy zone [24] where taxonomists might decide to split or merge species.
196 Hence, according to our working dataset, it seems that 50% of 15-mers is a good proxy for
197 estimating *Pseudomonas* clique. Despite the diverse range of habitats colonized by different
198 *Pseudomonas* populations [19], it is likely that the percentage of shared k -mers has to be
199 adapted when investigating other bacterial genera. Indeed, since population dynamics,
200 lifestyle and location impact molecular evolution, it is somewhat illusory to define a fixed
201 threshold for species delineation [25]. While 15-mers is a good starting point for investigating
202 infra-specific to infra-generic relationships between genome sequences, the computational
203 speed of KI-S offers the possibility to perform large scale genomic comparisons at different k
204 sizes to select the most appropriate threshold.

205 Genomic relatedness using whole genome sequences becomes a standard for
206 bacterial strain identification and bacterial taxonomy [24,26]. This proposition is primarily
207 motivated by fast and inexpensive sequencing of bacterial genome together with the limited
208 availability of cultured specimen for performing classical polyphasic approach. Whether full
209 genome sequences should represent the basis of taxonomic classification is an ongoing
210 debate between systematians [27]. While this consideration is well beyond the objectives of

211this work, obtaining a classification of bacterial genome sequences into coherent groups is of
212general interest. Indeed, the number of misidentified genomes sequences is exponentially
213growing in public databases. A number of initiatives such as Digital Protologue Database
214(DPD [28]), Microbial Genomes Atlas (MiGA [29]), Life Identification Numbers database
215(LINbase [30]) or the Genome Taxonomy Database (GTDB [26]) proposed services to
216classify and rename bacterial strains based ANIb values or single copy marker proteins.
217Using the percentage of shared *k*-mers between unknown bacterial genome sequences and
218reference genome sequences associated to these databases could provide a rapid
219complementary approach for bacterial classification. Moreover, KI-S tool, provides a friendly
220visualization interface that could help systematians to curate whole genome databases.
221Indeed, zoomable circle packing could be employed for highlighting (*i*) misidentified strains,
222(*ii*) bacterial taxa that possess representative type strains or (*iii*) bacterial taxa that contain
223few genomes sequences.

224 Association between a taxonomic group and its distribution across a range of habitats
225is useful for inferring the role of this taxa on its host or environment. For instance, community
226profiling approaches based on molecular marker such as hypervariable regions of 16S rRNA
227gene have been helpful for highlighting correlations between host fitness and microbiome
228composition. Finer-grained taxonomic resolution of microbiome composition could be
229achieved with metagenomics through *k*-mers based classification of reads. In this study we
230demonstrate that employing a database with a classification of strains reflecting their
231genomic relatedness greatly improve taxonomic assignments of reads. Therefore,
232investigating relationship between bacterial genome sequences not only benefits bacterial
233taxonomy but also deserves microbial ecology.

234 **Competing interests**

235 The authors declare that they have neither competing interests nor conflict of interest.

236 **Funding**

237 This research was supported in part by grant awarded by the Region des Pays de la Loire
238 (metaSEED, 2013 10080).

239 **Acknowledgements**

240 The authors wish to thank Claire Lemaitre and Guillaume Rizk for their assistance with the
241 SIMKA software.

242

243References

1. Parte AC. LPSN - List of Prokaryotic names with Standing in Nomenclature (bacterio.net), 20 years on. *Int J Syst Evol Microbiol.* 2018;68:1825–9.
2. Amann R, Rosselló-Móra R. After All, Only Millions? *MBio.* 2016;7:e00999-16.
3. Locey KJ, Lennon JT. Scaling laws predict global microbial diversity. *PNAS.* 2016;113:5970–5.
4. Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM. DNA-DNA hybridization values and their relationship to whole-genome sequence similarities. *Int J Syst Evol Microbiol.* 2007;57:81–91.
5. Richter M, Rosselló-Móra R. Shifting the genomic gold standard for the prokaryotic species definition. *Proc Natl Acad Sci USA.* 2009;106:19126–31.
6. Lee I, Ouk Kim Y, Park S-C, Chun J. OrthoANI: An improved algorithm and software for calculating average nucleotide identity. *Int J Syst Evol Microbiol.* 2016;66:1100–3.
7. Yoon S-H, Ha S-M, Lim J, Kwon S, Chun J. A large-scale evaluation of algorithms to calculate average nucleotide identity. *Antonie Van Leeuwenhoek.* 2017;110:1281–6.
8. Varghese NJ, Mukherjee S, Ivanova N, Konstantinidis KT, Mavrommatis K, Kyripides NC, et al. Microbial species delineation using whole genome sequences. *Nucleic Acids Res.* 2015;43:6761–71.
9. Benoit G, Peterlongo P, Mariadassou M, Drezen E, Schbath S, Lavenier D, et al. Multiple Comparative Metagenomics using Multiset k-mer Counting. *PeerJ Computer Science.* 2016 ; 2:e94
10. Déraspe M, Raymond F, Boisvert S, Culley A, Roy PH, Laviolette F, et al. Phenetic Comparison of Prokaryotic Genomes Using k-mers. *Mol Biol Evol.* 2017;34:2716–29.
11. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biology.* 2014;15:R46.
12. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome Res.* 2016; 26: 1721-1729
13. Ounit R, Wanamaker S, Close TJ, Lonardi S. CLARK: fast and accurate classification of metagenomic and genomic sequences using discriminative k-mers. *BMC Genomics.* 2015;16:236.
14. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical Assessment of Metagenome Interpretation – a benchmark of computational metagenomics software. *Nat Methods.* 2017;14:1063–71.
15. Peix A, Ramírez-Bahena M-H, Velázquez E. Historical evolution and current status of the taxonomy of genus *Pseudomonas*. *Infect Genet and Evol.* 2009;9:1132–47.
16. Berge O, Monteil CL, Bartoli C, Chandeysson C, Guilbaud C, Sands DC, et al. A user's guide to a data base of the diversity of *Pseudomonas syringae* and its application to classifying strains in this phylogenetic complex. *PLoS ONE.* 2014;9:e105547.
17. Gomila M, Busquets A, Mulet M, García-Valdés E, Lalucat J. Clarification of Taxonomic Status within the *Pseudomonas syringae* Species Group Based on a Phylogenomic Analysis. *Front Microbiol.* 2017;8:2422.

18. Gomila M, Peña A, Mulet M, Lalucat J, García-Valdés E. Phylogenomics and systematics in *Pseudomonas*. *Front Microbiol*. 2015;6:214.
19. Peix A, Ramírez-Bahena M-H, Velázquez E. The current status on the taxonomy of *Pseudomonas* revisited: An update. *Infect Genet Evol*. 2018;57:106–16.
20. Garrido-Sanz D, Meier-Kolthoff JP, Göker M, Martín M, Rivilla R, Redondo-Nieto M. Genomic and Genetic Diversity within the *Pseudomonas fluorescens* Complex. *PLoS ONE*. 2016;11:e0150183.
21. Hesse C, Schulz F, Bull CT, Shaffer BT, Yan Q, Shapiro N, et al. Genome-based evolutionary history of *Pseudomonas* spp. *Environ Microbiol*. 2018; doi: 10.1111/1462-2920.14130
22. Torres-Cortés G, Bonneau S, Bouchez O, Genthon C, Briand M, Jacques M-A, et al. Functional Microbial Features Driving Community Assembly During Seed Germination and Emergence. *Front Plant Sci*. 2018;9:902.
23. Patrick A Grimont. Use of DNA reassociation in bacterial classification. *Canadian J Microbiol*; 1988; 34:541-6.
24. Rosselló-Móra R, Amann R. Past and future species definitions for Bacteria and Archaea. *Syst Appl Microbiol*. 2015;38:209–16.
25. Bromham L. Why do species vary in their rate of molecular evolution? *Biol Lett*. 2009;5:401–4.
26. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarszewski A, Chaumeil P-A, et al. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nature Biotech*. 2018;36:996–1004.
27. Garrity GM. A New Genomics-Driven Taxonomy of Bacteria and Archaea: Are We There Yet? *J Clin Microbiol*. 2016;54:1956–63.
28. Rossello-Mora R, Sutcliffe IC. Reflections on the introduction of the Digital Protologue Database — A partial success? *System Appl Microbiol*. 2019; 42:1-2.
29. Rodriguez-R LM, Gunturu S, Harvey WT, Rosselló-Mora R, Tiedje JM, Cole JR, et al. The Microbial Genomes Atlas (MiGA) webserver: taxonomic and gene diversity analysis of Archaea and Bacteria at the whole genome level. *Nucleic Acids Res*. 2018;46:W282–8.
30. Vinatzer BA, Tian L, Heath LS. A proposal for a portal to make earth's microbial diversity easily accessible and searchable. *Antonie van Leeuwenhoek*. 2017;110:1271–9.

244 **Figures and Supplemental files**

245 **Figure 1: Distribution of shared *k*-mers values.** Relatedness between genome sequences
246 were estimated with ANIb (green) or shared *k*-mers (blue). The x axis represents ANIb or
247 percentage of shared *k*-mers while the y axis represents the number of values by class in the
248 subset of 934 *Pseudomonas* genomic comparison.

249 **Figure 2: Comparison of various *k*-mers length and ANIb values.** Pairwise similarities
250 between genome sequences were assessed with average nucleotidic identity based on
251 BLAST (ANIb, x-axis) and percentage of shared *k*-mers of length 10 (**A**), 15 (**B**) and 20 (**C**).
252 The red line corresponds to ANIb of 0.95, a threshold commonly employed for delineating
253 species level.

254 **Figure 3: Hierarchical clustering of *Pseudomonas* genome sequences.** Zoomable circle
255 packing representation of *Pseudomonas* genome sequences ($n = 3,623$). Similarities
256 between genome sequences were assessed by comparing the percentage of shared 15-
257 mers. Each dot represents a genome sequence, which is colored according to its group of
258 species [16,21]. These genome sequences have been grouped at three distinct thresholds
259 for assessing intraspecific (0.75), species-specific (0.5) and interspecies relationships (0.25).

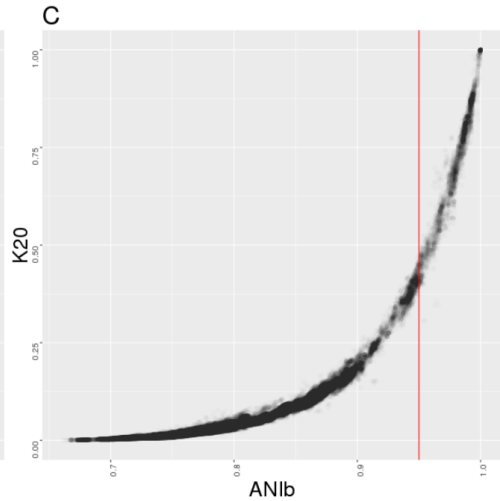
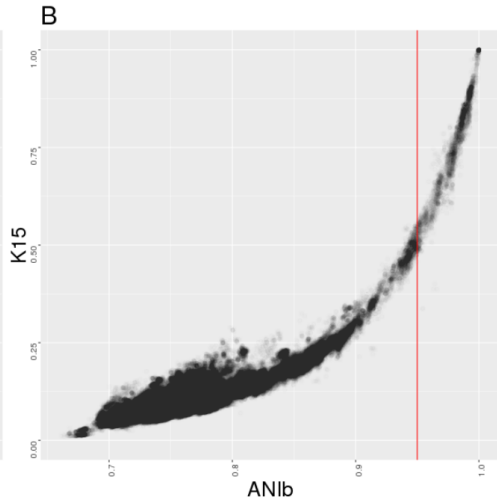
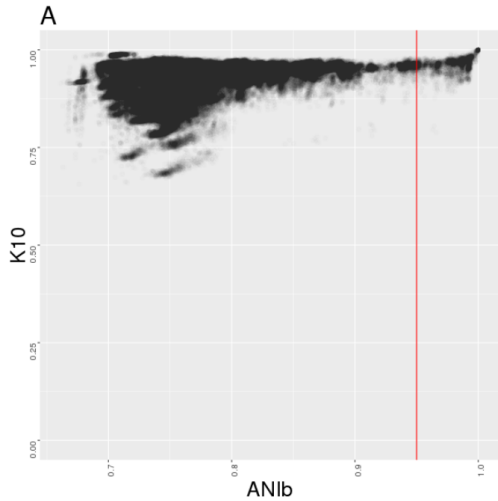
260 **Figure 4: Percentage of classified reads.** Classification of metagenome read sets derived
261 from bean seeds, germinating seeds and seedlings with Clark [13]. Three distinct databases
262 were employed for read classification: the original Clark database (red), Clark database
263 supplemented with 3,623 *Pseudomonas* genome sequences (green) and the Clark database
264 supplemented with 3,623 *Pseudomonas* genome sequences that were classified according
265 to their percentage of shared *k*-mers (blue).

266 **TableS1.csv : *Pseudomonas* cliques.** Description of the 350 cliques obtained after
267 clustering at 50% of shared 15-mers. For each clique, the *Pseudomonas* group [21] and
268 subgroup [16,21] are displayed.

269 **FigureS1.html: Zoomable circle packing representation of *Pseudomonas* genome
270 sequences.** Similarities between genome sequences were assessed by comparing the

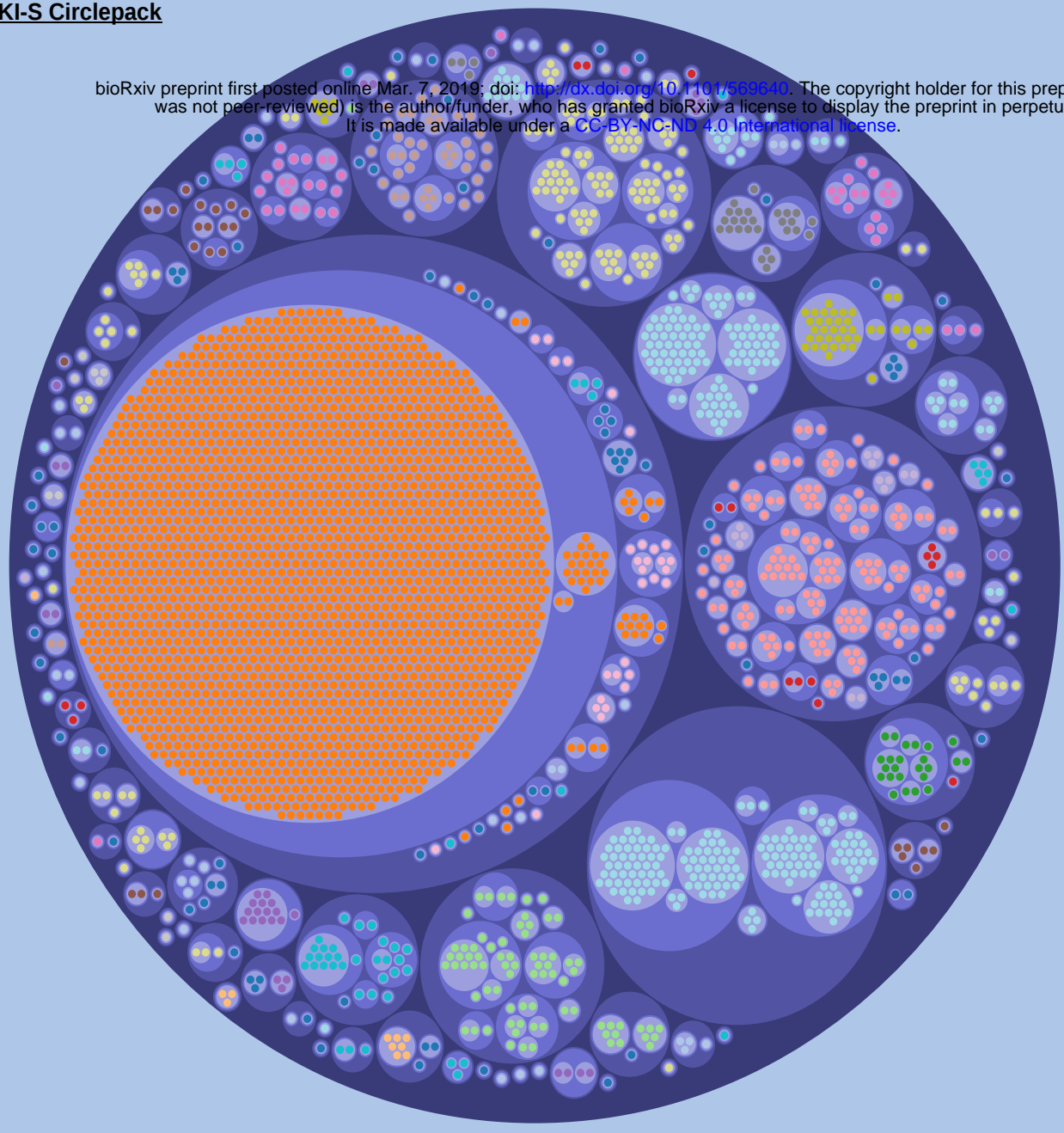
271percentage of shared 15-mers. Each dot represents a genome sequence, which is colored
272according to its group of species [16,21]. These genome sequences have been grouped at
273three distinct thresholds for assessing intraspecific (0.75), species-specific (0.5) and
274interspecies relationships (0.25).

275



KI-S Circlepack

bioRxiv preprint first posted online Mar. 7, 2019; doi: <https://doi.org/10.1101/389840>. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).



Thresholds

■ 0.25 ■ 0.5 ■ 0.75

Groups

- P_putida_group
- P_gessardii_subgroup
- P_fragi_subgroup
- P_corrugata_group
- P_pertucinogena_group
- Other
- P_fluorescens_group
- P_aeruginosa_group
- P_jessenii_subgroup
- P_syringae_group
- P_koreensis_subgroup
- P_oryzihabitans_group
- P_protegens_subgroup
- P_stutzeri_group
- P_chlororaphis_group
- P_asplenii_subgroup
- P_fluorescens_subgroup
- P_mandelii_subgroup
- P_oleovorans_group
- Unknown

