

## Supplementary appendix

### A dual barcoding approach to bacterial strain nomenclature: Genomic taxonomy of *Klebsiella pneumoniae* strains

Melanie Hennart <sup>a,b</sup>, Julien Guglielmini <sup>c</sup>, Sébastien Bridel <sup>a</sup>, Martin M.C. Maiden <sup>e</sup>, Keith A. Jolley <sup>e</sup>,  
Alexis Criscuolo <sup>c</sup> and Sylvain Brisse <sup>a</sup>

<sup>a</sup> Institut Pasteur, Université Paris Cité, Biodiversity and Epidemiology of Bacterial Pathogens, Paris,  
France

<sup>b</sup> Sorbonne Université, Collège doctoral, F-75005 Paris, France

<sup>c</sup> Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, F-75015 Paris, France

<sup>d</sup> Department of Zoology, University of Oxford, 11a Mansfield Road, Oxford, OX1 3SZ, United  
Kingdom

#### **Contents:**

*Detection of inter-phylogroup hybrids*

*Outbreak datasets to estimate variation of shallow-level classification groups*

*Minimum Spanning tree-based clustering of cgMLST profiles*

*Nomenclature inheritance algorithm*

*Impact of strains input order on LIN codes, and use of Prim's algorithm*

*List of supplementary figure legends*

*List of supplementary tables*

*Supplementary references*

### 23 *Detection of inter-phylogroup hybrids*

24 To define hybrid genomes (*i.e.*, arising from multiple ancestral populations), the 45 reference  
25 genomes representative of the seven phylogroups Kp1-Kp7 (**Table S9**) were used as models. First, for  
26 each of the 7,388 other genome sequences, the closest reference genome was determined by  
27 estimating the average nucleotide identity (ANI) using FastANI v1.1 <sup>1</sup>. Every genome *x* with ANI  
28 percentage > 99% against its closest reference genome *y* was then classified without ambiguity into  
29 the same phylogroup as the one of *y*. Second, for each genome classified into a phylogroup (Kp1-  
30 Kp7), all its cgMLST alleles were labelled with this phylogroup. Third, for each of the 629 scgMLSTv2  
31 loci, every distinct allele associated to more than one phylogroup labels was unlabeled, given that  
32 such an allele cannot be considered as a reliable representative of a unique phylogroup (*e.g.*, it was  
33 too conserved, or involved in horizontal transfer between phylogroups). Fourth, for each locus, every  
34 unlabeled sequence identical to one of the remaining labelled alleles (*i.e.*, sequence belonging to a  
35 genome that was not assigned to a phylogroup during step 1) was labelled accordingly. Such a  
36 procedure enabled the characterization of a large set of alleles that are each representative of one of  
37 the seven phylogroups Kp1-Kp7.

38 As a result, almost all cgMLST profiles were mostly made up by alleles belonging to only one  
39 phylogroup label (see **Figure S15**). However, notable exceptions were observed, with some cgMLST  
40 profiles being composed of alleles belonging to two phylogroup labels (see *e.g.*, **Figure S15**). To  
41 define putative hybrid profiles, a phylogroup homogeneity index was determined for each profile,  
42 defined as the proportion of loci labelled with the predominant phylogroup (normalized by the  
43 number of non-missing alleles called in the profile). As expected, for each phylogroup Kp1-Kp7, most  
44 profiles are associated with high phylogroup homogeneity indices (see distributions in **Figure S12**).  
45 However, a total of 138 cgMLST profiles (1.9%; mainly within Kp1, Kp2 and Kp4) were associated with  
46 atypically smaller homogeneity indices, and mapping of allele phylogroup labels along the  
47 chromosome showed that many of these 138 cgMLST profiles appeared to result from large-scale  
48 inter-phylogroup recombination, while a few were made up by many unlabeled alleles (**Figure S1**).

49 Large recombination events were detected in 1.9% (138/7198) genomes and mainly involved  
50 phylogroups Kp1, Kp2 and Kp4; we found 42 Kp4 genomes resulting from large-scale recombination  
51 of Kp1 out of the 50 hybrid Kp4 genomes (**Figure S1**), while others form a multitude of small-scale  
52 recombination events. Next, 17 Kp1 genomes were observed with a large-scale recombination (12  
53 with a Kp2 insertion, 4 with a Kp4 insertion and 1 with a Kp3 insertion), as well as 3 Kp2 genomes  
54 resulting from a large Kp1 insertion. In addition, 42 profiles of phylogroup Kp3 resulted from  
55 horizontal gene transfer (but not large-scale recombination events) from non-KpSC donors (**Table 1**;

56 **Figure S2; Table S8**). The recombination breakpoints were non-randomly distributed along the  
57 genomes: most (109/126, 86.5%) were localized in the second half of the genome (3 Mb – 5.2 Mb of  
58 NTUH-K2044 genome coordinates), whereas in the first part (0 – 3 Mb) accounted for only 15  
59 breakpoints.

60 These 138 genomes presenting hybrid profiles (or with multiple alleles of undefined origins) were  
61 therefore discarded during our initial population structure analyses and classification steps, which  
62 were based on the remaining 7,060 genomes that likely arose from vertical evolution.

### 63 *Outbreak datasets to estimate variation of shallow-level classification groups*

64 We searched for previously published genomic epidemiology studies. These genomic investigations  
65 of outbreaks (or clusters of cases) together comprised 9 sets of isolates defined as related based on  
66 epidemiologically and genomic evidence (**Table S7**). Distribution of the cgMLST pairwise distances  
67 among isolates within each outbreak cluster was investigated (**Tables S6, S7**).

68

### 69 *Minimum Spanning tree-based clustering of cgMLST profiles: building and assessment*

70 A pairwise dissimilarity between two cgMLST profiles can be defined by the proportion of loci with  
71 two distinct alleles among the loci where alleles are defined in both profiles. A pairwise dissimilarity  
72 matrix can be computed from  $n$  cgMLST profiles, and can be used to build a minimum spanning tree  
73 (MStree; *e.g.*, Kruskal, 1956; Prim, 1957a; Dijkstra, 1959), allowing to infer a clustering of the cgMLST  
74 profiles, defined by the  $k$  different connected components obtained by removing from the MStree all  
75 edges of length larger than a specified threshold  $t$ . Such an MStree-based clustering is closely related  
76 to the single-linkage classification of the  $n$  cgMLST profiles (*e.g.*, Gower and Ross, 1969; Johnson,  
77 1967).

78 In order to determine optimal thresholds  $t$ , several criteria can be used. Among these criteria, the  
79 average silhouette coefficient  $S_t$  assesses the ability of an MStree-based clustering to consistently  
80 represents in  $k$  class(es) the 'natural' grouping of the cgMLST profiles<sup>7,8</sup>. When  $S_t$  is close to 1, the  
81 clustering can be considered as accurate. A confidence interval for  $S_t$  can be also obtained by  
82 considering the distribution of the average silhouette coefficients of different clustering computed  
83 from the distance matrix with 'noised' entries.

84 Further, in order to assess whether an MStree-based clustering  $C_t$  (using threshold  $t$ ) is robust to any  
85 subsampling biases, a simple approach is to build another MStree-based clustering  $C_t'$  from a  
86 subsample of cgMLST profiles, and to measure the agreement between the cgMLST profile partitions

87 induced by  $C_t'$  and  $C_t$ . When the level of agreement remains high for different subsampling rates, the  
88 corresponding threshold  $t$  can be considered as being leading to stable clustering. Among different  
89 agreement metrics between partitions, the second adjusted Wallace coefficient  $w^{9,10}$  estimates the  
90 probability of observing a pair of profiles in the same class in  $C_t$  when they are clustered in the same  
91 class in  $C_t'$ . In order to derive a single coefficient  $W_t$  from a range of different subsampling rates  $r$  (=   
92 10% to 90%), different coefficients  $w$  were estimated and averaged for each rate  $r$ ; the area under  
93 the resulting curve (*i.e.*, rates  $r$  on X-axis; coefficient  $w$  on Y-axis) was computed and normalized  
94 (using its maximum expected value). Such a normalized area  $W_t$  is close to 1 when the different  
95 adjusted Wallace coefficients  $w$  (*i.e.*, derived from varying subsampling rates) are all close to their  
96 maximum value, therefore showing that the corresponding MStree-based clustering (based on the  
97 threshold  $t$ ) is robust to any subsampling biases. A confidence interval for  $W_t$  can be also obtained  
98 using the same approach as for  $S_t$  (see above).

99 The MStree-based clustering of cgMLST profiles, as well as the two consistency and stability indices  $S_t$   
100 and  $W_t$ , respectively, were implemented in the MSTclust tool  
101 (<https://gitlab.pasteur.fr/GIPhy/MSTclust>). For more details, see  
102 <https://gitlab.pasteur.fr/GIPhy/MSTclust/-/blob/0.21b/Technical.Notes.pdf>.

### 103 *Nomenclature inheritance algorithm*

104 In order to attribute to each clonal group (CG), an identifier that would maximally reflect the widely  
105 adopted 7-gene ST identifier of the corresponding isolates, we developed a set of naming rules that  
106 prioritize the most abundant ST observed among isolates of each CG, as well as some supplementary  
107 rules in case of ties. This algorithm is summarized below, and its implementation as a Python script is  
108 provided at <https://gitlab.pasteur.fr/BEBP/inheritance-algorithm>. **Figure S15** illustrates the process  
109 for an example.

110 Here the process for the CG level is described, but the algorithm was also applied to the SL level.  
111 Briefly, the data (*e.g.*, a list of CG-ST pairs) can be formalized as a bipartite graph, in which each CG  
112 and ST are nodes, and each non-empty CG-ST intersection is an edge. The weight of each edge is  
113 equal to the number of isolates sharing the corresponding CG and ST identifiers. Based on this  
114 representation, the algorithm will consist of following all edges in the input graph, in the order of  
115 decreasing weight. The approach prioritizes the most frequent ST/CG pairs of isolates, *i.e.*, those that  
116 are predominant in the dataset and thus naturally transfers to the CG nomenclature, the identifiers  
117 of the highest frequency STs. Rules were implemented to treat the cases of equality of  
118 representation of two or more STs connected to the same CG. Once all edges were removed from  
119 the graph, it may be that some CGs were not named, for example, because the identifier of their  
120 unique corresponding ST was already attributed to another CG. For these orphan CGs, iteratively, the  
121 attributed identifier corresponds to the maximal CG identifier already attributed, plus one (**Figure**  
122 **S15**).

### 123 Definitions and notations

124 Let  $G = (U, V, E)$  be a weighted bipartite graph where:

- 125 •  $U$  is a set of clonal groups (CG) inferred from a cgMLST scheme
- 126 •  $V$  is the set of sequence types (ST) induced by a MLST scheme
- 127 •  $E$  is the set of edges  $\{u, v\}$  with  $u \in U$  and  $v \in V$
- 128 •  $w(\{u, v\})$  is the weight of the edge  $\{u, v\}$ , *i.e.*, the number of isolates inside  $u \cap v$

129 Let  $L(v)$  be the label associated to node  $v$  (*i.e.*, the ST identifier), and  $L(u)$  the one to determine for  $u$   
130 (*i.e.*, the CG identifier).

131 Let  $d_G(v)$  be the degree of a node  $v$  inside the graph  $G$ , *i.e.*, the number of edges incident to node  $v$ .

132 Let  $s(u) := \sum_{v \in V} w(\{u, v\})$  be the size of  $u$ , *i.e.*, the number of strains belonging to the CG  $u$ .

133 Let  $\Gamma(G)$  be the edge-induced subgraph of a graph  $G$  defined by the edge(s) of maximal weight in  $G$ .

134 Let  $\Delta_G(u)$  be the set of nodes  $v$  that are joined to  $u$  inside the graph  $G$  and of minimum degree, *i.e.*,

$$135 \Delta_G(u) := \{v' \in V : \{u, v'\} \in E, d_G(v') = \min_{\{u, v\} \in E} d_G(v)\}.$$

### 136 Algorithm

```

137
138 (a) ◦  $\lambda := \max_{v \in V} L(v)$ 
139 (b) ◦ while  $E \neq \emptyset$ 
140     do
141     ◦ for each connected component  $G = (U', V', E')$  of  $\Gamma(G)$ 
142         do
143         ◦ if  $U' = \{\mu\}$ 
144             then
145         ◦  $v := \operatorname{argmin}_{v' \in V'} L(v')$ 
146         else
147         ◦  $U'' := \operatorname{argmin}_{u' \in U'} s(u')$ 
148         ◦ if  $U'' \neq \{\mu\}$ 
149             then
150         ◦  $\mu := \operatorname{argmin}_{u'' \in U''} L(u'')$ 
151         ◦  $v := \operatorname{argmin}_{v' \in \Delta_G(\mu)} L(v')$ 
152         ◦  $L(\mu) := L(v)$ 
153         ◦ removing  $\mu$  and  $v$  from  $G$ , as well as nodes  $v$  such that  $w(\{\mu, v\}) = w(\{\mu, v\})$ 
154 (l) ◦ for each  $\mu \in U$ 
155     do
156     (l) ◦  $\lambda := \lambda + 1$ 
157     (m) ◦  $L(\mu) := \lambda$ 

```

158 *Impact of strains input order on LIN codes, and use of Prim's algorithm*

159 By design, the input order of genomes into the cgLINcode nomenclature system influences their  
160 attributed code, as is the case for the original LIN code system <sup>12</sup>. We evaluated this impact by  
161 quantifying the variation in the number of partitions at a given threshold, as defined by the number  
162 of distinct prefixes: for each threshold varying from 1% to 99%, a LIN encoding was defined using this  
163 threshold, and the 7,060 high-quality, non-hybrid cgMLST profiles were encoded 500 times with  
164 random input orders. This experiment made it possible to determine (i) the threshold values  
165 associated with a stronger variability in the final number of values; and (ii) the magnitude of this  
166 variability. In the example illustrated in **Figure S10**, we observed that the number of distinct prefixes  
167 was affected by the order of encoding, especially in the 450 - 530 mismatches range. Note that this  
168 experiment can help to select position thresholds, for example, by favoring those that minimize the  
169 variance of the number of partitions (*i.e.*, are less affected by input order).

170 We next sought to minimize this problem by defining an optimal input order. The one that answered  
171 our expectations is the input order guided by a Prim's algorithm <sup>3</sup>. More precisely, the number of  
172 categories in a given LIN encoding bin is minimal (*i.e.*, identical to the number of groups created by a  
173 single-linkage clustering using the threshold associated to the bin) when the profiles are encoded  
174 following the order induced by the traversal of an MStree. Indeed, when following such an order,  
175 when a new profile is considered for encoding, then its closest profile is already encoded (by  
176 definition of a tree traversal). The optimal order we suggest is therefore verified by noting that the  
177 Prim's (1957) algorithm to infer a MStree induces such an MStree traversal. A comparison between  
178 the MLSL approach and the cgLIN codes was performed (cgLIN codes in optimal *versus* arbitrary  
179 order). We found that the partitioning created by the MLSL approach and that created by the cgLIN  
180 codes according to the optimal order, were identical (**Table S1**).

181 We generated 500 random input orders and then generated cgLIN codes with two bins (the first  
182 varies from 1 to 100% with a step of one allelic difference *i.e.* 100/629, and the second fixed at  
183 100%). Then we counted the number of prefixes, up to the first bin, that were created. We observed  
184 that the 10 identifier bins differ in their sensitivity to input order (**Figure S10**). The most affected bins  
185 correspond to regions of the pairwise distance distribution with high density; in particular around  
186 485 mismatches, before the mode that corresponds to inter-sublineage distances.

187 The algorithm below was used to define the input order, without even having to construct an  
188 MStree. Indeed, thanks to a simple traversal of the matrix of dissimilarities between profiles, the  
189 algorithm makes it possible to quickly determine the optimal order for LIN encoding.

190 Algorithm

191

192 (a) ◦ Create a set "mstSet" that keeps track of vertices already included in MST

193 (b) ◦ Assign a key value to all vertices in the input graph. Initialize all key  
194 values as  $\infty$ . Assign key value as 0 for the first vertex so that it is  
195 picked first.

196 (c) ◦ **while** "mstSet" doesn't include all vertices

197 **do**

198 (d) ◦ Pick a vertex  $u$  which is not there in "mstSet" and has minimum key value.

199 (e) ◦ Include  $u$  to "mstSet".

200 (f) ◦ Update key value of all adjacent vertices of  $u$ . To update the key values,  
201 iterate through all adjacent vertices. For every adjacent vertex  $v$ , if  
202 weight of edge  $u - v$  is less than the previous key value of  $v$ , update  
203 the key value as weight of  $u - v$ .

204

205 Note that using key values enables to pick the minimum weight edge from cut. The key values are

206 used only for vertices which are not yet included in MStree; the key value for these vertices indicate

207 the minimum weight edges connecting them to the set of vertices included in MStree. The time

208 complexity required by Prim's (1957) algorithm is  $O(E \log V)$  where  $E$  is the number of edges and  $V$  is

209 the number of vertices.

210



211 *List of supplementary figures*

212

213 Figure S1. cgMLST profile painting illustrates large recombinations

214 Figure S2. Genomes inclusion flowchart

215 Figure S3. Characteristics of the 629 loci of the cgMLST scheme

216 Figure S4. The distribution of pairwise distances based on Average Nucleotide Identity (ANI) and  
217 cgMLST. ANI values were calculated using the entire genomic sequence (not just the cgMLST allele  
218 sequences) using FastANI v1.1.

219 Figure S5. Details of cgMLST pairwise distances distributions

220 Figure S6. Correspondence of ST, sublineage and clonal group classifications for 9 major K.  
221 pneumoniae sublineages

222 Figure S7. The distribution of pairwise distances based on Average Nucleotide Identity (ANI) and  
223 cgMLST, with hybrid genomes

224 Figure S8. Impact of inter-phylogroup hybrid genomes on cgMLST classification groups

225 Figure S9. Virulence and resistance scores in major sublineages

226 Figure S10. Impact of input order on the number of partitions in the resulting LIN codes

227 Figure S11. Relationships between ST, cgMLST and cgLIN codes, and their behavior upon novel  
228 genomes inclusion

229 Figure S12. Distribution of the phylogroup homogeneity index

230 Figure S13. Principle of cgLIN code implementation

231 Figure S14. cgLIN codes implementation for nearly-identical cgMLST profiles

232 Figure S15. Step-by-step illustration of the taxonomic inheritance algorithm

233 *List of supplementary tables*

234

235 Table S1. Dataset of 7,433 genomes

236 Table S2. Hybrid genomes breakdown by phylogroup

237 Table S3. Characteristics of the 629 loci of the cgMLST scheme

238 Table S4. Correspondence between SLs, CGs and STs

239 Table S5. Characteristics of the clonal groups

240 Table S6. Outbreak dataset

241 Table S7. Within-outbreak variation

242 Table S8. Correspondence between ANI and cgMLST distance thresholds

243 Table S9. Reference genomes

244

245 *Supplementary references*

- 246 1. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput ANI  
247 analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* **9**, 5114  
248 (2018).
- 249 2. Kruskal, J. B. On the shortest spanning subtree of a graph and the traveling salesman problem.  
250 *Proc. Amer. Math. Soc.* **7**, 48–48 (1956).
- 251 3. Prim, R. C. Shortest Connection Networks And Some Generalizations. *Bell System Technical*  
252 *Journal* **36**, 1389–1401 (1957).
- 253 4. Dijkstra, E. W. A note on two problems in connexion with graphs. *Numer. Math.* **1**, 269–271  
254 (1959).
- 255 5. Gower, J. C. & Ross, G. J. S. Minimum Spanning Trees and Single Linkage Cluster Analysis. *Applied*  
256 *Statistics* **18**, 54 (1969).
- 257 6. Johnson, S. C. Hierarchical clustering schemes. *Psychometrika* **32**, 241–254 (1967).
- 258 7. Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster  
259 analysis. *Journal of Computational and Applied Mathematics* **20**, 53–65 (1987).
- 260 8. Lengyel, A. & Botta-Dukát, Z. Silhouette width using generalized mean—A flexible method for  
261 assessing clustering efficiency. *Ecol Evol* **9**, 13231–13243 (2019).
- 262 9. Wallace, D. L. A Method for Comparing Two Hierarchical Clusterings: Comment. *Journal of the*  
263 *American Statistical Association* **78**, 569–576 (1983).
- 264 10. Severiano, A., Pinto, F. R., Ramirez, M. & Carriço, J. A. Adjusted Wallace Coefficient as a Measure  
265 of Congruence between Typing Methods. *J. Clin. Microbiol.* **49**, 3997–4000 (2011).
- 266 11. Holt, K. E. *et al.* Genomic analysis of diversity, population structure, virulence, and antimicrobial  
267 resistance in *Klebsiella pneumoniae*, an urgent threat to public health. *Proc Natl Acad Sci U S A*  
268 **112**, E3574-81 (2015).

269 12. Marakeby, H. *et al.* A system to automatically classify and name any individual genome-  
270 sequenced organism independently of current biological classification and nomenclature. *PLoS*  
271 *One* **9**, e89142 (2014).