

1 **An update to the database for *Acinetobacter baumannii* capsular**
2 **polysaccharide locus typing extends the extensive and diverse repertoire**
3 **of genes found at and outside the K locus**

4
5 **1.1 Author names**

6 Sarah M. Cahill,¹ Ruth M. Hall,² Johanna J. Kenyon^{1*}

7
8 **1.2 Affiliation**

9 ¹ *Centre for Immunology and Infection Control, School of Biomedical Sciences, Faculty of*
10 *Health, Queensland University of Technology, Brisbane, Australia*

11 ² *School of Life and Environmental Sciences, The University of Sydney, Sydney, Australia*

12
13 **1.3 Corresponding author**

14 * email: johanna.kenyon@qut.edu.au; Ph: +61 7 3138 2552

15
16 **1.4 Keyword**

17 *Acinetobacter baumannii*; Kaptive; capsular polysaccharide; K locus; KL

18
19
20
21
22
23
24
25
26
27
28

29 **2. Abstract**

30 Several novel non-antibiotic therapeutics for the critical priority bacterial pathogen,
31 *Acinetobacter baumannii*, rely on specificity to the cell-surface capsular polysaccharide
32 (CPS). Hence, prediction of CPS type deduced from genes in whole genome sequence data
33 underpins the development and application of these therapies. In this study, we provide a
34 comprehensive update to the *A. baumannii* K locus reference sequence database for CPS
35 typing (available in *Kaptive v. 2.0.1*) to include 145 new KL, providing a total of 237 KL
36 reference sequences. The database was also reconfigured for compatibility with the updated
37 *Kaptive v. 2.0.0* code that enables prediction of ‘K type’ from special logic parameters
38 defined by detected combinations of KL and additional genes outside the K locus. Validation
39 of the database against 8994 publicly available *A. baumannii* genome assemblies from NCBI
40 databases identified the specific KL in 73.45% of genomes with perfect, very high or high
41 confidence. Poor sequence quality or the presence of insertion sequences were the main
42 reasons for lower confidence levels. Overall, 17 KL were overrepresented in available
43 genomes, with KL2 the most common followed by the related KL3 and KL22. Substantial
44 variation in gene content of the central portion of the K locus, that usually includes genes
45 specific to the CPS type, included 34 distinct groups of genes for synthesis of various
46 complex sugars and >400 genes for forming linkages between sugars or adding non-sugar
47 substituents. A repertoire of 681 gene types were found across the 237 KL, with 88.4% found
48 in <5% of KL.

49

50 **3. Significance as a BioResource to the community**

51 New therapies that target the bacterial polysaccharide capsule (CPS) show promise as
52 effective tools to curb the high mortality rates associated with extensively resistant *A.*
53 *baumannii*; one of the world’s most troublesome Gram-negative pathogens. As important
54 information about the CPS structure produced by an isolate can be extracted from Whole
55 Genome Sequences (WGS), simple bioinformatic tools and definitive sequence databases are
56 needed to facilitate robust prediction of CPS type from WGS data. Here, we provide a
57 comprehensive update to the international CPS sequence typing database for *A. baumannii*,
58 increasing the utility of this resource for prediction of CPS type from WGS to assist with
59 clinical surveillance, and/or the design and application of CPS-targeted therapies. This study
60 is expected to further inform epidemiological tracking efforts, as well as the design of

61 therapeutics targeting the CPS, enhancing global efforts to identify, trace and treat infections
62 caused by this pathogen.

63 **4. Data summary**

- 64 1. The updated *A. baumannii* KL reference sequence database including 241 fully
65 annotated gene clusters is available for download under *Kaptive v. 2.0.1* at
66 <https://github.com/katholt/Kaptive>.
- 67 2. Genome assemblies, short read data, or GenBank records used as representative
68 reference sequence for each K locus are listed in Supplementary Table S1, and are
69 referenced within each entry in the *A. baumannii* KL reference sequence database.

70

71 **The authors confirm all supporting data, code and protocols have been provided within**
72 **the article or through supplementary data files.**

73

74 **5. Introduction**

75 Failure of antibiotic therapy due to the emergence of pan-resistant bacteria is a growing
76 global health crisis. *Acinetobacter baumannii* is ranked as one of six leading bacterial species
77 responsible for nearly three quarters of deaths associated with antibiotic resistance
78 worldwide, with an estimated 80% of circulating isolates in many low- and middle-income
79 countries resistant to last-line carbapenems [1]. Hence, new therapeutic options are urgently
80 needed for treatment of carbapenem-resistant *A. baumannii*. Promising strategies include
81 monoclonal antibodies or bacteriophage [2]. Both strategies involve binding to cells via
82 interaction with exposed structures on the bacterial cell surface, and can display specificity
83 for structural epitopes of the polysaccharide capsule (CPS). However, in this species, even
84 closely related isolates can produce different forms of CPS, making knowledge of the specific
85 CPS type in the infection to be treated critical. Hence, the ability to determine CPS type is
86 needed to underpin the design and application of these therapies. The genetics underlying the
87 CPS type has also proven valuable as an epidemiological marker [3-7]. Finally, recent studies
88 have associated some specific CPS types [8, 9] or alterations in the CPS structure [10] with
89 increased virulence. Hence, the determination of the specific type produced in problem
90 strains is important in several areas.

91 As information about CPS type can be deduced from genes in bacterial genomes,
92 whole genome sequencing (WGS) is an attractive approach for CPS typing that is more
93 readily accessible than traditional laboratory-based serological typing methods. For *A.*
94 *baumannii*, most of the genes responsible for CPS biosynthesis are clustered together in the
95 chromosome between *fkpA* and *lldP* genes [11]. However, many different sets of genes have
96 been found at this 'K locus' (KL). To facilitate their identification, 92 fully annotated KL
97 reference sequences were recently compiled into a curated database and released publicly
98 [12]. The database is compatible with the bioinformatics search tool, *Kaptive* [13] and
99 *Kaptive-Web* [14]. This database was validated against 3415 genome sequences available in
100 the NCBI non-redundant and WGS databases at that time and 642 genomes assembled from
101 reads available in NCBI SRA database. However, it was noted that additional KL
102 configurations were known, and that there may be many more KL yet to be documented [12].
103 In fact, more than 128 distinct K loci were known at the time and an additional 78 KL have
104 since been identified as additional sequence data became available [15, 16, 17 and Kenyon,
105 unpublished data].

106 All known gene clusters at the *A. baumannii* K locus follow a general pattern that
107 includes 3 'regions' [11, 12]. Region 1 always consists of essential CPS export genes (*wza*,
108 *wzb* and *wzc*) that are transcribed in the opposite direction to the remainder of the locus.
109 Region 2, the central portion, includes many different sets of genes and these determine the
110 composition and structure of the K unit making up the specific CPS type. Region 3 flanks the
111 other side of Region 2, and always includes genes for the synthesis of simple sugar precursors
112 (*galU*, *ugd*, *gpi*, *pgm*), though genes can be variably inserted between *gpi* and *pgm*. The *gneI*
113 gene for D-Galp or D-GalpNAc synthesis is often present between *gpi* and *pgm* [11] but has
114 been found to be absent from some KL that do not include D-Galp or D-GalpNAc in the
115 corresponding CPS [17-22]. Other genes can also be found in this position giving rise to
116 some variation in Region 3 [19, 23-26].

117 In general, a unique KL identifier is assigned to a sequence when there is a detectable
118 difference in gene content with genes identified based on a product sequence identity cut-off
119 of 85%. However, as chemical structures of CPS produced by 70 distinct KL have been
120 determined [e.g. 17, 19, 20, 27-34], it is now known that differences in the 'conserved' genes
121 in Region 3 (i.e. *galU*, *ugd*, *gpi*, *pgm*) do not influence the type of CPS produced. In addition,
122 variation in the genes in Region 1 (*wza*, *wzb*, *wzc*) does not affect their essential role in
123 capsule export. Therefore, a new KL number is only assigned to a sequence when there is a

124 detectable difference in genes in Region 2 and/or the variable portion of Region 3 (between
125 *gpi* and *pgm*).

126 In most cases, complete correlation between the genetic content of specific KL and
127 the structural features of the corresponding CPS type have been reported. However, for a few
128 strains, the *wzy* gene has been shown to be missing from Region 2 of the K locus or
129 interrupted by an insertion sequence (IS), and a replacement *wzy* gene was found to the left of
130 Region 1 as in KL8 [11] or in a defined genomic island outside of the K locus [15, 35, 36]. In
131 a recent study, an additional Wzy polymerase gene was identified in prophage sequence
132 integrated elsewhere in the chromosome, and was found to alter the linkage between
133 oligosaccharide K-units that make up the CPS structure [37]. Acetyltransferase genes with
134 encoded products that have been shown to modify the CPS by acetylation have also been
135 found in integrated phage genomes [38]. Therefore, detection of these additional genetic
136 determinants in the genome will be essential to achieve robust prediction of CPS type from
137 WGS data.

138 Recently, the *Kaptive* code was updated (*Kaptive v. 2.0.0*) to include an additional
139 function that was designed specifically for discrimination of O-antigen serotypes in the
140 *Klebsiella pneumoniae* species complex [39]. For this function, determination of serotype or
141 'type' is based on the detection of either the O-antigen locus (OL) type alone or a defined
142 combination of OL and 'extra genes' in a genome assembly known to be involved in the
143 determination of a specific serotype. As an active CPS serotyping scheme does not exist for
144 *A. baumannii*, 'type' has previously been used to refer to the chemical structure of the CPS
145 produced by an isolate as defined by the KL number i.e. the KL2 sequence produces the K2
146 type CPS [32]. In cases where genes outside the K locus have been shown to modify the CPS
147 type, a suffix is now added to the K type name. For example, this was recently done for
148 K127-Wzy_{Ph}, which is defined as the structure formed by KL127 and Wzy encoded in
149 prophage (Ph) [37]. Other examples, such as K19 and K24 modified by Wzy proteins
150 encoded by genomic islands, GI-1 and GI-2, respectively, have also now been renamed to
151 indicate the role of extra-KL genes. Hence, the new *Kaptive v. 2.0.0* function can be
152 harnessed to predict *A. baumannii* CPS 'type' as defined by structural data where this is
153 available.

154 In this study, we provide a comprehensive update to the *A. baumannii* CPS reference
155 sequence database to include the known KL not included in the original version and new KL
156 sequences detected in 8994 publicly available *A. baumannii* genome sequences. Special logic
157 parameters to enable prediction of the CPS type based on KL or the detected combination of

158 a specific KL with 'extra genes' have also been included as well as information relating to K
159 type where structures have been determined. The updated database was validated against the
160 same large genome set and a smaller set of complete genomes. A detailed assessment of gene
161 repertoire at the chromosomal K locus was also conducted.

162

163 **6. Methods**

164 ***A. baumannii* genome assemblies**

165 A total of 9065 genome assemblies listed under the *Acinetobacter baumannii* taxonomic
166 classification in the NCBI non-redundant and WGS databases (10th June, 2021) were
167 downloaded for local analysis. Assemblies were first assessed for the presence of the *A.*
168 *baumannii*-specific *oxaAb* gene (also known as *bla*_{OXA-51-like}; available in *A. baumannii* strain
169 A1 complete genome sequence under GenBank accession number CP010781.1, base
170 positions 1753305 to 1754129) with BLASTn using a cut-off of >90% combined coverage
171 with >95% nucleotide sequence identity to confirm the *A. baumannii* species assignment (as
172 previously defined in [12]). Only confirmed *oxaAb*-positive genomes were used for
173 downstream analyses.

174

175 **Identification and annotation of novel K locus sequences**

176 *A. baumannii* genome assemblies (n = 8994) were screened against the original version of the
177 *A. baumannii* KL reference database [12] available in the *Kaptive versions 0.7.0-2.0.0*
178 (<https://github.com/katholt/Kaptive>) and then using an extended in-house database of 206 KL
179 (*unpublished*) with the command-line version of *Kaptive v 0.7.0* [13]. The search was
180 conducted using a parameter defining the “minimum gene identity” cut-off as 85% as is
181 standard for *A. baumannii* KL typing [11]. Output results from the in-house screen for
182 matches with a reported confidence level less than ‘perfect’ were examined. Matches with
183 length discrepancies, additional or missing genes, or those with <95% coverage and/or <95%
184 nucleotide sequence identity to the best matched reference sequence were manually inspected
185 to identify novel gene clusters.

186 Novel KL were annotated using the established nomenclature system for *A.*
187 *baumannii* K locus typing [11,12]. Briefly, KL were assigned a new number if any
188 differences were detected in gene presence/absence in Region 2 and between *gpi* and *pgm*
189 genes in Region 3, and standard gene names were used to indicate enzyme function. For

190 glycosyltransferase (*gtr*), acetyl or acyltransferase (*atr*), and pyruvyltransferase (*ptr*) genes
191 where sequence differences may result in a change of substrate preference, new numbers
192 were assigned when the product of the gene had <85% amino acid (aa) sequence identity to
193 the closest match.

194 KL comparisons were generated using EasyFig v 2.2.2 [40], and genome comparisons
195 assembled using Mauve v 2.4.0 [41] to order contigs, followed by BRIG [42] to generate a
196 circular comparison. Where necessary, read quality was assessed using FASTQC v 0.11.9
197 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) and assembly quality examined
198 using QUAST v. 5.02 [43]. For some cases, multi locus sequence typing (MLST) was
199 performed on the assembly using both Oxford and Institute Pasteur schemes established for
200 *A. baumannii* in PubMLST with the MLST package (<https://github.com/tseemann/mlst>).

201

202 **Curation of updated *A. baumannii* KL reference sequence database**

203 A GenBank format file (.gbk) for each new locus sequence was prepared, which included the
204 nucleotide sequence and annotations of all coding sequences in the locus. Where the only
205 available representative of a KL included an insertion sequence (IS), we substituted the
206 sequence with a manually generated version where the IS and target site duplication were
207 removed in order to include a KL that represents the presumptive ancestral, non-modified
208 sequence as is required for accurate typing by *Kaptive*. This was the case for 23 different KL,
209 which are indicated in Supplementary Table S1. An additional note field was added for all
210 reference loci to define the K type where structural data for the CPS was available for the
211 specific KL sequence as indicated in the reference record in each entry in the database. In
212 cases where no structural data is available, the note field specifies the type as unknown.

213 For CPS structures known to be modified by additional genes outside the K locus, the
214 note field indicates ‘special logic’ to be applied by *Kaptive*. This directs the tool to perform
215 an additional tBLASTn search for ‘extra genes’ supplied in the database, and an additional
216 “*Acinetobacter_baumannii_k_locus_primary_reference.logic*” file then specifies ‘type’ when
217 a specific combination of KL and extra genes is found. GenBank records for six ‘extra genes’
218 added to the database include: *wzy* and *atr25* genes in ‘genomic island 1’ (GI-1) involved in
219 K19 synthesis [35]; *wzy* in ‘genomic island 2’ (GI-2) for K24 synthesis [36]; *atr29* and *atr30*
220 in prophage (Ph) found to modify the K46 and K5 structures, respectively [38]; and *wzy* in
221 prophage that has recently been found to modify the K127 CPS [37].

222 All GenBank-format .gbk files were concatenated into a multi-record file to produce
223 an updated KL reference database for release as *Kaptive v. 2.0.1*. The updated database was

224 integrated with the *Kaptive-Web* platform (<http://kaptive.holtlab.net/>), and is available for
225 download from <https://github.com/katholt/Kaptive> for use with the command-line *Kaptive v.*
226 *2.0.0*. The database was validated on the same genome pool using *Kaptive v. 2.0.0* with a
227 parameter defining the minimum gene product identity cut-off as 85%.

228

229 **Analysis of sequence features and gene frequency**

230 To generate an overview of the sequence lengths and total repertoire of genes found at the *A.*
231 *baumannii* K locus, Prokka *v. 1.13* [44] was used to generate gff3 files for each individual K
232 locus sequence using manual annotations available from each gbk record. The complete
233 length of each gene cluster, along with the total number of open reading frames, and number
234 of *gtr* genes was manually tabulated. The gff3 files were then used as an input for the pan
235 genome tool Roary *v. 3.13.0* [45] to generate a gene presence/absence matrix using a cut-off
236 parameter of 85% aa sequence identity. The matrix was used to determined homology groups
237 defined as genes encoding products with >85% aa sequence identity and >330 bp in length.
238 Summarised data was visualised using the ggplot2 package in RStudio *v. 1.2.5033* [46].

239

240 **7. Results**

241 **Screening for novel CPS biosynthesis gene clusters at the K locus**

242 A total of 9,065 genome assemblies available under the *Acinetobacter baumannii* taxonomy
243 classification (Taxonomy ID: 470) were downloaded from NCBI GenBank and WGS
244 databases. The intrinsic *oxaAb* gene could not be identified in 71 assemblies, hence these
245 were excluded from further analysis. Confirmed *A. baumannii* genome assemblies (n=8994)
246 were screened against the original *A. baumannii* KL reference sequence database [12]
247 included in *Kaptive* versions (*v*) *0.7.0-2.0.0* (hereafter referred to as database *v 0.7.0-2.0.0*).
248 The confidence levels (categorical measure of match quality) called by *Kaptive v 0.7.0* using
249 this database were: 794 (perfect), 4794 (very high), 443 (high), 1776 (good), 192 (low) and
250 995 (none) (Supplementary Table S2; summarised in Table 1). This revealed that 62.13% of
251 genome assemblies could be confidently assigned a match indicated by a confidence level of
252 ‘Perfect’ (the identified locus is in a single contiguous sequence that shares 100% coverage
253 and 100% nucleotide sequence identity) or ‘Very high’ (a single contiguous sequence sharing
254 $\geq 99\%$ coverage and $\geq 95\%$ nucleotide sequence identity with the best match reference
255 sequence with no additional and/or missing coding sequences).

256 Since additional KL reference sequences have been characterised following the
257 release of the first database in 2020, the genome pool was reanalysed using an extended in-
258 house database that includes a total of 206 KL made up of the 92 KL reference sequences in
259 the original database plus the 36 previously characterised KL that were not included in the
260 original database, and an additional 78 distinct KL characterised since this time [15, 16, 17
261 and Kenyon, unpublished data]. Confidence levels obtained were: 922 (perfect), 5124 (very
262 high), 449 (high), 1726 (good), 135 (low) and 638 (none) (Table 1; Supplementary Table S3).
263 This secondary screen revealed a shift in the number of assemblies in each confidence level
264 with the proportion of matches scored ‘perfect’ or ‘very high’ confidence rising slightly to
265 67.22%. Given that a match with ‘perfect’ confidence indicates identity to the reference
266 sequence, only matches assigned with a confidence level of ‘very high’ or less were further
267 examined to identify novel locus sequences.

268

269 **‘Very high’ confidence matches are close relatives or IS variants of reference sequences**

270 Manual inspection of the output data from the screen using the in-house database of 206 KL
271 revealed that 5124 assemblies had a match assigned with ‘very high’ confidence. Of these
272 5124 assignments, 5094 (99.4%) were considered very close relatives of the best match locus
273 with single nucleotide polymorphisms (SNPs). However, 30 (0.6%) had a discrepancy in the
274 total length of the locus match with >700 bp of additional sequence. For 28 of these 30
275 assignments, the additional sequence was found to be an insertion sequence (IS), and thus
276 these loci were deemed variants of the archetypal KL reference sequence in the database. The
277 remaining two KL with discrepant lengths had variable-length insertions of ‘N’ bases, which
278 indicated possible issues with sequence or assembly quality.

279

280 **‘High’ confidence matches include IS variants and novel KL sequences**

281 A total of 449 assemblies were assigned a match with a confidence level of ‘high’ indicating
282 that the KL was found in a single piece with $\geq 99\%$ coverage but less than three missing gene
283 products and no extra genes. For 434 of 449 assignments (96.66%), the locus sequence had
284 been correctly identified but with detectable problems. For 415 of these, the locus match had
285 $\geq 99\%$ coverage, $\geq 89\%$ nucleotide sequence identity and no more than ± 101 bp of
286 discrepancy in sequence length to a KL reference sequence in the database, though SNPs or
287 base insertions/deletions resulting in frameshifts in known coding sequences were found. For
288 the other 19 assignments, significant discrepancies in the total sequence length of the match
289 (greater than ± 101 bp) were reported by *Kaptive*. For 15 of these, the additional sequence

290 was confirmed to be one or more IS insertions. Three others had strings of missing bases,
291 whereas one had a string of additional 'N' bases. Hence, 19 were considered variants or
292 possible variants of the best match KL reference sequence. Therefore, the specific KL had
293 been correctly identified in 72.05% of genomes assigned with either 'perfect', 'very high' and
294 'high' matches.

295 For 11 of the remaining 15 'high' assemblies, nucleotide sequence identity to a best
296 match locus of KL33 was <94%. However, the expected *psaD* and *psaE* genes that encode a
297 nucleotidase and an acetyltransferase involved in the synthesis of 5,7-di-*N*-acetylpsseudaminic
298 acid, respectively, were reported missing by *Kaptive*. Analysis of these genome assemblies
299 revealed that they all carried the same sequence at the K locus, which differed from KL33
300 only in a small segment where the *psaD* and *psaE* genes of KL33 are replaced by two related
301 but novel genes, designated here as *psaI* and *psaJ* (Fig. 1A). Both the encoded PsaI and PsaJ
302 products share 78.9% aa sequence identity to their PsaD and PsaE homologues. Therefore,
303 the KL sequence was considered novel and designated KL235. The *psaI* and *psaJ* genes were
304 also identified in place of *psaD/psaE* in a further 2 of the 15 assemblies, which both had
305 >95% identity to KL121. These were also considered novel and designated KL218 (Fig. 1A).
306 Though PsaI and PsaJ may produce the same sugar product as PsaD and PsaE, it is possible
307 that the difference in sequence could result in a new acylated derivative of pseudaminic acid,
308 and structural studies will be needed to confirm this.

309 The final two assemblies (NCBI assembly accession numbers GCA_005280695.1 and
310 GCA_013305465.1) had been assigned a best match locus of KL12 with 100% sequence
311 coverage and 94.7% nucleotide sequence identity by *Kaptive*. However, the output indicated
312 an additional 1083 bp of sequence, and the expected gene coding for the Wzy polymerase
313 was missing. Manual inspection of the two genome assemblies revealed that they carried the
314 same K locus sequence, and direct comparison of these DNA sequences to KL12 (Fig. 1B)
315 revealed that the *wzy_{KL12}* gene was present but interrupted by an *IS_{Aba125}* insertion
316 sequence. One of these assemblies, GCA_013305465.1, had been reported in a clinical isolate
317 from Australia [47], and an additional gene sharing 100% identity with *wzy* from KL183 in
318 the database was identified in the *Kaptive* output field, 'Other genes outside locus'. However,
319 this *wzy* gene is in fact in the locus (i.e. between *fkpA* and *lldP*) but in an unusual position
320 between *fkpA* and *wzc* at the 5'-end of the locus (Fig. 1B). The location of *wzy* at the
321 beginning of the locus adjacent to *fkpA* was previously reported for KL8 [11], and the *wzy*
322 gene from KL183 is identical to the KL8 *wzy*. A similar configuration was also found for

323 KL217 (see below). Hence, this gene was assigned the name wzy_{KL8} , and KL234 was assigned
324 to this novel locus.

325

326 **‘Good’, ‘Low’ or ‘None’ confidence matches**

327 A total of 2499 genome assemblies were assigned a match with a ‘good’ (1726), ‘low’ (135)
328 or ‘none’ (638) confidence level. These included only 231 (9.24%) with a best match locus
329 found in a single contiguous sequence, and 2268 genomes (90.75%) found in two or more
330 pieces (indicated by a ‘?’ problem score in the *Kaptive* output). As detectable breaks in KL
331 loci often suggest that a genome sequence or an assembly is poor quality or that loci are
332 variants in which an IS has interrupted the KL sequence, the 2268 assemblies with loci found
333 in more than 1 piece were not further investigated. For the 231 contiguous sequences, 135
334 loci (76 ‘good’, 24 ‘low’ and 35 ‘none’) were found to include numerous SNPs relative to the
335 assigned reference sequence or insertions of ‘N’ bases suggesting problems with sequence or
336 assembly quality. A further 46 ‘good’ matches included an IS indicating these were variants
337 of the best match reference sequence.

338 Of the 50 remaining contiguous matches, four genome assemblies were found to be
339 missing significant portions of the K locus sequence. One of these had a match confidence of
340 ‘good’ and was missing 13% of the assigned KL124 locus sequence, while a second locus
341 had a match confidence of ‘none’ and was missing 42% of the assigned KL13 locus. Another
342 two genome assemblies (GCA_001862175.1 and GCA_001862305.1) were assigned a best
343 match to KL92 with a ‘none’ confidence level and 0 of 22 expected genes identified. Manual
344 inspection of the two genome assemblies and comparison to the complete genome sequence
345 of a related strain revealed that both assemblies had a ~150 kb deletion that included the K
346 locus (Fig. 1C). The associated NCBI assembly data indicated these were clinical isolates
347 sequenced using Illumina Hiseq 2000 and assembled using CLC Genomic workbench v.
348 8.5.1. To assess if the deletion may be due to poor assembly or read quality, genome
349 assemblies were subjected to QUAST, and their short reads (SRR3381523 and SRR3381529)
350 to FastQC. Results outputs suggested good read and assembly quality (<50 contigs, length=
351 3.6-3.7 Mbp, 38.9% GC), and these genome assemblies were not further investigated. These
352 are surprising findings that arise from the fact that the current version of *Kaptive* still assigns
353 a best match KL even when the sequence is not present, and this will need to be addressed in
354 a future update to the *Kaptive* code.

355 The remaining 46 genome assemblies (29 ‘good’, 2 ‘low’, 15 ‘none’) were found to have
356 <95% DNA sequence coverage, <95% DNA sequence identity, significant length

357 discrepancies (>400 bp), missing expected genes and/or presence of unexpected genes in the
358 locus sequence. Amongst the 46 assemblies, 28 novel KLs were identified by manual
359 inspection. 27 of these KL were found to follow the same general pattern as for other gene
360 clusters described at the *A. baumannii* K locus to date, in that they consisted of three defined
361 regions with one *wzx* gene and one *wzy* gene in Region 2 of each gene cluster. The exception,
362 KL217, included the *wzy_{KL8}* gene in the location at the start of the K locus as described for
363 KL234 (see Fig. 1A). Therefore, together with KL218, KL234 and KL235 described above, a
364 total of 31 novel KL were identified amongst the 8994 genomes studied, bringing the total
365 number of known KL to 237.

366

367 **Annotation of novel genes**

368 Annotations were manually curated in accordance to the standard nomenclature system for *A.*
369 *baumannii* CPS biosynthesis genes [11,12] for 145 KL, which included the novel 31 KL
370 detected above, as well as the 114 not included in the previous version of the database or
371 identified since the first release. Several novel genes and gene modules were identified across
372 the 145 types and are described in further detail below.

373 Amongst the 145 additional KLs to be included in the new iteration of the database, a
374 total of 75 genes were predicted to encode novel glycosyltransferases (defined as <85% aa
375 identity to known types) not seen in the previous database. The products of three of these
376 were found to be homologues of glycosyltransferases previously annotated as KpsS1 and
377 KpsS2, and hence the genes were named *kpsS3-kpsS5* consistent with the nomenclature used
378 previously for this Gtr type [11, 32, 48]. All other predicted glycosyltransferases were
379 assigned new *gtr* numbers. Similarly, 19 genes were predicted to encode new acetyl-/acyl-
380 transferases and were assigned new *atr* names, while four new putative pyruvyltransferase
381 genes were found and assigned *ptr* names. In addition, 12 genes of unknown function (*orf*)
382 were also identified and further work will be needed to determine if these play a role in CPS
383 biosynthesis.

384 Several novel genes likely to be involved in the synthesis of a monosaccharide were
385 also found. For 6 KL (KL62, KL79, KL97, KL110, KL183, KL192), a homologue of the
386 *elaA* gene, designated *elaA2*, was identified adjacent to *elaBC* and in place of *elaA*. *ElaA* is a
387 putative oxidoreductase involved in the synthesis of 8-epilegionaminic acid (8eLeg) in the
388 K49 CPS structure [31]. As *ElaA* and *ElaA2* share 84% aa sequence identity, they likely
389 catalyse the same reaction. However, structural studies of the CPS produced by these 6 loci
390 will be needed to assess if *ElaA2* also produces 8eLeg or a related sugar. Other additional

391 homologues of sugar synthesis genes were identified for *mnaA* and *dmaA*, and these genes
392 are predicted to be involved in the synthesis of UDP-D-ManpNAc and UDP-2,3-diacetamido-
393 2,3-dideoxy-D-mannuronic acid, respectively. Gene homologues (encode proteins <85%
394 identical) were assigned numbers (*mnaA1-mnaA4* and *dmaA1-dmaA4*) and further structural
395 studies of the CPS will also be needed to confirm the type of sugar(s) produced.

396 Two genes located adjacent to each other in seven different KL (KL126, KL207,
397 KL208, KL209, KL219, KL228, KL236; Fig. 2A) were found to encode homologues (>85%
398 aa identity) of WeeE and WeeF from *Acinetobacter venetianus* RAG-1 ‘emulsan’ gene
399 cluster (GenBank accession number AJ243431.1). WeeE and WeeF have previously been
400 postulated to be involved in the synthesis of UDP-N-acyl-L-galactosaminuronic acid (UDP-L-
401 GalpNAcA), which is present in the RAG-1 CPS [49]. As in the nomenclature system for *A.*
402 *baumannii* CPS biosynthesis gene clusters are designated names after the putative sugar
403 product or function of the encoded enzyme in the synthesis pathway [11], the two genes were
404 named *gnlA* and *gnlB* for UDP-L-GalpNAcA, rather than *weeE* and *weeF*. The *gnlB* gene was
405 found without *gnlA* in an eighth gene cluster, KL215.

406 Three novel genes found in KL166 and KL224 (Fig. 2B) encode proteins sharing 33-
407 66% aa identity with WeiS, WeiP, and WeiQ from the *Escherichia coli* O109 O-antigen gene
408 cluster (GenBank accession number HM485572.1). WeiS, WeiP, and WeiQ have been
409 predicted to be involved in the synthesis of 2,3-diacetamido-2,3,6-trideoxy-L-mannose (L-
410 RhaNAc3NAc) found in the O109 structure [50]. As for *gnlAB*, the three genes were assigned
411 new names, *rhnA*, *rhnB* and *rhnC* for L-RhaN. Hence, the *gnl* and *rhn* names were added to
412 the nomenclature scheme for *A. baumannii* CPS biosynthesis genes and an updated list of
413 name descriptors can be found in Table 2.

414

415 **Updating the *A. baumannii* KL reference sequence database**

416 A representative of each novel KL sequence was chosen for addition to the reference
417 database. Fully annotated locus sequences were curated into GenBank format files (.gbk),
418 then added to the multi-record database to create a new iteration that includes a total of 237
419 KL. Two reference sequences, KL38 and KL78, were retained in the updated database
420 although the isolates carrying them have since been reassigned to other *Acinetobacter* species
421 and are no longer considered *A. baumannii* isolates. The KL38 and KL78 representative
422 reference sequences will be replaced with an *A. baumannii* sequence if one is later identified
423 in the species. Information on the representative sequences chosen for the database are
424 available in Supplementary Table S1.

425 As the latest release of the *Kaptive* code (v. 2.0.0) includes a new function that
426 enables ‘type’ to be inferred from locus sequence, an additional note field was added for all
427 reference loci to define the ‘type’ of CPS structure produced by a given KL. Citations for
428 associated structural data were also integrated into the database in the reference section of
429 each record, so that users are referred to the relevant publication(s). Where a structure for a
430 specific KL has not yet been determined, the integrated note field defines the ‘type’ as
431 unknown.

432 For CPS that have been found to require or be modified by additional genes located
433 outside the K locus, the added note indicates that ‘special logic’ needs to be applied by
434 *Kaptive*. This directs the tool to perform an additional tBLASTn search for ‘extra genes’
435 supplied in the database file. If detected, an additional
436 “*Acinetobacter_baumannii_k_locus_primary_reference.logic*” file then specifies the ‘type’
437 when a specific combination of KL and extra genes are found. This feature will be further
438 developed in future updates.

439

440 **Validation of the updated *A. baumannii* KL reference sequence database**

441 To assess the utility of the updated KL reference sequence database (hereafter referred to as
442 v. 2.0.1) with the constructed special logic file, the pool of 8994 genome assemblies were re-
443 screened using *Kaptive* v. 2.0.0 using the same parameters. The output confidence levels
444 were: 1012 (perfect), 5158 (very high), 436 (high), 1647 (good), 97 (low) and 644 (none)
445 (Table 1; Supplementary Table S4). Hence, the percentage of genome assemblies that could
446 be typed with a confidence score of ‘perfect’, ‘very high’, or ‘high’ rose to 73.45% (Fig. 3A).
447 As a greater number of KL were confidently assigned with the updated database, there was an
448 observed decrease in the overall number of ‘problem’ scores called by *Kaptive*. This included
449 a decrease of 627 assemblies with missing expected genes (denoted in the output by a ‘-’
450 symbol), 375 assemblies with additional genes (+), and 418 assemblies with one or more
451 expected genes with low identity (*). An unexpected decrease was observed for matches
452 found in more than one piece (?), with 30 fewer assemblies reported with this problem score
453 using *Kaptive* 2.0.1 (Fig. 3B).

454

455 **Assessment of complete genome sequences**

456 To further evaluate the impact of sequence quality on the number of problem scores called by
457 *Kaptive*, results for 264 completed genome sequences with the chromosome available in a
458 single contig were extracted for more detailed analyses. The match confidence scores, length

459 discrepancies, and problem scores were summarised for the 264 complete genomes, and are
460 shown in Table 3. More than 80% of the completed genomes were assigned a match
461 confidence score less than ‘perfect’ (Fig. 3A). For these, all matches with a detected length
462 discrepancy of >500 bp (20 total) were found to include an IS insertion. A total of 71
463 complete genomes were assigned at least one problem score (Fig. 3B). However, the majority
464 were due to low detectable identity of some genes, frameshifts or better sequence matches of
465 the same gene found in a different KL reference sequence in the database. Interestingly, one
466 complete genome received a ‘?’ problem score suggesting that the K locus was not in a single
467 piece. For this genome, the sequence had been opened within the K locus rather than at the
468 origin of replication, leading to detection of the K locus at both the start and the end of the
469 opened chromosome.

470 Seven of the 264 completed genomes were records for the *A. baumannii* reference
471 strain, ATCC 17978 (Table 4), known to carry the KL3 locus [11]. The ATCC 17978 genome
472 was originally sequenced via pyrosequencing [51] and first made available in 2007 under
473 GenBank accession number CP000521 (NCBI assembly number GCA_000015425.1).
474 However, this assembly was later found to include errors [52] and/or gene frameshifts [11],
475 and was subsequently re-sequenced using a combination of Illumina short read and PacBio
476 long read data (GenBank accession number CP012004.1; NCBI assembly number
477 GCA_001077675.1). Here, the re-sequenced genome was assigned to KL3 with ‘perfect’
478 confidence, whereas the original sequence has a confidence score of ‘high’ with the
479 previously reported KL gene frameshifts resulting in missing genes (‘-’) detected.

480 Three further sequenced versions of the ATCC 17978 genome were assigned a
481 ‘perfect’ match to KL3, whereas another that was sequenced and assembled using only
482 PacBio technology was assigned a ‘high’ KL3 match (Table 4). The remaining assembly
483 (NCBI assembly number GCA_011067065) was a ‘very high’ match to KL48. As this is
484 inconsistent with the previous finding of KL3 in ATCC 17978, the assembly was aligned to
485 the ATCC 17978 reference sequence (GCA_001077675) and found to have only 87%
486 sequence coverage with 98% nucleotide sequence identity. Further inspection using MLST
487 revealed that the genome belongs to sequence type ST2 in the Institut Pasteur scheme, rather
488 than ST437 as for all other ATCC 17978 genome sequences. This suggests that
489 GCA_011067065 is incorrectly named in the GenBank record as ATCC 17978.

490

491 **General features of *A. baumannii* K locus sequences**

492 Characterisation of the 237 distinct CPS biosynthesis gene clusters affords the opportunity to
493 re-examine common features of sequences found at the K locus in *A. baumannii* genomes.
494 Amongst the 237 KL, sequence lengths varied between 18.5 kb and 36.8 kb with a mean
495 length of 25 kb (Fig. 4A). The total number of open reading frames (ORFs) per KL also
496 varied, ranging between 16 and 31, with the majority of KL (~65%) including 20 to 23 ORFs
497 (Fig. 4B). The size of the locus correlated with the number of ORFs present, and the smallest
498 carried no modules for the synthesis of complex sugars. The larger gene clusters generally
499 included larger or more gene modules for complex sugar biosynthesis rather (see Fig. 1 and
500 Fig. 2). For example, large gene modules are required for synthesis of non-2-ulosonic acids
501 such as 5,7-di-*N*-acetylacetaminic acid that requires 10 genes [53].

502 A small group of five KL (KL92, KL99, KL142, KL143 and KL145), have a
503 configuration considered unusual for *A. baumannii* as they include a novel segment in Region
504 2 that includes *itrA4* and *wzi_{KL}* genes. Previously, this segment was suggested to have been
505 acquired from a source outside of *A. baumannii* [54].

506 All KL sequences included between 1 and 6 genes predicted to encode
507 glycosyltransferases, with most KL carrying 3 or 4 *gtr* genes (Fig. 4C) suggesting
508 tetrasaccharide and pentasaccharide K-units are common in *A. baumannii* CPS. The
509 correlation between the number of Gtr encoded and the number of sugars in the K unit has
510 been supported by structural studies with exceptions only for CPS containing L-rhamnose
511 [16, 17, 54].

512

513 **Repertoire of genes included in the database**

514 To understand the diversity and distribution of CPS biosynthesis genes across the 237 KL, all
515 genes were grouped into clusters of homologous gene groups using Roary with a cut-off
516 parameter of 85% aa minimum identity for the products, and the groups were used to
517 calculate frequency. This revealed a total of 681 different gene homology groups found
518 across the 237 KL, of which 42.6% of genes were found only in one KL and a further 45.82%
519 found in 2-12 KL (i.e. <5%; Fig. 4D). Nine gene groups occurred in 164 or more gene
520 clusters (>69.2% of KL) and only 1 gene, *pgm*, was found to encode products of >85% aa
521 identity for all 237 KL (100%). This finding was unexpected, as all CPS gene clusters
522 described for *A. baumannii* to date include the same eight genes: *wza*, *wzb*, and *wzc* genes in
523 'Region 1', *gna* in 'Region 2, and *galU*, *ugd*, *gpi*, and *pgm* genes in 'Region 3' [11]. Hence,
524 further assessment of gene product homology groups at the K locus was undertaken.

525

526 **Variation in the eight genes always present at the *A. baumannii* K locus**

527 With the exception of *pgm*, 2-4 homology groups were detected for each of the other seven
528 genes that are always present at the K locus, indicating these genes are not completely
529 conserved. The occurrence of >1 homology group for common genes may be due to multiple
530 imports of the same genes into the species via homologous recombination resulting in a
531 change of KL sequence. Sequence diversity in *wza*, *wzb* and *wzc* genes had been observed
532 previously [12]. However, the level of variation detected here for these genes, as well as for
533 other common genes, is significant and suggests a complex evolutionary history for the *A.*
534 *baumannii* K locus. For example, two homology groups were found for both the *ugd* and *gpi*
535 genes; one group is present in 97.9% of KL and the smaller group (in 2.1% of KL) occurs in
536 the KL92, KL99, KL142, KL143 and KL145 group described above. However, while 85% aa
537 identity is the cut-off used to define new gene types in *A. baumannii* K loci, variants of these
538 common genes are not currently numbered. Nor are they considered in assignment of new
539 numbers to KL.

540 For *gna*, three sequence groups had previously been reported in the same position at
541 the beginning of Region 2 [11]. Two of the three types were shown to form a module with
542 either *gne2* (*gna1*) for synthesis of D-GalNAcA or *dgaABC* (*gna3*) for synthesis of D-
543 GlcNAc3NAcA. A third type (*gna2*) is present in all other KL, though its role in CPS
544 production is still unknown. Here, the same three *gna* homology groups were found.

545

546 **Further variation in Region 3**

547 When present, the *gne1* gene is located between *gpi* and *pgm*. This gene is often present and
548 is required for synthesis of UDP-D-GalNAc and/or UDP-D-Gal [11]. Though the presence of
549 a small variable portion of Region 3 between the *gpi* and *pgm* genes has been previously
550 reported [11, 12], diversity in this region had not been further investigated. Here, a total of
551 217 of 237 KL (91.56%) were found to carry additional coding sequence(s) between *gpi* and
552 *pgm*. A *gne1* gene was present on its own in 95 KL (40.08%), though a further 97 KL
553 (40.93%) included *gne1* adjacent to additional *pgt1*, *pgt2*, *pet1* or *atr* genes (Table 5). A
554 further 23 KL include only *pgt1* between *gpi* and *pgm*. Besides *gne1*, a role for all other
555 genes found between *gpi* and *pgm* in CPS biosynthesis has not yet been established.

556 As the variation at this position in Region 3 is known to affect CPS structure only if
557 D-Galp or D-GalpNAc are present, some groups of KL are likely to produce the same CPS
558 structure. This has been the case for KL2/KL81 and KL3/KL22 pairs that are known to
559 produce the same structure [23]. The gene clusters in these two pairs differ from each other

560 only in the presence of a *pgt1* gene between *gpi* and *pgm*, which has no defined role in CPS
561 biosynthesis. A further 21 examples of pairs or groups of KL that differ in the
562 presence/absence of *pgt1* and/or other genes found between *gpi* and *pgm* in Region 3 were
563 detected amongst the 237 KL (listed in Table 6). Comparisons of some of these pairs or
564 groups can be seen in published KL compilations [15, 18, 25]. Hence, it is possible that
565 further examples of KL that produce the same K type may be found as more CPS structures
566 are determined. Of these 21, 15 groups include one KL with associated structural data (bold
567 in Table 6), possibly representing the structure produced by the other KL of the same group.
568

569 **Genes for biosynthesis of sugars and addition of substituents**

570 The number of homology groups relating to specific modules of genes for sugar biosynthesis
571 or functional categories of gene products were manually curated to gain insights into the
572 possible diversity in sugars and non-sugar substituents that can be incorporated into *A.*
573 *baumannii* CPS. A total of 34 possible modules of gene(s) for the synthesis of complex
574 sugars (described in Table 7) were found. Of these modules, 29 have been reported
575 previously and three are variants of known modules. Two modules, *rhnABC* and *gnLAB*, are
576 described here for the first time (see above). While most KL include at least one gene module
577 for complex sugar biosynthesis, 30 KL did not include any sugar synthesis gene module(s) in
578 Region 2 and these are likely to produce CPS with neutral sugars [23, 26, 29, 37] synthesised
579 by common genes in Region 3 [11]. The *rmlBDAC* module for L-Rhamnose synthesis is the
580 most common sugar gene module across the 237 KL types (found in 15.6%), though gene
581 modules for synthesis of sugars belonging to the non-2-ulosonic acid family (i.e. *psa*, *lga*,
582 *aci*, *ela*, and *neu*) were collectively found in 29.11% of KL.

583 In addition to complex sugar biosynthesis genes/modules, 40 different genes for
584 modifying a monosaccharide in the K unit via the addition of a substituent were identified in
585 Region 2. These included 31 *atr* genes for the transfer of acyl-/acetyl groups, 8 *ptr* genes for
586 the addition of pyruvate, and 1 *alt* gene for transfer of L-alanine. The presence of 34 sugar
587 biosynthesis gene modules and 40 different genes for the addition of non-sugar substituents
588 to the CPS therefore indicates significant potential for diversity in sugar composition and K-
589 unit decoration of *A. baumannii* CPS structures.

590

591 **Genes for initiating transferases and glycosyltransferases**

592 Biosynthesis of *A. baumannii* CPS in the cytoplasm is known to begin with the transfer of a
593 'first' sugar to a lipid carrier in the inner membrane, and six possible first sugars transferred

594 by one of six distinct initiating transferases (Itrs) belonging to one of two families (ItrA or
595 ItrB) are known [55]. A seventh non-functional Itr type, ItrB2, has also been reported [11].
596 Here, the same seven *itr* genes (*itrA1-itrA4* and *itrB1-itrB3*) were found across the 237 KL as
597 expected with no new types identified. The *itrA2* and *itrA3* genes were found to be the most
598 common, present in 91 and 81 different KL, respectively. This indicates that UDP-D-GalNAc
599 (ItrA2) and UDP-D-GlcNAc (ItrA3) are common first sugar substrates for *A. baumannii*
600 CPSs, consistent with what has been observed with available structural data.

601 Following the addition of the first sugar to the lipid carrier, CPS biosynthesis
602 progresses with the addition of further monosaccharides to the first sugar to build a complete
603 K-unit oligosaccharide [11]. The glycosidic linkages between sugars are formed by
604 glycosyltransferase enzymes that are encoded by either *gtr* or *kpsS* genes at the K locus. As
605 the specificity of Gtr/KpsS enzymes for their sugar donor and acceptor substrates can vary,
606 numbering new types is important. While new numbers are currently assigned to genes using
607 a cut-off of 85% aa identity, evidence has emerged that some Gtrs that share <85% aa
608 identity can form the same linkage as one another [17], while others sharing >85% aa identity
609 can form different linkages [48, 56]. Nonetheless, differentiation between different Gtr/KpsS
610 types can provide insights into diversity in the linkages possible between sugars in the K-unit.
611 Using the cut-off of 85% aa identity, a total of 272 homology groups (267 *gtr* and 5 *kpsS*)
612 were found across the 237 KL. Further work will be needed to analyse Gtr sequences to
613 identify relationships between them and how these relate to linkages formed.

614

615 **Genes for K-unit and CPS processing**

616 Region 2 usually also harbours *wzx* and *wzy* genes required for processing oligosaccharide
617 units to form long chain polymers as part of the Wzy-dependent pathway for CPS
618 biosynthesis [12]. Wzx is the translocase that flips oligosaccharide units into the periplasm
619 for polymerisation into chains by the Wzy polymerase. Across the 237 KL, a total of 81 *wzx*
620 and 137 *wzy* gene groups were found. It is unclear what effect *wzx* sequence diversity has on
621 CPS biosynthesis, though the variety of *wzy* genes may reflect many different linkages
622 possible between oligosaccharide units in the CPS polymer. With the exception of the *wzy*_{KL8}
623 type (Fig. 1B) which is located to the left of the *wza-wzb-wzc* genes, and the three (KL19,
624 KL24 and KL39) that do not contain any *wzy*, all *wzy* groups were found in Region 2. The
625 absence of *wzy* in these three gene clusters has been described previously, where Wzy
626 function is supplemented by a *wzy* gene located in genomic islands present elsewhere in the
627 chromosome [35, 36]. However, a relative of KL24, KL146, that does include a *wzy* gene was

628 recently identified [4] indicating that the original *wzy* gene may have been lost. Interestingly,
629 two KL (KL67 and KL134) include 2 distinct *wzy* genes, though it is not known if both
630 contribute to CPS assembly. As *wzx* and *wzy* genes can be shared by K loci, to distinguish
631 *wzx* and *wzy* groups for future typing, a suffix was added to each *wzx* and *wzy* type (defined
632 by the 85% aa identity cutoff) indicating the name of the first KL a group was identified in.

633 Finally, CPS assembly on the cell surface is mediated by a Wzi outer membrane
634 protein encoded by a *wzi* gene located outside the K locus [55]. However, in rare cases, a
635 second *wzi* type has been found at the K locus [54], and is referred to as *wzi_{KL}*. This *wzi* gene
636 is co-located with *itrA4* and was found only in the 5 KL that carry different *ugd* and *gpi* genes
637 (see above).

638

639 **Abundance of KL and K types in the genome pool**

640 To examine the ability of the updated database to detect diversity, the frequency of KL and K
641 types detected across the 8994 genome assemblies was calculated. An assessment of best
642 match loci reported by *Kaptive* with a confidence score of ‘good’ or above, regardless of
643 detected ‘problems’ or possible discontinuous locus sequence, revealed that 19 KL were not
644 found amongst the 8994 genome assemblies studied. Two of these were KL38 and KL78 that
645 are no longer considered *A. baumannii* sequences, and 17 KL that are currently only available
646 as short reads or extracted locus sequences in the GenBank non-redundant database. Some
647 KL were overrepresented, with 17 KL each found in >1% of genome assemblies (Fig. 5A)
648 and collectively representing 74.6% of the total genome pool. The other 25.4% of genome
649 assemblies included 201 KL. KL2 was found to be the most common K locus sequence found
650 in 16.5% of genomes, while KL3 and KL22, both of which produce a K3 type CPS, together
651 represent 20% of genomes.

652 A similar assessment of K type demonstrated that 66 K types could be predicted from
653 the assemblies reported with a confidence score of ‘good’ or above. Eight K types were
654 represented in >1% of genomes each (Fig. 5B), whereas 58 K types represented 19.1% and
655 the remaining 20.9% of matches were listed as ‘unknown’ because relevant structures are not
656 available. As the KL3 and KL22 gene clusters are known to produce the same K3 type
657 structure [23], K3 was found to be the most common consistent with the high proportion of
658 KL3 and KL22 in the genome pool followed by K2 (Fig. 5A).

659

660 8. Discussion

661 In this study, we provide a comprehensive update to the *A. baumannii* CPS gene reference
662 sequence database, providing 237 fully annotated K loci sequences and six 'extra' gene loci
663 outside of the K locus. Each KL record has linked structural information to enable 'type' to
664 be inferred wherever possible based on detection of specific KL with or without 'extra' genes
665 based on special logic parameters. However, since the conclusion of this work, an additional
666 4 novel KL have been identified in recently reported *A. baumannii* genomes. Hence, a total of
667 241 distinct KL sequences have been released into the *Kaptive v. 2.0.1 A. baumannii* CPS
668 gene reference sequence database. As we did not further investigate genomes with K loci
669 detected in more than 1 piece in this study, we expect there may be more sequences yet to be
670 documented amongst the 8994 genome assemblies included in this study.

671 Of the 6726 genome assemblies with a KL match found in a single piece, 98.22%
672 were able to be assigned matches with either 'perfect', 'very high' or 'high' confidence
673 demonstrating the utility of the database to type KL in *A. baumannii* genomes. However, for
674 the majority of 'high' matches, *Kaptive* had detected missing gene products that appeared to
675 be the result of frameshifts in gene sequence(s). Further work will be needed to determine if
676 these frameshifts are the result of sequence/assembly errors, or whether frameshifts are real
677 and if they give rise to changes in the CPS composition therefore warranting a new KL
678 designation to distinguish mutants from wildtype sequences. Given that our inspection of 264
679 complete genome sequences suggests that the quality of both the sequence and the assembly
680 influences the confidence score, the remaining genome assemblies that were not assigned
681 matches with either 'perfect', 'very high' or 'high' confidence, may be of poorer quality.
682 Nonetheless, KL variants were found in the genomes with confidence scores of 'very high' or
683 lower, and novel KL were found amongst assemblies with 'high', 'good', 'low' or even 'none'
684 confidence levels. Hence, users are encouraged to check sequence/assembly quality prior to
685 KL typing, and manually inspect any assembly assigned a *Kaptive* match less than 'perfect'.

686 With the addition of 'type' to the update, users are now also able to infer structural
687 properties of the CPS from their genome data when this is available or can be reasonably
688 inferred. Whilst the inclusion of >237 KL significantly enhances the ability to predict CPS
689 type using WGS, structural studies are currently the only definitive way to determine if genes
690 outside the K locus contribute to the determination of a specific CPS structure. Hence, further
691 examples of genes elsewhere in the genome are likely to be found as more structural data
692 coupled with KL gene content analysis becomes available. Though further work is needed to

693 improve the predictive power of WGS for determination of CPS type, the database is a
694 valuable tool for epidemiological studies. Our analysis showed that 17 KL represent a large
695 proportion of the genome assemblies included in this study. Though this may be due to bias
696 in the dataset associated with strain sampling for outbreak studies, particularly the over-
697 representation of genome for GC2 isolates, new KL described here were each found in <1%
698 of genomes in the pool, indicating that the previous iteration of the database had captured the
699 most common KL.

700 Additional analysis of K locus gene repertoire revealed a total of 681 gene types
701 amongst 237 KL, including 601 (88.3%) genes found in <5% of KL and 286 of these (42%)
702 present in only one KL. Interestingly, 95 of 237 KL include one or more unique genes,
703 suggesting that new gene clusters may arise by acquisition of novel genes, as well as by the
704 reassortment or exchange of genes between already established KL. With the K locus known
705 to undergo recombination [57], sequences of shared genes are likely to exhibit a degree of
706 sequence variation. However, the finding that 7 of the 8 genes always present at the K locus
707 included 2 or more product homology groups of <85% aa identity was unexpected,
708 suggesting that multiple imports of these genes into the species has occurred over time. While
709 new KL numbers are usually assigned to any new combination of genes at the K locus
710 defined by an 85% aa identity cut-off in one or more gene products, new numbers are not
711 warranted for KL that differ from a reference sequence only in one or more of these 8 genes.
712 Hence, we continue to assign new numbers only for any new combination of genes in Region
713 2 (between *gna* and *galU*) and the variable portion of Region 3 (between *gpi* and *pgm*),
714 regardless of whether a difference in CPS structure is expected. As structural data for
715 corresponding CPS continues to grow, examples of KL that produce the same CPS structure
716 will likely be uncovered.

717 Gene repertoire analysis further identified genes for 34 different complex sugars,
718 >400 for glycosidic linkages (*gtr*, *kpsS* and *wzy*) and >40 for K-unit modifications (*atr*, *alt*,
719 and *ptr*) across 237 KL sequences, predicting substantial variation in CPS sugars, sugar
720 linkages, and also acetyl-, acyl-, pyruvyl, and L-alanine decorations. This extraordinary
721 diversity observed in CPS biosynthesis genes complicates next-generation therapies,
722 including vaccines and bacteriophage strategies. However, fine-scale analysis of individual
723 CPS structures and their specific biosynthesis genes can inform tailored approaches to
724 alternate patient treatments. The increase in the number of KL and inclusion of K type
725 information in the updated reference database will significantly enhance epidemiological

726 tracking efforts and assist with building a comprehensive understanding of the circulation of
727 important strains both locally and globally.

728

729 **9. Author statements**

730 **9.1 Authors and contributors**

731 Conceptualization, JJK; Data curation, JJK; Formal analysis, SMC and JJK; Funding
732 acquisition SMC, RMH and JJK; Investigation, SMC, RMH and JJK; Methodology, SMC
733 and JJK; Visualization, SMC and JJK; Supervision, RMH and JJK; Writing – original draft,
734 SMC and JJK; Writing – review & editing, RMH and JJK

735

736 **9.2 Conflicts of interest**

737 The authors declare that there are no conflicts of interest.

738

739 **9.3 Funding information**

740 This work was supported by an Australia Government student stipend to SMC, a National
741 Health and Medical Research Council (NHMRC) of Australia Investigator grant
742 (GNT1194978) to RMH, and an Australian Research Council (ARC) DECRA Fellowship
743 (DE180101563) to JJK.

744

745 **9.4 Ethical approval**

746 N/A

747

748 **9.5 Consent for publication**

749 N/A

750

751 **9.6 Acknowledgements**

752 We thank Kelly Wyres and Ryan Wick from Monash University, Australia, and Kathryn Holt
753 from the London School of Hygiene & Tropical Medicine, UK, for their assistance with
754 updating the database on the *Kaptive* platforms.

755

756 **10. References**

- 757 1. Murray CJL, Ikuta KS, Sharara F, Swetschinski L, Aguilar GR, Gray A, et al. Global
758 burden of bacterial antimicrobial resistance in 2019: a systematic analysis. *The Lancet*.
759 2022;399(10326):629-55.
- 760 2. Isler B, Doi Y, Bonomo RA, Paterson DL. New treatment options against
761 carbapenem-resistant *Acinetobacter baumannii* infections. *Antimicrob Ag Chemother*.
762 2019;63(1):e01110-18.
- 763 3. Holt KE, Hamidian M, Kenyon JJ, Wynn M, Hawkey J, Pickard DJ, et al. Genome
764 sequence of *Acinetobacter baumannii* strain A1, an early example of antibiotic-resistant
765 global clone 1. *Genome Announcements*. 2015;3(2):e00032-15.
- 766 4. Koong J, Johnson C, Rafei R, Hamze M, Myers GS, Kenyon JJ, et al. Phylogenomics
767 of two ST1 antibiotic-susceptible non-clinical *Acinetobacter baumannii* strains reveals
768 multiple lineages and complex evolutionary history in global clone 1. *Microb Genom*.
769 2021;7(12).
- 770 5. Wright M, Haft D, Harkins D, Perez F, Hujer K, Bajaksouzian S, et al. New insights
771 into dissemination and variation of the health care-associated pathogen *Acinetobacter*
772 *baumannii* from genomic analysis. *mBio*. 2014;5(1):e00963-13
- 773 6. Luo YC, Hsieh YC, Wu JW, Quyen TLT, Chen YY, Liao WC, et al. Exploring the
774 association between capsular types, sequence types, and carbapenemase genes in
775 *Acinetobacter baumannii*. *Int J Antimicrob Ag*. 2022;59(1):106470.
- 776 7. Loraine J, Heinz E, Soontarach R, Blackwell GA, Stabler RA, Voravuthikunchai SP,
777 et al. Genomic and phenotypic analyses of *Acinetobacter baumannii* isolates from three
778 tertiary care hospitals in Thailand. *Front Microbiol*. 2020;11:548.
- 779 8. Yang JL, Yang CJ, Chuang YC, Sheng WH, Chen YC, Chang SC. Association of
780 capsular polysaccharide locus 2 with prognosis of *Acinetobacter baumannii* bacteraemia.
781 *Emerg Microbes & Infect*. 2021;11(1):83-90.
- 782 9. Deng Q, Zhang J, Zhang M, Liu Z, Zhong Y, Liu S, et al. Rapid identification of
783 KL49 *Acinetobacter baumannii* associated with clinical mortality. *Infect Drug Resist*.
784 2020;13:4125.

- 785 10. Talyansky Y, Nielsen TB, Yan J, Carlino-Macdonald U, Di Venanzio G, Chakravorty
786 S, et al. Capsule carbohydrate structure determines virulence in *Acinetobacter baumannii*
787 PLoS Pathog. 2021;17:e1009291.
- 788 11. Kenyon J, Hall R. Variation in the complex carbohydrate biosynthesis loci of
789 *Acinetobacter baumannii* genomes. PLoS One. 2013;8(4):e62160.
- 790 12. Wyres KL, Cahill SM, Holt KE, Hall RM, Kenyon JJ. Identification of *Acinetobacter*
791 *baumannii* loci for capsular polysaccharide (KL) and lipooligosaccharide outer core (OCL)
792 synthesis in genome assemblies using curated reference databases compatible with Kaptive.
793 Microb Genom. 2020;6(3):e000339.
- 794 13. Wyres K, Wick R, Gorrie C, Jenney A, Follador R, Thomson N, et al. Identification
795 of *Klebsiella* capsule synthesis loci from whole genome data. Microb Genom.
796 2016;2(12):e000102.
- 797 14. Wick R, Heinz E, Holt K, Wyres K. Kaptive Web: User-friendly capsule and
798 lipopolysaccharide serotype prediction for *Klebsiella* genomes. J Clin Microbiol.
799 2018;56(6):e00197-18.
- 800 15. Kenyon JJ, Hall RM. Updated analysis of the surface carbohydrate gene clusters in a
801 diverse panel of *Acinetobacter baumannii* isolates. Antimicrob Ag Chemother.
802 2021;66(1):e01807-21.
- 803 16. Kenyon JJ, Kasimova AA, Sviridova AN, Shpirt AM, Shneider MM, Mikhaylova
804 YV, et al. Correlation of *Acinetobacter baumannii* K144 and K86 capsular polysaccharide
805 structures with genes at the K locus reveals the involvement of a novel multifunctional
806 rhamnosyltransferase for structural synthesis. Int J Biol Macromol. 2021;193:1294-300.
- 807 17. Kenyon JJ, Arbatsky NP, Sweeney EL, Zhang Y, Senchenkova SN, Popova AV, et al.
808 Involvement of a multifunctional rhamnosyltransferase in the synthesis of three related
809 *Acinetobacter baumannii* capsular polysaccharides, K55, K74 and K85. Int J Biol Macromol.
810 2021;166:1230-7.
- 811 18. Kenyon JJ, Senchenkova SN, Shashkov AS, Shneider MM, Popova AV, Knirel YA,
812 et al. K17 capsular polysaccharide produced by *Acinetobacter baumannii* isolate G7 contains
813 an amide of 2-acetamido-2-deoxy-D-galacturonic acid with D-alanine. Int J Biol Macromol.
814 2019;144:857-62.
- 815 19. Kenyon JJ, Shashkov AS, Senchenkova SN, Shneider MM, Liu B, Popova AV, et al.
816 *Acinetobacter baumannii* K11 and K83 capsular polysaccharides have the same 6-deoxy-l-
817 talose-containing pentasaccharide K units but different linkages between the K units. Int J
818 Biol Macromol. 2017;103:648-55.

- 819 20. Kasimova AA, Arbatsky NP, Timoshina OY, Shneider MM, Shashkov AS, Chizhov
820 AO, et al. The K26 capsular polysaccharide from *Acinetobacter baumannii* KZ-1098:
821 Structure and cleavage by a specific phage depolymerase. *Int J Biol Macromol.*
822 2021;191:182-91.
- 823 21. Shashkov AS, Kenyon JJ, Arbatsky NP, Shneider MM, Popova AV, Knirel YA, et al.
824 Genetics of biosynthesis and structure of the K53 capsular polysaccharide of *Acinetobacter*
825 *baumannii* D23 made up of a disaccharide K unit. *Microbiology.* 2018;164:1289-92.
- 826 22. Arbatsky NP, Popova AV, Shneider MM, Shashkov AS, Hall RM, Kenyon JJ, et al.
827 Structure of the K87 capsular polysaccharide and KL87 gene cluster of *Acinetobacter*
828 *baumannii* LUH5547 reveals a heptasaccharide repeating unit. *Carbohydr Res.*
829 2021;509:108439.
- 830 23. Arbatsky NP, Shneider MM, Kenyon JJ, Shashkov AS, Popova AV, Miroshnikov
831 KA, et al. Structure of the neutral capsular polysaccharide of *Acinetobacter baumannii*
832 NIPH146 that carries the KL37 capsule gene cluster. *Carbohydr Res.* 2015;413:12-5.
- 833 24. Kasimova AA, Kenyon JJ, Arbatsky NP, Shashkov AS, Popova AV, Shneider MM, et
834 al. *Acinetobacter baumannii* K20 and K21 capsular polysaccharide structures establish roles
835 for UDP-glucose dehydrogenase Ugd2, pyruvyl transferase Ptr2 and two
836 glycosyltransferases. *Glycobiology.* 2018;28(11):876-84.
- 837 25. Kasimova AA, Kenyon JJ, Arbatsky NP, Shashkov AS, Popova AV, Knirel YA, et al.
838 Structure of the K82 capsular polysaccharide from *Acinetobacter baumannii* LUH5534
839 containing a D-galactose 4,6-pyruvic acid acetal. *Biochem (Mos).* 2018;83(7):831-5.
- 840 26. Shashkov AS, Kenyon JJ, Arbatsky NP, Shneider MM, Popova AV, Miroshnikov
841 KA, et al. Related structures of neutral capsular polysaccharides of *Acinetobacter baumannii*
842 isolates that carry related capsule gene clusters KL43, KL47, and KL88. *Carbohydr Res.*
843 2016;435:173-9.
- 844 27. Kasimova AA, Cahill SM, Shpirt AM, Dudnik AG, Shneider MM, Popov AV, et al.
845 The K139 capsular polysaccharide produced by *Acinetobacter baumannii* MAR17-1041
846 belongs to a group of related structures including K14, K37 and K116. *Int J Biol Macromol.*
847 2021;193:2297-303.
- 848 28. Arbatsky NP, Kenyon JJ, Kasimova AA, Shashkov AS, Shneider MM, Popova AV, et
849 al. K units of the K8 and K54 capsular polysaccharides produced by *Acinetobacter*
850 *baumannii* BAL 097 and RCH52 have the same structure but contain different di-N-acyl
851 derivatives of legionaminic acid and are linked differently. *Carbohydr Res.* 2019;483:107745.

- 852 29. Shashkov AS, Kenyon JJ, Arbatsky NP, Shneider MM, Popova AV, Miroshnikov
853 KA, et al. Structures of three different neutral polysaccharide of *Acinetobacter baumannii*,
854 NIPH190, NIPH201, and NIPH615, assigned to K30, K45, and K48 capsule types,
855 respectively, based on capsule biosynthesis gene clusters. *Carbohydr Res.* 2015;417:81-8.
- 856 30. Kenyon JJ, Kasimova AA, Notaro A, Arbatsky NP, Speciale I, Shashkov AS, et al.
857 *Acinetobacter baumannii* K13 and K73 capsular polysaccharides differ only in K-unit side
858 branches of novel non-2-ulosonic acids: di-N-acetylated forms of either acinetaminic acid or
859 8-epiacinetaminic acid. *Carbohydr Res.* 2017;452:149-55.
- 860 31. Vinogradov E, MacLean L, Xu HH, Chen W. The structure of the polysaccharide
861 isolated from *Acinetobacter baumannii* strain LAC-4. *Carbohydr Res.* 2014;390:42-5.
- 862 32. Kenyon JJ, Marzaioli AM, Hall RM, De Castro C. Structure of the K2 capsule
863 associated with the KL2 gene cluster of *Acinetobacter baumannii*. *Glycobiology.*
864 2014;24(6):554-63.
- 865 33. Senchenkova SN, Shashkov AS, Shneider MM, Arbatsky NP, Popova AV,
866 Miroshnikov KA, et al. Structure of the capsular polysaccharide of *Acinetobacter baumannii*
867 ACICU containing di-N-acetylpsseudaminic acid. *Carbohydr Res.* 2014;391:89-92.
- 868 34. Russo TA, Beanan J, Olson R, MacDonald U, Cox A, St. Michael F, et al. The K1
869 capsular polysaccharide from *Acinetobacter baumannii* is a potential therapeutic target via
870 passive immunization. *Infect Immun.* 2013;81(3):915-22.
- 871 35. Kenyon JJ, Shneider MM, Senchenkova SN, Shashkov AS, Siniagina MN, Malanin
872 SY, et al. K19 capsular polysaccharide of *Acinetobacter baumannii* is produced via a Wzy
873 polymerase encoded in a small genomic island rather than the KL19 capsule gene cluster.
874 *Microbiology.* 2016;162(8):1479-89.
- 875 36. Kenyon JJ, Kasimova AA, Shneider MM, Shashkov AS, Arbatsky NP, Popova AV, et
876 al. The KL24 gene cluster and a genomic island encoding a Wzy polymerase contribute genes
877 needed for synthesis of the K24 capsular polysaccharide by the multiply antibiotic resistant
878 *Acinetobacter baumannii* isolate RCH51. *Microbiology.* 2017;163:355-63.
- 879 37. Arbatsky NP, Kasimova AA, Shashkov AS, Shneider MM, Popova AV, Shagin DA,
880 et al. Involvement of a Phage-Encoded Wzy Protein in the Polymerization of K127 Units To
881 Form the Capsular Polysaccharide of *Acinetobacter baumannii* Isolate 36-1454. *Microbiol*
882 *Spectr.* 2022;27:e0150321.
- 883 38. Kenyon JJ, Arbatsky NP, Shneider MM, Popova AV, Dmitrenok AS, Kasimova AA,
884 et al. The K46 and K5 capsular polysaccharides produced by *Acinetobacter baumannii* NIPH

- 885 329 and SDF have related structures and the side-chain non-ulosonic acids are 4-O-acetylated
886 by phage-encoded O-acetyltransferases. *PLoS One*. 2019;14(6):e0218461.
- 887 39. Lam MMC, Wick RR, Judd LM, Holt KE, Wyres KL. Kaptive 2.0: updated capsule
888 and LPS locus typing for the *Klebsiella pneumoniae* species complex. *Microb Genom*.
889 2022;8(3):000800.
- 890 40. Sullivan MJ, Petty NK, Beatson SA. Easyfig: a genome comparison visualizer
891 *Bioinformatics*. 2011;27(1):1009-10.
- 892 41. Darling AC, Mau B, Blattner FR, Perna NT. Mauve: multiple alignment of conserved
893 genomic sequence with rearrangements. *Genome Res*. 2004;14(7):1394-403.
- 894 42. Alikhan NF, Petty, N.K., Ben Zakour, N.L. and Beatson, S.A. BLAST Ring Image
895 Generator (BRIG): simple prokaryote genome comparisons. *BMC Genom*. 2011;12(1):1-10.
- 896 43. Gurevich A, Saveliev, V., Vyahhi, N. and Tesler, G. QUAST: quality assessment tool
897 for genome assemblies. *Bioinformatics*. 2013;29(8):1072-5.
- 898 44. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*.
899 2014;30(14):2068–9.
- 900 45. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, et al. Roary:
901 rapid large-scale prokaryote pan genome analysis. *Bioinformatics*. 2015;31(22):3691–3.
- 902 46. Wickham H. Data analysis. In *ggplot2* Springer, Cham.; 2016.
- 903 47. Roberts L, Forde BM, Hurst T, Ling W, Nimmo GR, Bergh H, et al., Genomic
904 surveillance, characterization and intervention of a polymicrobial multidrug-resistant
905 outbreak in critical care. *Microb Genom*. 2021;7, 000530.
- 906 48. Senchenkova S, Popova A, Shashkov A, Shneider M, Mei Z, Arbatsky N, et al.
907 Structure of a new pseudaminic acid-containing capsular polysaccharide of *Acinetobacter*
908 *baumannii* LUH5550 having the KL42 capsule biosynthesis locus. *Carbohydr Res*.
909 2015;407:154-7.
- 910 49. Nakar D, Gutnick D. Analysis of the *wee* gene cluster responsible for the biosynthesis
911 of the polymeric bioemulsifier from the oil-degrading strain *Acinetobacter lwoffii* RAG-1.
912 *Microbiology*. 2001;147:1937-46.
- 913 50. Perepelov AV, Ni Z, Wang Q, Shevelev SD, Senchenkova SN, Shashkov AS, et al.
914 Structure and gene cluster of the O-antigen of *Escherichia coli* O109; chemical and genetic
915 evidences of the presence of L-RhaN3N derivatives in the O-antigens of *E. coli* O109 and
916 O119. *FEMS Immunol Med Microbiol*. 2011;61(1):47-53.

- 917 51. Smith M, Gianoulis T, Pukatzki S, Mekalanos J, Ornston L, Gerstein M, et al. New
918 insights into *Acinetobacter baumannii* pathogenesis revealed by high-density pyrosequencing
919 and transposon mutagenesis. *Genes Dev.* 2007;21(5):601-14.
- 920 52. Weber BS, Ly PM, Irwin JN, Pukatzki S, Feldman MF. A multidrug resistance
921 plasmid contains the molecular switch for type VI secretion in *Acinetobacter baumannii*.
922 *Proc Natl Acad Sci USA.* 2015;112(30):9442-7.
- 923 53. Kenyon JJ, Marzaioli A, De Castro C, Hall RM. 5,7-Di-*N*-acetylacinetaminic acid - a
924 novel non-2-ulosonic acid found in the capsule of an *Acinetobacter baumannii* isolate.
925 *Glycobiology.* 2015;25(6):644-54.
- 926 54. Senchenkova SN, Shashkov AS, Shneider MM, Popova AV, Balaji V, Biswas I, et al.
927 A novel ItrA4 d-galactosyl 1-phosphate transferase is predicted to initiate synthesis of an
928 amino sugar-lacking K92 capsular polysaccharide of *Acinetobacter baumannii* B8300. *Res*
929 *Microbiol.* 2021;172(3):103815.
- 930 55. Tickner J, Hawas S, Totsika M, Kenyon JJ. The Wzi outer membrane protein
931 mediates assembly of a tight capsular polysaccharide layer on the *Acinetobacter baumannii*
932 cell surface. *Sci Rep.* 2021;11(1):1-12.
- 933 56. Kenyon JJ, Arbatsky NP, Sweeney EL, Shashkov AS, Shneider MM, Popova AV, et
934 al. Production of the K16 capsular polysaccharide by *Acinetobacter baumannii* ST25 isolate
935 D4 involves a novel glycosyltransferase encoded in the KL16 gene cluster. *Int J Biol*
936 *Macromol.* 2019;128:101-6.
- 937 57. Holt KE, Kenyon JJ, Hamidian M, Schultz M, Pickard DJ, Dougan G, et al. Five
938 decades of genome evolution in the globally distributed, extensively antibiotic resistant
939 *Acinetobacter baumannii* global clone 1. *Microb Genom.* 2016;2(2):e000052.

940 **11. Figures and tables**

941 **Table 1. Summary of *Kaptive* search results**

	Match Confidence Scores						Total
	Perfect	Very high	High	Good	Low	None	
Database v <i>0.7.0-2.0.0</i>	794	4794	443	1776	192	995	8994
Database; in-house	922	5124	449	1726	135	638	8994
Database v <i>2.0.1</i>	1012	5158	436	1647	97	644	8994

942

943 **Table 2. Updated gene nomenclature key for *A. baumannii* K loci**

Gene name	Predicted reaction product	Predicted protein
<i>aci</i>	CMP- <u>a</u> cinetaminic <u>a</u> cid derivative	Multiple
<i>atr</i>	-	<u>A</u> cyl- or <u>A</u> cetyl- <u>t</u> ransferase
<i>alt</i>	-	D- <u>A</u> lanine <u>t</u> ransferase
<i>dga</i>	UDP-2,3- <u>d</u> iacetamido-2,3-dideoxy-D-glucuronic <u>a</u> cid	Multiple
<i>dmaA</i>	UDP-2,3-diacetamido-2,3-dideoxy-D-mannuronic acid	2-epimerase
<i>ela</i>	CMP-8- <u>e</u> pi <u>l</u> egionaminic <u>a</u> cid derivative	Multiple
<i>fdt</i>	dTDP-D-Fucp3NAc	Multiple
<i>fnl</i>	dTDP-L-FucpNAc	Multiple
<i>fnr</i>	UDP-D-FucpNAc	UDP-6-deoxy-4-keto-D-GalpNAc 4- <u>r</u> eductase
<i>galU</i>	UDP-D-Glcp	UTP-glucose-1-phosphate uridylyltransferase
<i>gdr</i>	UDP-4-keto-6-deoxy-D-GlcpNAc	UDP-GlcpNAc 4,6- <u>d</u> ehydratase
<i>glf</i>	D-Gal <u>f</u> (D-galactofuranose)	UDP-galactopyranose mutase
<i>gna</i>	UDP-D-GlcpNAc <u>A</u>	UDP-D-GlcpNAc dehydrogenase
<i>gne1</i>	UDP-D-GalpNAc, UDP-D-Galp	UDP-D-Glcp/UDP-D-GlcpNAc <u>e</u> pimerase
<i>gne2</i>	UDP-D-GalpNAcA	UDP-D-GlcpNAcA <u>e</u> pimerase
<i>gnl¹</i>	UDP-L-GalpNAcA	Multiple
<i>gpi</i>	L-Fructose-6-phosphate	<u>g</u> lucose-6- <u>p</u> hosphate <u>i</u> somerase
<i>gtr</i>	-	<u>G</u> lycosyl <u>t</u> ransferase
<i>itr</i>	-	<u>I</u> nitiating <u>t</u> ransferase
<i>lga</i>	CMP- <u>L</u> egionaminic <u>a</u> cid derivative	Multiple
<i>man</i>	GDP-D-mannose	Multiple
<i>mna</i>	UDP-D-ManpNAc	Multiple
<i>neu</i>	CMP-N-acetylneuraminic acid	Multiple
<i>pet</i>	-	<u>P</u> hosphoethanolamine <u>t</u> ransferase
<i>pgm</i>	D-Glucose-1-phosphate	Phosphoglucomutase
<i>pgt</i>	-	<u>P</u> hosphoglycerol <u>t</u> ransferase
<i>psa</i>	CMP- <u>P</u> seudaminic <u>a</u> cid derivative	Multiple
<i>ptr</i>	-	<u>P</u> yruvyl <u>t</u> ransferase
<i>qdt</i>	dTDP-D-Quip3NAc	Multiple
<i>qhb</i>	UDP-D-QuipNAc4NH <u>b</u>	Multiple
<i>qnr</i>	UDP-D-QuipNAc	UDP-6-deoxy-4-keto-D-GlcpNAc 4- <u>r</u> eductase
<i>rhn¹</i>	dTDP-L-RhaNAc	Multiple
<i>rml</i>	dTDP-L-Rhamnose	Multiple
<i>tle</i>	dTDP-6-deoxy-L-talose	dTDP-L-Rhamnose <u>e</u> pimerase
<i>ugd</i>	UDP-D-GlcpA	<u>U</u> DP-D- <u>G</u> lcp <u>d</u> ehydrogenase
<i>vio</i>	dTDP-4-acetamido-4,6-dideoxy-D-glucose	Multiple
<i>wza</i>	-	Outer membrane protein
<i>wzb</i>	-	Protein tyrosine phosphatase
<i>wzc</i>	-	Protein tyrosine kinase
<i>wzx</i>	-	Repeat unit translocase
<i>wzy</i>	-	Repeat unit polymerase

944 ¹ New gene annotations (this study)

945 **Table 3. Summary of *Kaptive v 2.0.1* results for complete genomes only**

Confidence score	Total number of assemblies	Matches with length discrepancy (> +/- 500bp)	Matches including genes with low identity to best match (**)	Matches including missing genes ('-')	Matches including additional genes ('+')	Matches with KL in >1 piece ('?')
Perfect	48 (18.2%)	0	0	0	0	0
Very high	170 (64.4%)	13 ¹	25	0	0	0
High	37 (14%)	7 ¹	1	37 ²	0	0
Good	2 (0.8%)	0	0	2 ²	1 ³	1
Low	1 (0.4%)	0	0	1 ²	1 ³	0
None	6 (2.3%)	0	3	6 ²	4 ³	0
Total	264	20	29	46	6	1

946

947 ¹ all length discrepancies are confirmed as additional IS insertion(s)

948 ² all genes reported 'missing' are due to SNPs predicting frameshifts or due to better

949 sequence matches of the same gene found a different KL reference sequence in the database

950 ³ all genes reported 'extra' are better sequence matches of the same gene found a different KL

951 reference sequence

952 **Table 4. Complete genome sequences for *A. baumannii* strain ATCC 17978**

Genome assembly accession number	Best match locus	Match confidence	Problems	Sequence coverage	DNA identity	Length discrepancy	Sequencing platform(s)	Assembly and/or polishing software	Read coverage	Upload date
GCA_000015425 ¹	KL3	High	- ³	100.00%	99.98%	0 bp	high-density pyrophosphate DNA sequencing	Pyrosequencing assembly	21.1X	2007
GCA_001077675 ²	KL3	Perfect		100.00%	100.00%	0 bp	Illumina; PacBio	SPAdes v. 2.5.0; HGAP v. 2.2.0.133377-patch-3	153X	2015
GCA_001593425	KL3	Perfect		100.00%	100.00%	0 bp	Illumina MiSeq	Geneious v. 9.1.5	300.0x	2016
GCA_004797155	KL3	Perfect		100.00%	100.00%	0 bp	PacBio	PacBio SMRT Analysis v. 2.3.0	247.19x	2019
GCA_014672775	KL3	High	- ³	100.00%	100.00%	-1 bp	PacBio RSII	HGAP v. 3.0	399.24x	2020
GCA_013372085	KL3	Perfect		100.00%	100.00%	0 bp	Illumina HiSeq; Oxford Nanopore MiniION	Unicycler v. 0.4.2	80.0x	2020
GCA_011067065	KL48	Very high		100.00%	96.76%	-2 bp	PacBio	Pacbio v. 20K	231.08x	2020

953

954 ¹ Genome sequence excluded from RefSeq database due to poor sequence quality

955 ² Re-sequenced genome reported in [52]

956 ³ ‘-‘ is Kaptive problem score indicating missing gene(s)

957 **Table 5. Number of KL with different gene combinations between *gpi* and *pgm* in Region**
958 **3.**

Gene combination	Number of KL
<i>gne1</i> only	95
<i>gne1/pgt1</i>	76
<i>gne1/pgt2</i>	1
<i>pgt1</i> only	23
<i>gne1/pet1/orf/orf</i>	3
<i>gne1/orf/atr32/atr33</i>	3
<i>gne1/atr15</i>	2
<i>gne1/orf/atr20</i>	1
<i>gne1/atr42/atr43</i>	6
<i>gne1/atr24</i>	1
<i>gne1/orf/atr5-like</i>	1
<i>gne1/atr12</i>	3
None	22
TOTAL	237

959 **Table 6. K loci predicted to produce the same CPS structure**

Pair	Gene module between <i>gpi</i> and <i>pgm</i>										
	None	<i>gne1</i>	<i>gne1</i> / <i>pgt1</i>	<i>pgt1</i>	<i>gne1</i> / <i>atr20</i>	<i>gne1</i> / <i>orf</i> / <i>atr32</i> / <i>atr33</i>	<i>gne1</i> / <i>atr15</i>	<i>gne1</i> / <i>atr12</i>	<i>gne1</i> / <i>pet1</i>	<i>gne1</i> / <i>orf</i> / <i>atr20</i>	<i>gne1</i> / <i>atr32</i>
1		KL2	KL81								
2		KL3	KL22			KL159					
3	KL1 ¹	KL107									
4	KL17	KL18	KL237								
5	KL109	KL9	KL149		KL168	KL173					
6		KL150					KL50				
7		KL64	KL160								
8	KL147	KL15									
9		KL33			KL77						
10		KL170	KL225								
11		KL42	KL216								
12		KL155	KL210								
13			KL27					KL130			
14		KL196	KL52								
15	KL201	KL25									
16	KL91	KL40									
17			KL195	KL11							
18		KL34	KL199						KL20		
19		KL161							KL118		
20										KL124	KL82
21		KL152				KL151				KL133	
22		KL231	KL47								
23		KL32	KL200		KL164		KL100				

960

961 ¹ Bold face text indicates KL has a known CPS structure. Citations supplied in database file.

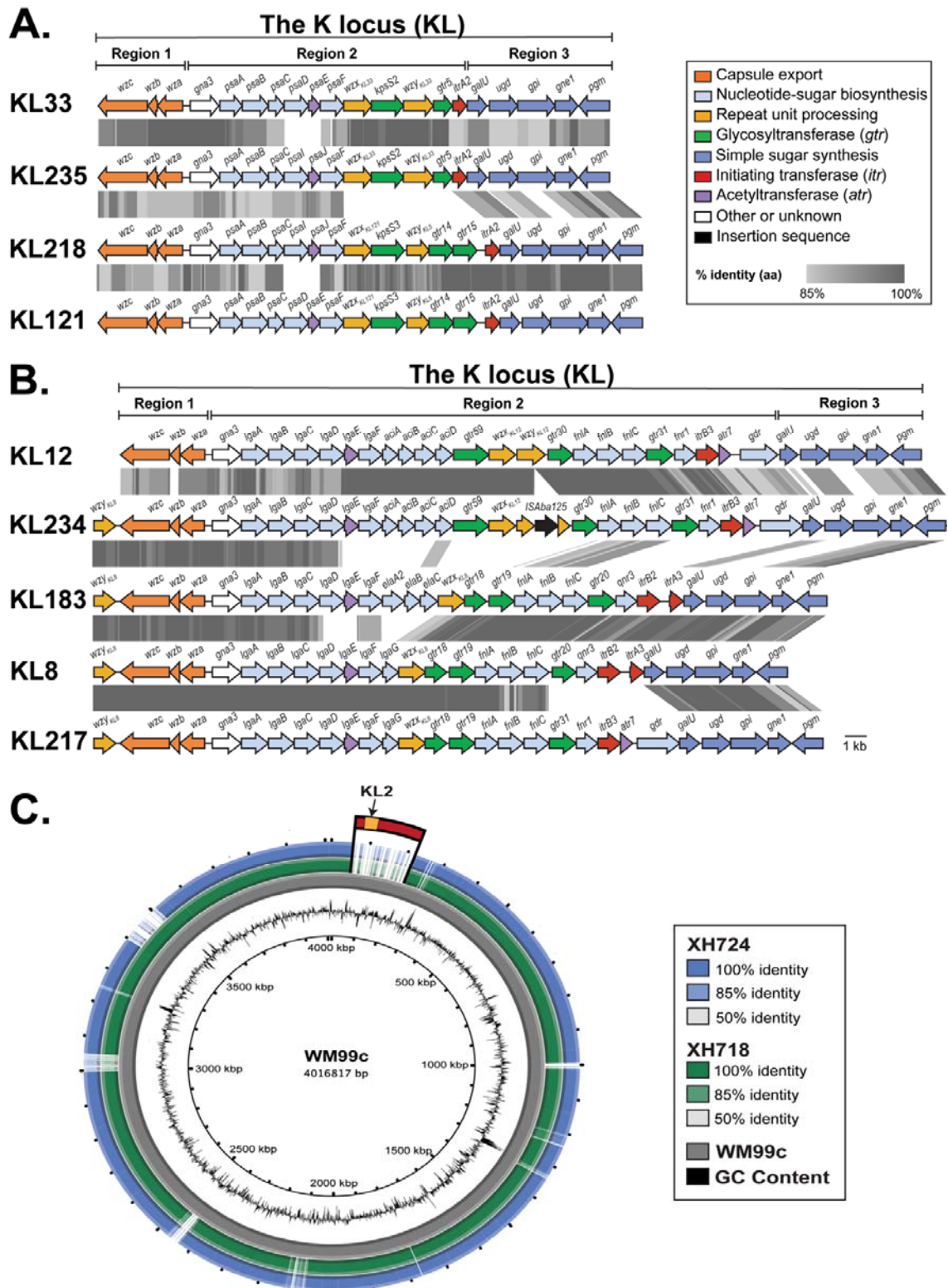
962 **Table 7. Different sugar synthesis gene modules identified in Region 2 across 237 KL.**

Sugar synthesis gene module	Predicted sugar	Number of KL
<i>lgaABCDEF/aciABCD</i>	Ac5A7Ac (5,7-di- <i>N</i> -acetylacinetaminic acid)	5
<i>lgaABCDEF/aciAECD</i>	8eAc5A7Ac (5,7-di- <i>N</i> -acetyl-8-epiacinetaminic acid)	1
<i>lgaABCDEF</i>	Leg5Ac7Ac (5,7-di- <i>N</i> -acetyllegionaminic acid)	18
<i>lgaABCHIFG</i>	Leg5Ac7R (5- <i>N</i> -acetyl-7- <i>N</i> -[(<i>R</i>)-3-hydroxybutanoyl]legionaminic acid)	12
<i>lgaABCDEF/elaABC</i>	8eLeg5Ac7Ac (5,7-di- <i>N</i> -acetyl-8-epilegionaminic acid)	2
<i>lgaABCDEF/ela2BC</i>	(8-epilegionaminic acid derivative) ¹	6
<i>psaABCDEF</i>	Pse5Ac7Ac (5,7-di- <i>N</i> -acetylpsseudaminic acid)	24
<i>psaABC</i> <i>G</i> <i>H</i> <i>F</i>	Pse5Ac7R (5- <i>N</i> -acetyl-7- <i>N</i> -[(<i>R</i>)-3-hydroxybutanoyl]psseudaminic acid)	13
<i>psaABC</i> <i>I</i> <i>J</i> <i>F</i>	(Pseudaminic acid derivative) ¹	2
<i>neuAB</i>	(Neuraminic acid derivative) ¹	1
<i>gna/dgaABC</i>	D-GlcNAc3NAcA (2,3-diacetamido-2,3-dideoxy-D-glucuronic acid)	26
<i>ugd2-ugd5</i>	D-GlcA (D-glucuronic acid) ²	24
<i>gna/gne2</i>	D-GalNAcA (<i>N</i> -acetyl-D-galactosaminuronic acid)	21
<i>gnlAB</i>	(<i>N</i> -acetyl-L-galactosaminuronic acid) ¹	7
<i>gnlB</i>	(<i>N</i> -acetyl-L-galactosamine) ¹	1
<i>glf</i>	(D-galactofuranose) ¹	1
<i>dmaA</i>	(2,3-diacetamido-2,3-dideoxy-D-mannuronic acid) ^{1,2}	12
<i>mnaA</i>	D-ManNAc (<i>N</i> -acetyl-D-mannosamine) ²	9
<i>mnaAB</i>	D-ManNAcA (<i>N</i> -acetyl-D-mannosaminuronic acid)	17
<i>manC</i>	D-Man (D-mannose)	2
<i>fnlABC</i>	L-FucNAc (<i>N</i> -acetyl-L-fucosamine)	32
<i>fnr/gdr</i>	D-FucNAc (<i>N</i> -acetyl-D-fucosamine)	21
<i>rmlBA/fdtACDB</i>	D-Fuc3NAc (3-acetamido-3,6-dideoxy-D-galactose)	4
<i>rmlBA/fdtEB</i>	D-Fuc3NAc (3-acetamido-3,6-dideoxy-D-galactose)	5
<i>qnr/gdr</i>	D-QuiNAc (<i>N</i> -acetyl-D-quinosamine)	3
<i>rmlBA/qdtACDB</i>	D-Qui3NAc (3-acetamido-3,6-dideoxy-D-glucose)	1
<i>rmlBA/qdtEB</i>	D-Qui3NAc (3-acetamido-3,6-dideoxy-D-glucose)	5
<i>rmlBA/vioAB</i>	D-Qui4NAc (4-acetamido-4,6-dideoxy-D-glucose)	4
<i>qhbAB/gdr</i>	D-QuiNAc4NAc (2,4-diacetamido-2,4,6-trideoxy-D-glucose)	14
<i>qhbCB/gdr</i>	D-QuiNAc4NHb (2-acetamido-4-[(<i>S</i>)-3-hydroxybutanoyl]amino-2,4,6-trideoxy-D-glucose)	22
<i>rhnABC</i>	(2,3-diacetamido-2,3,6-trideoxy-L-mannose) ¹	2
<i>rmlBDAC</i>	L-Rha (L-rhamnose)	37
<i>tle</i>	6d-L-Tal (6-deoxy-L-talose)	12
None	-	30

963 ¹ Structural data not available to indicate sugar formed by this gene module

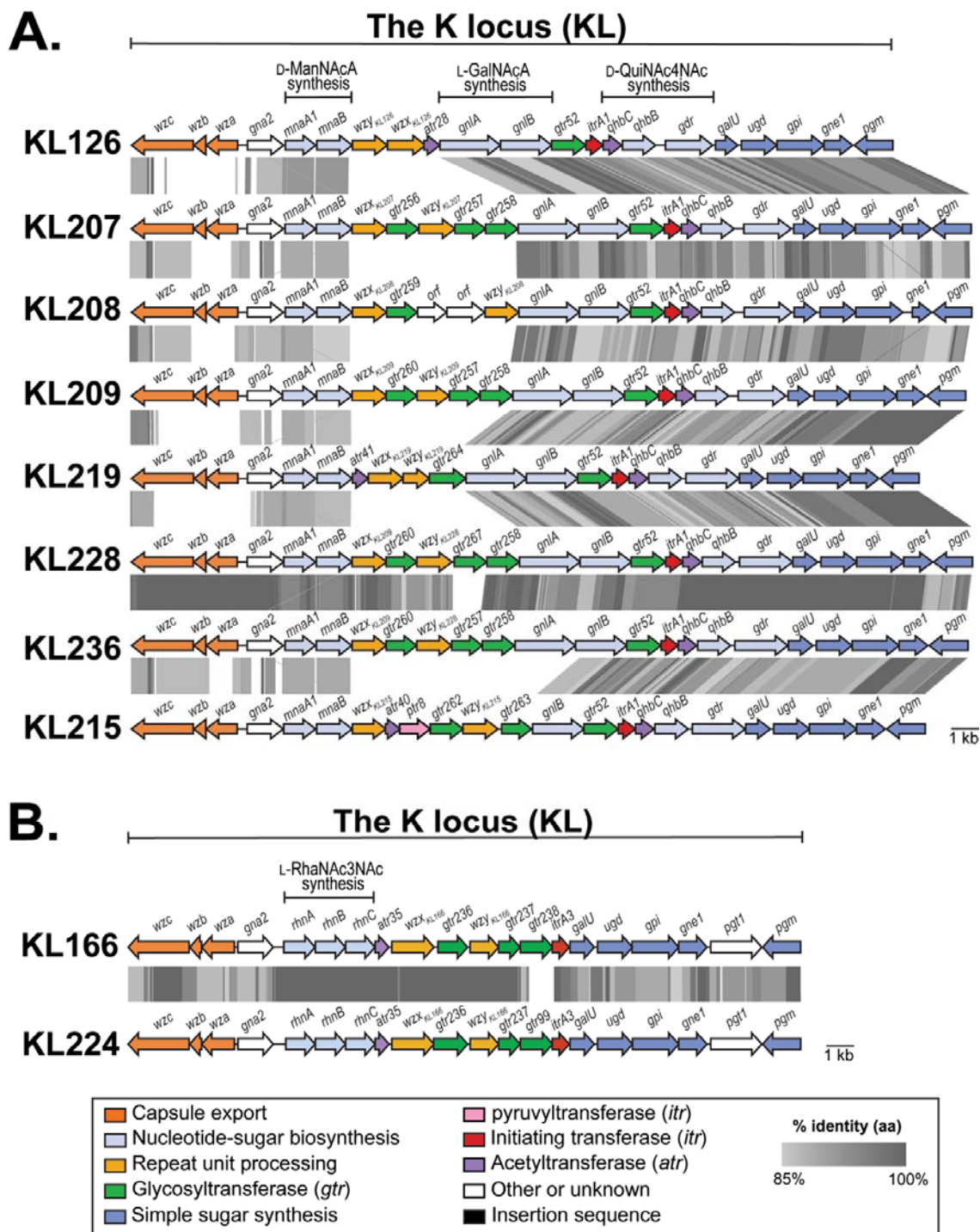
964 ² Structural data not available to confirm sugar made by all homology groups of specific gene

965



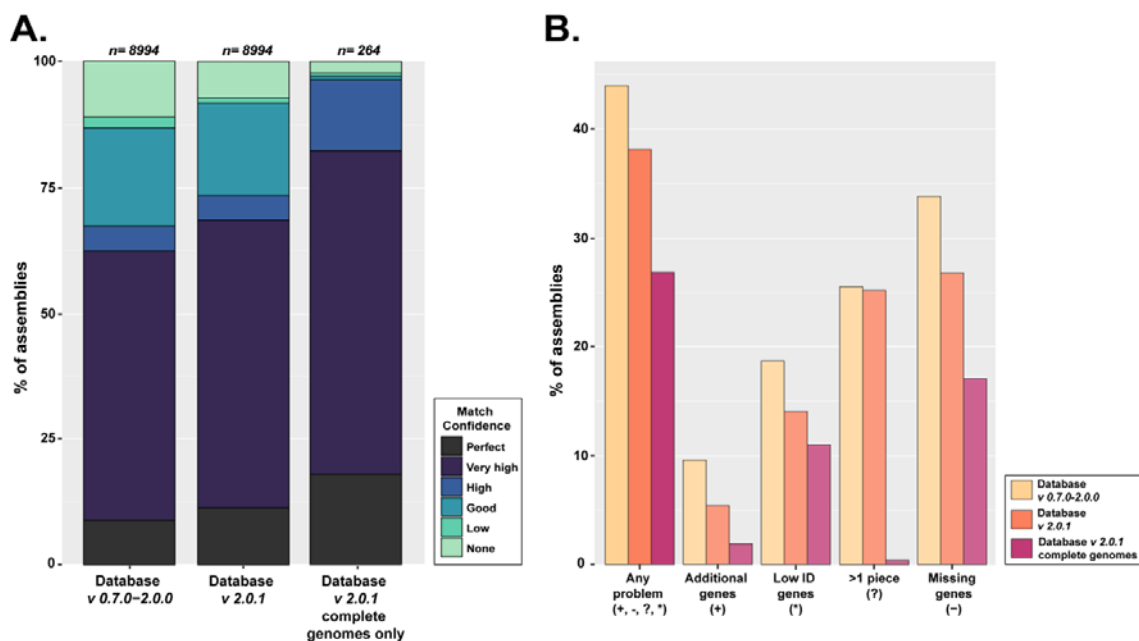
966
 967 **Figure 1. (A)** Comparison of KL33 to the best match reference sequence, KL235, and related
 968 sequences, KL218 and KL121. **(B)** Comparison of KL234 to the best match reference
 969 sequence, KL12, and to related sequences with *wzy* downstream of *wzc*. Genes are coloured
 970 to the functional category of their gene products with legend shown top right. Grey shading

971 between gene clusters indicates amino acid sequence identities with scale shown in the
972 legend. Figures drawn to scale using Easyfig [40] and annotated/coloured in Adobe
973 Illustrator. (C) BRIG multiple sequence alignment of genomes from strains XH724 and
974 XH718 aligned to the WM99c reference genome (NCBI accession number CP032055.1).
975 Contigs were reordered using MAUVE [41] prior to BRIG [42]. Location of the KL2 locus in
976 the WM99c reference genome is marked orange.



977
978

979 **Figure 2.** Comparison of *A. baumannii* KL that include (A) novel *gnaI* and/or *gnaB* genes,
980 and (B) novel *rhn* genes. Genes are coloured to the functional category of their gene products
981 with colour legend shown below. Grey shading between gene clusters indicates amino acid
982 sequence identities with scale shown in the legend below. Figures drawn to scale using
983 Easyfig [40] and annotated/coloured in Adobe Illustrator.

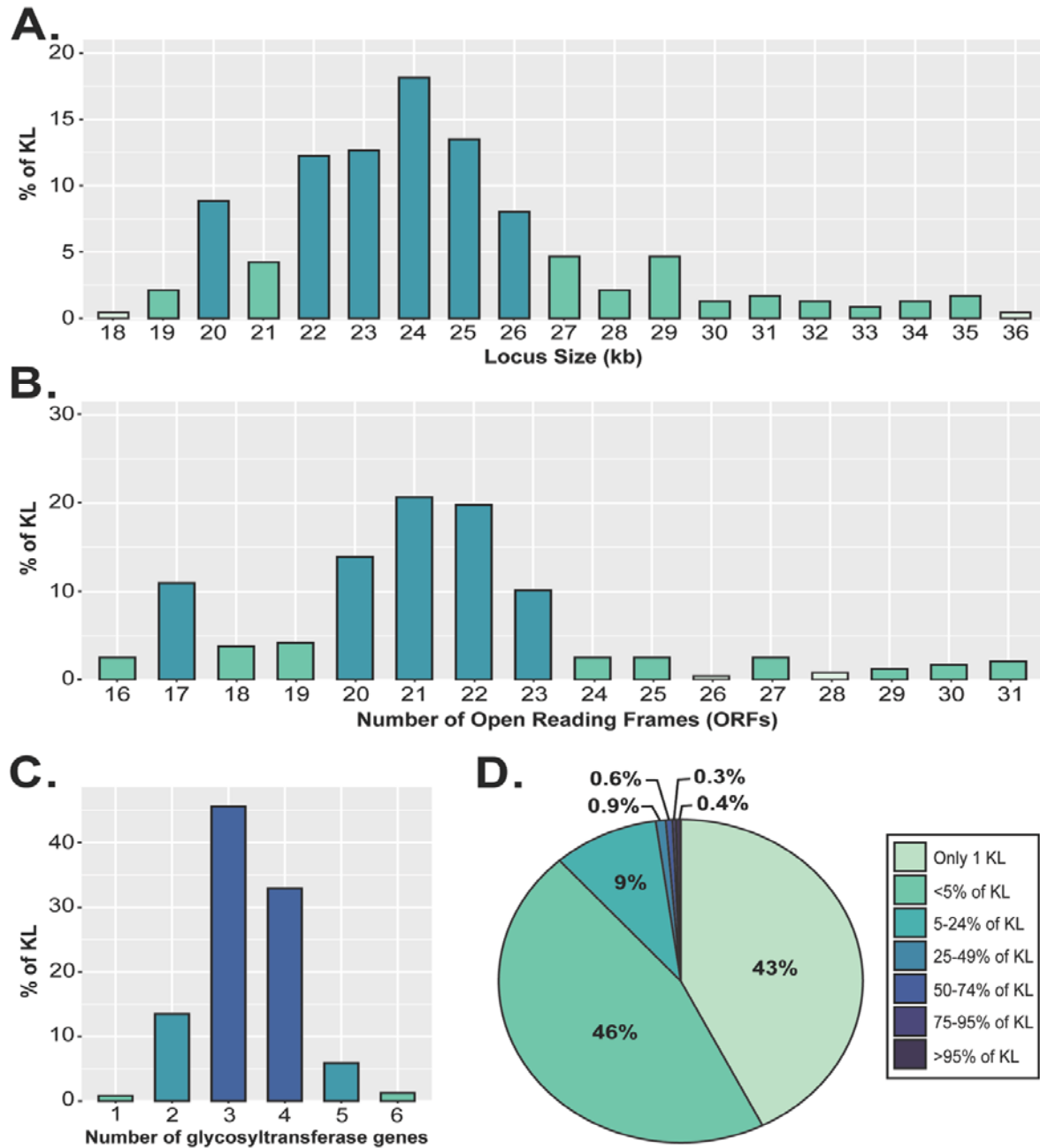


984

985 **Figure 3. Comparison of *Kaptive* performance using the original and updated *A.***

986 ***baumannii* KL reference sequence databases. (A)** Comparison of match confidence scores
987 visualised as a stack bar blot. Database version is indicated below, and the number of
988 genomes assessed per column is shown above. Colours indicate match confidence scores, and
989 the key is shown on the right. Match confidence is a categorical measure of match quality
990 between query and reference sequence in the database, and definitions for each category can
991 be found in [12]. **(B)** Comparison of ‘problem’ scores visualised as a bar plot. Database
992 version is indicated by the colour key shown on the right, and type of problem score is shown
993 below. Figures were created using ggplot2 package in RStudio [46].

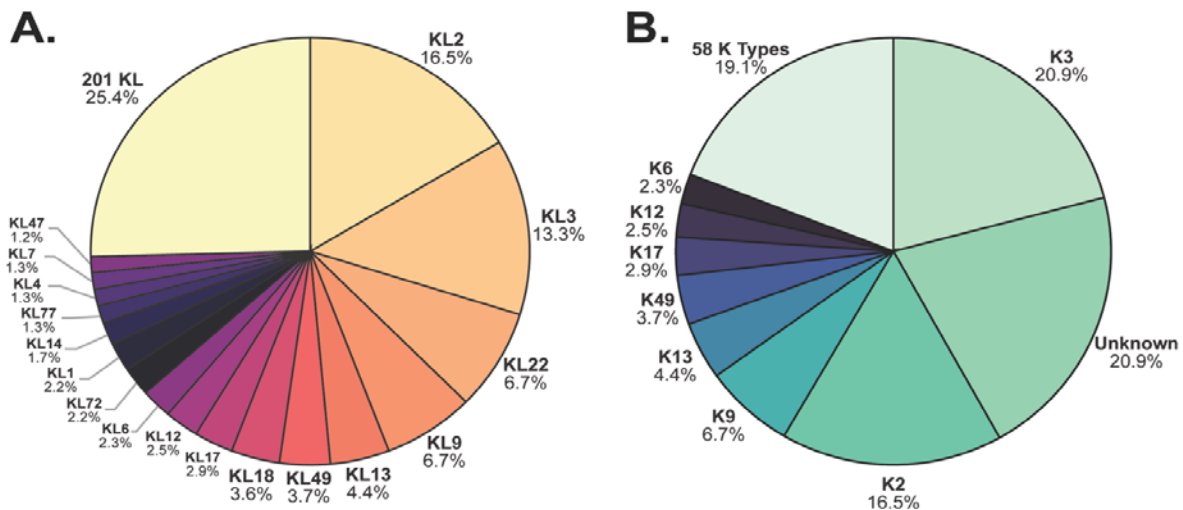
994



995

996

997 **Figure 4.** General features of the 237 K loci identified in *A. baumannii* genomes. **(A)**
 998 Percentage of 237 KL with specific total sequence length (kb). **(B)** Percentage of 237 KL
 999 with specific number of Open Reading Frames (ORFs). **(C)** Percentage of 237 KL with
 1000 specific number of glycosyltransferase genes (*gtrs*). **(D)** Frequency of 681 gene types
 1001 (homology groups) found at the K locus across 237 KL. Figures were created in RStudio
 1002 [46].



1003
1004 **Figure 5.** Distribution of (A) best match locus and (B) best match type amongst the 8994
1005 genome assemblies assigned a confidence score of 'good' or above by *Kaptive*.
1006 Figures were created in RStudio [46].