

LONDON  
SCHOOL of  
HYGIENE  
& TROPICAL  
MEDICINE



LSHTM Research Online

Hamidian, Mohammad; Wick, Ryan R; Hartstein, Rebecca M; Judd, Louise M; Holt, Kathryn E; Hall, Ruth M; (2019) Insights from the revised complete genome sequences of *Acinetobacter baumannii* strains AB307-0294 and ACICU belonging to global clones 1 and 2. *Microbial genomics*. ISSN 2057-5858 DOI: <https://doi.org/10.1099/mgen.0.000298>

Downloaded from: <http://researchonline.lshtm.ac.uk/id/eprint/4654549/>

DOI: <https://doi.org/10.1099/mgen.0.000298>

**Usage Guidelines:**

Please refer to usage guidelines at <https://researchonline.lshtm.ac.uk/policies.html> or alternatively contact [researchonline@lshtm.ac.uk](mailto:researchonline@lshtm.ac.uk).

Available under license: <http://creativecommons.org/licenses/by/2.5/>

<https://researchonline.lshtm.ac.uk>

# Insights from the revised complete genome sequences of *Acinetobacter baumannii* strains AB307-0294 and ACICU belonging to global clones 1 and 2

Mohammad Hamidian<sup>1,\*</sup>, Ryan R. Wick<sup>2</sup>, Rebecca M. Hartstein<sup>3</sup>, Louise M. Judd<sup>2</sup>, Kathryn E. Holt<sup>2,4</sup> and Ruth M. Hall<sup>3</sup>

## Abstract

The *Acinetobacter baumannii* global clone 1 isolate AB307-0294, recovered in the USA in 1994, and the global clone 2 (GC2) isolate ACICU, isolated in 2005 in Italy, were among the first *A. baumannii* isolates to be completely sequenced. AB307-0294 is susceptible to most antibiotics and has been used in many genetic studies, and ACICU belongs to a rare GC2 lineage. The complete genome sequences, originally determined using 454 pyrosequencing technology, which is known to generate sequencing errors, were re-determined using Illumina MiSeq and MinION (Oxford Nanopore Technologies) technologies and a hybrid assembly generated using Unicycler. Comparison of the resulting new high-quality genomes to the earlier 454-sequenced versions identified a large number of nucleotide differences affecting protein coding sequence (CDS) features, and allowed the sequences of the long and highly repetitive *bap* and *blp1* genes to be properly resolved for the first time in ACICU. Comparisons of the annotations of the original and revised genomes revealed a large number of differences in the protein CDS features, underlining the impact of sequence errors on protein sequence predictions and core gene determination. On average, 400 predicted CDSs were longer or shorter in the revised genomes and about 200 CDS features were no longer present.

## DATA SUMMARY

The corrected complete genome sequence of *Acinetobacter baumannii* AB307-0294 has been deposited in GenBank under accession number CP001172.2 (chromosome; url - <https://www.ncbi.nlm.nih.gov/nucore/CP001172.2>). The corrected complete genome sequence of *A. baumannii* ACICU has been deposited in GenBank under accession numbers CP031380 (chromosome; url - <https://www.ncbi.nlm.nih.gov/nucore/CP031380>), CP031381 (pACICU1; url - <https://www.ncbi.nlm.nih.gov/nucore/CP031381>) and CP031382 (pACICU2; url - <https://www.ncbi.nlm.nih.gov/nucore/CP031382>).

## INTRODUCTION

*Acinetobacter baumannii* is a Gram-negative bacterium that has emerged as an important opportunistic pathogen, and is a research priority because of its high levels of resistance

to antibiotics [1–3], desiccation and heavy metals [4, 5]. On a global scale, members of two clinically important clones, known as global clone 1 (GC1) and global clone 2 (GC2), have been responsible for the majority of outbreaks caused by multiply antibiotic-resistant *A. baumannii* strains [1–3, 6–8]. Whole-genome sequencing technologies have revolutionized the study of bacterial pathogens, allowing the entire gene repertoire of bacterial strains to be determined; hence, enabling the study of the relationships between outbreak strains with an unprecedented high resolution [9]. However, accuracy is important.

The first 10 complete genomes of *A. baumannii* strains were reported between 2006 and 2012 (Table 1), and are still used as a baseline in many studies of this micro-organism [10–12]. Except for three strains (AYE, TCDC-AB0715 and TYTH-1), all of the early *A. baumannii* complete genomes

Received 18 May 2019; Accepted 05 September 2019; Published 26 September 2019

**Author affiliations:** <sup>1</sup>The ithree Institute, University of Technology Sydney, Ultimo, NSW, Australia; <sup>2</sup>Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria, Australia; <sup>3</sup>School of Life and Environmental Sciences, University of Sydney, Sydney, Australia; <sup>4</sup>London School of Hygiene and Tropical Medicine, London WC1E 7HT, UK.

**\*Correspondence:** Mohammad Hamidian, mohammad.hamidian@uts.edu.au; mohamidian@gmail.com

**Keywords:** *Acinetobacter baumannii*; AB307-0294; ACICU; global clone 1; global clone 2; complete genome sequence.

**Abbreviations:** CDS, coding sequence; GC1, global clone 1; GC2, global clone 2; MLST, multilocus sequence typing; ONT, Oxford Nanopore Technologies; PacBio, Pacific Biosciences; SND, single nucleotide difference; ST, sequence type.

The GenBank/EMBL/DBJ accession numbers for the complete genome sequences of the *Acinetobacter baumannii* strains and plasmids are CP001172.2 (AB307-0294 chromosome), CP031380 (ACICU chromosome), CP031381 (pACICU1) and CP031382 (pACICU2).

**Data statement:** All supporting data, code and protocols have been provided within the article or through supplementary data files.

000298 © 2019 The Authors

This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

were sequenced using 454 pyrosequencing technology and assembled using PCR. Pyrosequencing is known to generate frequent systematic sequencing errors, especially errors in the length of homopolymeric runs [13]; and these errors lead to erroneous protein coding sequence (CDS) prediction, often associated with fragmentation of genuine ORFs.

An additional problem in *A. baumannii* genomes determined using short-read sequence data followed by PCR gap closure arises from the many short internal repeats present in the very large *bap* gene (~8–25 kbp), which is hard to assemble accurately. This gene encodes the biofilm-associated protein Bap [14–17]. The *bap* gene was originally cloned from AB307-0294 (GC1), and found to be 25 863 bp with a complex configuration of internal repeats [15]. However, the size of the *bap* gene from a GC2 isolate was estimated at approximately 16 kbp [16]. In another study, the length of Bap proteins predicted from *A. baumannii* genomes available in GenBank appeared to be highly variable, mainly due to different numbers of copies of the various repeated segments and the ORF was often fragmented [17]. The *blp1* gene, which is 9–10 kbp, encodes a further very large protein that also has internal repeats and is associated with biofilm formation [17].

Newer sequencing technologies such as PacBio (Pacific Biosciences) and MinION (Oxford Nanopore Technologies; ONT) can generate much longer sequencing reads [9], allowing gaps to be spanned. MinION-only assemblies are also prone to errors [18], but can be combined with high-accuracy Illumina short-read data to produce very-high-quality finished genome assemblies [19]. Long-read sequence data have enabled a re-assessment of early completed *A. baumannii* genomes, including several of the first 10 to be sequenced (Table 1). For example, in 2016, ATCC 17978 was re-sequenced using PacBio. This revealed the presence of a 148 kb conjugative plasmid, pAB3, fragments of which were erroneously merged into the chromosome in the original 454-based assembly [20]. This plasmid sequence brought together the parts of *GIsul2*, fragmented pieces of which had been randomly distributed in the chromosome in the original sequence [21]. In 2017, we revised the 454-based genome sequence of the GC1 strain AB0057 using Illumina HiSeq technology, and found that hundreds of single base additions or deletions changed >200 protein CDS features [22]. An additional copy of the *oxa23* carbapenem-resistance gene, located in Tn2006, was also found in the revised sequence of the chromosome (GenBank accession no. CP001182.2) [22, 23].

A recent revision of the 454-based genome of the GC2 strain MDR-ZJ06 using PacBio sequencing led to the correction of hundreds of CDS features and allowed reassessment of the localization of important antimicrobial-resistance regions [24]. The position of transposon Tn2009, which carries the *oxa23* gene, was revised; and a region originally reported as a plasmid, but that had been predicted to be a chromosomally located AbGRI3-type resistance island [25], was incorporated into the chromosome (CP001937.2) [24]. In the revised genome, the two arrays of gene cassettes carrying

### Impact Statement

The genomes of the first 10 *Acinetobacter baumannii* strains to be completely sequenced underpin a large amount of published genetic and genomic analysis. However, most of their genome sequences contain substantial numbers of errors as they were sequenced using 454 pyrosequencing, which is known to generate errors particularly in homopolymer regions; and employed manual PCR and capillary sequencing steps to bridge contig gaps and repetitive regions in order to finish the genomes. Assembly of the very large and internally repetitive genes for the biofilm-associated proteins Bap and Blp1 was a recurring problem. As these strains continue to be used for genetic studies and their genomes continue to be used as references in phylogenomics studies, including core gene determination, there is value in improving the quality of their genome sequences. To this end, we re-sequenced two such strains that belong to the two major globally distributed clones of *A. baumannii*, using a combination of highly accurate short-read and gap-spanning long-read technologies. Annotation of the revised genome sequences eliminated hundreds of incorrect coding sequence (CDS) feature annotations and corrected hundreds more. Given that these revisions affected hundreds of non-existent or incorrect CDS features currently cluttering GenBank protein databases, it can be envisaged that similar revision of other early bacterial genomes that were sequenced using error-prone technologies will affect thousands of CDSs currently listed in GenBank and other databases. These corrections will impact the quality of predicted protein sequence data stored in public databases. The revised genomes will also improve the accuracy of future genetic and comparative genomic analyses incorporating these clinically important strains.

antibiotic-resistance genes in class 1 integrons are now in the correct resistance islands. These revisions exemplify the challenges encountered when relying solely on short-read data to assemble bacterial genomes, and highlight the extent and impact of pyrosequencing errors particularly on CDS predictions.

Two further *A. baumannii* strains for which only early 454-based genome sequences are available are the largely antibiotic-susceptible isolate AB307-0294, recovered from the blood of a patient hospitalized in Buffalo, NY, USA, in 1994 [26], and the extensively antibiotic-resistant isolate ACICU recovered in 2005 from the cerebrospinal fluid of a patient in San Giovanni Addolorata Hospital in Rome, Italy (GenBank accession no. CP000863) [27]. AB307-0294 was one of the first GC1 strains to be completely sequenced [26] and has been extensively used in genetic studies [28–32]. It belongs to CC1 (clonal complex 1) [sequence type 1 (ST1)] in the

**Table 1.** Properties of early *A. baumannii* completed genomes

Strain/plasmid	Country	Isolation date	GC	Original sequence			Revised						
				Length (bp)	GenBank no.	Sequencing technology	Reference	Length (bp)	GenBank no.	Sequencing technology	Reference		
ATCC 17978	NK	1951	-										
Chromosome				3976747	CP000521.1	454	[41]	Yes	3857743	CP012004.1	PacBio		[20]
pAB1				13408	CP000522.1	454		No	Not present	NA	NK		
pAB2				11302	CP000523.1	454		No	Not present	NA	NK		
pAB3				Not present	-	454		Yes	148955	CP012005.1	PacBio		
AYE	France	2001	1										
Chromosome				3936291	CU459141.1	454	[42]	No	-	-	-		-
p1ABAYE				5644	CU459137.1	454		No	-	-	-		-
p2ABAYE				9661	CU459138.1	454		No	-	-	-		-
p3ABAYE				94413	CU459140.1	454		No	-	-	-		-
p4ABAYE				2726	CU459139.1	454		No	-	-	-		-
AB307-0294	USA	1994	1										
Chromosome				3760981	CP001172.1	454	[26]	Yes	3759495	CP001172.2	MinION and Illumina		This study
AB0057	USA	2004	1										
Chromosome				4050513	CP001182.1	454	[26]	Yes	4055148	CP001182.2	Illumina		[22]
pAB0057				8729	CP001183.1	454		Yes	8731	CP001183.2	Illumina		
1656-2	South Korea	2011*	2										
Chromosome				3940614	CP001921.1	454	[43]	No	-	-	-		-
ABkp1				74451	CP001922.1	454		No	-	-	-		-
ABkp2				8041	CP001923.1	454		No	-	-	-		-
AGICU	Italy	2005	2										
Chromosome				3904116	CP000863.1	454	[27]	Yes	3919274	CP031380.1	MinION and Illumina		This study
pAC1CU1				28279	CP000864.1	454		Yes	24268	CP031381.1	MinION and Illumina		This study
pAC1CU2				64366	CP000865.1	454		Yes	70101	CP031382.1	MinION and Illumina		This study
TDCC-AB0715†	Taiwan	2007	2										
Chromosome				4130792	CP002522.1	454	[44]	Yes	4138388	CP002522.2	Illumina		-
p1ABTDC0715				8731	CP002523.1	454		No	-	-	-		-
p2ABTDC0715				70894	CP002524.1	454		No	-	-	-		-

Continued

Table 1. Continued

Strain/plasmid	Country	Isolation date	GC	Original sequence			Revised sequence						
				Length (bp)	GenBank no.	Sequencing technology	Reference	Revised	Length (bp)	GenBank no.	Sequencing technology	Reference	
MDR-ZJ06	China	2006	2										
Chromosome				3991133	CP001937.1	454	[45]	Yes	4022275	CP001937.2	PacBio	[24]	
pMDR-ZJ06†				20301	CP001938.1	454	[46]	Yes	Not present†	NA	NA	NA	
TYTH-1	Taiwan	2008	2										
Chromosome				3957368	CP003856	Illumina		No					
MDR-TJ	China	2012‡	2										
Chromosome				3964912	CP003500.1	454	[47]	No					
pABTJ1				77528	CP003501.1	454		No					
pABTJ2				110967	CP004359.1	454		No					

\*Genome submission date; isolation date is not known.

†pMDR-ZJ06 is not present in the revised genome

‡Recovered between 2007 and 2012.

NA, not applicable; NK, not known.

Institut Pasteur multilocus sequence typing (MLST) scheme and to ST231 in the Oxford MLST scheme, and carries the KL1 capsule genes and OCL1 at the outer core locus [33]. Compared to other GC1 strains characterized to date, AB307-0294 is relatively susceptible to antibiotics [26], exhibiting resistance only to chloramphenicol (intrinsic) and nalidixic acid (acquired). It contains no plasmids.

ACICU was the first GC2 isolate to be sequenced [27]. It belongs to ST2 in the Institut Pasteur MLST scheme, ST437 in the Oxford MLST scheme, and carries the KL2 capsule genes and OCL1 at the outer core locus [33]. ACICU is carbapenem resistant and also resistant to multiple antibiotics, including third-generation cephalosporins, sulfonamides, tetracycline, amikacin, kanamycin, netilmicin and ciprofloxacin [27]. It contains two plasmids [27]. However, we previously showed that the largest plasmid, pACICU-2, which was reported to include no resistance genes, is larger and contains the amikacin-resistance gene *aphA6* in transposon *TnaphA6*. The central segment of *TnaphA6*, including the *aphA6* gene and one of the ISAb125 copies as well as a 4.7 kb backbone segment, were missing in the original 454-based whole-genome sequence [34].

Here, we report revised complete genome sequences for *A. baumannii* strains AB307-0294 (GC1) and ACICU (GC2), generated using MiSeq (Illumina) and MinION (ONT) sequence data. The new genome sequences correct hundreds of protein CDS features generated by the presence of single nucleotide differences (SNDs) and small insertion/deletions of mainly 1–3 bases in the earlier 454 genome sequences.

## METHODS

### Whole-genome sequencing, assembly and annotation

Whole-cell DNA was isolated and purified using a protocol described elsewhere [1, 35]. Libraries were prepared from whole-cell DNA isolated from AB307-0294 and ACICU, and were sequenced using Illumina MiSeq and ONT MinION. Paired-end reads of 150 bp and MinION reads of up to 20 kb were used to assemble each genome using Unicycler software (v0.4.0) [19] with default parameters.

Protein CDS, rRNA and tRNA genes were annotated using the automatic annotation program Prokka v1.13 [36]. Regions containing antibiotic-resistance genes and the polysaccharide biosynthesis loci, biofilm-associated proteins and genes used in the MLST schemes were annotated manually.

To compare previous CDS ( $\geq 25$  aa CDS features) annotations with our new results, we wrote a script ([github.com/rrwick/Compare-annotations](https://github.com/rrwick/Compare-annotations)) to quantify the differences. This script classifies CDSs in the annotations as either exact matches, inexact matches, only present in the first annotation or only present in the second annotation. We also used the Ideel pipeline of Dr Mick Watson ([github.com/mw55309/ideel](https://github.com/mw55309/ideel)) to assess the completeness of CDSs annotated in each genome, by comparing the length of each CDS to that of its longest

**Table 2.** Comparison of *bap* and *blp1* genes in early *A. baumannii* complete genomes and their revisions

Genome	Revision technology	<i>bap</i>		<i>blp1</i>	
		Size (bp)	Locus ID	Size (bp)	Locus ID
<b>ATCC 17978</b>					
Original		6306	A1S_2696	–*	–
Revised	PacBio	6225	ACX60_04030	–	–
<b>AB307-0294</b>					
Original		22 920	ABBFA_000776	10 071	ABBFA_000810
Revised	Nanopore	25863	ABBFA_00771	10 089	ABBFA_00802
<b>ACICU</b>					
Original		6420	ACICU_02938-46	9510	ACICU_02910
Revised	Nanopore	22 212	DMO12_08904	9813	DMO12_08811
<b>MDR-ZJ06</b>					
Original		2115	ABZJ_03124	9135	ABZJ_03096
Revised	PacBio	7947	ABZJ_03955	9813	ABZJ_03096

\*ATCC 17978 does not contain the *blp1* gene.

BLAST hit in the UniProt database (as described in <http://www.opiniomics.org/a-simple-test-for-uncorrected-insertions-and-deletions-indels-in-bacterial-genomes/>).

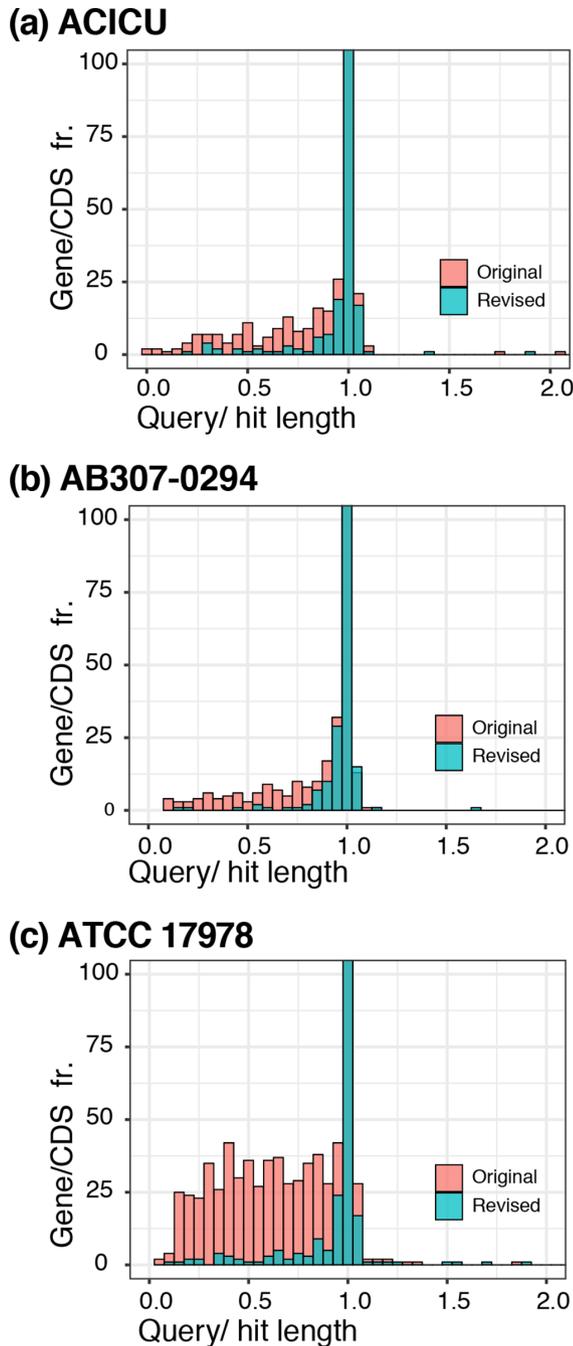
## RESULTS AND DISCUSSION

### Revised genome of ACICU

ACICU, the first GC2 strain to be completely sequenced, contains AbaR2 in the chromosomal *comM* gene [27]. As this AbaR resistance island type is more usually found in this location in GC1 strains [37] with an AbGRI1 type island in GC2 isolates [38], ACICU may represent a rare GC2 lineage. Here, the ACICU genome was re-sequenced using a combination of Illumina (MiSeq, 58× depth) and ONT (MinION, 253× depth) data. The new contiguous ACICU chromosomal sequence comprised 3 919 274 bp (GenBank accession no. CP001172.2), compared to 3 904 116 bp in the original submission (GenBank accession no. CP000863), making the revised chromosome 15 158 bp longer (Table 1). Most of the additional length in the revised chromosome was found to be due to a 11.2 kbp longer *bap* gene (Table 2), which is just over 11 kbp and in nine smaller ORFs in the original sequence (locus\_ids ACICU\_02938 to ACICU\_2946), as noted previously [17]. In the revised genome sequence, the *bap* gene is 22.2 kbp (locus\_id DMO12\_08904), mainly due to a large number of short strings of repeated sequences missing previously. Hence, some of the variation in the length of *bap* reported elsewhere [17] may be due to sequencing and assembly issues rather than genuine length variation in the *A. baumannii* population. The *blp1* gene in the original sequence (locus\_id ACICU\_02910) is 9510 bp and 9813 bp (locus\_id DMO12\_08811) in the revised genome (Table 2).

The revised chromosome of ACICU differs from the original at 281 positions, including 40 SNDS and 241 insertions or deletions of 1–3 bases (mostly in homopolymeric runs of As or Ts). The original annotation included 3677 protein-encoding features (CDS features are  $\geq 25$  aa), whereas the revised genome annotation contains 3605 CDS features. Comparison of the CDS features indicated that only 3129 CDSs are identical between the two versions. The differences are mostly due to the correction of ORFs that were interrupted or fused due to errors in the 454 sequence, and include 80 CDSs unique to the revised version and 142 CDS features in the original sequence that could not be found in the corrected chromosome. A further 396 CDSs that are present in both versions are altered: of these, 8 have the same length, 285 are longer in the revised chromosome and 103 are shorter. Overall, 98.8 % of all genes ( $n=3568$ ) in the new assembly are within 5 % of the maximum length of homologous proteins in UniProt (i.e. the expected length), calculated using the Ideel pipeline (see Methods). In the old assembly, only 95.8 % ( $n=3494$ ) of all genes were within 5 % of this expected length. The distribution of length ratios is shown in Fig. 1(a), highlighting a substantial population of CDSs annotated in the old assembly that have lengths well below those of homologous proteins in UniProt.

ACICU carries two plasmids (Table 1), pACICU1 and pACICU2 [27], which encode the RepAci1 and RepAci6 replication initiation proteins [39]. The original pACICU1 sequence (GenBank accession no. CP000864) is 28 279 bp long and contains two copies of the carbapenem-resistance gene *oxa58*, while the revised pACICU1 (GenBank accession no. CP031381) is 24 268 bp long and includes only a single *oxa58* copy. It lacks the region between the two IS26 and one



**Fig. 1.** Histograms of CDS lengths relative to the length of the top hit in UniProt, in the original versus revised genomes. (a) ACICU GenBank accession no. CP000863.1 (original) and CP031380 (revised), (b) AB307-0294 GenBank accession no. CP001172.1 (original) and CP001172.2 (revised), and (c) ATCC 17978 GenBank accession no. CP000521.1 (original) and CP012004.1 (revised). The x-axis shows the ratio of CDS length to the length of the closest hit in the UniProt TrEMBL database. The y-axis shows gene frequency and is truncated at 100 (the centre bar extends to ~3000 genes). A tight distribution around 1.0 indicates that the assembly's CDSs match known proteins, supporting few indel errors in the assembly. A left-skewed distribution is characteristic of an assembly with indel errors that lead to premature stop codons.

copy of IS26 in the original sequence. The IS26-mediated duplication may have been generated during growth in selective media. The original and revised pACICU1 sequences also differed by three SNDS, six single bp insertions, and one single bp and two 2 bp deletions. We previously used a PCR mapping strategy [34] to show that the *aphA6* gene and an additional ISAbA125, as well as a 4.7 kb long backbone segment, located between two copies of a ~420 bp repeated segment, are missing from the original sequence of pACICU2, the larger plasmid of ACICU [34]. Here, the long-read sequences generated for pACICU2 (GenBank accession no. CP031382) confirmed this. The revised plasmid sequence differs by six SNDS from pAb-G7-2 (GenBank accession no. KF669606.1), a conjugative plasmid from a GC1 isolated in Australia in 2003 reported previously [40].

### Revised genome of AB307-0294

The AB307-0294 genome was also sequenced using a combination of Illumina (MiSeq, 63× depth) and ONT (MinION, 120× depth) technologies. The hybrid assembly resulted in a single 3 759 495 bp chromosome (GenBank accession no. CP001172.2) compared with 3 760 981 bp in the original genome (GenBank accession no. CP001172.1), making the revised genome 1486 bp shorter (Table 1). As with AB0057, the majority of differences were found to be additions or deletions of 1–3 bases, usually in ‘A’ or ‘T’ in homopolymeric runs of these nucleotides. The original annotation included 3427 CDSs, while the revised annotation contains 3458 (≥25 aa), of which 2937 CDSs are identical in the two versions. Corrections of insertion/deletion errors changed 354 ORFs leading to merging and splitting of CDS regions. Amongst these 354 CDS features, 286 CDSs in the revised genome are longer and 65 are shorter than the corresponding CDSs in the original annotation, and 3 have the same length but differ internally. The revised genome also includes 136 novel CDS features, compared to the original sequence, while there are also 167 CDSs in the old sequence that no longer exist in the revised genome, again indicating the high impact of the errors caused by the use of 454 pyrosequencing technology. Overall, 98.9 % of all genes ( $n=3387$ ) in the new assembly are within 5 % of the expected length, calculated using the Ideel pipeline, versus just 96.4 % ( $n=3336$ ) in the old assembly (Fig. 1b).

The *bap* gene was 25 863 bp (locus\_id ABBFA\_00771), the same length as reported originally [15] but 1067 bp shorter than the 26 930 bp *bap* gene in the original genome sequence where it is split into two ORFs (locus\_id ABBFA\_000776 and locus\_id ABBFA\_000777). The revised genome was found to contain a 10 089 bp *blp1* gene (ABBFA\_00802), only 18 bp longer than that in the original sequence (Table 2). Interestingly, both the original and revised genomes appear to be devoid of any insertion sequences.

### Revised genomes affect many predicted protein sequences

To date, six early *A. baumannii* genome sequences, including AB307-0294 and ACICU reported here, have been corrected, and in each case the revised genome has resulted in the

correction of ~600 CDS features on average [20, 22]. In each comparison of revised and original genome sequences, 100–150 new CDS features appeared, 150–200 CDSs disappeared and 150–200 CDSs changed. As the extent of errors had not been reported previously [20], we also compared the original (GenBank accession no. CP000521.1) and revised (GenBank accession no. CP012004.1) genomes of *A. baumannii* ATCC 17978. This revealed that the revised sequence has extensively re-ordered parts of the chromosome, correcting a large number of inversions, insertion/deletions and other misassemblies. A striking difference between the two genomes is the inclusion in the original chromosome assembly of several large segments that in fact make up a 148 kb plasmid (pAB3) carrying the *sul2* sulfonamide-resistance gene (GenBank accession no. CP012005). The misassembly issues precluded a simple alignment of the two chromosome sequences, but alignment of 14 separate chromosomal segments totalling 3 843 892 bp, revealed 334 SNPs as well as 635 deletions and 754 insertions of 1–3 bases, mainly As or Ts in runs of As or Ts. Overall, 3503 genes (98.2 % of all genes) in the new assembly are within 5 % of the expected length, calculated using the Ideel pipeline, versus 3381 (86.4 %) in the old assembly (see Fig. 1c). Hence, the original assembly was substantially flawed and should not be used in future. However, although the original study reported that ATCC 17978 contains two cryptic plasmids of 13 kb, pAB1 (GenBank accession no. CP000522.1) and 11 kb, pAB2 (GenBank accession no. CP000523.1) [41], the revised genome does not include either of these plasmids. This may be due to an assembly parameter setting to filter out the small contigs, which would remove pAB1 and pAB2 from the final assembly.

Granted the large effects observed on the length of *bap* and *blp* in ACICU using long-read data, their sizes in original and revised genomes in the remainder of the first set of 10 sequenced *A. baumannii* (Table 1) were compared and significant differences were observed only where long-read data were used in the revision. In the GC2 strain MDR-ZJ06 (GenBank accession no. CP001937), *blp1* (locus tag ABZJ\_03096) is 9812 bp in the revised genome (CP001937.2) versus 9134 bp in the original sequence (locus tag ABZJ\_03096; Table 2). Further, *bap*, which is 7946 bp in the revised genome (locus\_id ABZJ\_03955), was split into three ORFs, ranging in size from 2 to 2.5 kb, in the original sequence. In ATCC 17978, the *blp1* gene is not present in either the original or the revised genome. However, the *bap* gene, which was split into two ORFs (locus\_id A1S\_2696, 6306 bp; and A1S\_2724, 1161 bp) and separated by 41 kbp in the original sequence is now in a single ORF (locus\_id ACX60\_04030; 6225 bp) in the revised genome and 842 bp shorter compared to the original genome (Table 2).

## Conclusions

The revised genome sequences of AB307-0294 and ACICU will underpin more accurate studies of the genetics and genomic evolution of related *A. baumannii* strains belonging to GC1 and GC2. This work highlights the need to review

and revise early bacterial genomes sequenced using short-read data and assembled with (or sometimes without) PCR to join contigs. Special attention needs to focus on the genomes determined using the 454 pyrosequencing technology in order to correct predicted protein sequences.

Long-read data, such as those generated by PacBio and ONT (MinION) technologies, allow for complete genome assembly without manual intervention. While assembling long-read data alone can result in sequence errors and failure to detect small plasmids, hybrid assembly (using both short and long reads) can produce assemblies that are both complete and accurate. However, repetitive sequences in the genome, such as the genes encoding Bap and Blp1, are difficult to perfect even with hybrid assembly, so variations in these regions should be interpreted with caution.

Finally, as the original GenBank entries are replaced by revised genomes, there is a need to eliminate non-existent and incorrect predicted protein sequences in order to simplify the already complex task of protein sequence searches. It can be assumed that this problem is not only limited to *A. baumannii* genomes as many bacterial species so far have been sequenced using 454 pyrosequencing technology.

## Funding information

This work was supported by National Health and Medical Research Council (NHMRC) grant GNT1079616 to R.M.H. M.H. is supported by a Chancellor's Postdoctoral Research Fellowship (CPDRF PR017-4005) from the University of Technology Sydney, Australia. K.E.H. is supported by a Senior Medical Research Fellowship from the Viertel Foundation of Australia.

## Acknowledgements

We would like to thank Professor Thomas A. Russo, State University of New York, Buffalo, NY, USA, for kindly providing AB307-0294, and Professor Alessandra Carratoli, Istituto Superiore di Sanità, Rome, Italy, for supplying ACICU.

## Author contributions

Conceptualization, R.M.H., M.H.; data curation, M.H., R.W.; formal analysis, M.H., R.W., K.E.H., R.M.H.; funding, R.M.H., K.E.H., M.H.; investigation, M.H., R.W., L.J.; resources, K.E.H.; visualization, M.H., R.W., K.E.H.; manuscript preparation – original draft, M.H. and R.M.H.; manuscript preparation – review and editing, R.M.H., M.H., R.W., K.E.H.

## Conflicts of interest

The authors declare that there are no conflicts of interest.

## Ethical statement

No human nor animal experimentation is reported.

## Data bibliography

- Adams MD, Goglin K, Molyneux N, Hujer KM, Lavender H *et al.* NCBI GenBank accession no. CP012952, *A. baumannii* AB307-0294 (2008).
- Carattoli A, Villa L, Fortini D, Cassone A. NCBI GenBank accession no. CP000863.1, *A. baumannii* ACICU, complete genome (2007).
- Hamidian M, Wick R, Judd L, Russo TA, Holt KE *et al.* NCBI GenBank accession no. CP012952, *A. baumannii* AB307-0294 (2017).
- Hartstein RM, Hamidian M, Nigro SJ, Wick R, Judd L *et al.* NCBI GenBank accession no. CP031380.1, *A. baumannii* isolate ACICU (2019).
- Hua X. NCBI GenBank accession no. CP001937.2, *A. baumannii* MDR-ZJ06, complete genome (2018).
- Smith MG, Gianoulis TA, Pukatzki S, Mekalanos JJ, Ornston LN *et al.* NCBI GenBank accession no. CP000521.1, *A. baumannii* ATCC 17978, complete genome (2008).

7. Weber BS, Ly PM, Irwin JN, Pukatzki S, Feldman MF. NCBI GenBank accession no. CP012004.1, *A. baumannii* ATCC 17978-mff, complete genome (2013).

## References

- 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.
- Post V, Hall RM. AbaR5, a large multiple-antibiotic resistance region found in *Acinetobacter baumannii*. *Antimicrob Agents Chemother* 2009;53:2667–2671.
- Post V, White PA, Hall RM. Evolution of AbaR-type genomic resistance islands in multiply antibiotic-resistant *Acinetobacter baumannii*. *J Antimicrob Chemother* 2010;65:1162–1170.
- Eijkelkamp BA, Hassan KA, Paulsen IT, Brown MH. Investigation of the human pathogen *Acinetobacter baumannii* under iron limiting conditions. *BMC Genomics* 2011;12:126.
- Giannouli M, Antunes LCS, Marchetti V, Triassi M, Visca P et al. Virulence-related traits of epidemic *Acinetobacter baumannii* strains belonging to the international clonal lineages I-III and to the emerging genotypes ST25 and ST78. *BMC Infect Dis* 2013;13:282.
- Adams MD, Chan ER, Molyneaux ND, Bonomo RA. Genomewide analysis of divergence of antibiotic resistance determinants in closely related isolates of *Acinetobacter baumannii*. *Antimicrob Agents Chemother* 2010;54:3569–3577.
- 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.
- 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
- Quainoo S, Coolen JPM, van Hijum SAFT, Huynen MA, Melchers WJG et al. Whole-genome sequencing of bacterial pathogens: the future of nosocomial outbreak analysis. *Clin Microbiol Rev* 2017;30:1015–1063.
- Sahl JW, Gillece JD, Schupp JM, Waddell VG, Driebe EM et al. Evolution of a pathogen: a comparative genomics analysis identifies a genetic pathway to pathogenesis in *Acinetobacter*. *PLoS One* 2013;8:e54287.
- Sahl JW, Johnson JK, Harris AD, Phillipy AM, Hsiao WW et al. Genomic comparison of multi-drug resistant invasive and colonizing *Acinetobacter baumannii* isolated from diverse human body sites reveals genomic plasticity. *BMC Genomics* 2011;12:291.
- Chan AP, Sutton G, DePew J, Krishnakumar R, Choi Y et al. A novel method of consensus pan-chromosome assembly and large-scale comparative analysis reveal the highly flexible pan-genome of *Acinetobacter baumannii*. *Genome Biol* 2015;16:143.
- Balzer S, Malde K, Jonassen I. Systematic exploration of error sources in pyrosequencing flowgram data. *Bioinformatics* 2011;27:i304–i309.
- Brossard KA, Campagnari AA. The *Acinetobacter baumannii* biofilm-associated protein plays a role in adherence to human epithelial cells. *Infect Immun* 2012;80:228–233.
- Loehfelm TW, Luke NR, Campagnari AA. Identification and characterization of an *Acinetobacter baumannii* biofilm-associated protein. *J Bacteriol* 2008;190:1036–1044.
- Goh HMS, Beatson SA, Totsika M, Moriel DG, Phan M-D et al. Molecular analysis of the *Acinetobacter baumannii* biofilm-associated protein. *Appl Environ Microbiol* 2013;79:6535–6543.
- De Gregorio E, Del Franco M, Martinucci M, Roschetto E, Zarrilli R et al. Biofilm-associated proteins: news from *Acinetobacter*. *BMC Genomics* 2015;16:933.
- Watson M, Warr A. Errors in long-read assemblies can critically affect protein prediction. *Nat Biotechnol* 2019;37:124–126.
- 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.
- Weber BS, Ly PM, Irwin JN, Pukatzki S, Feldman MF. A multi-drug resistance plasmid contains the molecular switch for type VI secretion in *Acinetobacter baumannii*. *Proc Natl Acad Sci USA* 2015;112:9442–9447.
- Nigro SJ, Hall RM. Glsul2, a genomic island carrying the sul2 sulphonamide resistance gene and the small mobile element CR2 found in the *Enterobacter cloacae* subspecies cloacae type strain ATCC 13047 from 1890, *Shigella flexneri* ATCC 700930 from 1954 and *Acinetobacter baumannii* ATCC 17978 from 1951. *J Antimicrob Chemother* 2011;66:2175–2176.
- Hamidian M, Venepally P, Hall RM, Adams MD. Corrected genome sequence of *Acinetobacter baumannii* strain AB0057, an antibiotic-resistant isolate from lineage 1 of global clone 1. *Genome Announc* 2017;5:e00836-17
- 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:mgen.0.000242.
- Hua X, Xu Q, Zhou Z, Ji S, Yu Y. Relocation of Tn2009 and characterization of an ABGR13-2 from re-sequenced genome sequence of *Acinetobacter baumannii* MDR-ZJ06. *J Antimicrob Chemother* 2019;74:1153–1155.
- Blackwell GA, Holt KE, Bentley SD, Hsu LY, Hall RM. Variants of AbGR13 carrying the *armA* gene in extensively antibiotic-resistant *Acinetobacter baumannii* from Singapore. *J Antimicrob Chemother* 2017;72:1031–1039.
- Adams MD, Goglin K, Molyneaux N, Hujer KM, Lavender H et al. Comparative genome sequence analysis of multidrug-resistant *Acinetobacter baumannii*. *J Bacteriol* 2008;190:8053–8064.
- Iacono M, Villa L, Fortini D, Bordoni R, Imperi F et al. Whole-genome pyrosequencing of an epidemic multidrug-resistant *Acinetobacter baumannii* strain belonging to the European clone II group. *Antimicrob Agents Chemother* 2008;52:2616–2625.
- 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.
- Russo TA, Manohar A, Beanan JM, Olson R, MacDonald U et al. The response regulator BfmR is a potential drug target for *Acinetobacter baumannii*. *mSphere* 2016;1:e00082-16
- Vallejo JA, Beceiro A, Rumbo-Feal S, Rodríguez-Palero MJ, Russo TA et al. Optimisation of the *Caenorhabditis elegans* model for studying the pathogenesis of opportunistic *Acinetobacter baumannii*. *Int J Antimicrob Agents* 2015;S0924-8579(15)00241-1
- Wang-Lin SX, Olson R, Beanan JM, MacDonald U, Balthasar JP et al. The capsular polysaccharide of *Acinetobacter baumannii* is an obstacle for therapeutic passive immunization strategies. *Infect Immun* 2017;85:e00591-17
- Hamidian M, Ambrose SJ, Hall RM. A large conjugative *Acinetobacter baumannii* plasmid carrying the sul2 sulphonamide and strAB streptomycin resistance genes. *Plasmid* 2016;87-88:43–50.
- Kenyon JJ, Hall RM. Variation in the complex carbohydrate biosynthesis loci of *Acinetobacter baumannii* genomes. *PLoS One* 2013;8:e62160.
- Hamidian M, Hall RM. pACICU2 is a conjugative plasmid of *Acinetobacter* carrying the aminoglycoside resistance transposon TnaphA6. *J Antimicrob Chemother* 2014;69:1146–1148.
- Wilson K. Preparation of genomic DNA from bacteria. *Curr Protoc Mol Biol* 2001;56:2.4.1–2.4.5.
- Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014;30:2068–2069.
- Hamidian M, Hall RM. The AbaR antibiotic resistance islands found in *Acinetobacter baumannii* global clone 1 - structure, origin and evolution. *Drug Resist Updat* 2018;41:26–39.

38. Nigro SJ, Hall RM. Tn6167, an antibiotic resistance island in an Australian carbapenem-resistant *Acinetobacter baumannii* GC2, ST92 isolate. *J Antimicrob Chemother* 2012;67:1342–1346.
39. Bertini A, Poirel L, Mugnier PD, Villa L, Nordmann P et al. Characterization and PCR-based replicon typing of resistance plasmids in *Acinetobacter baumannii*. *Antimicrob Agents Chemother* 2010;54:4168–4177.
40. Hamidian M, Holt KE, Pickard D, Dougan G, Hall RM. A GC1 *Acinetobacter baumannii* isolate carrying AbaR3 and the aminoglycoside resistance transposon TnaphA6 in a conjugative plasmid. *J Antimicrob Chemother* 2014;69:955–958.
41. Smith MG, Gianoulis TA, Pukatzki S, Mekalanos JJ, Ornston LN et al. New insights into *Acinetobacter baumannii* pathogenesis revealed by high-density pyrosequencing and transposon mutagenesis. *Genes Dev* 2007;21:601–614.
42. Vallenet D, Nordmann P, Barbe V, Poirel L, Mangenot S et al. Comparative analysis of acinetobacters: three genomes for three lifestyles. *PLoS One* 2008;3:e1805.
43. Park JY, Kim S, Kim S-M, Cha SH, Lim S-K et al. Complete genome sequence of multidrug-resistant *Acinetobacter baumannii* strain 1656-2, which forms sturdy biofilm. *J Bacteriol* 2011;193:6393–6394.
44. Chen C-C, Lin Y-C, Sheng W-H, Chen Y-C, Chang S-C et al. Genome sequence of a dominant, multidrug-resistant *Acinetobacter baumannii* strain, TCDC-AB0715. *J Bacteriol* 2011;193:2361–2362.
45. Zhou H, Zhang T, Yu D, Pi B, Yang Q et al. Genomic analysis of the multidrug-resistant *Acinetobacter baumannii* strain MDR-ZJ06 widely spread in China. *Antimicrob Agents Chemother* 2011;55:4506–4512.
46. Liou M-L, Liu C-C, Lu C-W, Hsieh M-F, Chang K-C et al. Genome sequence of *Acinetobacter baumannii* TYTH-1. *J Bacteriol* 2012;194:6974.
47. Gao F, Wang Y, Liu Y-J, Wu X-M, Lv X et al. Genome sequence of *Acinetobacter baumannii* MDR-TJ. *J Bacteriol* 2011;193:2365–2366.

#### Five reasons to publish your next article with a Microbiology Society journal

1. The Microbiology Society is a not-for-profit organization.
2. We offer fast and rigorous peer review – average time to first decision is 4–6 weeks.
3. Our journals have a global readership with subscriptions held in research institutions around the world.
4. 80% of our authors rate our submission process as 'excellent' or 'very good'.
5. Your article will be published on an interactive journal platform with advanced metrics.

Find out more and submit your article at [microbiologyresearch.org](http://microbiologyresearch.org).