

Identification of *Acinetobacter baumannii* loci for capsular polysaccharide (KL) and lipooligosaccharide outer core (OCL) synthesis in genome assemblies using curated reference databases compatible with *Kaptive*

Kelly L. Wyres¹, Sarah M. Cahill², Kathryn E. Holt^{1,3}, Ruth M. Hall⁴ and Johanna J. Kenyon^{2,*}

Abstract

Multiply antibiotic-resistant *Acinetobacter baumannii* infections are a global public health concern and accurate tracking of the spread of specific lineages is needed. Variation in the composition and structure of capsular polysaccharide (CPS), a critical determinant of virulence and phage susceptibility, makes it an attractive epidemiological marker. The outer core (OC) of lipooligosaccharide also exhibits variation. To take better advantage of the untapped information available in whole genome sequences, we have created a curated reference database of 92 publicly available gene clusters at the locus encoding proteins responsible for biosynthesis and export of CPS (K locus), and a second database for 12 gene clusters at the locus for outer core biosynthesis (OC locus). Each entry has been assigned a unique KL or OCL number, and is fully annotated using a simple, transparent and standardized nomenclature. These databases are compatible with *Kaptive*, a tool for *in silico* typing of bacterial surface polysaccharide loci, and their utility was validated using (a) >630 assembled *A. baumannii* draft genomes for which the KL and OCL regions had been previously typed manually, and (b) 3386 *A. baumannii* genome assemblies downloaded from NCBI. Among the previously typed genomes, *Kaptive* was able to confidently assign KL and OCL types with 100% accuracy. Among the genomes retrieved from NCBI, *Kaptive* detected known KL and OCL in 87 and 90% of genomes, respectively, indicating that the majority of common KL and OCL types are captured within the databases; 13 of the 92 KL in the database were not detected in any publicly available whole genome assembly. The failure to assign a KL or OCL type may indicate incomplete or poor-quality genomes. However, further novel variants may remain to be documented. Combining outputs with multilocus sequence typing (Institut Pasteur scheme) revealed multiple KL and OCL types in collections of a single sequence type (ST) representing each of the two predominant globally distributed clones, ST1 of GC1 and ST2 of GC2, and in collections of other clones comprising >20 isolates each (ST10, ST25, and ST140), indicating extensive within-clone replacement of these loci. The databases are available at <https://github.com/katholt/Kaptive> and will be updated as further locus types become available.

DATA SUMMARY

1. Databases including fully annotated gene cluster sequences for *A. baumannii* K loci and OC loci are available for download at <https://github.com/katholt/Kaptive>
2. Details of the *Kaptive* search results validating *in silico* serotyping of K and OC loci using our approach are provided

as supplementary files, available in the online version of this article: Dataset 1 (92 KL reference sequences and 12 OCL reference sequences), Dataset 2 (642 genomes assembled from reads available in NCBI SRA) and Dataset 3 (3415 genome assemblies downloaded from NCBI GenBank).

Received 09 December 2019; Accepted 28 January 2020; Published 02 March 2020

Author affiliations: ¹Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Australia; ²Institute of Health and Biomedical Innovation, School of Biomedical Sciences, Faculty of Health, Queensland University of Technology, Brisbane, Australia; ³Department of Infection Biology, London School of Hygiene and Tropical Medicine, London, UK; ⁴School of Life and Environmental Sciences, The University of Sydney, Sydney, Australia.

***Correspondence:** Johanna J. Kenyon, johanna.kenyon@qut.edu.au

Keywords: *Acinetobacter baumannii*; *Kaptive*; capsular polysaccharide; K locus; outer-core oligosaccharide; OC locus.

Abbreviations: CPS, capsular polysaccharide; GC, global clone; IS, insertion sequence; KL, K locus; LOS, lipooligosaccharide; LPS, lipopolysaccharide; MLST, multilocus sequence typing; OCL, outer-core locus; OPS, O-antigen polysaccharide; ST, sequence type.

Data statement: All supporting data, code and protocols have been provided within the article or through supplementary data files. Supplementary material is available with the online version of this article.

000339 © 2020 The Authors



This is an open-access article distributed under the terms of the Creative Commons Attribution License.

INTRODUCTION

One of the most imminent global health crises is the increasing prevalence and global dissemination of highly resistant bacterial pathogens that are able to persist in hospital environments despite infection control procedures. In 2017, the World Health Organization identified carbapenem-resistant strains of the opportunistic Gram-negative bacterium *Acinetobacter baumannii* as a critical priority for therapeutics development due to alarming levels of resistance against nearly all clinically suitable antibiotics [1]. The success of extensively antibiotic-resistant *A. baumannii* isolates can be attributed, in part, to the evolution and expansion of well-adapted clonal lineages [2–5], including the two major globally disseminated clones, Global Clone 1 (GC1) and Global Clone 2 (GC2), and other lineages that are found less frequently (e.g. sequence type 25; ST25) or on only one or two continents (e.g. ST78) [6]. Hence, precise epidemiological tracking methods for *A. baumannii* isolates, in particular those from important clonal lineages, are urgently needed to enhance surveillance and improve our understanding of how *A. baumannii* circulates both locally and globally.

Traditionally, epidemiological studies tracing important bacterial lineages associated with human and animal infections have used serological typing of the polysaccharides produced on the cell surface [7, 8], as there can be significant variation in structures observed on different isolates of the same species [9–12]. The cell-surface polysaccharides targeted in these schemes included capsular polysaccharide (CPS, K, or capsule) and/or O-antigen polysaccharide (OPS or O) that is attached to lipooligosaccharide (LOS) forming a lipopolysaccharide (LPS). In early studies, an *A. baumannii* serological typing scheme was developed for a major immunogenic polysaccharide, believed at the time to be the O antigen [13, 14], and 38 different serovars were included in the last update to the scheme nearly two decades ago [15]. However, this system is no longer used.

In the last decade, it has been shown that the major immunogenic polysaccharide produced by the species is CPS not O antigen [16–18]. The CPS of *A. baumannii* is a major virulence determinant as isolates lacking CPS do not cause infections [17]. CPS is also a key target of potential novel control strategies including phage therapy [19, 20] and vaccinations [21, 22]. Unfortunately, the current lack of knowledge about capsule diversity and epidemiology in the broader *A. baumannii* population, and a lack of tools to readily detect changes in the population distribution hinder effective design of these controls.

Most of the genes that direct the synthesis of the CPS are clustered at the K locus (KL) that is located between the *fkpA* and *lldP* genes in the *A. baumannii* chromosome [16, 23]. The general arrangement of the K locus features three main regions (Fig. 1a). On one side, a module of genes for CPS export machinery (*wza*–*wzb*–*wzc*) are in a separate operon, divergently transcribed from the remainder of the gene cluster. On the other side lies a module of genes involved in the synthesis of simple sugar substrates. However, the *gneI*

Significance as a BioResource to the community

The ability to identify and track closely related isolates is key to understanding, and ultimately controlling, the spread of multiply antibiotic-resistant *A. baumannii* causing difficult to treat infections, which are an urgent public health threat. Extensive variation in the KL and OCL gene clusters responsible for biosynthesis of capsule and the outer core of lipooligosaccharide, respectively, are potentially highly informative epidemiological markers. However, clear, well-documented identification of each variant and simple-to-use tools and procedures are needed to reliably identify them in genome sequence data. Here, we present curated databases compatible with the available web-based and command-line Kaptive tool to make KL and OCL typing readily accessible to assist epidemiological surveillance of this species. As many bacteriophages recognize specific properties of the capsule and attach to it, capsule typing is also important in assessing the potential of specific phages for therapy on a case by case basis.

gene can be lost (e.g. Fig. 2b) if N-acetyl-D-galactosamine (D-GalpNAc) is not present in the CPS, and various other genes have been found between *gne* (or *gpi*) and *pgm* in some K loci [24–26]. The genetic content of the central region is specific to the CPS structure produced. It includes genes for the required number of glycosyltransferases, and the capsule processing genes (*wzx* and *wzy*). If complex sugars (pseudaminic acid, legionaminic acid, acinetaminic acid, bacillosamine, etc.) are included in the CPS, the central region will also contain genes for the synthesis and modification of these sugars [16, 27–30]. Each distinct gene cluster, defined by a difference in gene content between *fkpA* and *lldP*, is assigned a unique identifying number (KL1, KL2, etc.). To date, more than 128 KL gene clusters (KL types) have been identified at the KL in *A. baumannii* genomes [31].

A clear nomenclature system for CPS biosynthesis genes in *A. baumannii* was developed in 2013 to identify the specific function of KL-encoded proteins for the non-expert [16]. Where possible, gene names indicate enzyme function (i.e. Gtr assigned to GlycosylTTransferases and Itr to the transferases initiating K unit synthesis). For enzymes (e.g. Gtrs and Itrs) where sequence differences probably result in a change of substrate preference, a number indicating the different types (cut-off value of 85% amino acid sequence identity) is included in the name as a suffix. The current gene names are listed in Table 1. Most published annotations use this system (e.g. [26, 29–36]). However, sometimes other nomenclature systems have been used [23, 37].

A second locus with variable gene content involved in the production of a surface polysaccharide [16] has been shown to be responsible for synthesis of the outer-core (OC) component of the LOS [38]. The OC locus (OCL) is located in the

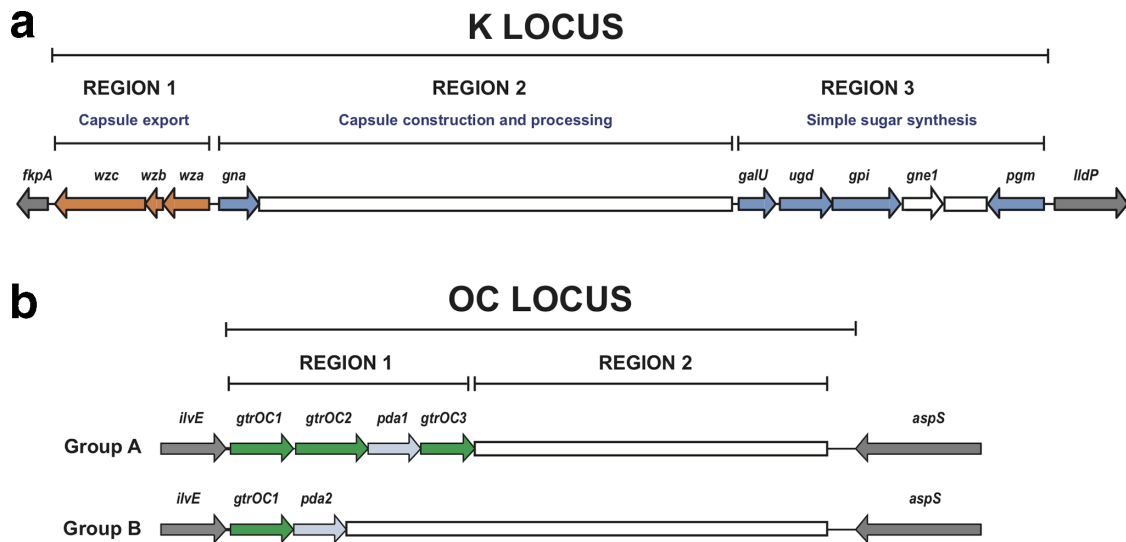


Fig. 1. General arrangement of the surface polysaccharide synthesis loci in *A. baumannii*. KL and OCL boundaries are shown and flanking locus genes are coloured grey. Variable sequence portions are indicated by white boxes, and conserved genes at each locus are represented by coloured arrows. (a) Organization of the KL with marked regions defining the roles of common modules. CPS export genes are in orange, and genes in dark blue are involved in the synthesis of common sugar substrates. *gne1* is not always present but is often critical to the synthesis of many CPS structures. (b) Organization of the two groups (A and B) of the OC locus with marked regions defining conserved or variable portions. Genes in green encode conserved glycosyltransferases and genes in light blue are those involved in complex sugar synthesis.

chromosome between the *aspS* and *ilvE* genes [16, 39]. Each distinct gene cluster found between the flanking genes is assigned a unique number identifying the locus type (OCL1, OCL2, etc.), and to date, 14 different gene clusters (OCL1–12 [39] and OCL15 and 16 [40]) have been reported. The nomenclature for OCL genes is also shown in Table 1, and Gtrs encoded at the OCL are differentiated from KL-encoded Gtrs by the addition of OC to the name (GtrOC#). Generally, OC gene clusters fall into two broad families (Fig. 1b), designated Group A and Group B, defined by the presence of *pda1* and *pda2* genes, respectively [39].

Several studies have highlighted the extremely plastic nature of the *A. baumannii* genome, revealing very poor correlation between KL and OCL types and other genomic features including ST [2, 16, 23, 41–43]. Therefore, the most valuable framework for tracing important genetic lineages of *A. baumannii* currently involves a combination approach, including phylogenetic analysis with multilocus sequence typing (MLST) using both Institut Pasteur and Oxford schemes, resistance and virulence gene mapping, and KL and OCL typing [2, 40–43]. Bioinformatics tools and databases currently exist for MLST and resistance gene typing, allowing multiple genomes to be processed quickly. However, the lack of computational tools and databases to rapidly extract interpretable, actionable information about K and OC loci from large data sets is a current bottleneck.

Recently, a computational tool, named *Kaptive*, was developed to rapidly identify reference K and O loci in *Klebsiella pneumoniae* species complex genome sequences taking as input a curated database of reference sequences and a query

genome assembly [44, 45]. Although the computational tool can be used to type loci in any species, a complete and curated compendium of appropriate, species-specific KL, OL or OCL sequences is needed. In *A. baumannii*, such databases are not currently available.

Here, we present curated databases of annotated reference sequences for *A. baumannii* K and OC loci that are compatible with *Kaptive*, enabling rapid typing of genomes for this clinically significant pathogen. We evaluate the accuracy of this approach by comparison of K and OC locus calls for >630 genomes typed previously using manual methods. Additionally, we apply this approach to type >3300 *A. baumannii* genomes retrieved from the NCBI database, highlighting the extent of K and OC locus variability in the broader population and among clinically important clonal complexes, and confirming that the vast majority of genomes harbour loci matching those in our reference databases.

METHODS

K and OC reference sequences

Nucleotide sequences for reference isolates carrying each KL and OCL type were downloaded from NCBI non-redundant or WGS databases (accession numbers are listed in Tables S1 and S2). Where possible, whole genome sequences were assessed for the presence of the *A. baumannii*-specific *oxaAb* gene (also referred to as *bla*_{OXA-51}; GenBank accession number CP010781.1, base positions 1753305–1754129) to confirm the sequences were obtained from an *A. baumannii* isolate. A GenBank format file (.gbk) for each distinct locus type was

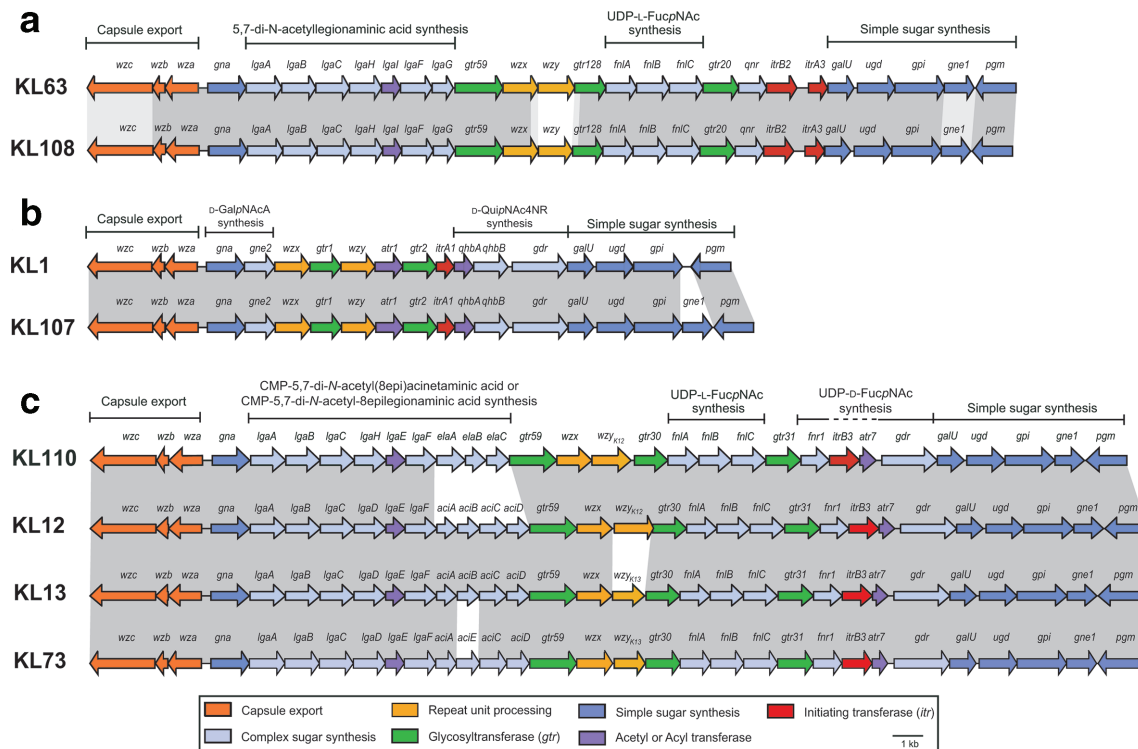


Fig. 2. Closely related capsule biosynthesis gene clusters demonstrating cases of small genetic replacements. Genes are represented by arrows orientated in the direction of transcription that are coloured according to the scheme shown below. Shading between gene clusters indicates regions of >95% nucleotide sequence similarity (dark grey) or 90–95% nucleotide sequence similarity (light grey). Figure drawn to scale using the GenBank accession numbers listed in Table S1. (a) KL63 and KL108 gene clusters differing in *wzy* sequence. (b) KL1 and KL107 are examples of *gne1* presence vs. absence. (c) KL13, KL73, KL12 and KL110 are examples of several closely related gene clusters with small sequence replacements altering the synthesis pathway of a complex sugar substrate, or topology of the CPS structure.

prepared. This file includes the nucleotide reference sequence for the locus without flanking sequence, the annotations of all coding sequences in the locus, and citation(s) for the annotations and/or polysaccharide structural data, if available.

Curated *Kaptive* databases

The individual KL files were concatenated into a multirecord GenBank-format file to produce a data set containing annotated KL reference sequences. Likewise, the OCL files were compiled to generate a separate data set. Both reference databases were integrated with the *Kaptive-Web* platform (<http://kaptive.holtlab.net/>), which enables users to submit their genome sequence queries to a browser and receive the output in a visual format, as described in detail previously [45]. The KL and OCL databases have also been made freely available for download from the *Kaptive* github repository (<https://github.com/katholt/Kaptive>) for use with the command-line version of *Kaptive* [44], or other tools.

Genome sequence collections

Acinetobacter genome assemblies from our collection for which the KL and OCL types had been previously determined via manual or automated sequence inspection ([2, 41]; and

unpublished data) were used to assess the level of typing accuracy that could be achieved through the use of our novel databases with *Kaptive*. Paired-end Illumina read data (described in [2, 41] and available under BioProject accession PRJEB2801) were *de novo* assembled using SPAdes v. 3.13.1 [46] and optimized with Unicycler v. 0.4.7 [47]. High-quality genome assemblies ($n=719$) with a maximum contig number of 300 and minimum assembly length of 3.6 Mbp were included in the analysis (cut-offs determined empirically by manual inspection of the contig number and assembly length distributions, respectively). These assemblies were assessed for *oxaAb* presence using BLASTN (>95% nucleotide sequence similarity and >90% combined coverage) to confirm the *A. baumannii* species assignment. Confirmed *A. baumannii* sequences ($n=642$) were analysed using both KL and OCL reference databases with command-line *Kaptive* v. 0.7 [44] with default parameters.

The same method was used to test databases against 3412 genome sequences available in the NCBI non-redundant and WGS databases as of February 2019. These genome assemblies were bulk downloaded from NCBI as a compressed .tar file for local analysis. Genomes lacking *oxaAb* were removed prior to typing but quality control (QC) analysis as described above was applied to this data set only after typing was complete.

Table 1. Gene nomenclature key for *A. baumannii* K and OC loci

Gene name	Predicted reaction product	Predicted protein
K locus		
<i>aci</i>	CMP- <u>A</u> cetaminic acid 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-di <u>a</u> cetamido-2,3-dideoxy-D-glucuronic <u>a</u> cid	Multiple
<i>dmaA</i>	UDP-2,3-di <u>a</u> cetamido-2,3-dideoxy-D-mannuronic acid	2-epimerase
<i>ela</i>	CMP-8-epi <u>l</u> egionaminic acid 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>gna</i>	UDP-D-GlcpNAcA	UDP-D-GlcpNAc dehydrogenase
<i>gne1</i>	UDP-D-GalpNAc	UDP-D-GlcpNAc epimerase
<i>gne2</i>	UDP-D-GalpNAcA	UDP-D-GlcpNAcA epimerase
<i>gpi</i>	L-Fructose-6-phosphate	glucose-6-phosphate isomerase
<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 acid 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> hospho <u>t</u> hanolamine transferase
<i>pgm</i>	D-Glucose-1-phosphate	Phosphoglucosyltransferase
<i>pgt</i>	–	<u>P</u> hosphoglycerol transferase
<i>psa</i>	CMP- <u>P</u> seudaminic acid derivative	Multiple
<i>ptr</i>	–	<u>P</u> yruvyl <u>t</u> ransferase
<i>qdt</i>	dTDP-D-Quip3NAc	Multiple

Continued

Table 1. Continued

Gene name	Predicted reaction product	Predicted protein
<i>qhb</i>	UDP-D-QuipNAc4NHb	Multiple
<i>qnr</i>	UDP-D-QuipNAc	UDP-6-deoxy-4-keto-D-GlcpNAc 4- <u>r</u> eductase
<i>rml</i>	dTDP-L-rhamnose	Multiple
<i>tle</i>	dTDP-6-deoxy-L-talose	dTDP-L-rhamnose epimerase
<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
OC locus		
<i>ahy</i>	–	Predicted acylhydrolase
<i>gtrOC</i>	–	<u>G</u> lycosyl <u>t</u> ransferase (outer core)
<i>pda</i>	UDP-D-GlcN	Polysaccharide deacetylase
<i>ptrOC</i>	–	<u>P</u> yruvyl transferase (outer core)
<i>wecB</i>	UDP-D-ManpNAc	UDP-D-GlcpNAc C2 epimerase

Interpretation of *Kaptive* output

The *Kaptive* output is described in detail elsewhere [45]. Briefly, *Kaptive* uses a combination of BLASTN and TBLASTN searches to identify the best matching reference locus for each query genome and indicates a corresponding confidence level. The latter is dependent on the BLASTN coverage and identity for the full-length reference locus, the number of reference locus genes (expected genes) or other genes (unexpected genes) found within the locus region of the query genome (determined by TBLASTN, default coverage cut-off $\geq 90\%$, identity $\geq 80\%$), and whether the locus is found on a single or multiple assembly contigs. A ‘perfect’ confidence match indicates that the locus was found in the query genome on a single contig with 100% coverage and 100% nucleotide identity to the best-match reference locus. ‘Very high’ confidence matches are those for which the locus is present in the query genome in a single assembly contig with $\geq 99\%$ coverage and $\geq 95\%$ nucleotide sequence identity to the best-match reference locus, and no missing or unexpected genes within the locus. ‘High’ confidence matches are defined as those for which the locus was found on a single contig with $\geq 99\%$ coverage to the best-match reference locus, ≤ 3 missing genes

and no unexpected genes within the locus. ‘Good’ confidence matches indicate that the locus was found on a single contig or split across multiple assembly contigs with $\geq 95\%$ coverage to the best-match locus, ≤ 3 missing genes and ≤ 1 unexpected gene within the locus. ‘Low’ confidence matches indicate that the locus was found on a single contig or split across multiple assembly contigs with $\geq 90\%$ coverage to the best-match locus, ≤ 3 missing genes and ≤ 2 unexpected genes within the locus. A confidence level of ‘None’ indicates that the match does not meet the criteria for any other confidence level.

Distribution of K and OC loci

For NCBI genome assemblies, STs were assigned with the mlst script (github.com/tseeman/mlst) using the Insitut Pasteur scheme for *A. baumannii* (abaumannii_2 scheme) available at https://pubmlst.org/bigdb?db=pubmlst_abaumannii_pasteur_seqdef. KL and OCL variation were visualized for STs with ≥ 20 isolate representatives with ‘good’ or better confidence matches called by *Kaptive*.

RESULTS

KL and OCL numbering and nomenclature

The development of curated databases for numbered and fully annotated *A. baumannii* K and OC loci relies on the consistent application of a standardized nomenclature and numbering system for these loci. Here, the system developed for transparent annotation of both the K and the OC loci [16] has been used. As new KL and OCL types with additional gene families have been discovered since 2013, the gene nomenclature has been extended and is summarized in Table 1. For consistency, K loci that were originally published using other nomenclatures or typing systems have been re-annotated, and where possible the corresponding GenBank entries have been updated with the permission of the original authors (see Table S1).

In several cases, KL types that differ only by a small portion of the locus have been found (e.g [16, 48]), and examples are shown in Fig. 2. In cases where structures have been determined, the locus difference is associated with changes in the composition or structure of the CPS [26, 27, 29, 31, 35, 49–54] but some locus differences are now known to have no effect on CPS structure [24, 55]. As all differences in genetic content are relevant in epidemiological studies, all K loci comprising a unique combination of genes were distinguished with a new KL number.

The curated KL reference database

Curated annotations already existed in NCBI GenBank for 78 KL types, three of which were submitted as third-party annotations (TPAs) (see Table S1). Here, an additional 14 sequences were extracted from genome sequences available in the NCBI WGS database (see Table S1) giving a total of 92. Although we have identified 128 KL types in total (J. J. Kenyon, unpublished), sequences for the remaining 36 KL

types (128 minus 92) are not currently available in the public domain.

Complete annotations for the 92 publicly available *A. baumannii* KL reference sequences spanning the full length of each gene cluster (between *fkpA* and *lldP*) were therefore compiled into a KL reference database for use with *Kaptive*. Where the only available representative of a KL type included an insertion sequence (IS), we substituted the sequence with a manually generated version with the IS and target site duplication removed in order to include a KL that represents the presumptive ancestral, non-modified sequence as is required for accurate typing by *Kaptive* [44]. This was the case for KL types KL27, KL44, KL82, KL87, KL93, KL114 and KL118 (Table S1).

The curated OCL reference database

The annotations for 12 different OCL types have been described in the literature [39]. A complete list is found in Table S2. However, only six of them were available in GenBank. The remaining six OCL sequences were identified in the WGS database, and the WGS accession numbers are given in Table S1. Complete annotations for the 12 publicly available OCL spanning the full length of the gene clusters (between *ilvE* and *aspS*) were combined into a single OCL reference database for use with *Kaptive*.

Compatibility of the KL and OCL databases with *Kaptive*

To confirm the compatibility of the KL and OCL databases for *Kaptive*-based typing, we created two query sequence sets comprising FASTA sequences of the reference KL and OCL, respectively. *Kaptive* was applied to each of these query sets, and was able to successfully identify the correct locus in all cases (Dataset 1).

Comparison of *Kaptive* assignments with previous KL assignments

We assessed the accuracy of *Kaptive*-based KL typing using our curated KL database by application to a collection of 642 *A. baumannii* genome assemblies (see Dataset 2), which had been typed previously using BLASTN plus manual inspection ([2, 41]; and unpublished data). For these assemblies, the confidence levels called by *Kaptive* were: 176 (perfect), 385 (very high), 28 (high), 53 (good), 0 (low) and 0 (none) (Fig. 3a; Dataset 2). Notably, 561 matches were assigned ‘perfect’ or ‘very high’ confidence calls, demonstrating that *Kaptive* could very confidently assign a KL type to the majority (87.4%) of the 642 genome assemblies provided.

The 28 ‘high’ confidence matches each included one or more single base deletions within the locus leading to the interruption of a coding sequence, which *Kaptive* reports as one or more missing genes when the resulting TBLASTN matches have $< 90\%$ coverage to the reference gene sequence. Such deletions may represent sequencing and/or assembly errors but may also represent true sequence variations with the potential to result in altered CPS structure. Because *Kaptive* is

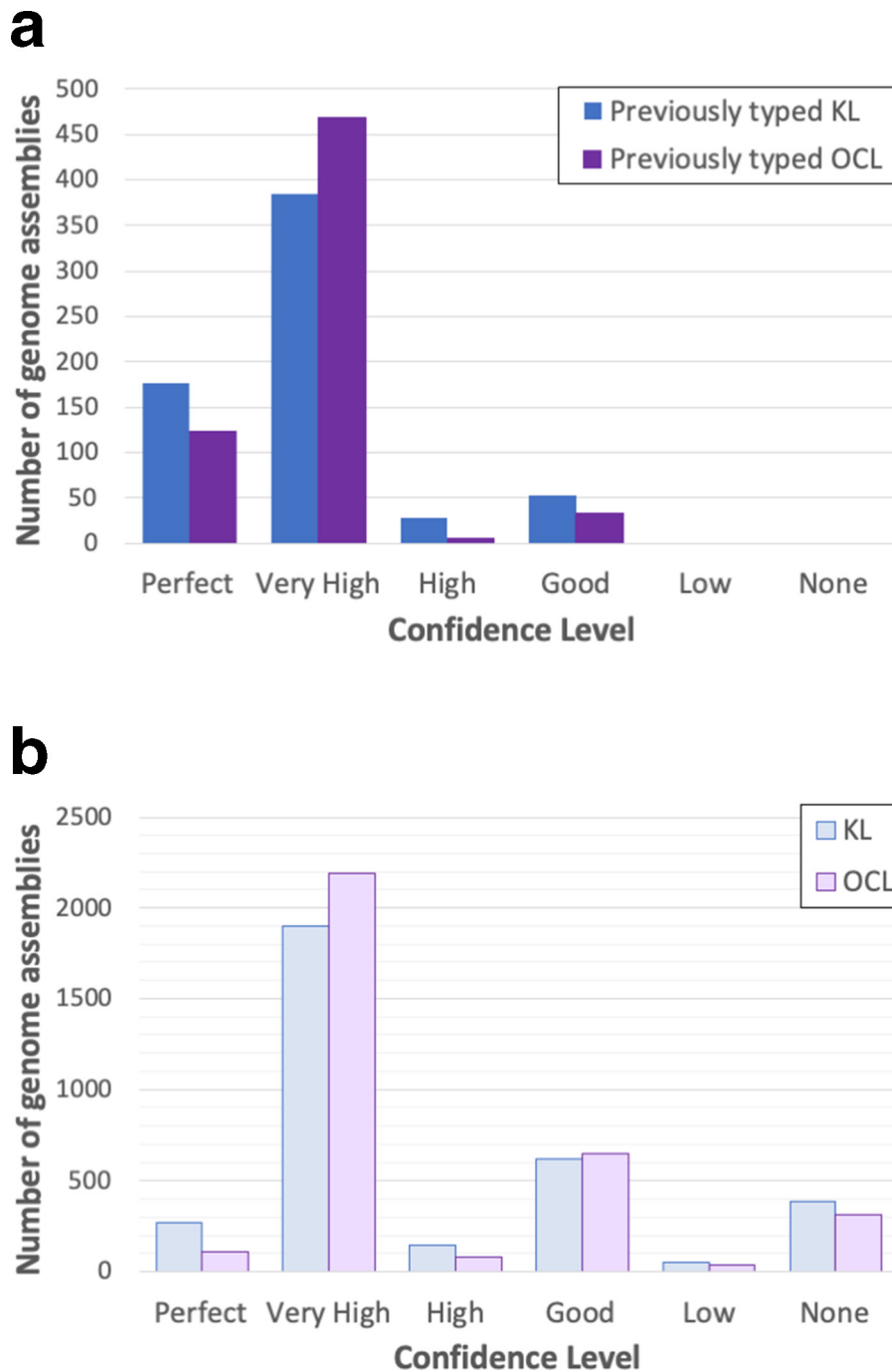


Fig. 3. Breakdown of confidence levels for *Kaptive* locus calls using the *A. baumannii* KL and OCL databases. (a) Results following database quality checking using a private collection of 680 *A. baumannii* genome assemblies (Dataset 2). Colour key is shown in the top right corner. (b) Results of applying the databases to 3412 genome assemblies available in NCBI databases (Dataset 3). Colour key is shown in the top right corner.

unable to distinguish these possibilities it reports the ‘missing’ gene and lowered confidence score in order to alert the user and facilitate further investigation.

Manual inspection of the relevant assembly graphs showed that 50 of 53 (94.4%) assignments with a ‘good’ confidence

level were locus variants in which an IS had interrupted the KL gene cluster, breaking it into two or more contigs in the query genome. The three remaining assemblies that were typed with a ‘good’ confidence level were also broken into multiple contigs that represented dead-ends in the assembly

graphs, and hence it was not possible to determine if these also represented IS variants or were simply the result of assembly problems (e.g. due to low sequencing depth in the KL region of the genome).

Of the 642 assemblies with a KL type that was assigned previously, 641 (99.8%) were concordant and one (0.2%) was discrepant. The KL of *A. baumannii* isolate BAL_266 had previously been described as KL63 [41]. However, *Kaptive* assigned it to KL108 with a 'very high' confidence level (99.98% nucleotide sequence similarity; 100% coverage). The sequence of this isolate was manually checked again and confirmed to be KL108. The KL63 and KL108 gene clusters are 97.96% identical across 95% of the locus, differing from each other only in a ~1.3kb segment in the central region that includes the *wzy* gene (Fig. 2a). This small difference between the two gene clusters was missed in the original manual typing but probably alters the linkage between the K units. This highlights the need to look for any regions of sequence difference when manually typing.

Comparison of *Kaptive* assignments with previous OCL assignments

We also assessed the accuracy of OCL identification using our curated OCL database applied to the same collection of *A. baumannii* genomes. The OCL region of 631 of these had previously been typed using BLASTN plus manual inspection ([2, 41]; and unpublished data). The confidence levels for the OCL matches for the 631 typed genomes were: 124 (perfect), 469 (very high), 5 (high), 33 (good), 0 (low) and 0 (none) (Fig. 3a; Dataset 2). As for the KL database, the large number of 'perfect' and 'very high' confidence matches (593, 94.0%) demonstrates the capacity of the OCL database to type the majority of genome assemblies provided as a query. Manual inspection confirmed that the five 'high' confidence matches included those with one or more base deletions in coding sequences, and the 33 'good' matches represented variants of the corresponding reference sequences interrupted by one or more ISs. In this set, there were no discrepancies between the previous locus assignments and those determined by *Kaptive*.

Application of KL and OCL databases for *A. baumannii* genome typing

As the KL and OCL regions in the majority of NCBI genome sequences have not yet been examined, the publicly available genomes provide a large dataset to begin to explore KL and OCL diversity in the species. Available genome assemblies of 3412 isolates annotated as *A. baumannii* in the NCBI non-redundant and WGS databases were first checked for the presence of the *oxaAb* gene to ensure correct assignment to the *A. baumannii* species. The *oxaAb* gene was absent from 34 assemblies (0.99%), and these were removed from the analysis, bringing the total number of assemblies examined to 3378.

For the KL database, the confidence levels of the matches called by *Kaptive* were: 272 (perfect), 1901 (very high), 149 (high), 622 (good), 51 (low) and 383 (none). Among the 2944

genomes with KL confidence matches 'good' or better, there were 79 distinct KL types, 36 (45.6%) of which were identified in five or fewer genomes. Notably 13 of the loci included in the KL reference database were not identified among any of the genome assemblies retrieved from the NCBI database. The most common KL types were KL2 (713 of 2948 genomes, 24.2%), KL9 (343, 11.6%), KL22 (330, 11.2%), KL3 (294, 10.0%) and KL13 (155, 5.3%).

For the OCL database, the confidence levels were as follows: 108 (perfect), 2192 (very high), 80 (high), 645 (good), 39 (low) and 314 (none) (Fig. 3b; Dataset 3). All 12 of the reference OC loci were identified among the 3029 genomes with OCL confidence matches 'good' or better. Among these genomes the most common OCL types were OCL1 (2086, 68.9%), OCL3 (272, 9.0%), OCL6 (157, 5.2%), OCL2 (150, 5.0%) and OCL5 (125, 4.1%).

Therefore, among the *A. baumannii* genomes retrieved from NCBI, KL and OCL calls were obtained for 87 and 90% of the assemblies, respectively. However, the 'low' and 'none' confidence levels may result from poor-quality sequence assembly and/or may indicate that a novel locus is present in the query assembly [44]. Indeed, the application of the same quality control cut-off used for inclusion in our own data set (see above) revealed that 13/51 'low' and 174/387 'none' confidence matches for the KL assignments may be assemblies of poor quality. Similarly, 12/39 'low' and 76/314 'none' confidence matches for the OCL assignments did not meet the same quality control cut-off. Hence, it is recommended that users perform additional investigations to confirm the quality of their assemblies before excluding 'low' and/or 'none' confidence matches from their analyses.

KL and OCL variation in clonal lineages

Variations in the KL and OCL in the major multidrug-resistant clonal lineages have largely been examined using small data sets (e.g. [2, 16, 38, 39]). For the GC2 lineage, these studies assessed diversity amongst isolates predominantly recovered from the same outbreak or region [41–43] or sporadic isolates [26, 56], limiting the ability to gain a complete picture of surface polysaccharide variation in this clone. Across these studies, at least 14 KL and five OCL have been reported in GC2. Of the 3386 genome assemblies we analysed here, 2016 belonged to ST2 in the Institut Pasteur scheme, representing the most common ST in GC2 and the largest group of isolates belonging to a single ST [6]. Among the 2016 ST2 genomes, *Kaptive* identified 30 KL and three OCL (Fig. 4) in those with confidence matches 'good' or better. The most common KL arrangements were KL2 (32.2%) and KL22 (14.4%), whereas OCL1 represented the most predominant OCL type (78.6%). Only one KL, KL63, was found in a single ST2 genome. For the remaining assemblies, 107 (5.3%) and 256 (12.7%) were assigned 'low' and 'none' confidence matches against the KL and OCL databases, respectively. These assemblies may be of poor quality or they may carry novel types, but this was not further investigated.

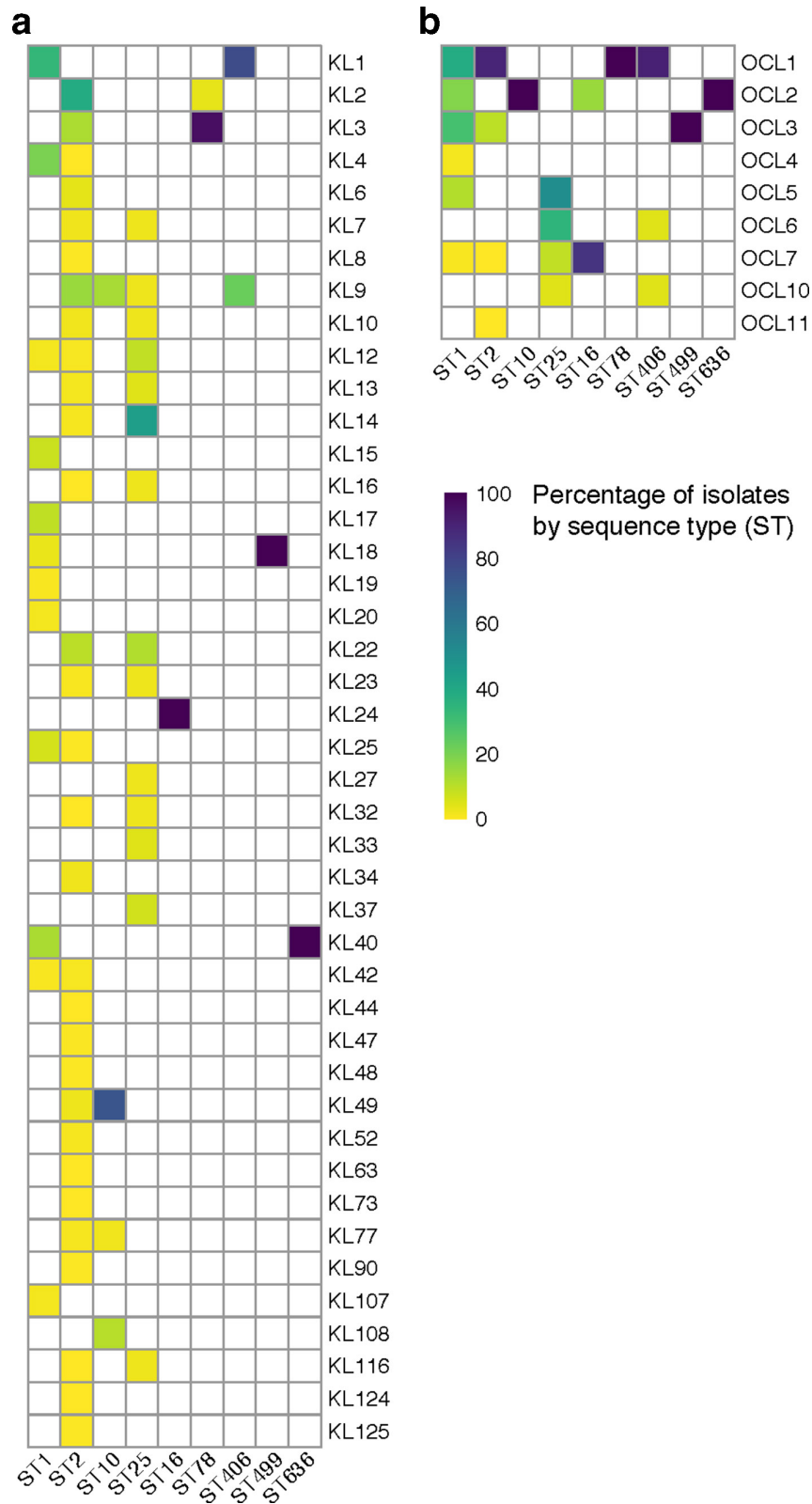


Fig. 4. Distribution of K and OC loci by sequence type. Heat maps show the distribution of distinct K (a) and OC (b) loci among genomes assigned to nine common multilocus sequence types (STs). Coloured shading indicates the percentage of isolates belonging to a given ST that were assigned a given KL or OCL type, as indicated by the colour legend. A. *baumannii* genome assemblies were retrieved from the NCBI database; only confirmed A. *baumannii* for which both K and OC loci were assigned by *Kaptive* with confidence level 'good' or better are shown ($n=2002$; 125 ST1, 1669 ST2, 46 ST10, 20 ST16, 43 ST25, 28 ST78, 22 ST406, 29 ST499, 20 ST636).

KL and OCL diversity have also previously been reported for the other major clonal lineage, GC1. An in-depth study of 45 *A. baumannii* GC1 isolates identified eight KL and five OCL types in this clone [2], with one additional KL type found in a subsequent study [57]. In the set of 3386 genome assemblies, we found 134 that belong to ST1, which represents the most common GC1 ST. *Kaptive* identified a total of 10 KL and six OCL types in the ST1 lineage (Fig. 4), expanding the number of distinct types observed previously. Among these ST1 genomes, the most common KL types were KL1 (31.3%) and KL4 (18.7%), while the most common OCL types were OCL1 (36.6%), OCL2 (17.2%) and OCL3 (29.1%). KL19 and KL42 and also OCL7 were found in single isolates.

We also examined a further seven STs for which there were ≥ 20 isolate representatives with confident *Kaptive* matches ('good' or better). Of these STs, ST10 included the largest number of genome assemblies (47 of 3386 assemblies), and four KL types and one OCL type were found in this group. ST25, the second largest group with 46 assemblies, had very high variation with 14 KL and four OCL types. ST406 (22 assemblies) also included four OCL types but only two KL types. However, one of two KL types and one or two OCL types were found in ST16 (20 assemblies), ST78 (29), ST499 (29) and ST636 (20). Notably, specific KL and OCL types were not confined to single STs, with several locus types found in more than one ST.

DISCUSSION

In this study, we present *Kaptive*-compatible databases of annotated reference sequences for *A. baumannii* K and OC loci, extending the utility of *Kaptive* and broadening the ability of researchers, clinicians and public health professionals to analyse genome data sets. Using these databases, *Kaptive* was able to confidently and accurately assign KL and OCL types to the majority of *A. baumannii* genome assemblies examined. Among >630 *A. baumannii* genomes typed previously using manual methods, only a single discrepancy between the previous KL assignment and that of *Kaptive* was identified. This was traced to an error in the previous manual assignment, which had overlooked a small genetic replacement within the locus. As sequence replacements of <2 kb are common in *A. baumannii* KL and OCL regions (examples shown in Fig. 2), the ability of *Kaptive* to correctly identify the KL type demonstrated the stringent nature of the tool and the quality of the databases described herein. The KL and OCL databases were also used to probe the collection of *A. baumannii* genome assemblies available through the NCBI GenBank and WGS databases. *Kaptive* was able to confidently assign locus types to more than 87% of these genome assemblies, indicating that the databases capture the majority of common KL and OCL types. However, to confirm the locus calls, all *Kaptive* assignments should be checked for length discrepancies that would reveal missing expected genes, and/or the presence of additional genes or ISs in the locus.

The remaining genomes that could not be confidently assigned a locus type (13% KL and 10% OCL unassigned) may include genomes with low coverage and/or poor assembly quality in the KL and/or OCL genome regions. Alternatively, these genomes may carry loci that are not represented in the current reference databases. In these cases, users are encouraged to undertake further investigations, such as by manual inspection of the assembly and/or assembly graphs and comparisons to the best-matching reference loci using visualization tools such as the Artemis Comparison Tool [58] and Bandage [59]. Further work will be needed to identify and include further novel loci, and the databases will be continuously updated as sequences and annotations for further KL and OCL types become available. We encourage users to contact us via the *Kaptive-Web* website and/or the *Kaptive* github page to submit novel loci for the assignment of KL and OCL numbers and addition to the publicly available databases.

The typing system and the databases have been designed strictly for use in *A. baumannii* and therefore users are encouraged to check the origin of their sequences to ensure reliable results. The presence of the intrinsic *oxaAb* gene in the genome sequence can be applied as a simple check to confirm a sequence is from an *A. baumannii* isolate prior to use of the databases, bearing in mind that it may be missing from poor-quality assemblies. However, this does not preclude the use of the *A. baumannii* KL and OCL databases on other species of *Acinetobacter*. Although not all locus types found in other species will be represented in the databases, K or OC loci with high similarity to those found among *A. baumannii* can be easily identified (see examples in Dataset 3). Hence, the *A. baumannii* databases may assist identification and annotation of the specific genetic content of loci in other *Acinetobacter* species.

It should be noted that the KL does not predict the structure of the CPS, although it does include information about the possible number and identity of sugars present. The CPS structure for each KL must be determined directly because in a number of cases additional genes involved in capsule synthesis are found outside the locus [28, 51, 54]. Hence the KL type is only a starting point for predicting if a particular isolate might be susceptible to a particular phage. However, the potential power of KL and OCL typing as epidemiological tools is highlighted by the analysis of K and OC loci found in single STs. KL and OCL typing have previously proven valuable in dissecting the evolution of two major global clones [2, 41–43]. However, in most studies the genomes were typed using a time-intensive manual process, which imposed a considerable limitation on the scale of datasets that could be explored. In contrast, in this study we were able to use the automated method implemented in *Kaptive* to type the K and OC loci of thousands of genomes, including 134 GC1 and 2016 GC2, revealing even more extensive variation, which is likely to be driven by exchange of locus sequences via recombination in both clones. Given that the available genomes are drawn from a biased, convenient sample of

genomes deposited in NCBI [6], they still may not reflect the true variation in these clones. Similar high levels of variation were found in two other clones (ST10 and ST25), suggesting that they are subject to similar molecular evolutionary processes. In contrast, there appeared to be limited KL and OCL variation among ST16, ST78, ST406, ST499 and ST636.

The findings reported here clearly demonstrate the utility of our novel KL and OCL databases to facilitate rapid and accurate typing of *A. baumannii* surface polysaccharide synthesis loci. This information can be used to distinguish lineages within the global clonal complexes [2, 41, 57] and hence provide valuable information for epidemiological studies, as well as essential information to guide the design of novel treatment or control strategies targeting *A. baumannii* capsules and lipooligosaccharides.

Funding information

This work was supported by an Australian Research Council (ARC) DECRA Fellowship DE180101563 to J.J.K. K.E.H. was supported by a Senior Medical Research Fellowship from the Viertel Foundation of Australia.

Conflicts of interest

The authors declare that there are no conflicts of interest.

Data Bibliography

1. Wyres KL and Holt KE. The Kaptive software, which can be used to screen new genomes against the K and O locus database is available at <https://github.com/katholt/Kaptive> (command-line code) and <http://kaptive.holtlab.net/> (interactive web service).

References

- World Health Organisation (WHO). 2017. Global priority list of antibiotic-resistant bacteria to guide research, discovery, and development of new antibiotics. https://www.who.int/medicines/publications/WHO-PPL-Short_Summary_25Feb-ET_NM_WHO.pdf
- Holt K, Kenyon JJ, Hamidian M, Schultz MB, Pickard DJ et al. Five decades of genome evolution in the globally distributed, extensively antibiotic-resistant *Acinetobacter baumannii* global clone 1. *Microb Genom* 2016;2:e000052.
- Diancourt L, Passet V, Nemec A, Dijkshoorn L, Brisse S. The population structure of *Acinetobacter baumannii*: expanding multiresistant clones from an ancestral susceptible genetic pool. *PLoS One* 2010;5:e10034.
- Sahl JW, Del Franco M, Pournaras S, Colman RE, Karah N et al. Phylogenetic and genomic diversity in isolates from the globally distributed *Acinetobacter baumannii* ST25 lineage. *Sci Rep* 2015;5:15188.
- Zarrilli R, Pournaras S, Giannouli M, Tsakris A. Global evolution of multidrug-resistant *Acinetobacter baumannii* clonal lineages. *Int J Antimicrob Agents* 2013;41:11–19.
- Hamidian M, Nigro SJ. Emergence, molecular mechanisms and global spread of carbapenem-resistant *Acinetobacter baumannii*. *Microb Genom* 2019;5.
- Orskov I, Orskov F, Jann B, Jann K. Serology, chemistry, and genetics of O and K antigens of *Escherichia coli*. *Bacteriol Rev* 1977;41:667–710.
- Ørskov I, Ørskov F. Serotyping of *Klebsiella*. *Method. Microbiol* 1984;14:143–164.
- Liu B, Knirel YA, Feng L, Perepelov AV, Senchenkova S et al. Structure and genetics of *Shigella* O antigens. *FEMS Microbiol Rev* 2008;32:627–653.
- Liu B, Knirel YA, Feng L, Perepelov A, Senchenkova S et al. Structural diversity in *Salmonella* O antigens and its genetic basis. *FEMS Microbiol Rev* 2014;38:56–89.
- Kenyon JJ, Cunneen MM, Reeves PR. Genetics and evolution of *Yersinia pseudotuberculosis* O-specific polysaccharides: a novel pattern of O-antigen diversity. *FEMS Microbiol Rev* 2017;41:200–217.
- Stenutz R, Weintraub A, Widmalm G. The structures of *Escherichia coli* O-polysaccharide antigens. *FEMS Microbiol Rev* 2006;30:382–403.
- Traub WH. *Acinetobacter baumannii* serotyping for delineation of outbreaks of nosocomial cross-infection. *J Clin Microbiol* 1989;27:2713–2716.
- Pantophlet R. Lipopolysaccharides of *Acinetobacter*. In: Gerischer U (editor). *Acinetobacter Molecular Microbiology*. Norfolk, UK: Horizon Scientific Press; 2008.
- Traub WH, Bauer D. Surveillance of nosocomial cross-infections due to three *Acinetobacter* genospecies (*Acinetobacter baumannii*, genospecies 3 and genospecies 13) during a 10-year observation period: serotyping, macrorestriction analysis of genomic DNA and antibiotic susceptibilities. *Chemotherapy* 2000;46:282–292.
- Kenyon JJ, Hall RM. Variation in the complex carbohydrate biosynthesis loci of *Acinetobacter baumannii* genomes. *PLoS One* 2013;8:e62160.
- Russo TA, Luke NR, Beanan JM, Olson R, Sauberman SL et al. The K1 capsular polysaccharide of *Acinetobacter baumannii* strain 307-0294 is a major virulence factor. *Infect Immun* 2010;78:3993–4000.
- Fregolino E, Gargiulo V, Lanzetta R, Parrilli M, Holst O et al. Identification and structural determination of the capsular polysaccharides from two *Acinetobacter baumannii* clinical isolates, MG1 and SMAL. *Carbohydr Res* 2011;346:973–977.
- Oliveira H, Costa AR, Ferreira A, Konstantinides N, Santos SB et al. Functional analysis and antivirulence properties of a new depolymerase from a Myovirus that infects *Acinetobacter baumannii* capsule K45. *J Virol* 2019;93:e01163–18.
- Oliveira H, Costa AR, Konstantinides N, Ferreira A, Akturk E et al. Ability of phages to infect *Acinetobacter calcoaceticus*-*Acinetobacter baumannii* complex species through acquisition of different pectate lyase depolymerase domains. *Environ Microbiol* 2017;19:5060–5077.
- Russo TA, Beanan JM, Olson R, MacDonald U, Cox AD et al. The K1 capsular polysaccharide from *Acinetobacter baumannii* is a potential therapeutic target via passive immunization. *Infect Immun* 2013;81:915–922.
- Yang F-L, Lou T-C, Kuo S-C, Wu W-L, Chern J et al. A medically relevant capsular polysaccharide in *Acinetobacter baumannii* is a potential vaccine candidate. *Vaccine* 2017;35:1440–1447.
- Hu D, Liu B, Dijkshoorn L, Wang L, Reeves PR. Diversity in the major polysaccharide antigen of *Acinetobacter baumannii* assessed by DNA sequencing, and development of a molecular serotyping scheme. *PLoS One* 2013;8:e70329.
- Kenyon JJ, Senchenkova SYN, Shashkov AS, Shneider MM, Popova AV et al. K17 capsular polysaccharide produced by *Acinetobacter baumannii* isolate G7 contains an amide of 2-acetamido-2-deoxy-D-galacturonic acid with D-alanine. *Int J Biol Macromol* 2020;144:857–862.
- Kenyon JJ, Kasimova AA, Shashkov AS, Hall RM, Knirel YA. *Acinetobacter baumannii* isolate BAL_212 from Vietnam produces the K57 capsular polysaccharide containing a rarely occurring amino sugar N-acetylviosamine. *Microbiology* 2018;164:217–220.
- Kasimova AA, Kenyon JJ, Arbatsky NP, Shashkov AS, Popova AV et al. *Acinetobacter baumannii* K20 and K21 capsular polysaccharide structures establish roles for UDP-glucose dehydrogenase Ugd2, pyruvyl transferase Ptr2 and two glycosyltransferases. *Glycobiology* 2018;28:876–884.
- Kenyon JJ, Shashkov AS, Senchenkova Sof'ya N, Shneider MM, Liu B et al. *Acinetobacter baumannii* K11 and K83 capsular polysaccharides have the same 6-deoxy-L-talose-containing

- pentasaccharide K units but different linkages between the K units. *Int J Biol Macromol* 2017;103:648–655.
28. Kenyon JJ, Kasimova AA, Shneider MM, Shashkov AS, Arbatsky NP et al. The KL24 gene cluster and a genomic island encoding a Wzy polymerase contribute genes needed for synthesis of the K24 capsular polysaccharide by the multiply antibiotic resistant *Acinetobacter baumannii* isolate RCH51. *Microbiol* 2017;163:355–363.
 29. Kenyon JJ, Kasimova AA, Notaro A, Arbatsky NP, Speciale I et al. *Acinetobacter baumannii* K13 and K73 capsular polysaccharides differ only in K-unit side branches of novel non-2-ulosonic acids: di-N-acetylated forms of either acinetaminic acid or 8-epiacinetaminic acid. *Carbohydr Res* 2017;452:149–155.
 30. Kenyon JJ, Marzaioli AM, Hall RM, De Castro C. Structure of the K2 capsule associated with the KL2 gene cluster of *Acinetobacter baumannii*. *Glycobiology* 2014;24:554–563.
 31. Arbatsky NP, Kasimova AA, Shashkov AS, Shneider MM, Popova AV et al. Structure of the K128 capsular polysaccharide produced by *Acinetobacter baumannii* KZ-1093 from Kazakhstan. *Carbohydr Res* 2019;485:107814.
 32. Arbatsky NP, Shneider MM, Dmitrenok AS, Popova AV, Shagin DA et al. Structure and gene cluster of the K125 capsular polysaccharide from *Acinetobacter baumannii* MAR13-1452. *Int J Biol Macromol* 2018;117:1195–1199.
 33. Kasimova AA, Shneider MM, Arbatsky NP, Popova AV, Shashkov AS et al. Structure and gene cluster of the K93 capsular polysaccharide of *Acinetobacter baumannii* B11911 containing 5-N-Acetyl-7-N-[(R)-3-hydroxybutanoyl]pseudaminic acid. *Biochemistry Moscow* 2017;82:483–489.
 34. Senchenkova SN, Shashkov AS, Popova AV, Shneider MM, Arbatsky NP et al. Structure elucidation of the capsular polysaccharide of *Acinetobacter baumannii* AB5075 having the KL25 capsule biosynthesis locus. *Carbohydr Res* 2015;408:8–11.
 35. Shashkov AS, Kenyon JJ, Senchenkova SN, Shneider MM, Popova AV et al. *Acinetobacter baumannii* K27 and K44 capsular polysaccharides have the same K unit but different structures due to the presence of distinct wzy genes in otherwise closely related K gene clusters. *Glycobiology* 2016;26:501–508.
 36. Kenyon JJ, Hall RM, De Castro C. Structural determination of the K14 capsular polysaccharide from an ST25 *Acinetobacter baumannii* isolate, D46. *Carbohydr Res* 2015;417:52–56.
 37. Lees-Miller RG, Iwashkiw JA, Scott NE, Seper A, Vinogradov E et al. A common pathway for O-linked protein-glycosylation and synthesis of capsule in *Acinetobacter baumannii*. *Mol Microbiol* 2013;89:816–830.
 38. Kenyon JJ, Holt KE, Pickard D, Dougan G, Hall RM. Insertions in the OCL1 locus of *Acinetobacter baumannii* lead to shortened lipooligosaccharides. *Res Microbiol* 2014;165:472–475.
 39. Kenyon JJ, Nigro SJ, Hall RM. Variation in the OC locus of *Acinetobacter baumannii* genomes predicts extensive structural diversity in the lipooligosaccharide. *PLoS One* 2014;9:e107833.
 40. Meumann EM, Anstey NM, Currie BJ, Piera KA, Kenyon JJ et al. Genomic epidemiology of severe community-onset *Acinetobacter baumannii* infection. *Microb Genom* 2019;5 [Epub ahead of print 26 02 2019].
 41. Schultz MB, Pham Thanh D, Tran Do Hoan N, Wick RR, Ingle DJ et al. Repeated local emergence of carbapenem-resistant *Acinetobacter baumannii* in a single hospital ward. *Microb Genom* 2016;2:e000050.
 42. Wright MS, Haft DH, Harkins DM, Perez F, Hujer KM et al. New insights into dissemination and variation of the health care-associated pathogen *Acinetobacter baumannii* from genomic analysis. *mBio* 2014;5:e00963–13.
 43. Adams MD, Wright MS, Karichu JK, Venepally P, Fouts DE et al. Rapid replacement of *Acinetobacter baumannii* strains accompanied by changes in lipooligosaccharide loci and resistance gene repertoire. *mBio* 2019;10:e00356–19.
 44. Wyres KL, Wick RR, Gorrie C, Jenney A, Follador R et al. Identification of *Klebsiella* capsule synthesis loci from whole genome data. *Microb Genom* 2016;2:e000102.
 45. Wick RR, Heinz E, Holt KE, Wyres KL. Kaptive Web: User-friendly capsule and lipopolysaccharide Serotype prediction for *Klebsiella* genomes. *J Clin Microbiol* 2018;56:e00197–18.
 46. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* 2012;19:455–477.
 47. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: resolving bacterial genome assemblies from short and long sequencing reads. *PLoS Comput Biol* 2017;13:e1005595.
 48. Kenyon JJ, Marzaioli AM, De Castro C, Hall RM. 5,7-di-N-acetyl-acinetaminic acid: a novel non-2-ulosonic acid found in the capsule of an *Acinetobacter baumannii* isolate. *Glycobiology* 2015;25:644–654.
 49. Arbatsky NP, Kenyon JJ, Shashkov AS, Shneider MM, Popova AV et al. The K5 capsular polysaccharide of the bacterium *Acinetobacter baumannii* SDF with the same K unit containing Leg5Ac7Ac as the K7 capsular polysaccharide but a different linkage between the K units. *Russ Chem Bull* 2019;68:163–167.
 50. Shashkov AS, Kenyon JJ, Arbatsky NP, Shneider MM, Popova AV et al. Structures of three different neutral polysaccharides of *Acinetobacter baumannii* NIPH190, NIPH201, and NIPH615, assigned to K30, K45, and K48 capsule types, respectively, based on capsule biosynthesis gene clusters. *Carbohydr Res* 2015;417:81–88.
 51. Kenyon JJ, Shneider MM, Senchenkova SN, Shashkov AS, Sinia-gina MN et al. K19 capsular polysaccharide of *Acinetobacter baumannii* is produced via a Wzy polymerase encoded in a small genomic island rather than the KL19 capsule gene cluster. *Microbiology* 2016;162:1479–1489.
 52. Shashkov AS, Kenyon JJ, Arbatsky NP, Shneider MM, Popova AV et al. Related structures of neutral capsular polysaccharides of *Acinetobacter baumannii* isolates that carry related capsule gene clusters KL43, KL47, and KL88. *Carbohydr Res* 2016;435:173–179.
 53. Shashkov AS, Cahill SM, Arbatsky NP, Westacott AC, Kasimova AA et al. *Acinetobacter baumannii* K116 capsular polysaccharide structure is a hybrid of the K14 and revised K37 structures. *Carbohydr Res* 2019;484:107774.
 54. Kenyon JJ, Arbatsky NP, Shneider MM, Popova AV, Dmitrenok AS et al. The K46 and K5 capsular polysaccharides produced by *Acinetobacter baumannii* NIPH 329 and SDF have related structures and the side-chain non-ulosonic acids are 4-O-acetylated by phage-encoded O-acetyltransferases. *PLoS One* 2019;14:e0218461.
 55. Arbatsky NP, Shneider MM, Kenyon JJ, Shashkov AS, Popova AV et al. Structure of the neutral capsular polysaccharide of *Acinetobacter baumannii* NIPH146 that carries the KL37 capsule gene cluster. *Carbohydr Res* 2015;413:12–15.
 56. Kenyon JJ, Notaro A, Hsu LY, De Castro C, Hall RM. 5,7-Di-N-acetyl-8-epiacinetaminic acid: A new non-2-ulosonic acid found in the K73 capsule produced by an *Acinetobacter baumannii* isolate from Singapore. *Sci Rep* 2017;7:11357.
 57. Hamidian M, Hawkey J, Wick R, Holt KE, Hall RM. Evolution of a clade of *Acinetobacter baumannii* global clone 1, lineage 1 via acquisition of carbapenem- and aminoglycoside-resistance genes and dispersion of ISAba1. *Microb Genom* 2019;5:e000242.
 58. Carver TJ, Rutherford KM, Berriman M, Rajandream M-A, Barrell BG et al. Act: the Artemis comparison tool. *Bioinformatics* 2005;21:3422–3423.
 59. Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics* 2015;31:3350–3352.



Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Wyres, KL; Cahill, SM; Holt, KE; Hall, RM; Kenyon, JJ

Title:

Identification of *Acinetobacter baumannii* loci for capsular polysaccharide (KL) and lipooligosaccharide outer core (OCL) synthesis in genome assemblies using curated reference databases compatible with Kaptive

Date:

2020-03-01

Citation:

Wyres, K. L., Cahill, S. M., Holt, K. E., Hall, R. M. & Kenyon, J. J. (2020). Identification of *Acinetobacter baumannii* loci for capsular polysaccharide (KL) and lipooligosaccharide outer core (OCL) synthesis in genome assemblies using curated reference databases compatible with Kaptive. *MICROBIAL GENOMICS*, 6 (3), <https://doi.org/10.1099/mgen.0.000339>.

Persistent Link:

<http://hdl.handle.net/11343/245861>

File Description:

published version

License:

CC BY