

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

## **2relatedness.**

3

4**Authors and affiliations** : Briand M<sup>1\*</sup>, Bouzid M<sup>1†</sup>, Hunault G<sup>2</sup>, Legeay M<sup>3</sup>, Fischer-Le Saux

5M<sup>1</sup>, Barret M<sup>1</sup>

6

7<sup>1</sup>IRHS, Agrocampus-Ouest, INRA, Université d'Angers, SFR4207 QuaSaV, 49071

8Beaucouzé, France.

9<sup>2</sup> Université d'Angers, Laboratoire d'Hémodynamique, Interaction Fibrose et Invasivité

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

11<sup>3</sup>Novo Nordisk Foundation Center for Protein Research, University of Copenhagen, Denmark

12

## **13Authors email addresses:**

14Martial Briand: [martial.briand@inra.fr](mailto:martial.briand@inra.fr)

15Mariam Bouzid : [mariam.bouzid@icloud.com](mailto:mariam.bouzid@icloud.com)

16Gilles Hunault : [gilles.hunault@univ-angers.fr](mailto:gilles.hunault@univ-angers.fr)

17Marc Legeay : [legeay.marc@free.fr](mailto:legeay.marc@free.fr)

18Marion Fischer-Le Saux : [marion.le-saux@inra.fr](mailto:marion.le-saux@inra.fr)

19Matthieu Barret: [matthieu.barret@inra.fr](mailto:matthieu.barret@inra.fr)

20

21\*Corresponding author

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

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

## 24Abstract

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

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

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

40

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

42

43

## 44Background

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

53 Prokaryotic species are primarily described as cohesive genomic groups and  
 54approaches based on similarity of whole genome sequence, also known as overall genome  
 55relatedness indices (OGRI), have been proposed for delineating species. Genome Blast  
 56Distance Phylogeny (GBDP [4]) and Average nucleotide identity (ANI) are currently the most  
 57frequently used OGRI for assessing relatedness between genomic sequences. Distinct ANI  
 58algorithms such as ANI based on BLAST (ANId [5]), ANI based on MUMmer (ANIm [6]) or  
 59ANI based on orthologous genes (OrthoANId [7]; OrthoANId [8]; gANI,AF [9]), which differ in  
 60their precision but more importantly in their calculation times [8], have been developed.  
 61Indeed, improvement of calculation time for whole genomic comparison of large datasets is  
 62an essential parameter. As of November 2018, the total number of prokaryotic genome  
 63sequences publicly available in the NCBI database is 170,728. Considering an average time  
 64of 1 second for calculating ANI values for one pair of genome sequence, it would take  
 65approximately 1,000 years to obtain ANI values for all pairwise comparisons.

66 The number of words of length  $k$  ( $k$ -mers) shared between read sets [10] or genomic  
 67sequences [11] is an alignment-free alternative for assessing the similarities between entities.  
 68Methods based on  $k$ -mer counts, such as SIMKA [10], can quickly compute pairwise  
 69comparison of multiple metagenome read sets with high accuracy. In addition, specific  $k$ -mer  
 70profiles are now routinely employed by multiple read classifiers for estimating the taxonomic

71structure of metagenome read sets [12–14]. While these *k*-mer based classifiers differ in term  
72of sensitivity and specificity [15], they rely on accurate genome databases for affiliating read  
73to a taxonomic rank.

74       The objective of the current work was to evaluate an alternative method based on *k*-  
75mer counts to study species delimitation on extensive genome datasets. We therefore  
76decided to employ *k*-mer counting to assess the similarity among genome sequences  
77belonging to the *Pseudomonas* genus. Indeed, this genus contains an important diversity of  
78species ( $n = 207$ ), whose taxonomic affiliation is under constant evolution [16–22], and  
79numerous genome sequences are available in public databases. We also proposed an  
80original visualization tool based on D3 Zoomable Circle Packing  
81(<https://gist.github.com/mbostock/7607535>) for assessing relatedness of thousands of  
82genome sequences. Finally, the benefit of taxonomic curation of reference database on the  
83taxonomic affiliation of metagenomics read sets was assessed. The developed workflow was  
84integrated in the user-friendly KI-S tool which is available in the galaxy toolbox of CIRM-  
85CFBP (<https://iris.angers.inra.fr/galaxypub-cfbp>).

86

## 87Methods

88

### 89Genomic dataset

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

93

### 94Calculation of Overall Genome Relatedness Indices

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

102

### 103Development of KI-S tool

104An integrative tool named KI-S was developed. The number of shared  $k$ -mers between  
105genome sequences was initially calculated with Simka [10]. A custom R script was then  
106employed to cluster the genome sequences according to their connected components at  
107different selected thresholds (e.g. 50% of shared 15-mers). The clustering result is visualized  
108with Zoomable 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  
110the following address: <https://sourcesup.renater.fr/projects/ki-s/>. A wrapper for accessing KI-S  
111in a user-friendly Galaxy tool is also available at the following address:  
112<https://iris.angers.inra.fr/galaxypub-cfbp>.

113

# 114Taxonomic inference of metagenomic read sets

115The taxonomic profiles of 9 metagenomic read sets derived from seed, germinating seeds  
 116and seedlings of common bean (*Phaseolus vulgaris* var. Flavert) were estimated with Clark  
 117version 1.2.4 [14]. These metagenomic datasets were selected because of the high relative  
 118abundance of reads affiliated to *Pseudomonas* [23]. The following Clark default parameters –  
 119k 31 –t <minFreqTarget> 0 and -o <minFreqObject> 0 were used for the taxonomic profiling.  
 120Three distinct Clark databases were employed: (i) the original Clark database from  
 121NCBI/RefSeq at the species level (ii) the original Clark database supplemented with the  
 1223,623 *Pseudomonas* genome sequences and their original NCBI taxonomic affiliation (iii) the  
 123original Clark database supplemented with the 3,623 *Pseudomonas* genome sequences  
 124whose taxonomic affiliation was corrected according to the reclassification based on the  
 125number of shared *k*-mers. For this third database, genome sequences were clustered at  
 126>50% of 15-mers.

127

# 128Results

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

130Using the percentage of shared *k*-mers as an OGRI for species delineation first required the  
 131determination of the optimal *k*-mer size. This was performed by comparing the percentage of  
 132shared *k*-mers to a widely employed OGRI, ANIb [5], among 934 *Pseudomonas* genome  
 133sequences. Since the species delineation threshold was initially proposed following the  
 134observation of a gap in the distribution of pairwise comparison values [24], the distribution  
 135profiles obtained with *k*-mer lengths ranging from 10 to 20 were compared to ANIb values.  
 136Short *k*-mers ( $k < 12$ ) were evenly shared by most strains and not discriminative (**Fig. 1**). As  
 137the length of the *k*-mer increased, a multimodal distribution based on four peaks was  
 138observed (**Fig. 1**). The first peak related to the genome sequences that do not belong to the  
 139same species. Then, depending on *k* length, the second and third peaks (e.g. 50% and 80%  
 140for  $k = 15$ ) corresponded to genome sequences associated to the same species and  
 141subspecies, respectively. The fourth peak at 100% of shared *k*-mers was related to identical  
 142genome sequences.

143 Fifty percent of 15-mers is close to ANIb value of 0.95 (**Fig. 2**), a threshold commonly  
 144employed for delineating bacterial species [5]. More precisely, the median percentage of  
 145shared 15-mers is 49% [34%-66%] for ANIb value ranging from 0.94 to 0.96. In addition, 15-  
 146mers allows the investigation of inter-and infra-specific relationship at lower and higher  
 147percentage of shared 15-mers, respectively.

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

151

## 152Classification of *Pseudomonas* genome sequences

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

155 identified 350 cliques. The clique containing the most abundant number of genome  
156 sequences was related to *P. aeruginosa* ( $n = 2,341$ ), followed by the phylogroups PG1 ( $n =$   
157 111), PG3 ( $n = 92$ ) and PG2 ( $n = 74$ ) of *P. syringae* species complex ([17]; **Table S1**). At the  
158 clustering threshold employed, 185 cliques were composed of a single genome sequences,  
159 therefore highlighting the high *Pseudomonas* strain diversity. Moreover, according to Chao1  
160 index, *Pseudomonas* species richness is estimated at 629 cliques [ $\pm 57$ ], which indicates that  
161 additional strain isolations and sequencing effort are needed to cover the whole diversity of  
162 this bacterial genus. Graphical representation of hierarchical clustering by dendrogram for a  
163 large dataset is generally not optimal. Here we employed Zoomable circle packing as an  
164 alternative to dendrogram for representing similarity between genome sequences (**Fig. 3** and  
165 **FigS1.html**). The different clustering thresholds that can be superimposed on the same  
166 graphical representation allow the investigation of inter- and intra- groups relationships (**Fig.**  
167 **3** and **FigS1.html**). This is useful for affiliating a specific clique to a group or subgroup of  
168 *Pseudomonas* species.

169

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

171 The taxonomic composition of metagenome read sets is frequently estimated with *k*-mer  
172 based classifiers. While these *k*-mer based classifiers differ in term of sensitivity and  
173 specificity, they all rely on accurate genome databases for affiliating reads to taxonomic rank.  
174 Here, we investigated the impact of database content and curation on taxonomic affiliation.  
175 Using Clark [14] as a taxonomic profiler with the original Clark database, we classified  
176 metagenome read sets derived from bean seeds, germinating seeds and seedlings [23].  
177 Adding the 3,623 *Pseudomonas* genome sequences with their original taxonomic affiliation  
178 from NCBI to the original Clark database did not increase the percentage of classified reads  
179 (**Fig. 4**). However, adding the same genome sequences reclassified in cliques according to  
180 their percentage of shared *k*-mers ( $k=15$ ; threshold= 50%) increased 1.4-fold on average the  
181 number of classified reads (**Fig. 4**).



182

183

# 184Discussion

185Classification of bacterial strains on the basis on their genome sequence similarities has  
186emerged over the last decade as an alternative to the cumbersome DNA-DNA hybridizations  
187[4, 25]. Although ANIb is one widely employed method for investigating genomic relatedness,  
188its intensive computational time prohibited its used for comparing large genome datasets [8].  
189In contrast, investigating the percentage of shared  $k$ -mers is scalable for comparing  
190thousands of genome sequences.

191 In a method based on  $k$ -mer counts, choosing the length of  $k$  is a compromise  
192between accuracy and speed. The distribution of shared  $k$ -mer values between genome  
193sequences is impacted by  $k$  length. For  $k = 15$ , four peaks were observed at 15%, 50%, 80%  
194and 100% of shared  $k$ -mers. The second peak is close to ANIb value of 0.95 and falls in the  
195so called grey or fuzzy zone [25] where taxonomists might decide to split or merge species.  
196Hence, according to our working dataset, it seems that 50% of 15-mers is a good proxy for  
197estimating *Pseudomonas* clique. Despite the diverse range of habitats colonized by different  
198*Pseudomonas* populations [20], it is likely that the percentage of shared  $k$ -mers has to be  
199adapted when investigating other bacterial genera. Indeed, since population dynamics,  
200lifestyle and location impact molecular evolution, it is somewhat illusory to define a fixed  
201threshold for species delineation [26]. While 15-mers is a good starting point for investigating  
202infra-specific to infra-generic relationships between genome sequences, the computational  
203speed of KI-S offers the possibility to perform large scale genomic comparisons at different  $k$   
204sizes to select the most appropriate threshold.

205 Genomic relatedness using whole genome sequences has become the standard  
206method for bacterial strain identification and bacterial taxonomy [4,25,27]. This is primarily  
207motivated by fast and inexpensive sequencing of bacterial genomes together with the limited  
208availability of cultured specimen for performing classical polyphasic approach. Whether full  
209genome sequences should represent the basis of taxonomic classification is an ongoing  
210debate between systematians [28]. While this consideration is well beyond the objectives of

this work, obtaining a classification of bacterial genome sequences into coherent groups is of general interest. Indeed, the number of misidentified genome sequences is exponentially growing in public databases. A number of initiatives such as Digital Protologue Database (DPD [29]), Microbial Genomes Atlas (MiGA [30]), Life Identification Numbers database (LINbase [31]) or the Genome Taxonomy Database (GTDB [27]) proposed services to classify and rename bacterial strains based ANI values or single copy marker proteins. Using the percentage of shared *k*-mers between unknown bacterial genome sequences and reference genome sequences associated to these databases could provide a rapid complementary approach for bacterial classification. Moreover, KI-S tool, provides a friendly visualization interface that could help systematians to curate whole genome databases. Indeed, zoomable circle packing could be employed for highlighting (i) misidentified strains, (ii) bacterial taxa that possess representative type strains or (iii) bacterial taxa that contain few genome sequences.

Association between a taxonomic group and its distribution across a range of habitats is useful for inferring the role of this taxa on its host or environment. For instance, community profiling approaches based on molecular marker such as hypervariable regions of 16S rRNA gene have been helpful for highlighting correlations between host fitness and microbiome composition. Higher taxonomic resolution of microbiome composition could be achieved with metagenomics through *k*-mer based classification of reads. In this study we demonstrate that employing a database with a classification of strains reflecting their genomic relatedness greatly improve taxonomic assignments of reads. Therefore, investigating the relationships between bacterial genome sequences not only benefits bacterial taxonomy but also 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 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 assistances with the  
241 SIMKA software and Jason Shiller for manuscript assessment and for editing the English.

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. Meier-Kolthoff JP, Auch AF, Klenk H-P, Göker M. Genome sequence-based species delimitation with confidence intervals and improved distance functions. *BMC Bioinformatics*. 2013;14:60.
5. 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.
6. Richter M, Rosselló-Móra R. Shifting the genomic gold standard for the prokaryotic species definition. *Proc Natl Acad Sci USA*. 2009;106:19126–31.
7. 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.
8. 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.
9. Varghese NJ, Mukherjee S, Ivanova N, Konstantinidis KT, Mavrommatis K, Kyrpides NC, et al. Microbial species delineation using whole genome sequences. *Nucleic Acids Res*. 2015;43:6761–71.
10. 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
11. 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.
12. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biology*. 2014;15:R46.
13. Kim D, Song L, Breitwieser FP, Salzberg SL. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome Res*. 2016; 26: 1721-1729
14. 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.
15. 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.
16. 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.
17. 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.

18. 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.
19. Gomila M, Peña A, Mulet M, Lalucat J, García-Valdés E. Phylogenomics and systematics in *Pseudomonas*. *Front Microbiol.* 2015;6:214.
20. 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.
21. 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.
22. 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
23. 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.
24. Patrick A Grimont. Use of DNA reassociation in bacterial classification. *Canadian J Microbiol.* 1988; 34:541-6.
25. Rosselló-Móra R, Amann R. Past and future species definitions for Bacteria and Archaea. *Syst Appl Microbiol.* 2015;38:209–16.
26. Bromham L. Why do species vary in their rate of molecular evolution? *Biol Lett.* 2009;5:401–4.
27. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski 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.
28. Garrity GM. A New Genomics-Driven Taxonomy of Bacteria and Archaea: Are We There Yet? *J Clin Microbiol.* 2016;54:1956–63.
29. Rossello-Mora R, Sutcliffe IC. Reflections on the introduction of the Digital Protologue Database — A partial success? *System Appl Microbiol.* 2019; 42:1-2.
30. 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.
31. 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 nucleotide 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 [17,22]. 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).

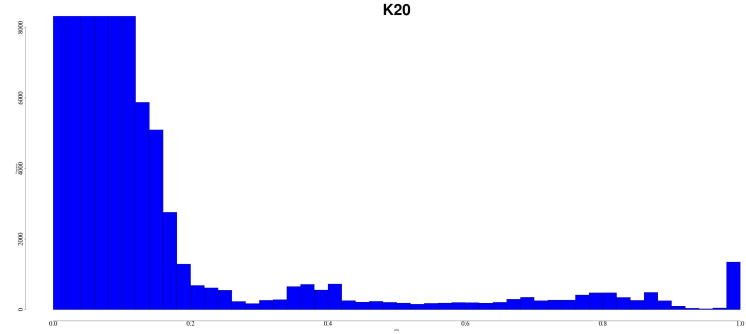
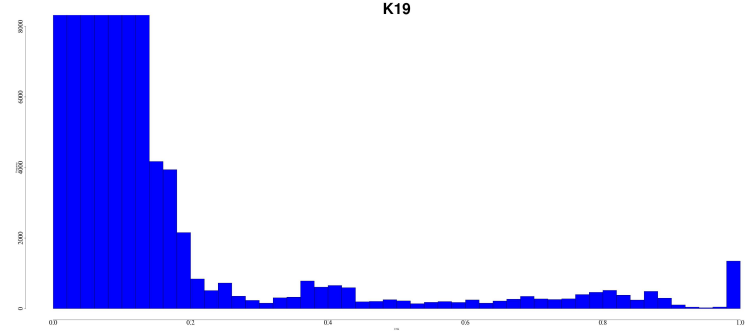
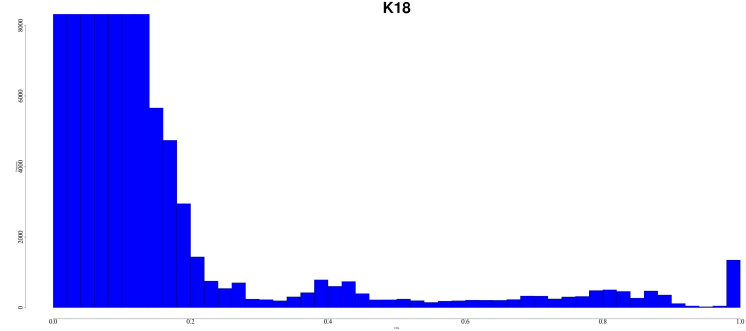
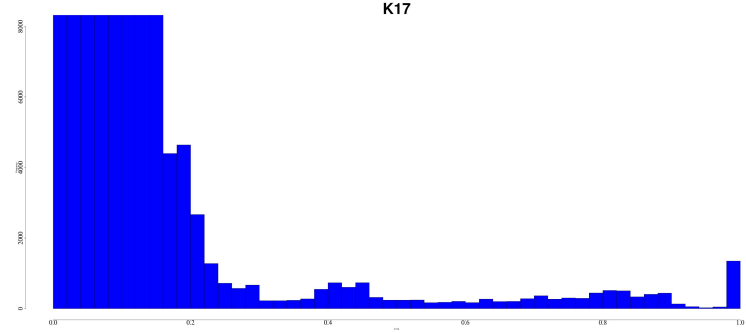
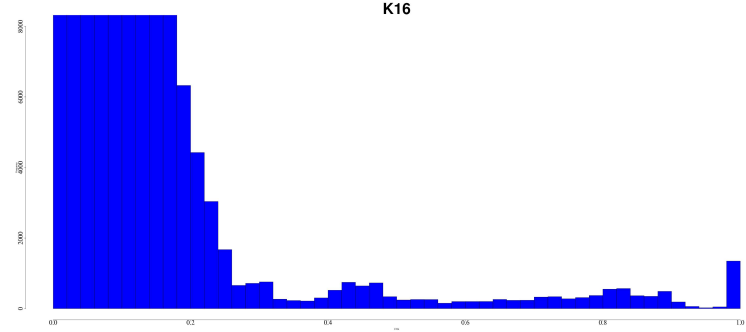
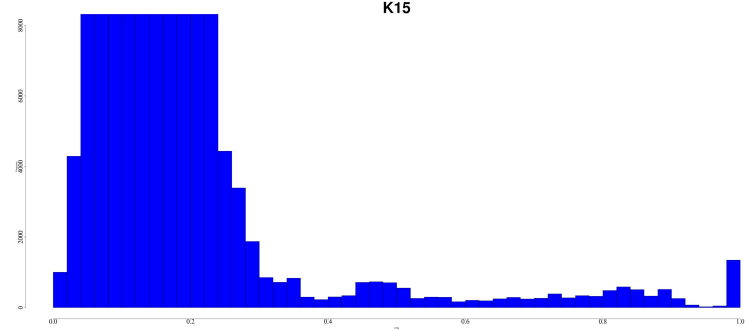
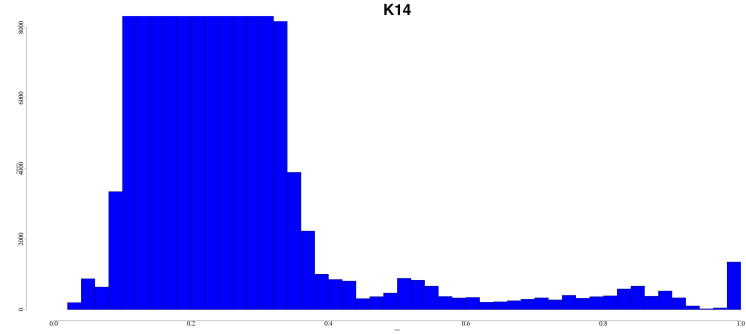
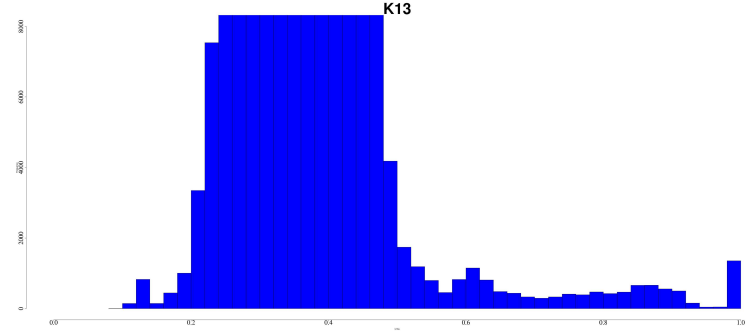
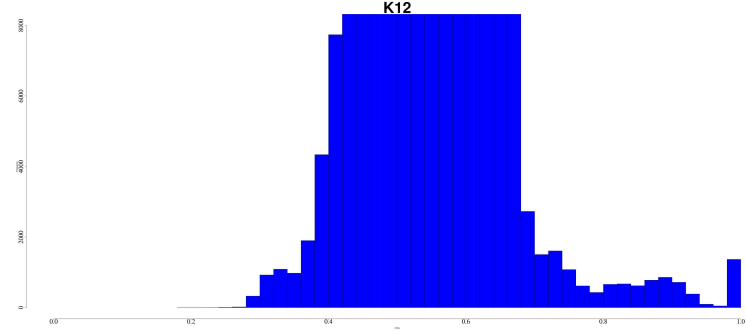
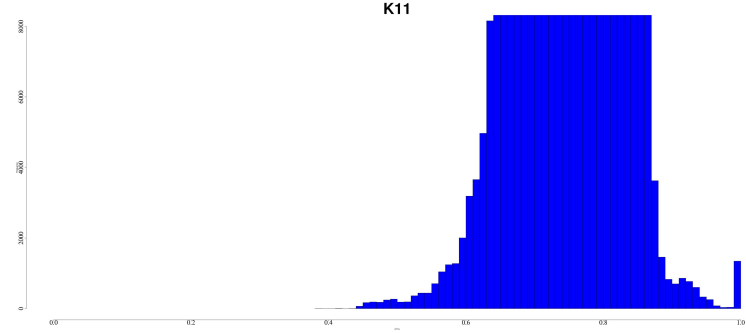
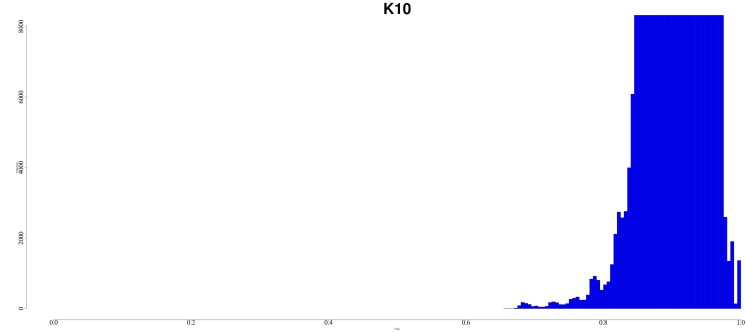
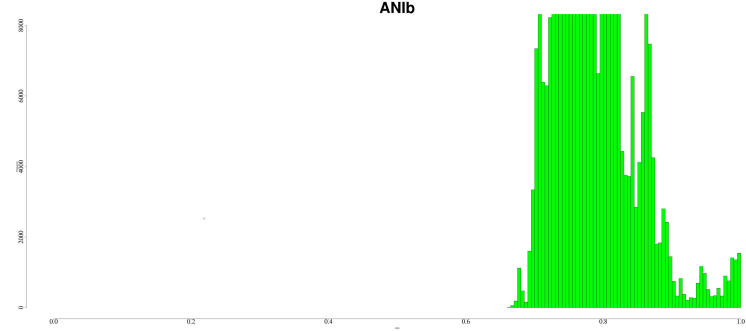
260 **Figure 4: Percentage of classified reads.** Classification of metagenome read sets derived  
261 from bean seeds, germinating seeds and seedlings with Clark [14]. Three distinct databases  
262 were employed for read classification: the original Clark database (red), Clark database  
263 supplemented 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).

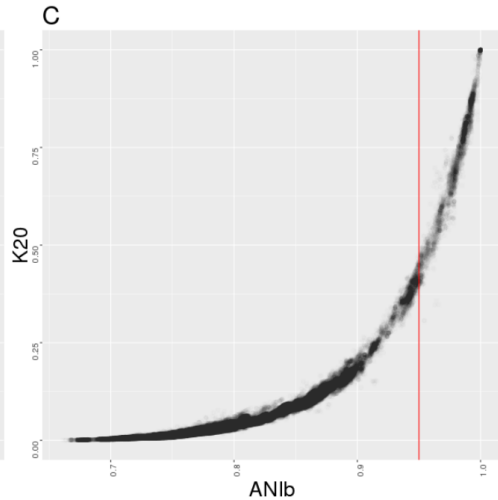
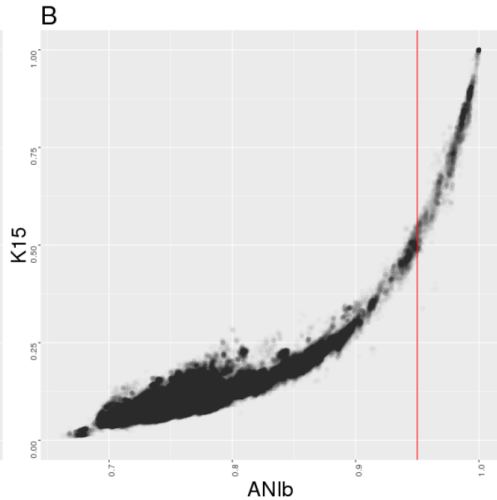
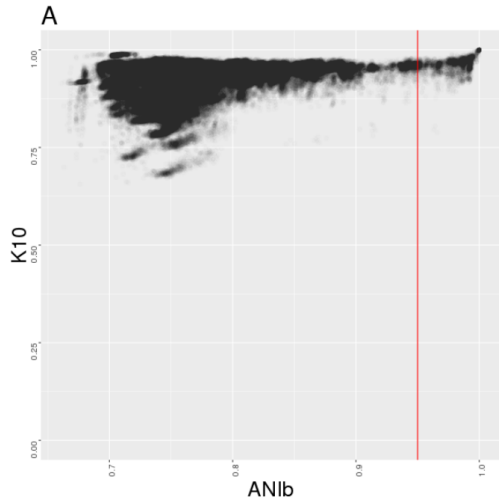
266 **TableS1.csv : *Pseudomonas* cliques.** Description of the 350 cliques obtained after  
267 clustering at 50% of shared 15-mers. For each clique, the *Pseudomonas* group [22] and  
268 subgroup [17,22] are displayed.

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

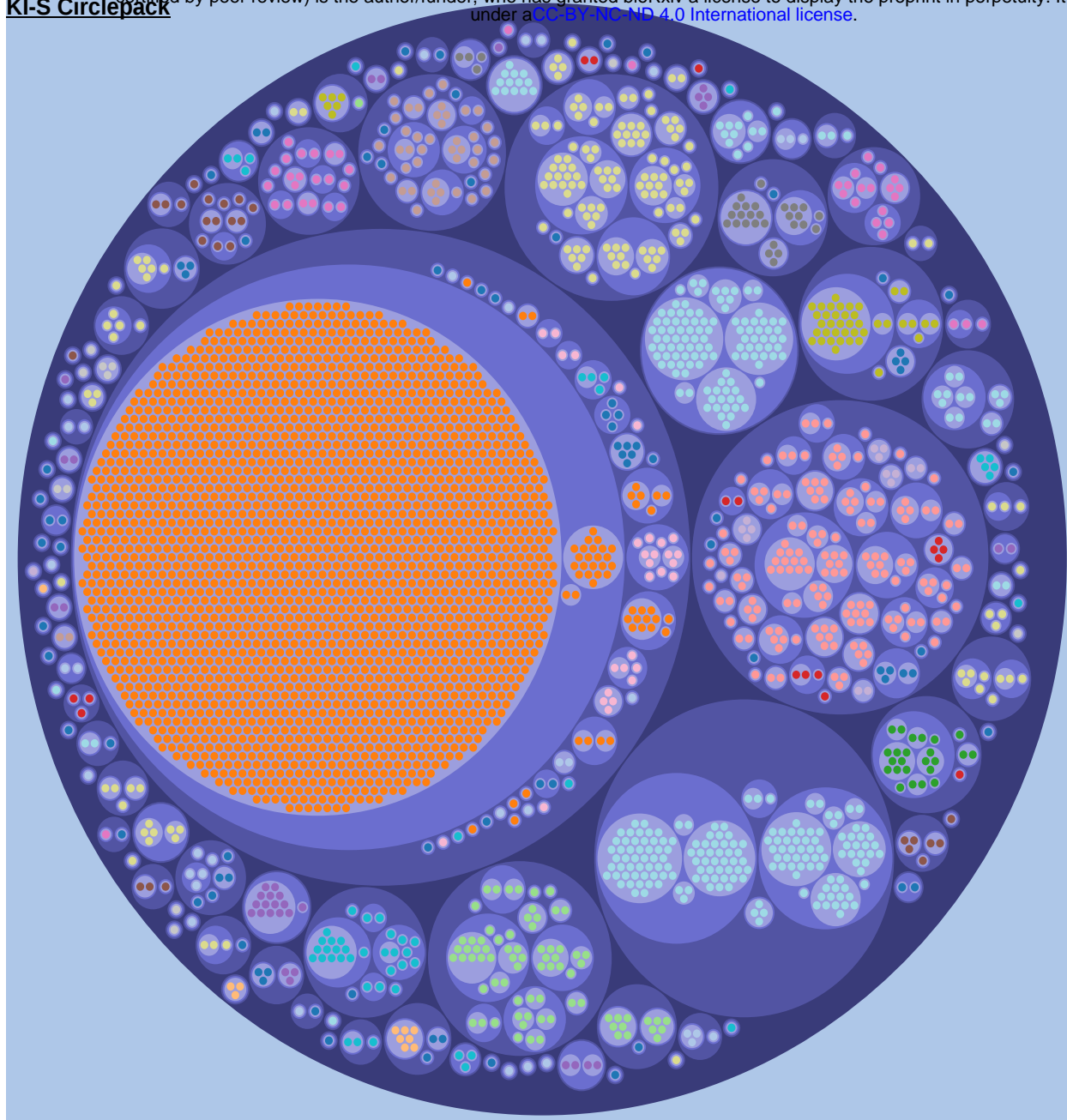
271percentage of shared 15-mers. Each dot represents a genome sequence, which is colored  
272according to its group of species [17,22]. These genome sequences have been grouped at  
273three distinct thresholds for assessing intraspecific (0.75), species-specific (0.5) and  
274interspecies relationships (0.25).







# KI-S Circlepack



## 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

