# Bacterial strain nomenclature in the genomic era:

# Life Identification Numbers using a gene-by-gene approach

Federica Palma[1,2,#], Melanie Hennart[1,#], Keith A. Jolley[3], Chiara Crestani[1], Kelly L. Wyres[4,5], Sebastien Bridel[1], Corin A. Yeats[6], Bryan Brancotte[7], Brice Raffestin[8], Sophia David[6], Margaret M. C. Lam[4], Radosław Izdebski[9], Virginie Passet[1], Carla Rodrigues[1], Martin Rethoret-Pasty[1], Martin C. J. Maiden[3], David M. Aanensen[6], Kathryn E. Holt[4,10], Alexis Criscuolo[2], Sylvain Brisse[1,2,*]

[#] Equal contribution

**Affiliations**

[1] Institut Pasteur, Université Paris Cité, Biodiversity and Epidemiology of Bacterial Pathogens, Paris, France

[2] Institut Pasteur, Université Paris Cité, Biological Resource Center of Institut Pasteur, Paris, France

[3] Department of Biology, University of Oxford, Oxford, United Kingdom

[4] Department of Infectious Diseases, School of Translational Medicine, Monash University, Melbourne, Australia

[5] Centre to Impact AMR, Monash University, Clayton, 3800, Australia

[6] Centre for Genomic Pathogen Surveillance, Pandemic Sciences Institute, University of Oxford, Oxford, United Kingdom

[7] Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, Paris, France

[8] Institut Pasteur, Université Paris Cité, HPC Core Facility, Paris, France

[9] Department of Molecular Microbiology, National Medicines Institute, Warsaw, Poland

[10] Department of Infection Biology, London School of Hygiene and Tropical Medicine, London, United Kingdom


[*] **Correspondence:** Prof Sylvain Brisse, Institut Pasteur, Biodiversity and Epidemiology of Bacterial Pathogens, 25-28 rue du Docteur Roux, F-75724, Paris, France; Phone: +33 1 45 68 83 34; E-mail: sylvain.brisse@pasteur.fr

## Abstract

Unified strain taxonomies are crucial for fostering international communication in microbiological research and for the epidemiological surveillance of bacterial pathogens. While multilocus sequence typing (MLST) has served as a foundation of strain taxonomy for two decades, whole genome sequencing enables more precise classifications and significantly improves discriminatory resolution. The core genome-wide extension of MLST (known as cgMLST) thus holds great promise for strain genotyping and classification, but its implementation faces challenges that include missing data, potential instability of cluster-based nomenclatures, and the necessity to ensure backwards compatibility with MLST identifiers. Life Identification Number (LIN) codes offer a solution by providing multi-level classification groups that are inherently stable. Here, we present, consolidate, and extend the cgMLST-based LIN code approach. We first develop a nicknaming system for LIN code prefixes, which enables flexible human-readable strain nomenclatures. Using *Klebsiella pneumoniae* (Kp) as an example, LIN code nicknames were attributed by inheritance from MLST identifiers, thus perpetuating the legacy of MLST nomenclatures in the genomic era. We show that while 7-gene MLST sometimes conflates unrelated sublineages into the same ST, cgMLST-based LIN codes are highly concordant with phylogenetic relationships. We implement this novel LIN code-based nomenclature in the BIGSdb platform, and illustrate, with Pathogenwatch, how it can also be used in other genomic epidemiology platforms. Finally, we demonstrate the value of LIN codes for tracking the strain diversity within high-risk internationally disseminated clonal groups of Kp and protracted outbreaks. Given its stability, precision, and flexibility, we recommend the adoption of the cgMLST-based LIN code taxonomic approach for Kp and suggest that this approach is widely applicable to other bacterial pathogens.

## Introduction

Taxonomies of bacterial strains responsible for infectious diseases are essential resources to ensure effective communication in population biology, epidemiological surveillance, and public health response to outbreaks. As illustrated by the SARS-CoV-2 variant nomenclature system, simple nicknames (*e.g.,* Alpha, Delta, Omicron) for pathogen variants can greatly improve communication between different public health sectors (Konings et al., 2021; Rambaut et al., 2020).

Currently, there are neither classification nor nomenclature standards to define sublineages, variants, types or clones (hereafter, collectively called "strains") within bacterial species ("International Code of Nomenclature of Prokaryotes," 2019). Ad-hoc phenotypic (*e.g.,* serotypes) and genotypic (*e.g.,* sequence types) approaches have long been used to define strains of particular species, but the advent of universally applicable whole genome sequencing (WGS) has the potential to refine and generalize strain taxonomy by providing the maximal discrimination needed for epidemiological surveillance, and a harmonized general approach across pathogen phyla (Maiden et al., 2013; Nadon et al., 2017; Struelens and Brisse, 2013). However, few attempts have been made to devise genomic taxonomies and evaluate their general applicability. With WGS implemented worldwide and in all sectors of microbiology (medical, veterinary, food, environmental), a precise and universal approach for describing strains of bacterial species becomes a key need to translate WGS data into relevant information that would support epidemiological surveillance, outbreak investigations, cross-niche or between host transmission detection, and public health actions that need international and cross-sectoral coordination.

Among the broad range of methods developed for bacterial strain typing and group naming (Struelens et al., 1998; van Belkum et al., 2007), multi-locus sequence typing (MLST), based on the analysis of a few (typically seven) conserved loci, was established over the last two decades as the method of choice for strain taxonomy of most bacterial species (Aanensen and Spratt, 2005; Maiden, 2006; Maiden et al., 1998). This gene-by-gene approach was logically extended to the genome scale, with core genome MLST (cgMLST) schemes encompassing thousands of loci (Bialek-Davenet et al., 2014; Maiden et al., 2013). Whether using the classical or the core genome MLST schemes, the "sequence type" (ST) nomenclature system is highly reproducible, portable, and easy to interpret (Feil, 2004). To recognize deeper phylogenetic associations, cgMLST allele profiles can be grouped at any level of similarity by single-linkage clustering or static aggregation to predefined groups or founder genotypes (Zhou et al., 2021).

A novel system for genome classification was proposed by Vinatzer and colleagues, using multi-position numerical codes attributed to each individual genome (Marakeby et al., 2014; Vinatzer et al., 2017). These codes, called Life Identification Numbers (LINs), were designed to encompass all domains of life in a single taxonomy, based on the Average Nucleotide Identity (ANI) metric (Goris et

86    al., 2007; Konstantinidis and Tiedje, 2005). However, the ANI-based genome similarity is imprecise
87    and non-reproducible for nearly identical strains, which are most often compared through sequences of
88    draft genomes. Leveraging the strengths of both approaches, some of us recently proposed combining
89    cgMLST and LIN codes to design taxonomies of bacterial strains within species (Hennart et al., 2022).
90    The use of cgMLST dissimilarities, rather than ANI-based similarities, provides robustness in
91    estimating small-scale genome relationships, which are efficiently summarized by cgMLST LIN codes
92    (hereafter, LIN codes for short).

93    In this article, we present further developments of the LIN code approach. We first design a
94    nicknaming approach for LIN codes, which can be used to recognize familiar groups that are
95    important in biological research or epidemiological surveillance. We further show the benefit of
96    inheriting these nicknames from MLST identifiers. We additionally describe practical
97    implementations of LIN codes in the widely used genotyping platforms BIGSdb (Argimón et al.,
98    2021; Jolley et al., 2018). We next illustrate the use and benefits of LIN code strain taxonomy using
99    the *Klebsiella pneumoniae* Species Complex (KpSC), a phenotypically and genetically diverse
100   ubiquitous pathogenic group (Wyres et al., 2020). We show that for this pathogen, classical (7-gene)
101   MLST classifications can be misleading, and that LIN codes can pinpoint these cases and mitigate
102   misclassifications. Lastly, we illustrate the benefit of LIN codes for defining and naming intraspecific
103   groups from epidemiologically important phylogenetic lineages down to outbreak strains in a stable
104   way.

### Section 1: LIN codes: definitions and practical implementation

**The principle of cgMLST-based LIN codes: an overview**

Here we explain in more detail how cgMLST-based LIN codes work, as originally proposed (Hennart et al., 2022), before describing new developments and applications of the system (see **Section 2: Novel developments and examples of applications)**. The core genome Life Identification Number classification code system combines the core genome MLST (cgMLST) approach with Life Identification Numbers (LIN) (Vinatzer et al., 2017). The LIN codes consist of multiple (for example, 10) predefined positions (or bins), each corresponding to a (range of) cgMLST profile similarity value, together representing a partition of the complete range [0%-100%]. From left to right, the positions of the code correspond to decreasing allele mismatch dissimilarity, *i.e.*, increasing similarity. The leftmost bins capture the lowest similarities reflective of deep phylogenetic divisions, whereas the rightmost bins capture the highest similarities. Each bin has a left border threshold (inclusive) that corresponds to a maximum number of pairwise allele differences between profiles and is delimited on the right by the next threshold (exclusive, as the threshold value corresponds to the left threshold of the downstream bin).

While any number of bins (up to the number of loci in the cgMLST scheme) can be chosen, in the case of the *Klebsiella pneumoniae* Species Complex (KpSC) used here as an example, 10 bins were determined to define their LIN codes (Hennart et al., 2022). The first four bins represent the deepest hierarchical levels of relatedness, corresponding to species, subspecies, sublineage and clonal group, respectively (Hennart et al., 2022). The last bins delineate six levels of high-resolution relatedness that might be useful for epidemiological surveillance. KpSC profiles are defined using a 629-loci cgMLST scheme; bins 1 to 4 have as right borders 610, 585, 190 and 43 allele mismatches, respectively, while bins 5 to 10 correspond to thresholds 10, 7, 4, 2, 1 and 0 mismatches, respectively. Thus, the first bin corresponds to the range [629-610[ of cgMLST mismatches (the '[' indicates the value 610 is excluded), whereas the last one corresponds to the range [1-0[ (note that it excludes complete identity, *i.e.,* 0 mismatch*,* 629 matches: in this case, the LIN code is simply copied from the reference, see below).

Formally, LIN codes are attributed to core genome Sequence Types (cgST) (Hennart et al., 2022). Therefore, before assigning LIN codes, cgMLST profiles must be assigned to cgSTs. Like the ST designation in classical 7-gene MLST, a cgST is defined for each unique cgMLST profile, characterized by a unique combination of alleles at all loci of the scheme. Profiles with too many missing loci can be filtered out at this stage. In practice, for the KpSC, cgMLST profiles are assigned to a cgST only when they comprise fewer than 30 missing alleles (*i.e.*, equal to or more than 600 called alleles). Profiles with 30 (4.77%) or more missing alleles (which are likely to correspond to poor quality genomes) are not considered further, and therefore not included in the KpSC LIN code

140 taxonomy. For any LIN code taxonomy, the proportion of tolerated missing data for cgST assignment

141 can be set to higher values (to increase the proportion of coded genomes) or lower values (to improve

142 the precision of LIN code classifications).

143 LIN codes are created for each distinct cgST. The formal process of LIN code assignment from

144 cgMLST data, first proposed in (Hennart et al., 2022), is presented in **Box 1** and summarized in

145 **Figure 1**. The system is initialized by creating, for an initial allelic profile, a LIN code with the integer

146 value 0 at every bin. This initial profile can be chosen randomly or based on a reference genome of the

147 species under consideration, as convenient. The next steps are the same for all subsequent individual

148 cgSTs.

149

150 **Box 1. The formal process of assigning LIN codes**

151 The LIN code of the first allelic profile is attributed 0 in every bin. Next, each new allele profile $j$ is

152 encoded from its closest already encoded profile $i$ (*i.e.*, that maximizes the allele similarity percentage

153 $s_{ij}$). After determining the pivot bin $p$, such that $s_{ij} \in [s_p, s_{p+1}[$ (*i.e.*, right threshold exclusive), the

154 encoding of the new profile $j$ is performed in three steps:

155 (i) the same prefix as code $i$ is attributed up to the bin $p-1$ (inclusive);

156 (ii) for the pivot bin $p$: the maximum value observed in this bin among the subset of codes sharing

157 the same prefix is incremented by 1;

158 (iii) 0 is attributed at each downstream bin from $p+1$ (inclusive).

159 Of note, when $s_{ij} = 100\%$, the LIN code of the new profile $j$ is given the complete LIN code of $i$

160 (including at the last bin).

161 Missing data, equal matches and input order of profiles are handled as explained in **Box 2**.

162

163 The process of assigning a LIN code to a cgMLST profile first involves matching it against all existing

164 defined LIN-encoded cgSTs to identify its closest neighbor (*i.e.*, the reference profile). If the two

165 profiles (new and reference) have no dissimilarity (*i.e.*, no allele mismatch among the loci called in

166 both profiles), the LIN code of the reference is simply assigned to the new profile. This will happen

167 when the new cgST differs from the reference only by its missing data pattern (see **Box 2**). Otherwise,

168 when the two profiles differ by at least one allele, a novel LIN code is created. For this, the pivot bin is

169 defined as the bin in which the observed allele dissimilarity falls, and the novel LIN code is created in

170 three steps (**Figure 1; Box 1**): (i) copying the LIN code prefix of the reference isolate, *i.e.* from the left

171 bin up to the pivot bin (excluded); (ii) incrementing by 1 the maximum integer value observed in the

172    pivot bin among the profile(s) sharing the same prefix used at step (i); (iii) attributing the integer value

173    0 at the bins downstream of the pivot, corresponding to initialization of the novel subdivision created

174    at the pivot bin level.



| | Bin number: | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Genome sequences ⇩ cgMLST profiles ⇩ cgST ⇩ | Max. allelic similarity*: | 19 | 44 | 439 | 586 | 619 | 622 | 625 | 627 | 628 | 629 |
| | Min. allelic difference*: | 610 | 585 | 190 | 43 | 10 | 7 | 4 | 2 | 1 | 0 |
| | Bins left thresholds: | | | | | | | | | | |
| | **Closest genome** (similarity %) | 0 | 3.02 | 6.99 | 69.79 | 93.16 | 98.41 | 98.88 | 99.36 | 99.68 | 99.84 | 100 |
| **Genome A** | Initialization | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Genome B** | A (3.50%) | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Genome C** | B (99.0%) | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| **Genome D** | B (7.00%) | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **…** | … | | | | … | | | | | | |
| **Genome X** | Y (5.00%) | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Genome Z** | X (98.90%) | 0 | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |

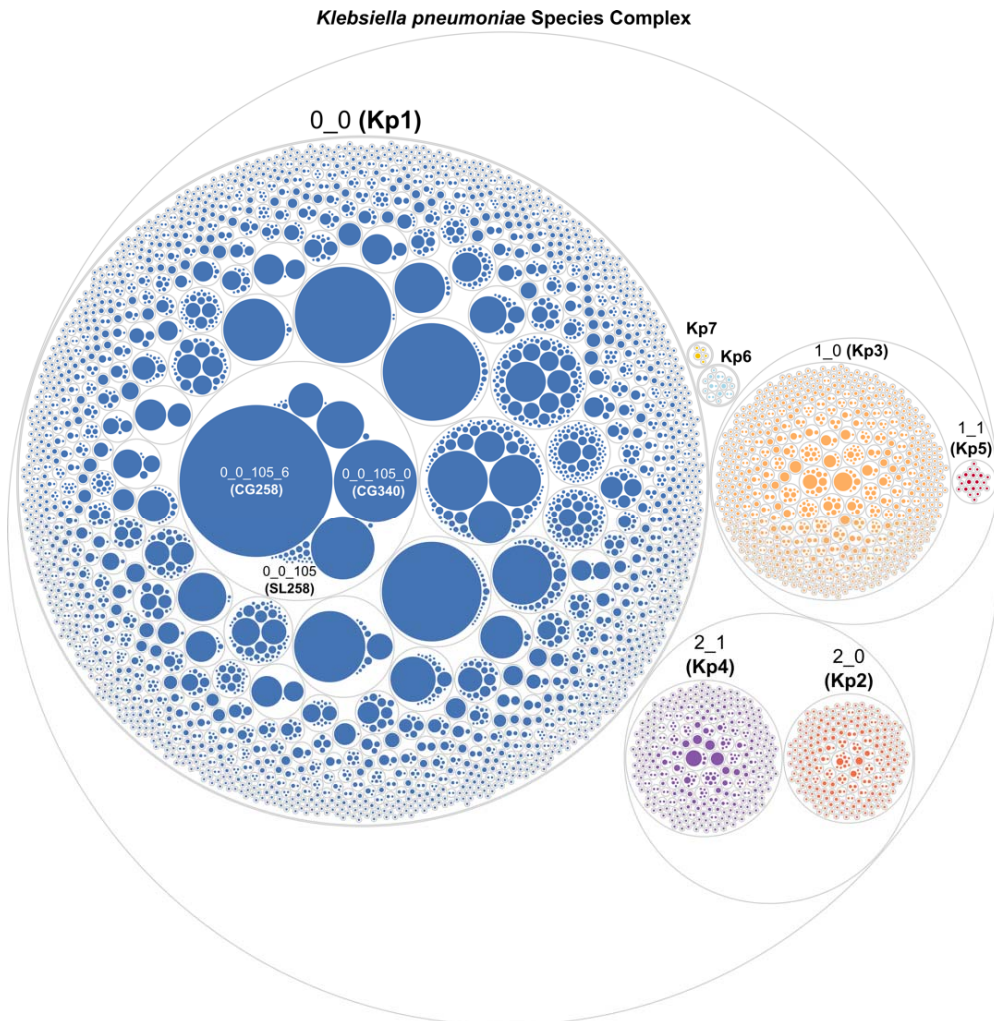Coding steps (genome Z):     (i)    (ii)    (iii)

175

176    **Figure 1**. **Overview of the process of cgMLST-based LIN code assignment.** The process starts with
177    assigning cgMLST profiles to genome sequences and classification of profiles into unique core
178    genome sequence types (cgST). After an initialization step (full-0 code for the first cgST), LIN codes
179    are created for each cgST using the similarity to its closest-related already encoded cgST (steps i, ii
180    and iii; see details in main text and Box 1). The bins and their threshold values are those chosen for the
181    KpSC. The asterisk (*) indicates that the values are for the right threshold of each bin, exclusive. Note
182    that there is no bin corresponding to complete similarity (gray column on the right), as in this case the
183    LIN codes are identical, *i.e.*, there is no need to create a novel LIN code.

184

185    A LIN code prefix can be defined as any bin subset that starts from the leftmost position of the

186    complete LIN code. The notion of prefix is important as it conveys a sense of genetic similarity among

187    profiles: the longer the common prefix of two LIN codes is, the more similar the two corresponding

188    profiles are. For a given cgST profile, its LIN code thus expresses how similar it is to other cgMLST

189    profiles. Very different profiles will show identity at few or no prefix positions of their LIN codes,

190    whereas nearly identical genomes will have LIN codes identical at most or all positions (see *e.g.*,

191    **Figure 1**, genomes Z *versus* X: shared prefix 0_2_0_0_0_0 implies a minimum similarity of 98.88%,

192   inclusive, and a maximum similarity of 99.36%, exclusive). Of note, our definition of LIN code prefix

193   is similar to the LINgroup concept proposed by Vinatzer and colleagues (Vinatzer et al., 2017).

194   An important particularity of LIN codes is that the numerical identifiers at a given bin position (except

195   the leftmost one) can only be interpreted in the context of the LIN code prefix preceding the

196   considered bin: the same integer value at a given bin position corresponds to group membership only

197   if the upstream prefixes are identical. In other words, groups at a given bin position are subdivisions of

198   the upstream prefixes and are numbered starting from zero independently for each prefix. This

199   particularity of LIN codes reduces the total number of integer identifiers observed in each position,

200   making them easier to read than systems in which a group identifier is created independently at each

201   level (for example, there are currently > 10,000 group identifiers at HierCC-1 level; (Achtman et al.,

202   2022)). Interestingly, the diversity observed within a group defined by a given prefix can immediately

203   be deduced from the maximal integer found among its members in the bin immediately downstream of

204   the prefix length (**Figure 2**).



205

206 **Figure 2**. **The hierarchical nature of LIN code positions**. Numbering starts from 0 for subdividing
207 each higher-level partition, characterized by a unique LIN code prefix. The hierarchical structure of
208 LIN codes is shown here with a circular packing plot obtained from the KpSC data from BIGSdb-
209 Pasteur. The circles correspond to LIN code prefixes of lengths 1 to 4 (an extra, all-encompassing
210 circle corresponds to the entire KpSC); the size of the circles is related to the number of genomes they
211 comprise. The first two bins in the LIN code prefix are used to identify phylogroups. Where for some
212 phylogroups the first bin is unique (*e.g.*, prefix 0 for Kp1), in other cases it is common to multiple
213 phylogroups (*e.g.*, prefix 2, which is associated with both Kp2 and Kp4), and therefore the second bin
214 is necessary to discriminate between them (*e.g.*, 2_0 and 2_1 for Kp2 and Kp4, respectively). The
215 hierarchical nature of LIN codes applies to subsequent levels of the prefix such as to those
216 corresponding to sublineages (third bin, e.g. Kp1 SL258 is identified with the LIN code prefix
217 0_0_105) and to clonal groups (fourth bin, e.g. Kp1 CG258 with the LIN code prefix 0_0_105_6).
218 Data was plotted in R v4.3.2 with ggplot2 and edited using Inkscape.

219

220 **Box 2. The particulars of LIN codes: handling of missing data, equal matches, input order and**
221 **computational precision**

222 **Missing data**. Whereas 7-gene MLST genotyping requires complete allelic profiles, cgMLST
223 approaches can tolerate the presence of missing alleles, as some core genes may not be essential, and
224 as genome assembly shortfalls occasionally result in the absence or incompleteness of some loci.
225 Therefore, the definition of cgSTs needs to accommodate missing data. Profiles may differ only by
226 loci where there is one or more missing allele(s) in one of the profiles, while otherwise identical at all
227 loci called in both profiles. Such profiles will be assigned to distinct cgSTs. We define as coincident
228 cgSTs, groups of cgST profiles that differ only by their missing data pattern. As the dissimilarity
229 between profiles is computed based solely on loci called in both profiles (Hennart et al., 2022),
230 coincident cgST profiles will have a 0 dissimilarity value between them, and therefore the same LIN
231 code.

232 Near-identical isolates or different WGS runs of the same isolate can lead to variable missing allele
233 calls but are otherwise identical in the called loci, and will as a consequence lead to the creation of two
234 or more coincident cgSTs. Each of these isolates' profiles will match with these multiple coincident
235 cgST. When a given profile matches two or more predefined coincident cgSTs, it will (by definition)
236 be attributed to all the coincident cgSTs. To minimize this phenomenon, a maximum number of
237 accepted missing data must be defined when implementing the cgST classification within BIGSdb.

238 **Equal matches and unicity of LIN codes**. As described above, an isolate's profile may match more
239 than one encoded cgST, due to missing loci. In this case, a unique LIN code will be defined (and
240 displayed) for the isolate. To choose between the different possibilities, the LIN code of the cgST with
241 the fewest missing allele(s) will be attributed. When two or more coincident cgSTs have the same
242 number of missing allele(s), the cgST with the smallest LIN code partition identifiers (considered from
243 left to right bin, *i.e.*, the lowest sort order) will be chosen. The same priority rule is applied to encode

244    every novel profile that is equidistant to two (or more) previously LIN-encoded non-coincidental

245    cgSTs.

246    **Input order**. The LIN code approach is dependent on input order, as the partition in a given bin may

247    vary slightly according to the order by which the genomes were encoded (Hennart et al., 2022). To

248    minimize this effect, BIGSdb uses the traversal of a minimum spanning tree (MStree; (Prim, 1957)) to

249    define the order by which the novel profiles are encoded. To code a novel batch of genomes, after

250    creating a MStree, the isolate chosen as the starting point for LIN encoding is the one that has the

251    closest similarity to an already encoded isolate in the database; next, the MStree is traversed from this

252    node. This approach (implemented since v1.36.1) maximizes reproducibility when adding a batch of

253    novel genomes. To minimize the number of resulting prefix-based partitions, novel genomes should be

254    encoded in batches as large as possible.

255    **Computational precision.** As for all categorizations that rely on thresholds, computational precision

256    is critical for reproducible results. For example, the pairwise dissimilarity between cgMLST profiles,

257    which is a ratio, may often have a higher number of decimals than can be handled by the computing

258    system, and its rounded value may lead to a slight underestimate (or overestimate) of the true value.

259    When the (true) dissimilarity between an incoming profile and its reference is exactly identical to the

260    left threshold of a bin (*i.e.*, the same ratio of distinct versus called alleles), a rounded value may

261    incorrectly correspond to the previous bin (**Figure 3**). Therefore, pairwise dissimilarity computations

262    should be performed in a way exactly identical to the bin thresholds themselves. In BIGSdb, ratios

263    corresponding to the thresholds are compared to the calculated dissimilarity values using Perl

264    platform-native floating point values (usually IEEE 754 double-precision).

265

**Figure 3**. **The effect of rounded cgMLST similarity values on LIN code assignment**. In this example, the use of a rounded value for the similarity between genome X and genome D leads to a slight underestimate, therefore creating a novel identifier in bin 7, instead of bin 8 when computing the similarity with the same precision as the threshold.

**LIN codes functionalities implemented within the BIGSdb platform**

The LIN code taxonomy of KpSC genomes was incorporated into the Institut Pasteur *K. pneumoniae* MLST and whole-genome MLST platform (https://bigsdb.pasteur.fr/klebsiella), using BIGSdb v1.34.0 and upwards (Hennart et al., 2022). For the KpSC, this database plays the role of the source database for the definitions of alleles, cgMLST profiles, cgSTs, and LIN codes.

In BIGSdb, LIN code schemes can be defined in the curator's interface of both the 'sequence definition' and 'isolates' databases. A LIN code taxonomy is created with reference to a defined indexed scheme, *e.g.*, cgMLST. An indexed scheme is a scheme with a unique identifier for each profile, *e.g.,* cgST here. To index a scheme, one needs to specify the maximum number of missing alleles accepted for profiles to be assigned to cgSTs. To create a LIN code taxonomy, allele mismatch thresholds that define the LIN code bins must simply be defined. In the case of KpSC, the 629-loci cgMLST scheme was selected, and ten thresholds were defined (**Figure 1**).

Users who wish to assign a novel LIN code for a KpSC isolate must submit the genome sequence(s) to the BIGSdb-Pasteur 'isolates and genomes' database. If all quality criteria are fulfilled

286  (https://bigsdb.pasteur.fr/klebsiella/genome-quality-check/), the genome(s) will be deposited in the

287  database for allele, cgMLST profile, cgST and LIN code definitions. The inferred cgMLST profiles, as

288  well as their cgST identifiers and LIN codes, will be made openly accessible through the sequence and

289  profile definition database ('seqdef'). To ensure confidentiality of users' data when requested, isolate

290  metadata and associated genome sequence(s) can be embargoed and released at a later stage.

291  Users can search *K. pneumoniae* isolates of interest using the LIN code matching functionalities

292  implemented in BIGSdb. A complete LIN code (or any prefix) can be used as a query. The nickname

293  nomenclature attached to LIN code prefixes can also be used to facilitate the query of groups of

294  interest (*e.g.,* SL258 members can be searched by using its attached prefix 0_0_105, or using the

295  SL258 nickname itself). The list of genomes from the query results can be further analyzed using the

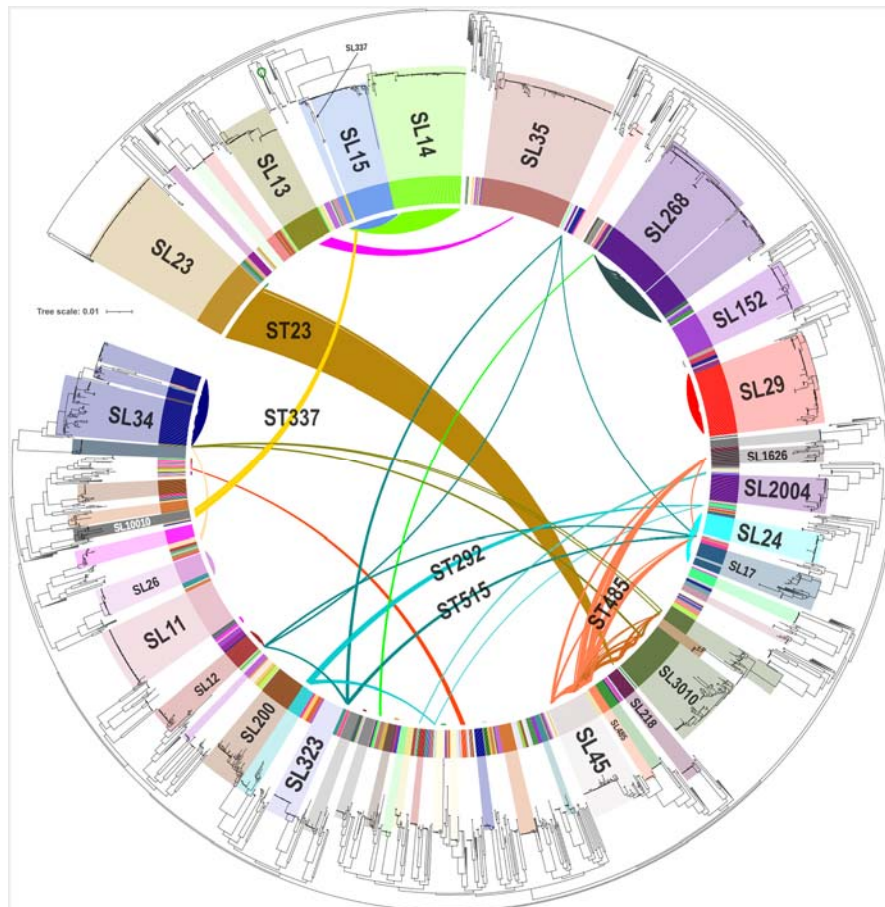296  available analytical tools within the BIGSdb platform, or exported for external use.

297

298

299 **Section 2: Novel developments and examples of applications**

300 **Multiple *Klebsiella pneumoniae* 7-gene MLST sequence types are polyphyletic**

301 Even though they are based on allelic profile comparisons rather than a sequence-based phylogenetic

302 analysis, LIN code prefixes of length 3 or 4 bins are compatible with phylogenetic classifications and

303 thus represent markers of their corresponding tree branches (Hennart et al., 2022). In contrast, 7-gene

304 MLST may conflate phylogenetically unrelated genomes in a single ST, for example through

305 recombination leading to the same ST being assigned to genomes from distinct parental lineages, or by

306 large recombinations affecting multiple cgMLST loci but leaving the 7-gene MLST loci unaffected

307 (Lam et al., 2023). Here we explore the extent of this phenomenon using 44,000 publicly available

308 genomes of *K. pneumoniae* (June 2023). We found that 113 STs are polyphyletic, defined here as

309 being observed in at least two unrelated LIN code sublineages (**Table S1**). We illustrate this

310 phenomenon for major STs in **Figure 4**. For example, ST485 was observed in four phylogenetically

311 unrelated sublineages: SL485 (0_0_157), SL45 (0_0_158), SL1626 (0_0_227) and SL11569

312 (0_0_1215). ST347 stands out as being observed in 8 distinct sublineages. This analysis also

313 confirmed the polyphyletic status of ST23 (Lam et al., 2023), which conflates isolates from distant

314 sublineages: SL23 (0_0_429) and SL218 (0_0_115).



315

316    **Figure 4**. **Phylogenetic tree of *K. pneumoniae* main sublineages**. A phylogenetic analysis of 5,665
317    *K. pneumoniae sensu stricto* genomes (LIN code prefix 0_0; see selection process in Methods) was
318    performed from the multiple sequence alignments of 629 cgMLST genes. Closely related leaves were
319    collapsed. The colored sectors in the inner circle correspond to the sublineages (SL) defined based on
320    their prefix of length 3 (*i.e.*, made of the three first bins); the major sublineages are highlighted by
321    lighter-colored sectors joining the circle to the tree leaves. The internal connectors between
322    sublineages represent frequent STs that were found in two or more sublineages. The full interactive
323    tree is available at: https://itol.embl.de/tree/1579917420525181688029926

324

325    **Nicknaming the LIN code prefixes enables carry-over of MLST identifiers into the genomic**
326    **taxonomy**

327    Whereas LIN code prefixes themselves can be used as canonical markers of groups of interest that are
328    easy to handle by computers, for humans, prefixes are not very easy to remember or pronounce. Here,
329    we propose to nickname the LIN code prefixes with simple denominations using a LIN code prefix
330    nicknaming system (newly implemented within BIGSdb;
331    https://bigsdb.readthedocs.io/en/latest/administration.html?highlight=prefix#setting-up-lincode-
332    definitions-for-cgmlst-schemes). It is thereby possible to nickname every prefix in any chosen way,
333    for example by incrementing an integer identifier for each novel prefix of a given length, analogous to
334    the numbering of 7-gene MLST STs. Other labels could be applied, such as Greek letters,
335    astronomical objects, or any other series of words that may be universally understandable and easy to
336    remember. This nicknaming process would be particularly useful for long prefixes, or prefixes of
337    particular relevance that subdivide the population at particularly informative levels.

338    For bacterial species where previous nomenclatures exist, a novel and unrelated naming system would
339    have the drawback of creating yet another nomenclature. Assigning nicknames to prefixes based on
340    the previous nomenclature system is therefore more meaningful. For *K. pneumoniae,* the classical
341    MLST nomenclature system is widely used, and knowledge has accumulated on the epidemiological
342    history and characteristics of predominant STs. We therefore aimed to create backward nomenclatural
343    compatibility of LIN codes with ST identifiers. We used a majority identifier inheritance rule that was
344    previously developed and applied to single-linkage cgMLST groups (Hennart et al., 2022). We applied
345    this approach to nickname LIN code prefixes of lengths 3 and 4 bins (which, for convenience, we have
346    defined as sublineages and clonal groups, respectively) by using ST identifiers as a source. In short,
347    for each LIN code prefix of length 3 or 4 bins, the identifier of the predominant ST among its genomes
348    was used as a label, wherever possible (*i.e.*, if not yet attributed). Following this approach, most SLs
349    and CGs were indeed labeled according to the ST identifiers of most of their isolates, whereas a
350    minority are nicknamed with incremental numbers (because the majority ST was already used for
351    another prefix). In **Figure 5,** we provide illustrative examples of correspondence between prefixes and
352    nicknames for major clonal groups. For example, ST258, and its derivative ST512, share the prefix
353    0_0_105, nicknamed SL258, and the 4-positions prefix 0_0_105_6, nicknamed CG258.

354

| LIN prefix | Phylo-group |
|---|---|
| 0_0 | Kp1 |
| 1_0 | Kp3 |
| 1_1 | Kp5 |
| 2_0 | Kp2 |
| 2_1 | Kp4 |
| 3_0 | Kp6 |
| 4_0 | Kp7 |

| LIN prefix | Main ST | Nickname |
|---|---|---|
| 0_0_0 | 15 | SL15 |
| 0_0_429 | 23 | SL23 |
| 0_0_105 | 258 | SL258 |
| 0_0_158 | 45 | SL45 |
| 0_0_197 | 147 | SL147 |
| 0_0_369 | 307 | SL307 |
| 0_0_750 | 6589 | SL10691 |

| LIN prefix | Main ST | Nickname |
|---|---|---|
| 0_0_105_6 | 258 | CG258 |
| 0_0_105_0 | 340 | CG340 |
| 0_0_105_2 | 11 | CG11 |
| 0_0_105_11 | 11 | CG3666 |
| 0_0_105_1 | 437 | GC10268 |
| 0_0_105_29 | 11 | CG12811 |
| 0_0_105_7 | 895 | CG895 |

355

**Figure 5**. **Nicknaming of LIN code prefixes enables inheritance of previous nomenclatures**. Nicknames of some LIN code prefixes of lengths 2 to 4 bins, inherited from phylogroup numbering or Linnaean taxonomy (2-bin prefix, left panel) or 7-gene MLST (prefixes of lengths 3 and 4 bins, central and right panels), are displayed.

Note that the MLST nickname inheritance rule was applied only using ST identifiers up to ST6500 (**Figure S1**). Given that the main sublineages of KpSC have long been sampled, the inheritance of MLST identifiers on SL and CG identifiers will apply to most of the extant diversity of the KpSC. For subsequent prefixes, SL and CG nicknames are numbered incrementally, starting with 10,000 (see example on **Figure 5**) in order to make clear that these new nicknames are not inherited from MLST nomenclature. In parallel, continual expansion of the MLST nomenclature will result in defining STs (incremented by one) upwards of 6500 (currently the highest ST is ST6859, January 21[st], 2024). Hence, a correspondence between ST identifiers >6500 and prefix nicknames >10,000 may exist but will not be immediately obvious. For novel sublineages and clonal groups that may emerge in the future, our recommendation is to prioritize their LIN code SL and CG nicknames, rather than their ST, when communicating on these groups. Note that the 2-bin prefixes of *Klebsiella* LIN codes each define a particular KpSC phylogroup, corresponding to the seven currently described species or subspecies (Hennart et al., 2022), and were thus nicknamed accordingly (**Figure 5**).

**From dual- to single-barcoding taxonomy of *Klebsiella pneumoniae* strains**

Previously, cgMLST groups were defined by the single-linkage (slink) clustering method using the same 10 thresholds as for the LIN codes, and the four highest-level groups were nicknamed by inheritance from Linnaean taxon names (for the two first) or MLST labels (for the levels defined by thresholds 190 and 43, dubbed Sublineage and Clonal Groups, respectively) (Hennart et al., 2022). Together with the LIN code taxonomy (which had no nickname in (Hennart et al., 2022)), this slink-based system formed a 'dual-barcoding approach'. However, because such slink groups suffered from

382     fusion of existing groups upon addition of subsequent genotypes, which occasionally had intermediate

383     distances between preexisting groups (*e.g.*, hybrid genotypes), the classification of cgMLST profiles

384     into slink groups was abandoned. Fortunately, when excluding the hybrid genotypes, a nearly

385     complete concordance was observed at the four first levels between slink clusters and LIN code

386     groupings optimized based on MStree (Hennart et al., 2022). As a result, the LIN code taxonomy

387     currently in use is nearly fully consistent with the one initially proposed (only SL10000 to SL10021,

388     and CG10000 to CG10276 correspond to groups that were renamed; table of correspondence available

389     upon request). The use of a single-barcoding taxonomic system based on LIN codes will stabilize and

390     simplify the way groups are defined and labeled.

391

**392     LIN code taxonomy usage in external genomic epidemiology platforms**

393     To make the LIN code taxonomy accessible for external tools, databases and analysis platforms, the

394     LIN code nomenclature components (alleles, profiles, cgSTs and LIN codes) can be extracted from

395     BIGSdb using an application programming interface (Jolley et al., 2017). This can be performed via a

396     single       query       using       the       following       link:

397     https://bigsdb.pasteur.fr/api/db/pubmlst_klebsiella_seqdef/schemes/18/profiles_csv. However it is

398     important to note that to be effective, external copies of the database need to be very frequently

399     synchronized with the primary nomenclature database. This is because, when genome sequences

400     (through their cgMLST profiles) are matched to the LIN code taxonomy, an incomplete LIN code may

401     be defined in many cases, as no identical cgMLST profile may be existing at this time in the source

402     LIN code taxonomy. In such cases, a new nomenclatural identifier must be defined and assigned, but

403     this is only possible within the source database otherwise consistency of nomenclature will be lost.

404     Inference of the query genome's LIN code in external resources can only be inferred up to the bin

405     preceding the pivot bin corresponding to the closest match. Notably though, when the LIN code prefix

406     up to the fourth bin (at least) can be defined for the query genome, information on species, subspecies,

407     SL and CG can be derived. If the query genome is closely related to one in the source database, its

408     LIN code will be almost completely defined. Therefore, although novel cgMLST alleles, cgST profiles

409     and LIN codes can only be defined in the source database of the nomenclature (BIGSdb-Pasteur for

410     the KpSC), the use of LIN codes in external databases or tools still has functional relevance. For any

411     genome (cgMLST profile) that has no complete LIN code, data submission to the source database is

412     encouraged, in order to update the LIN code taxonomy and define complete LIN codes for the novel

413     genomes.

414     To illustrate the external use of LIN codes, we implemented the KpSC LIN code taxonomy in the

415     Pathogenwatch platform, in which a KpSC database was set-up previously (Argimón et al., 2021).

416     First, on a regular basis, Pathogenwatch synchronizes from BIGSdb into its internal temporary

417  database, the defined alleles, cgSTs and associated LIN codes, using the API functionality of BIGSdb.

418  Second, the cgMLST allele sequences extracted from the query genome assembly are compared to

419  those in the temporary database, and the cgMLST profile is used to find the closest match in the

420  temporary database. If the query genome does not match completely with an existing source

421  nomenclature cgST, a provisional cgST is assigned, represented by the asterisk and a code (*e.g.,* cgST

422  *f26e). Pathogenwatch also indicates the closest cgST defined in the source taxonomy database and

423  provides a link to the list of all isolates within Pathogenwatch that have the same cgST genotype.

424  Third, an incomplete LIN code will be provided by Pathogenwatch based on the shared prefix with the

425  closest reference cgST (**Figure 6**). This process provides information about the relatedness of a query

426  Pathogenwatch genome compared to the existing taxonomy elements and can in most cases provide

427  sublineage and clonal group identification. In those cases where Pathogenwatch provides provisional

428  alleles, STs, cgSTs and/or LIN codes, the user is encouraged to submit the genomic sequence data to

429  the source BIGSdb-Pasteur database so that novel nomenclatural identifiers (alleles, STs, cgSTs, LIN

430  codes) can be created. Note that as Pathogenwatch uses its own algorithm to provide the species and

431  subspecies for KpSC genomes, this taxonomic information is not deduced from LIN codes in that

432  platform.

433



434  View all cgST *f26e

435

436  **Figure 6. Example of LIN code identification in Pathogenwatch**. Although the LIN code is
437  incomplete, the genome can be inferred to belong to clonal group 258 (defined as prefix 0_0_105_6),
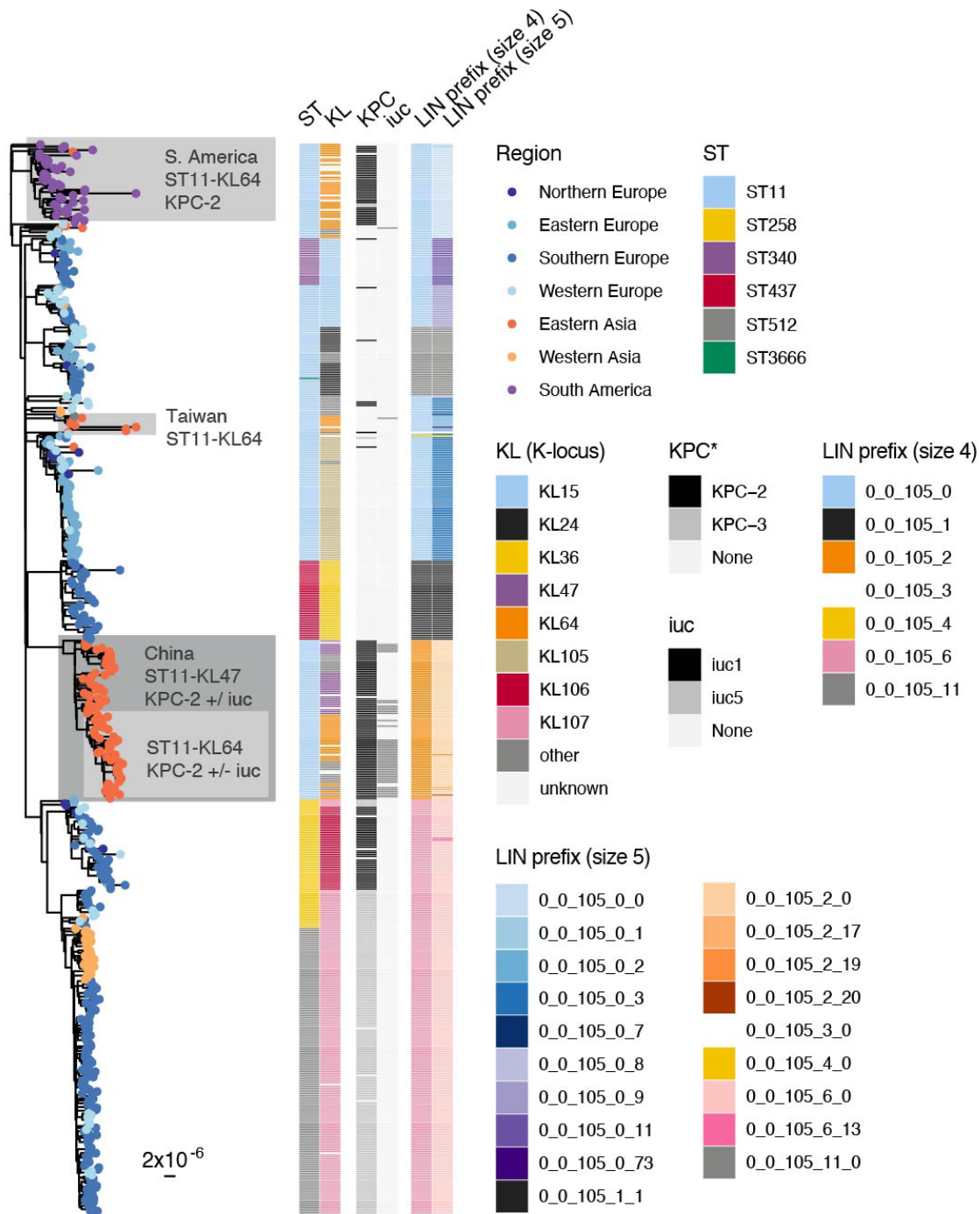438  which comprises ST258 and ST512 isolates (see **Figure 7**).

439

**Applications of LIN codes to subdivisions within high-risk Kp sublineages**

441  A number of *K. pneumoniae* sublineages, including SL258, SL147, SL307, SL17 and SL23, have been

442  recognized to cause a large burden of so-called hypervirulent or multidrug resistant infections. These

443  groups have been the subject of detailed studies, that have led to defining their geographical spread

444  and phylogenetic subgroups (Deleo et al., 2014; Hetland et al., 2023; Lam et al., 2018; Rodrigues et

445 al., 2022; Wyres et al., 2019). However, so far, a harmonized nomenclature of these subgroups has
446 been lacking, making it difficult to recognize them in subsequent studies. Here, we illustrate how LIN
447 codes can help track Kp dissemination at fine genetic scales within sublineages, using the example of
448 SL258, a major *Klebsiella pneumoniae* carbapenemase (KPC) producing sublineage of *K.*
449 *pneumoniae*.

450 SL258 is defined by its LIN code prefix, 0_0_105, and encompasses all isolates from 7-gene ST11,
451 ST258, ST340, ST512 and some others (**Figure 5**). Its phylogenetic structure shows that SL258 is
452 divided into several clades (**Figure 7**) that are labeled with their unique clonal group number. These
453 include CG258 (0_0_105_**6**), defined by LIN code position 4, which contains all ST258 and ST512
454 isolates. LIN code position 5 can further be used to distinguish major subclades within SL258,
455 including ST340 (0_0_105_**0_11**) and ST437 (0_0_105_**1_1**) and other subclades within ST11, some
456 of which appear to be associated with recombination events that include the capsule locus (KL column
457 in **Figure 7**). The LIN codes can also help distinguish between different subclades that are associated
458 with the same capsule locus. For example, they clearly distinguish 3 subclades that are all ST11-KL64
459 (grey shading on the tree branches, **Figure 7**). One of these is the major lineage circulating in China
460 (0_0_105_**2_0_0_2**, predominantly 0_0_105_2_0_0_2_17, 24/30 genomes) that carries KPC-2 and
461 often the *iuc1* aerobactin virulence locus, descended from ST11-KL47-KPC-2 (0_0_105_**2_0_0_2_***,
462 where * is not 17), as discussed broadly in the literature (Zhou et al., 2023, 2020). A second, unrelated
463 ST11 subclade carrying KL64 (0_0_105_**0_0**) is circulating in South America (encoding KPC-2, but
464 rarely *iuc*), while a third smaller clade (0_0_105_**0_2**) is detected primarily in Taiwan rather than in
465 mainland China (lacking KPC and with only one of eight genomes carrying *iuc*). The example of
466 SL258 illustrates how LIN code classification beneath the sublineage level can help recognize and
467 name subgroups of medical and epidemiological relevance, which should be the object of enhanced
468 surveillance.

**Figure 7. SL258 phylogenetic structure and LIN codes.** Maximum-likelihood phylogenetic tree of n=586 SL258 genomes inferred from a recombination-free variable site alignment (see Methods). Tips are coloured to indicate geographic region of origin as per the legend (United Nations region classifications). The distribution of 7-gene sequence types (STs), K-loci (KL), $bla_{KPC}$ (KPC) alleles, aerobactin locus lineages (*iuc*), LIN code prefixes of sizes 4 and 5, are indicated by colored blocks as per the legends (note that colors are independent to each column). Only K-loci identified with a Kaptive confidence score of 'Good' or better are shown (otherwise marked 'unknown'). Two isolates were detected with $bla_{KPC-30}$ and one with $bla_{KPC-12}$ but are not shown in the figure for brevity. Sub-clades described in the text are coloured and labeled accordingly.

479

### Application of LIN codes to outbreak strain identification

481 To illustrate the use of LIN codes to identify outbreak strains, and to track strain diversification during
482 protracted outbreaks, we explored the example of SL147. This is a prominent multidrug-resistant
483 international sublineage of *K. pneumoniae*, defined by its LIN code prefix 0_0_197. **Figure S2**
484 illustrates how the phylogenetic relationships within SL147 are captured by LIN codes, using a
485 previously described dataset (Rodrigues et al., 2022). SL147 comprises a single clonal group
486 (0_0_197_0) and three 7-gene STs (ST147, ST273 and ST392). At LIN code position 5, four partitions
487 (0_0_197_0_**0**, 0_0_197_0_**4**, 0_0_197_0_**17** and 0_0_197_0_**25**) correspond largely to ST273,
488 ST392 and two deep branches of ST147. In addition, both ST147 and ST273 are genetically
489 heterogeneous and structured phylogenetically into several minor branches, which were captured by
490 additional partitions of LIN code level 5 (**Figure S2, panel A**).

491 Protracted outbreaks often lead their investigators to define local clades (or subgroups) within the
492 closely related outbreak isolates. These clades are often attributed temporary placeholder names,
493 which are difficult to compare across studies *e.g.*, Clade A and Clade B (Martin et al., 2021). We
494 illustrate how LIN codes provide a way to define these clades definitively, using the diversity among
495 outbreak isolates from a metallo-β-lactamase (NDM)–producing carbapenem-resistant ST147 outbreak
496 in Tuscany (**Figure S2, panel B; Table S2**). The time span of the Tuscany outbreak is November
497 2018 - 2021. Most of the isolates in this outbreak have prefix 0_0_197_0_**4_1_0**, thus differing by no
498 more than 4 alleles out of 629 with another member of the group. The authors defined two clades, A
499 and B. Here, clade B corresponds to the set of LIN codes 0_0_197_0_4_1_0_8_x_x (*i.e.,* with prefix
500 0_0_197_0_4_1_0_**8**, with x meaning there may be variation at the two last positions). Clade A was
501 more diverse, and LIN codes classify this genetic variability in a definitive way, with six 8[th] position
502 prefixes (0_0_197_0_4_1_0_**7**, 0_0_197_0_4_1_0_**9**, 0_0_197_0_4_1_0_**10**, 0_0_197_0_4_1_0_**11**,
503 0_0_197_0_4_1_0_**12** and 0_0_197_0_4_1_0_**66**). This example highlights how *K. pneumoniae* LIN
504 codes can subdivide isolates from long-term outbreaks.

505 A search of the BIGSdb-Pasteur KpSC database (January 31[st], 2024) for prefix 0_0_197_0_4_1_**0**
506 identified n=395 *K. pneumoniae* genomes, isolated between 2014 and 2023 and coming from 20
507 countries from North America, Europe, Asia, Africa and Oceania, which indicate the global
508 dissemination of this particular subgroup of SL147. However, prefix 0_0_197_0_4_1_**0_8** was so far
509 only reported from the Italian outbreak. This example illustrates how LIN codes can facilitate the
510 tracking of strain dissemination, by enabling the identification of similar isolates from separate
511 studies. As an outbreak strain prefix can be easily discussed and shared among investigators and is
512 sufficient to exchange information on strain identity across countries, LIN codes enable genomic
513 surveillance investigations without the need to share genomic sequences, which may alleviate issues

514     around data confidentiality. Likewise, for the surveillance of particularly concerning strains, early

515     warnings could be triggered based on the detection of the specific LIN code of the strains under

516     surveillance.

517     Given that LIN codes are phylogenetically informative, they can be represented graphically as prefix

518     trees, which broadly approximate the phylogenetic relationships among isolates (Hennart et al., 2022).

519     Here, we introduce the tool LINtree to create prefix trees from LIN codes

520     (https://gitlab.pasteur.fr/GIPhy/LINtree). The input file contains a list of genome names and LIN

521     codes (one sample per row), with a header row indicating the level of similarity for each bin. LINtree

522     outputs a Newick-formatted tree showing the relationships between input genomes, based on the

523     hierarchy provided by the LIN codes and with branch lengths scaled using the similarity levels in the

524     header row. For example, the tree of the ST147 Italian outbreak shown in **Figure S2** was generated

525     using this tool, based on the input list of LIN codes. This example illustrates how the prefix tree

526     recapitulates the phylogenetic relationships of this outbreak strain with its ancestral relatives,

527     providing a useful aid in outbreak investigations.

528

## Discussion

530     Facilitating communication on the intraspecific diversity of bacterial strains is a key objective of strain

531     taxonomies, which entail classification and naming of groups within species. In the field of

532     epidemiological surveillance of pathogens, it has long been recognized that strain typing methods used

533     for long-term and global strain tracking should rely on an internationally standardized nomenclature

534     (Struelens, 1998). In turn, a robust and fine-grained strain taxonomy promotes the understanding of

535     the links between genotypes and clinical phenotypes, vaccine coverage and antimicrobial resistance

536     (Achtman et al., 2022; Maiden et al., 2013).

537     Here we have presented in detail the cgMLST-based LIN code approach and further developed this

538     novel strain taxonomy system. The stability of LIN code classification is a critical property, which has

539     been impossible to achieve with previous strain classification systems relying on single-linkage

540     clustering (such as MLST clonal complexes defined by BURST or cgMLST single-linkage groups).

541     cgMLST LIN codes are stable, as the incorporation of novel genomes has no effect on pre-existing

542     LIN codes (**Figure 1**). Here, we have presented important enhancements of our initial implementation,

543     by (i) improving the reproducibility of LIN encoding by addressing the dependency of this approach to

544     rounded genetic distance values; (ii) the implementation within the BIGSdb platform, of input order

545     rules for creating novel LIN codes, and (iii) implementing formal rules for handling missing data.

546     These improvements optimize the definition of LIN codes and have resulted in a robust strain

547     taxonomy system that is now in operation for *K. pneumoniae* since January 2023 and currently

548     comprises 37,070 cgSTs and 32,500 LIN codes, which correspond to 2,492 sublineages and 4,230

549     clonal groups (January 28[th], 2024).

550     In this work, we also extend the LIN code approach by proposing and implementing a nicknaming

551     system for LIN code prefixes. As shown previously (Hennart et al., 2022; Marakeby et al., 2014), LIN

552     codes are highly compatible with phylogenetic relationships, and their prefixes can therefore act as

553     markers of phylogenetic groups. Nicknaming was designed to be flexible, and can thus accommodate

554     any naming system of choice, either numerical or textual. To ensure continuity with 7-gene MLST

555     nomenclature, we had previously proposed to nickname cgMLST single-linkage groups (Hennart et

556     al., 2022). For *K. pneumoniae*, we had nicknamed the partitions within two special levels with

557     thresholds of 43 and 190 mismatches, defined as "sublineages" and "clonal groups", respectively.

558     However, because of the instability of the single-linkage clustering approach, we soon observed

559     fusions of previously defined (and nicknamed) groups, rendering the single-linkage-based

560     nomenclature unstable. Here, we instead nickname the LIN code prefixes of lengths 3 and 4 bins,

561     which correspond to the same thresholds as previously defined "sublineages" and "clonal groups",

562     respectively. Hence, we here redefined the "sublineages" and "clonal groups" as being based on LIN

563     code prefixes.

564     A key property of a novel nomenclature system is its continuity with previous nomenclatures, as it

565     minimizes confusion and facilitates its adoption by microbiologists and epidemiologists. Establishing

566     a dictionary of correspondence between novel and previous nomenclatures is a possibility but it

567     implies cumbersome handling of both series of identifiers. Here, we provide the possibility of

568     embedding any previous nomenclature(s) within the LIN code taxonomy. In the case of *K.*

569     *pneumoniae*, by using a previously described inheritance algorithm (Hennart et al., 2022) that has

570     mapped the 7-gene ST identifiers onto LIN code prefixes of lengths 3 and 4 bins, we provide

571     continuity between the novel nomenclature of sublineages and clonal groups with the widely used

572     MLST standard. Using LIN code prefix nicknames instead of MLST identifiers has the additional

573     benefit of enhancing the compatibility of the nomenclature with phylogenetic relationships: we have

574     shown here for *K. pneumoniae* that classical MLST profiles often conflate unrelated sublineages. Note

575     that we still recommend the maintenance and extension of the MLST nomenclature to classify future

576     *K. pneumoniae* isolates, in parallel to the novel genomic nomenclature. However, we suggest the

577     prioritization of LIN code nomenclature over MLST, which will be particularly important for

578     sublineage and clonal group designations above 10,000 that are not inherited from MLST.

579     Hierarchical clustering (HierCC) also provides stable classifications and is likewise implemented

580     based on cgMLST schemes (Zhou et al., 2021). Unlike for LIN codes, HierCC partition identifiers are

581     incremented independently across levels, necessitating the handling of large integers, particularly in

582     bins corresponding to the highest similarities, where over 100,000 partitions might be created. In

583    contrast, LIN codes re-initiate the numbering from 0 within a bin, for each subdivision of a partition in

584    the upper bin, resulting in a predominance of small integers, which are easier to handle for humans. By

585    design, HierCC is stable only in its production mode, whereas it relies on the unstable single-linkage

586    clustering approach in its development mode, implying an arbitrary decision on the switch from

587    development to production to achieve stability.

588    LIN codes, as well as HierCC, are multilevel classifications that provide proxies of strain

589    relationships. By conveying for each genome, its group membership and approximate degree of

590    relatedness at various phylogenetic depths simultaneously, they are phylogenetically informative. LIN

591    code prefixes are shared by genomes having at least the identity corresponding to the upper threshold

592    of the last prefix bin (exclusive). The LIN codes (or HierCC codes) can in fact themselves be

593    represented as a tree (formally, a prefix tree), with multifurcations corresponding to subdivisions of

594    each prefix (**Figure S2, panel C;** see also (Hennart et al., 2022)) and node height corresponding to bin

595    thresholds. This tree representation of LIN codes may serve as a proxy for the phylogenetic tree and

596    can be created with no need of initial sequences or cgMLST profiles.

597    A taxonomic system needs to be created and updated in a coordinated manner. For this purpose, the

598    cgMLST LIN code strain taxonomy approach was implemented in the BIGSdb platform. Its

599    integration in this widely used platform will make it publicly available, and will facilitate its

600    implementation for other bacterial species, as was recently illustrated for *Streptococcus pneumoniae*

601    (Brueggemann *et al.* bioRxiv 2023, doi: https://doi.org/10.1101/2023.12.19.571883). The applicability

602    to other bacterial species should be straightforward, provided that they comprise meaningful cgMLST

603    diversity, excluding the so-called monomorphic pathogens (Achtman, 2008), such as *Mycobacterium*

604    *tuberculosis* or *Salmonella enterica* serotype Typhi. Setting up LIN codes for other species will

605    require defining tailored bin thresholds based on population structures, which requires globally

606    representative genome datasets (**Figure S3, overview chart**). The approach could also be extended

607    with minor adaptations to other organisms with predominantly clonal reproduction, such as protozoan

608    parasites and fungi, even if they are not haploid (Bougnoux et al., 2004; Yeo et al., 2011). The wide

609    adoption of the standardized cgMLST LIN code strain taxonomy would result in a universal strain

610    nomenclature approach that could greatly enhance microbial biodiversity studies, international

611    genomic epidemiology and infectious disease surveillance.

612

613

## Methods

### Identification of MLST sequence types that are discordant with sublineage classification

We used the 44,000 public genomes available in BIGSdb in June 2023. To spot potential discordances between ST and LIN code prefixes, we first filtered out non-Kp1 phylogroup (prefix 0_0) genomes and removed nearly identical cgMLST profiles, by keeping a single representative of each partition at LIN code bin 5. STs observed only in a single isolate were filtered out. We next searched for all STs that were split in several clonal groups or sublineages (as defined by their prefix) and conversely, also looked for prefixes of length 3 or 4 which comprised several STs. We then placed these genomes in a phylogenetic tree built using IQtree v2.2.2.2 (Minh et al., 2020) using GTR+I+G model, from concatenated alignments of individual cgMLST gene alignments.

### SL258 phylogeny

Whole genome sequences representing SL258 were identified among the EuSCAPE collection (David et al., 2019) and two recent studies reporting 7-gene ST11 with K-locus (KL) 47 and KL64, for which multiple independent evolutions have been reported (Wang et al., 2023; Zhou et al., 2023). The ST11 genomes were subsampled to a manageable number as follows: (i) five randomly selected genomes per study year for each of ST11-KL47 and ST11-KL64 reported from sites across China, plus all ST11 with other K-loci reported in the same study included for context, total 92 genomes from this study (Zhou et al., 2023); (ii) 64 genomes representing ST11-KL64 clade 1 as defined in an analysis of public ST11-KL64 genomes (Wang et al., 2023). Genome assemblies were acquired from Pathogenwatch and those with >500 contigs and/or total assembly size < 4,969,898 or > 6,132,846 bp were removed (as per the KlebNET-GSP quality control definitions). Kleborate v2.3.2 (Lam et al., 2021) was used to determine 7-gene ST, $bla_{KPC}$ alleles, and $iuc$ lineages (aerobactin locus), and Kaptive v2.0.7 (Lam et al., 2022; Wyres et al., 2016) was used to identify KLs.

In order to infer a high-resolution phylogenetic tree, genome assemblies were used to simulate 100bp paired end reads with wgsim (without errors, https://github.com/lh3/wgsim). Reads were mapped against the NJST258-1 completed reference genome (NCBI accession: CP006923.1) and single nucleotide variants called using the RedDog pipeline (https://github.com/katholt/RedDog). The resultant allele table was converted to a pseudo-whole genome alignment and used as input for Gubbins v2.3.2 (Croucher et al., 2015), in order to detect and remove recombination (100 iterations). The final filtered alignment of 10,390 variable sites, representing 591 genomes, was used to infer a maximum likelihood (ML) phylogenetic tree using RAxML v8.2.9 with parameters: best of 5 runs, 1,000 bootstraps each, gamma model of rate variation (Stamatakis, 2014). Subsequently, five genomes were removed due to excessive branch lengths. The ML tree was visualized with R v4.3.1 and the

648    following packages: ape v5.7.1 (Paradis et al., 2004), phytools v 1.9-16 (Revell, 2024), and ggtree v

649    3.8.2 (Yu et al., 2017).

650

651    **Data availability**

652    There are currently 39,506 genomic sequences publicly available in BIGSdb. Genomes can be

653    downloaded from the sequence bin page, and LIN codes are available in the main table retrieved

654    following an isolates search page results.  The complete and up-to-date LIN code nomenclature

655    (comprising alleles, profiles, cgSTs and LIN codes) can be extracted from BIGSdb using a single

656    query                    at                    the                    following                    link:

657    https://bigsdb.pasteur.fr/api/db/pubmlst_klebsiella_seqdef/schemes/18/profiles_csv.

672    **Authors contributions**

673    *Klebsiella* genomic platform integration (KlebNET) conceptualization and coordination: SyB (Sylvain

674    Brisse), KEH, DMA. cgMLST LIN code conceptualization and developments: MH, AC, SB. BIGSdb-

675    Pasteur Platform maintenance: FP, BR, BB, SeB (Sebastien Bridel), SyB. Data acquisition and

676    curation: VP, RI, CR, FP, MH, CC, SeB, ML, KW, CAY, MRP, DA. Data analyses: MH, SeB, KW,

677    ML, CAY. PubMLST/BIGSdb platform software development: KAJ, MCJM. Pathogenwatch platform

678    maintenance and software development: DMA, CAY, SD. Data visualization: FP, MH, KW, ML, SeB,

679    CC. Writing of first draft: SyB, FP, MH. All authors contributed to, and approved, the final version of

680    the manuscript.

681     **Ethical statements**

682     Not relevant.

683     **Conflicts of interests**

684     The authors declare no conflict of interest.

## References

685 **References**

686  Aanensen, D.M., Spratt, B.G., 2005. The multilocus sequence typing network: mlst.net. Nucleic Acids
687  Res 33, W728-33.

688  Achtman, M., 2008. Evolution, population structure, and phylogeography of genetically monomorphic
689  bacterial pathogens. Annu Rev Microbiol 62, 53–70.

690  Achtman, M., Zhou, Z., Charlesworth, J., Baxter, L., 2022. EnteroBase: hierarchical clustering of 100
691  000s of bacterial genomes into species/subspecies and populations. Philos. Trans. R. Soc. Lond. B.
692  Biol. Sci. 377, 20210240. https://doi.org/10.1098/rstb.2021.0240

693  Argimón, S., David, S., Underwood, A., Abrudan, M., Wheeler, N.E., Kekre, M., Abudahab, K.,
694  Yeats, C.A., Goater, R., Taylor, B., Harste, H., Muddyman, D., Feil, E.J., Brisse, S., Holt, K., Donado-
695  Godoy, P., Ravikumar, K.L., Okeke, I.N., Carlos, C., Aanensen, D.M., NIHR Global Health Research
696  Unit on Genomic Surveillance of Antimicrobial Resistance, Fabian Bernal, J., Arevalo, A., Fernanda
697  Valencia, M., Osma Castro, E.C.D., Nagaraj, G., Shamanna, V., Govindan, V., Prabhu, A., Sravani,
698  D., Shincy, M.R., Rose, S., Ravishankar, K.N., Oaikhena, A.O., Afolayan, A.O., Ajiboye, J.J.,
699  Ewomazino Odih, E., Lagrada, M.L., Macaranas, P.K.V., Olorosa, A.M., Gayeta, J.M., Masim,
700  M.A.L., Herrera, E.M., Molloy, A., Stelling, J., 2021. Rapid Genomic Characterization and Global
701  Surveillance of *Klebsiella* Using Pathogenwatch. Clin. Infect. Dis. 73, S325–S335.
702  https://doi.org/10.1093/cid/ciab784

703  Bialek-Davenet, S., Criscuolo, A., Ailloud, F., Passet, V., Jones, L., Delannoy-Vieillard, A.-S., Garin,
704  B., Le Hello, S., Arlet, G., Nicolas-Chanoine, M.-H., Decré, D., Brisse, S., 2014. Genomic definition
705  of hypervirulent and multidrug-resistant Klebsiella pneumoniae clonal groups. Emerg. Infect. Dis. 20,
706  1812–1820. https://doi.org/10.3201/eid2011.140206

707  Bougnoux, M.-E., Aanensen, D.M., Morand, S., Théraud, M., Spratt, B.G., d'Enfert, C., 2004.
708  Multilocus sequence typing of Candida albicans: strategies, data exchange and applications. Infect.
709  Genet. Evol. J. Mol. Epidemiol. Evol. Genet. Infect. Dis. 4, 243–252.
710  https://doi.org/10.1016/j.meegid.2004.06.002

711  Croucher, N.J., Page, A.J., Connor, T.R., Delaney, A.J., Keane, J.A., Bentley, S.D., Parkhill, J., Harris,
712  S.R., 2015. Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome
713  sequences using Gubbins. Nucleic Acids Res. 43, e15. https://doi.org/10.1093/nar/gku1196

714  David, S., Reuter, S., Harris, S.R., Glasner, C., Feltwell, T., Argimon, S., Abudahab, K., Goater, R.,
715  Giani, T., Errico, G., Aspbury, M., Sjunnebo, S., EuSCAPE Working Group, ESGEM Study Group,
716  Feil, E.J., Rossolini, G.M., Aanensen, D.M., Grundmann, H., 2019. Epidemic of carbapenem-resistant
717  Klebsiella pneumoniae in Europe is driven by nosocomial spread. Nat. Microbiol. 4, 1919–1929.
718  https://doi.org/10.1038/s41564-019-0492-8

719  Deleo, F.R., Chen, L., Porcella, S.F., Martens, C.A., Kobayashi, S.D., Porter, A.R., Chavda, K.D.,
720  Jacobs, M.R., Mathema, B., Olsen, R.J., Bonomo, R.A., Musser, J.M., Kreiswirth, B.N., 2014.
721  Molecular dissection of the evolution of carbapenem-resistant multilocus sequence type 258 Klebsiella
722  pneumoniae. Proc. Natl. Acad. Sci. U. S. A. 111, 4988–4993.
723  https://doi.org/10.1073/pnas.1321364111

724  Feil, E.J., 2004. Small change: keeping pace with microevolution. Nat Rev Microbiol 2, 483–95.

725  Goris, J., Konstantinidis, K.T., Klappenbach, J.A., Coenye, T., Vandamme, P., Tiedje, J.M., 2007.
726  DNA-DNA hybridization values and their relationship to whole-genome sequence similarities. Int J
727  Syst Evol Microbiol 57, 81–91.

728    Hennart, M., Guglielmini, J., Bridel, S., Maiden, M.C.J., Jolley, K.A., Criscuolo, A., Brisse, S., 2022.
729    A Dual Barcoding Approach to Bacterial Strain Nomenclature: Genomic Taxonomy of Klebsiella
730    pneumoniae Strains. Mol. Biol. Evol. 39, msac135. https://doi.org/10.1093/molbev/msac135

731    Hetland, M.A.K., Hawkey, J., Bernhoff, E., Bakksjø, R.-J., Kaspersen, H., Rettedal, S.I., Sundsfjord,
732    A., Holt, K.E., Löhr, I.H., 2023. Within-patient and global evolutionary dynamics of Klebsiella
733    pneumoniae ST17. Microb. Genomics 9, mgen001005. https://doi.org/10.1099/mgen.0.001005

734    International Code of Nomenclature of Prokaryotes, 2019. . Int. J. Syst. Evol. Microbiol. 69, S1–S111.
735    https://doi.org/10.1099/ijsem.0.000778

736    Jolley, K.A., Bray, J.E., Maiden, M.C.J., 2018. Open-access bacterial population genomics: BIGSdb
737    software, the PubMLST.org website and their applications. Wellcome Open Res. 3, 124.
738    https://doi.org/10.12688/wellcomeopenres.14826.1

739    Jolley, K.A., Bray, J.E., Maiden, M.C.J., 2017. A RESTful application programming interface for the
740    PubMLST molecular typing and genome databases. Database 2017, bax060.
741    https://doi.org/10.1093/database/bax060

742    Konings, F., Perkins, M.D., Kuhn, J.H., Pallen, M.J., Alm, E.J., Archer, B.N., Barakat, A., Bedford,
743    T., Bhiman, J.N., Caly, L., Carter, L.L., Cullinane, A., de Oliveira, T., Druce, J., El Masry, I., Evans,
744    R., Gao, G.F., Gorbalenya, A.E., Hamblion, E., Herring, B.L., Hodcroft, E., Holmes, E.C., Kakkar,
745    M., Khare, S., Koopmans, M.P.G., Korber, B., Leite, J., MacCannell, D., Marklewitz, M., Maurer-
746    Stroh, S., Rico, J.A.M., Munster, V.J., Neher, R., Munnink, B.O., Pavlin, B.I., Peiris, M., Poon, L.,
747    Pybus, O., Rambaut, A., Resende, P., Subissi, L., Thiel, V., Tong, S., van der Werf, S., von Gottberg,
748    A., Ziebuhr, J., Van Kerkhove, M.D., 2021. SARS-CoV-2 Variants of Interest and Concern naming
749    scheme conducive for global discourse. Nat. Microbiol. 6, 821–823. https://doi.org/10.1038/s41564-
750    021-00932-w

751    Konstantinidis, K.T., Tiedje, J.M., 2005. Towards a genome-based taxonomy for prokaryotes. J
752    Bacteriol 187, 6258–64.

753    Lam, M.M.C., Holt, K.E., Wyres, K.L., 2023. Comment on: MDR carbapenemase-producing
754    Klebsiella pneumoniae of the hypervirulence-associated ST23 clone in Poland, 2009-19. J.
755    Antimicrob. Chemother. 78, 1132–1134. https://doi.org/10.1093/jac/dkad028

756    Lam, M.M.C., Wick, R.R., Judd, L.M., Holt, K.E., Wyres, K.L., 2022. Kaptive 2.0: updated capsule
757    and lipopolysaccharide locus typing for the Klebsiella pneumoniae species complex. Microb.
758    Genomics 8. https://doi.org/10.1099/mgen.0.000800

759    Lam, M.M.C., Wick, R.R., Watts, S.C., Cerdeira, L.T., Wyres, K.L., Holt, K.E., 2021. A genomic
760    surveillance framework and genotyping tool for Klebsiella pneumoniae and its related species
761    complex. Nat. Commun. 12, 4188. https://doi.org/10.1038/s41467-021-24448-3

762    Lam, M.M.C., Wyres, K.L., Duchêne, S., Wick, R.R., Judd, L.M., Gan, Y.-H., Hoh, C.-H., Archuleta,
763    S., Molton, J.S., Kalimuddin, S., Koh, T.H., Passet, V., Brisse, S., Holt, K.E., 2018. Population
764    genomics of hypervirulent Klebsiella pneumoniae clonal-group 23 reveals early emergence and rapid
765    global dissemination. Nat. Commun. 9, 2703. https://doi.org/10.1038/s41467-018-05114-7

766    Maiden, M.C., 2006. Multilocus sequence typing of bacteria. Annu Rev Microbiol 60, 561–88.

767    Maiden, M.C., Bygraves, J.A., Feil, E., Morelli, G., Russell, J.E., Urwin, R., Zhang, Q., Zhou, J.,
768    Zurth, K., Caugant, D.A., Feavers, I.M., Achtman, M., Spratt, B.G., 1998. Multilocus sequence typing:
769    a portable approach to the identification of clones within populations of pathogenic microorganisms.
770    Proc Natl Acad Sci U A 95, 3140–5.

771    Maiden, M.C., van Rensburg, M.J., Bray, J.E., Earle, S.G., Ford, S.A., Jolley, K.A., McCarthy, N.D.,
772    2013. MLST revisited: the gene-by-gene approach to bacterial genomics. Nat Rev Microbiol 11, 728–
773    36. https://doi.org/10.1038/nrmicro3093

774    Marakeby, H., Badr, E., Torkey, H., Song, Y., Leman, S., Monteil, C.L., Heath, L.S., Vinatzer, B.A.,
775    2014. A system to automatically classify and name any individual genome-sequenced organism
776    independently of current biological classification and nomenclature. PloS One 9, e89142.
777    https://doi.org/10.1371/journal.pone.0089142

778    Martin, M.J., Corey, B.W., Sannio, F., Hall, L.R., MacDonald, U., Jones, B.T., Mills, E.G., Harless,
779    C., Stam, J., Maybank, R., Kwak, Y., Schaufler, K., Becker, K., Hübner, N.-O., Cresti, S., Tordini, G.,
780    Valassina, M., Cusi, M.G., Bennett, J.W., Russo, T.A., McGann, P.T., Lebreton, F., Docquier, J.-D.,
781    2021. Anatomy of an extensively drug-resistant Klebsiella pneumoniae outbreak in Tuscany, Italy.
782    Proc. Natl. Acad. Sci. U. S. A. 118, e2110227118. https://doi.org/10.1073/pnas.2110227118

783    Minh, B.Q., Schmidt, H.A., Chernomor, O., Schrempf, D., Woodhams, M.D., von Haeseler, A.,
784    Lanfear, R., 2020. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the
785    Genomic Era. Mol. Biol. Evol. 37, 1530–1534. https://doi.org/10.1093/molbev/msaa015

786    Nadon, C., Van Walle, I., Gerner-Smidt, P., Campos, J., Chinen, I., Concepcion-Acevedo, J., Gilpin,
787    B., Smith, A.M., Man Kam, K., Perez, E., Trees, E., Kubota, K., Takkinen, J., Nielsen, E.M., Carleton,
788    H., FWD-NEXT Expert Panel, 2017. PulseNet International: Vision for the implementation of whole
789    genome sequencing (WGS) for global food-borne disease surveillance. Euro Surveill. Bull. Eur. Sur
790    Mal. Transm. Eur. Commun. Dis. Bull. 22. https://doi.org/10.2807/1560-7917.ES.2017.22.23.30544

791    Paradis, E., Claude, J., Strimmer, K., 2004. APE: Analyses of Phylogenetics and Evolution in R
792    language. Bioinformatics 20, 289–290. https://doi.org/10.1093/bioinformatics/btg412

793    Prim, R.C., 1957. Shortest connection networks and some generalizations. Bell Syst Tech J 1389–
794    1401.

795    Rambaut, A., Holmes, E.C., O'Toole, Á., Hill, V., McCrone, J.T., Ruis, C., du Plessis, L., Pybus,
796    O.G., 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic
797    epidemiology. Nat. Microbiol. 5, 1403–1407. https://doi.org/10.1038/s41564-020-0770-5

798    Revell, L.J., 2024. phytools 2.0: an updated R ecosystem for phylogenetic comparative methods (and
799    other things). PeerJ 12, e16505. https://doi.org/10.7717/peerj.16505

800    Rodrigues, C., Desai, S., Passet, V., Gajjar, D., Brisse, S., 2022. Genomic evolution of the globally
801    disseminated multidrug-resistant Klebsiella pneumoniae clonal group 147. Microb. Genomics 8.
802    https://doi.org/10.1099/mgen.0.000737

803    Stamatakis, A., 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
804    phylogenies. Bioinformatics 30, 1312–1313. https://doi.org/10.1093/bioinformatics/btu033

805    Struelens, M.J., Brisse, S., 2013. From molecular to genomic epidemiology: transforming surveillance
806    and control of infectious diseases. Euro Surveill 18, 20386.

807    Struelens, M.J., De Gheldre, Y., Deplano, A., 1998. Comparative and library epidemiological typing
808    systems: outbreak investigations versus surveillance systems. Infect Control Hosp Epidemiol 19, 565–
809    9.

810    van Belkum, A., Tassios, P.T., Dijkshoorn, L., Haeggman, S., Cookson, B., Fry, N.K., Fussing, V.,
811    Green, J., Feil, E., Gerner-Smidt, P., Brisse, S., Struelens, M., 2007. Guidelines for the validation and
812    application of typing methods for use in bacterial epidemiology. Clin Microbiol Infect 13 Suppl 3, 1–
813    46.

814    Vinatzer, B.A., Tian, L., Heath, L.S., 2017. A proposal for a portal to make earth's microbial diversity
815    easily accessible and searchable. Antonie Van Leeuwenhoek 110, 1271–1279.
816    https://doi.org/10.1007/s10482-017-0849-z

817    Wang, J., Feng, Y., Zong, Z., 2023. The Origins of ST11 KL64 Klebsiella pneumoniae: a Genome-
818    Based Study. Microbiol. Spectr. 11, e0416522. https://doi.org/10.1128/spectrum.04165-22

819    Wyres, K.L., Hawkey, J., Hetland, M.A.K., Fostervold, A., Wick, R.R., Judd, L.M., Hamidian, M.,
820    Howden, B.P., Löhr, I.H., Holt, K.E., 2019. Emergence and rapid global dissemination of CTX-M-15-
821    associated Klebsiella pneumoniae strain ST307. J. Antimicrob. Chemother. 74, 577–581.
822    https://doi.org/10.1093/jac/dky492

823    Wyres, K.L., Lam, M.M.C., Holt, K.E., 2020. Population genomics of Klebsiella pneumoniae. Nat.
824    Rev. Microbiol. 18, 344–359. https://doi.org/10.1038/s41579-019-0315-1

825    Wyres, K.L., Wick, R.R., Gorrie, C., Jenney, A., Follador, R., Thomson, N.R., Holt, K.E., 2016.
826    Identification of Klebsiella capsule synthesis loci from whole genome data. Microb. Genomics 2,
827    e000102. https://doi.org/10.1099/mgen.0.000102

828    Yeo, M., Mauricio, I.L., Messenger, L.A., Lewis, M.D., Llewellyn, M.S., Acosta, N., Bhattacharyya,
829    T., Diosque, P., Carrasco, H.J., Miles, M.A., 2011. Multilocus sequence typing (MLST) for lineage
830    assignment and high resolution diversity studies in Trypanosoma cruzi. PLoS Negl. Trop. Dis. 5,
831    e1049. https://doi.org/10.1371/journal.pntd.0001049

832    Yu, G., Smith, D.K., Zhu, H., Guan, Y., Lam, T.T.-Y., 2017. ggtree: an r package for visualization and
833    annotation of phylogenetic trees with their covariates and other associated data. Methods Ecol. Evol.
834    8, 28–36. https://doi.org/10.1111/2041-210X.12628

835    Zhou, K., Xiao, T., David, S., Wang, Q., Zhou, Y., Guo, L., Aanensen, D., Holt, K.E., Thomson, N.R.,
836    Grundmann, H., Shen, P., Xiao, Y., 2020. Novel Subclone of Carbapenem-Resistant Klebsiella
837    pneumoniae Sequence Type 11 with Enhanced Virulence and Transmissibility, China. Emerg. Infect.
838    Dis. 26, 289–297. https://doi.org/10.3201/eid2602.190594

839    Zhou, K., Xue, C.-X., Xu, T., Shen, P., Wei, S., Wyres, K.L., Lam, M.M.C., Liu, J., Lin, H., Chen, Y.,
840    Holt, K.E., BRICS Working Group, Xiao, Y., 2023. A point mutation in recC associated with
841    subclonal replacement of carbapenem-resistant Klebsiella pneumoniae ST11 in China. Nat. Commun.
842    14, 2464. https://doi.org/10.1038/s41467-023-38061-z

843    Zhou, Z., Charlesworth, J., Achtman, M., 2021. HierCC: A multi-level clustering scheme for
844    population assignments based on core genome MLST. Bioinforma. Oxf. Engl. btab234.
845    https://doi.org/10.1093/bioinformatics/btab234

846