

Nanopore-only assemblies for genomic surveillance of the global priority drug-resistant pathogen, *Klebsiella pneumoniae*

Ebenezer Foster-Nyarko^{1,*}, Hugh Cottingham², Ryan R. Wick², Louise M. Judd², Margaret M. C. Lam², Kelly L. Wyres², Thomas D. Stanton¹, Kara K. Tsang¹, Sophia David³, David M. Aanensen³, Sylvain Brisse⁴ and Kathryn E. Holt^{1,2}

Abstract

Oxford Nanopore Technologies (ONT) sequencing has rich potential for genomic epidemiology and public health investigations of bacterial pathogens, particularly in low-resource settings and at the point of care, due to its portability and affordability. However, low base-call accuracy has limited the reliability of ONT data for critical tasks such as antimicrobial resistance (AMR) and virulence gene detection and typing, serotype prediction, and cluster identification. Thus, Illumina sequencing remains the standard for genomic surveillance despite higher capital and running costs. We tested the accuracy of ONT-only assemblies for common applied bacterial genomics tasks (genotyping and cluster detection, implemented via Kleborate, Kaptive and Pathogenwatch), using data from 54 unique *Klebsiella pneumoniae* isolates. ONT reads generated via MinION with R9.4.1 flowcells were basecalled using three alternative models [Fast, High-accuracy (HAC) and Super-accuracy (SUP), available within ONT's Guppy software], assembled with Flye and polished using Medaka. Accuracy of typing using ONT-only assemblies was compared with that of Illumina-only and hybrid ONT+Illumina assemblies, constructed from the same isolates as reference standards. The most resource-intensive ONT-assembly approach (SUP basecalling, with or without Medaka polishing) performed best, yielding reliable capsule (K) type calls for all strains (100% exact or best matching locus), reliable multi-locus sequence type (MLST) assignment (98.3% exact match or single-locus variants), and good detection of acquired AMR genes and mutations (88–100% correct identification across the various drug classes). Distance-based trees generated from SUP+Medaka assemblies accurately reflected overall genetic relationships between isolates. The definition of outbreak clusters from ONT-only assemblies was problematic due to inflation of SNP counts by high base-call errors. However, ONT data could be reliably used to 'rule out' isolates of distinct lineages from suspected transmission clusters. HAC basecalling + Medaka polishing performed similarly to SUP basecalling without polishing. Therefore, we recommend investing compute resources into basecalling (SUP model), wherever compute resources and time allow, and note that polishing is also worthwhile for improved performance. Overall, our results show that MLST, K type and AMR determinants can be reliably identified with ONT-only R9.4.1 flowcell data. However, cluster detection remains challenging with this technology.

DATA SUMMARY

All supporting data and protocols have been provided within the article or as supplementary data files. The four supplementary figures are available with the online version of this article, and via figshare (DOI: 10.6084/m9.figshare.19745608). Illumina reads have been deposited in the NCBI SRA, under the BioProject ID: PRJEB6891, and ONT reads (SUP basecalled) have been

Received 12 July 2022; Accepted 21 November 2022; Published 08 February 2023

Author affiliations: ¹Department of Infection Biology, London School of Hygiene and Tropical Medicine, London, UK; ²Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, VIC, 3004, Australia; ³Centre for Genomic Pathogen Surveillance, Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, Nuffield Department of Medicine, Oxford University, Oxford OX3 7LF, UK; ⁴Institut Pasteur, Université Paris Cité, Biodiversity and Epidemiology of Bacterial Pathogens, Paris, France.

*Correspondence: Ebenezer Foster-Nyarko, ebenezer.foster-nyarko2@lshtm.ac.uk

Keywords: AMR; bacterial pathogens; basecalling; benchmarking; genomic surveillance; *Klebsiella pneumoniae*; MLST; phylogenetic clustering; serotyping; Nanopore sequencing.

Abbreviations: AMR, antimicrobial resistance; HAC, high-accuracy basecalling; MDR, multi-drug resistance; MLST, multi-locus sequence type; ONT, Oxford Nanopore Technologies; ORF, open reading frame; PW, Pathogenwatch; SNP, single nucleotide polymorphism; ST, sequence type; SUP, super-accuracy basecalling; WGS, whole genome sequencing; WHO, World Health Organization.

Data statement: All supporting data, code and protocols have been provided within the article or through supplementary data files. Four supplementary figures and one supplementary table are available with the online version of this article.

000936 © 2023 The Authors



This is an open-access article distributed under the terms of the Creative Commons Attribution License. This article was made open access via a Publish and Read agreement between the Microbiology Society and the corresponding author's institution.

Impact Statement

Oxford Nanopore Technologies (ONT) sequencing devices are increasingly adopted to generate high-quality hybrid short-plus-long-read assemblies for bacteria, due to their portability and low capital cost. However, high read-level error rates have precluded the adoption of ONT devices as a standalone platform for bacterial pathogen genomic epidemiology studies or public health investigations. The widespread adoption of ONT for virus sequencing and outbreak analysis during the COVID-19 pandemic – including in settings where Illumina platforms are not available and are unlikely to become so due to capital costs, reagent access and other logistical issues – has created an unparalleled opportunity to harness ONT sequencing for genomic epidemiology and surveillance of bacterial pathogens. This is particularly attractive for multi-drug resistance-associated pathogens such as *Klebsiella pneumoniae*, where the ability of long reads to resolve plasmids and other highly variable genomic components is highly beneficial and advantageous over short reads. However, it is unclear what levels of accuracy to expect with currently widespread flowcells (Mk9.4.1) and basecalling algorithms, when using ONT-only assemblies to perform common tasks such as identifying antimicrobial resistance and virulence traits and phylogenetic clustering. Here, we investigated the utility of ONT-only assemblies for genotyping and phylogenetic clustering analysis, using 54 clinical *K. pneumoniae* isolates and free open-access software tools for genomic surveillance of *K. pneumoniae* (Pathogenwatch, Kleborate and Kaptive), which are suitable for non-informatics specialists such as public health or clinical microbiologists and epidemiologists.

deposited under BioProject PRJNA646837. The SUP-basecalled, short-read-first, long-read polished assembly for each isolate has been deposited in GenBank (BioProject IDs: PRJNA646837 and PRJEB6891), while the reference assemblies (i.e. SUP-basecalled long-read-first, short-read polished) have been shared via figshare (DOI: 10.6084/m9.figshare.19745608). Individual accession numbers for these reads and Biosample IDs are provided in Table S1 (available with the online version of this article) and via figshare (DOI: 10.6084/m9.figshare.19745608). Additionally, all alternative assemblies discussed here, including 10 per isolate constructed using different long-read basecallers and with/without polishing, are available in figshare (DOI: 10.6084/m9.figshare.19745608) [1].

INTRODUCTION

Klebsiella pneumoniae is a highly versatile species widely linked with multi-drug-resistant (MDR) healthcare-associated infections and hypervirulent community-acquired infections [2, 3]. It contributes a third of all healthcare-associated Gram-negative infections globally and is a leading cause of Gram-negative invasive diseases, second to *Escherichia coli* [2–4]. Worryingly, the acquisition rate of antimicrobial resistance (AMR) determinants among *K. pneumoniae* has continued to rise rapidly over the last decades, particularly resistance to third-generation cephalosporins and carbapenems. The Global Research on AntiMicrobial resistance (GRAM) study estimated that *K. pneumoniae* was one of six pathogens that caused more than 250 000 deaths associated with AMR in 2019 [5]. In that study, *K. pneumoniae* was the leading Gram-negative pathogen in sub-Saharan Africa, contributing 20% of the deaths attributable to AMR [5]. At the same time, data aggregated from 151 studies encompassing 26 countries in sub-Saharan Africa indicate that *K. pneumoniae* is the top Gram-negative and second most important cause of neonatal sepsis [4]. *K. pneumoniae* infections are also characterised by a high case fatality rate (18–49%) [2, 4, 6, 7]. Consequently, *K. pneumoniae* has been flagged in the World Health Organisation's (WHO) highest category of priority pathogens, for which the development of new antibiotics is urgently needed [8].

Accumulating genomic evidence has highlighted the emergence of hybrid MDR-hypervirulent clones, driven by the acquisition of mobile genetic elements [6, 7, 9–11]. This phenomenon represents a significant challenge to treating infections caused by *K. pneumoniae* and signals the urgent need for appropriate surveillance systems to identify MDR-hypervirulent clones and their genetic elements. WHO recommends the robust surveillance of AMR as a critical part of the Global Action Plan on AMR [12]. However, classical approaches to surveillance have relied on antimicrobial susceptibility testing and low-resolution genotyping methods, which yield little information on the evolution and spread of AMR within bacterial populations.

AMR surveillance in critical pathogens such as *K. pneumoniae* has been dramatically enhanced by the recent democratisation of whole-genome sequencing (WGS), primarily by short-read (Illumina) sequencing. WGS is currently the gold standard approach for public health pathogen surveillance, facilitating cluster identification and genotyping of clinically relevant loci, including AMR, virulence traits and antigen prediction [13–16]. Furthermore, as few labs perform *K. pneumoniae* serology, prediction of capsular (K) and lipopolysaccharide (O) antigen serotypes and their population distribution relies on inference from K and O antigen biosynthesis loci captured in WGS data [17, 18]. Genomic surveillance for K and O types is of increasing interest, as these antigens are the main target for alternative control measures such as vaccines, phage therapy and monoclonal antibody therapy [19–22]. Thirteen O antigen locus (OL) and 162 K antigen locus (KL) types have been described to date [17]. Unfortunately, the gold standard Illumina short-read sequencing platforms result in fragmented genome assemblies, which cannot

accurately resolve plasmids and other mobile genetic elements that drive the dissemination of AMR determinants via horizontal and lateral gene transfer. Fragmentation of Illumina genome assemblies also complicates K and O typing, resulting in up to a third of Illumina-based genomes being untypeable [23]. Long-read sequencing platforms, such as those of Pacific Biosciences and Oxford Nanopore Technologies (ONT), can overcome these limitations, as they yield longer reads (often exceeding 10 kb, with current maximal read lengths reaching 1 Mb) that are capable of resolving structural variations, long repeat regions and genomic copy-number alterations [24–26]. ONT is increasingly adopted to generate high-quality hybrid Illumina-plus-long-read assemblies for bacteria, owing to its flexibility and low capital cost compared to PacBio [27, 28]; however, high read-level error rates have precluded the adoption of ONT devices as a standalone platform for pathogen genomic epidemiology studies or public health investigations [29]. The situation has shifted somewhat during the COVID-19 pandemic as ONT has been widely adopted for virus sequencing and outbreak analysis, including in settings where Illumina platforms are not available and are unlikely to become so due to capital costs, reagent access and other logistical issues [30]. There is, therefore, now an unparalleled opportunity to harness ONT sequencing for genomic epidemiology and surveillance of bacterial pathogens, which is particularly attractive for AMR-associated pathogens such as *K. pneumoniae*, where resolving plasmids and other highly variable genomic components is highly beneficial.

The current state-of-the-art ONT basecaller, Guppy, utilises the Fast, High-accuracy (HAC) and Super-accuracy (SUP) models. The Fast model, as the name indicates, is the fastest of the three algorithms but at the cost of accuracy. SUP is the most accurate and slowest; the speed of the basecaller being a factor of the number of parameters in the neural network model. HAC is intermediate between Fast and SUP and about six times slower than the Fast model (<https://github.com/rrwick/August-2019-consensus-accuracy-update#basecalling>). ONT polishing, utilising Medaka—a neural network-based tool that generates consensus sequences by aligning the individual sequence reads against a draft assembly [31]—is rapid, highly effective and believed to improve sequence accuracy significantly [32]. ONT continuously refines these algorithms to improve raw reads and consensus accuracy [33].

However, despite continuous improvement over the years, ONT is still susceptible to non-trivial numbers of base-call errors, which can impact the identification of clinically important features such as acquired virulence and AMR genes [34, 35]. For example, a mean global error rate of ~6% is reported for reads with quality scores of at least 10—corresponding to 94% raw read accuracy and >99.99% consensus accuracy with the primary nanopore in use currently (R9.4.1 chemistry), although this varies by the organism [36] and is improving with new chemistries such as R10 [37]. For a 5 Mb genome, such as *K. pneumoniae*, we previously reported roughly 3000 errors (corresponding to at least 337 substitutions per genome assembly) [38]; therefore, ONT-only assemblies are prone to erroneous SNP calls, which can potentially hamper outbreak investigations [39]. The performance of ONT-only assemblies in identifying AMR and virulence traits against the gold-standard (hybrid assemblies) is unclear but is gaining attention [40–42], and is pertinent to inform the research community working with long-read data on what levels of accuracy to expect with current basecalling algorithms.

Here, we use a collection of 54 clinical isolates of *K. pneumoniae* with matched Illumina and ONT data [43] to assess the performance of ONT-only data when used for the prediction of AMR determinants, virulence factors, K/O typing and phylogenetic clustering analysis. Current read-based genotyping tools for Illumina data (such as ARIBA [44] or SRST2 [45]) are not optimised to be used efficiently on noisy long reads such as those produced by ONT. However, assembly-based tools such as Pathogenwatch [46], Kleborate [18, 47] and Kaptive [17] can in principle be applied to any assembly regardless of the sequencing platform, allowing for direct comparisons and facilitating the straightforward application of common genotyping workflows to data generated from different sequencing instruments deployed for pathogen surveillance. Furthermore, Pathogenwatch provides a free and accessible online analysis platform, suitable for non-informatics specialists such as public health or clinical microbiologists and epidemiologists, which, together with the growing network of ONT-equipped labs, could facilitate the widespread acceleration of genomic surveillance for AMR pathogens. We therefore focused our investigation on the performance of these tools (Pathogenwatch, Kleborate and Kaptive) for genomic surveillance of *K. pneumoniae* using ONT-only assemblies.

METHODS

Bacterial isolates and sequencing

We utilised 54 *K. pneumoniae* that were previously isolated at the Alfred Health Clinical Microbiological Diagnostic Laboratory in Melbourne, Australia, as part of a year-long prospective study of isolates from clinical infections (hospital-wide) and screening swabs (in intensive care and geriatric wards) in the Alfred Health network [43, 48–50]. The accession numbers for all genome data analysed in this study are provided in Table S1.

Isolates were sequenced first via Illumina HiSeq 2500, generating 125 bp paired-end reads, as described previously [48]. Later, selected isolates were subcultured, and fresh DNA was extracted for long-read sequencing using GenFind (Beckman Coulter) kits (full protocol available under DOI: 10.17504/protocols.io.p5mdq46). ONT libraries were prepared using the ligation protocol with barcoding kits to multiplex 12–24 isolates per run (barcode kits EXP-NBD104 and EXP-NBD114) and sequenced on a MinION device with R9.4.1 flowcells as previously described [27].

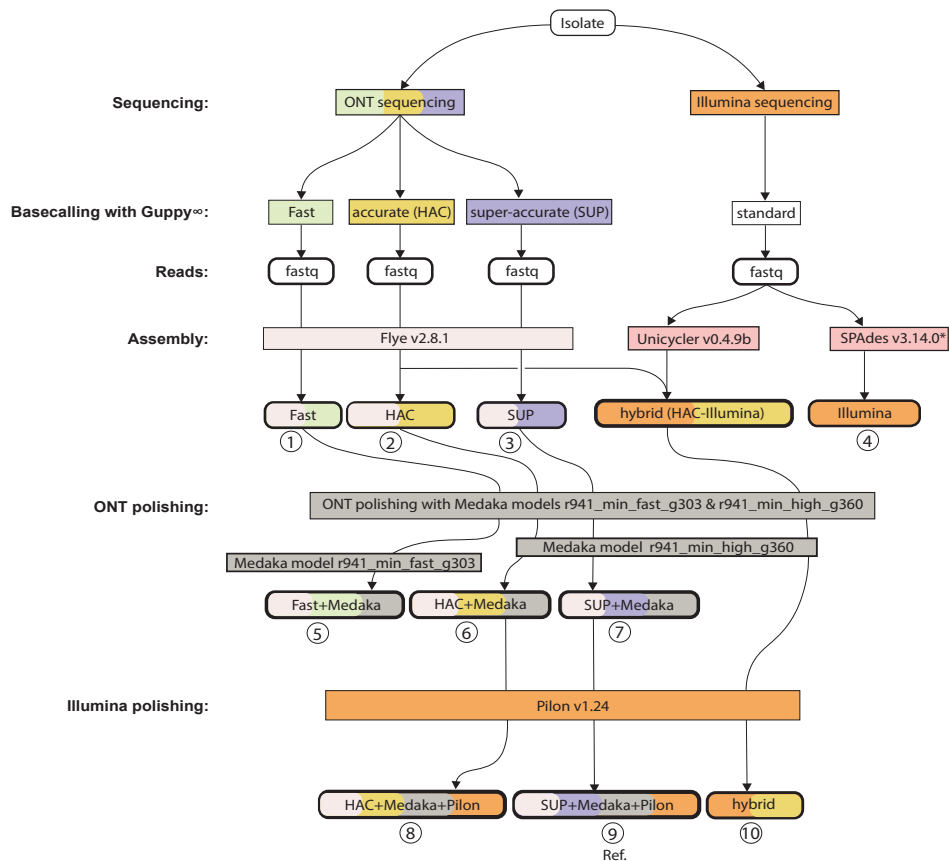


Fig. 1. Genome assembly flow diagram. We generated ten assemblies for each of 54 *Klebsiella pneumoniae* isolates, which were sequenced on both Illumina and ONT MinION (numbered 1–10). ONT reads were basecalled using three alternative models (via Guppy v4.0.14): fast, high-accuracy (HAC) and super-accuracy (SUP). The resulting basecalled read sets were each subjected to assembly using Flye (v2.8.1), generating three ONT-only assemblies per isolate. These assemblies were each then polished with the corresponding ONT read set (using Medaka models r941_min_fast_g303 and r941_min_high_g360) to generate a further three ONT-only polished assemblies per isolate, totalling six ONT-only assemblies (three polished, three unpolished) per isolate. For comparison, we used Illumina data to generate Illumina-only and ONT/Illumina hybrid assemblies for each isolate. Illumina reads were assembled using SPAdes v3.14.0 to generate an Illumina-only assembly. A short-read-first hybrid assembly was generated using Unicycler v0.4.9b and then polished using Pilon v1.24. Finally, two long-read-first hybrid assemblies were generated by polishing the ONT-only HAC and super-accuracy models using Illumina reads (via Pilon v1.24). *SPAdes was run via Unicycler v0.4.9b. ∞The Fast and HAC basecalling utilised Guppy v4.0.14 while SUP basecalling was done using Guppy v5.0.7.

Basecalling and assembly

Fast5 files generated by the ONT MinION were initially basecalled with Guppy v4.0.14 [51] on the Monash University M3 MASSIVE GPU cluster using the ‘Fast’ (r9.4.1_450bps_fast) and ‘High-accuracy’, aka HAC (r9.4.1_450bps_hac) models for each sample. The ‘Super-accuracy’ (r9.4.1_450bps_sup) model was later used when it was released in Guppy v5.0.7 [51, 52]. The basecalled FASTQ files were then concatenated into a single file per basecalled sequence run and demultiplexed into individual FASTQ files (one per sample) with the qcat command-line tool v1.1.0 [53] based on the barcode sequences. As a quality control step before assembly, we first filtered out poor-quality reads using fastp v0.20.1 [54] (short reads) and Filtlong v0.2.1 [55] (long reads) at a sequence similarity threshold of 90–95%, retaining only the reads with a minimum length of 1 kb and excluding the worst 5–10% of the reads (for FASTQ files less than 200 Mb, a keep_percent of 95% was used, while a keep_percent of 90% was applied to read files of size greater than 200 Mb).

Three separate approaches were employed to generate genome assemblies (Fig. 1). First, Unicycler v0.4.9b [56] was used to create Illumina-only assemblies for all the study isolates using default settings. When given short reads, Unicycler performs genome assembly using SPAdes v3.14.0 [57] but includes extra fine-tuning steps such as filtering out low-depth contigs and thus low-level contamination.

Next, we used Flye v2.8.1 [58] to generate ONT-only assemblies from basecalled reads (three read sets per genome, called with the three different basecallers). Flye was chosen based on our previous benchmarking of algorithms for ONT-only assemblies of

bacterial genomes [59]. Each ONT assembly was then polished with ONT reads using Medaka (v1.4.3) models r941_min_fast_g303 and r941_min_high_g360 [31], resulting in two ONT-only assemblies per basecalling model (polished vs unpolished); i.e., a total of six ONT-only assemblies per sample.

Finally, we produced two kinds of hybrid Illumina+ONT assemblies: short-read-first and long-read-first. Short-read-first hybrid assemblies were generated using Unicycler v0.4.9b, which starts by building a short-read assembly graph with SPAdes v3.14.0, then uses the corresponding long reads (HAC basecalled, in this case) to scaffold the genome, and finally runs Pilon v1.23 [60] in an attempt to fill gaps, correct bases and fix misassemblies using the short reads [27, 56] (HAC basecalled reads were used for the short-read-first assemblies as these were generated before the SUP basecalling model became available [52]). To generate long-read-first assemblies [61], we used Flye v2.8-1 to produce ONT-only assemblies for the set of reads basecalled with the HAC and SUP-accuracy basecalling models, followed by long-read polishing with Medaka [31] to repair any residual errors using ONT long reads, then finally short-read polishing using Illumina reads and Pilon v1.24 (following the recommendations noted in [62]). Thus, altogether, we produced ten assemblies per sample, encompassing reads derived from the three separate basecalling models (Fig. 1). QUAST v5.0.0 [de6973bb] [63] and CheckM v1.1.10 [64] were used to assess the quality and completion of the hybrid assemblies (see Supplementary File 1, available at <https://figshare.com/s/6026855223031e769d8a>; DOI: 10.6084/m9.figshare.19745608).

We designate the ONT-only assemblies, with and without Medaka ONT-read polishing, as 'Fast', 'Fast+Medaka', 'HAC', 'HAC+Medaka', 'SUP' and 'SUP+Medaka'. Based on the recommendations of Wick *et al.* [62], we consider the long-read-first hybrid assembly derived from SUP basecalled reads as the most accurate assembly and designate this as the reference sequence, against which the performance of other assemblies is benchmarked. We hereafter use the term 'hybrid' to refer to the short-read-first (Unicycler) hybrid assemblies, to distinguish them from long-read-first hybrid assemblies (otherwise referred to as 'SUP+Medaka+pilon' and 'HAC+Medaka+pilon').

Genotyping

Genome assemblies were uploaded to Pathogenwatch v2.3.1 [46] where Kleborate v2.2 [47] and Kaptive v2.0 [17] were automatically deployed to call multi-locus sequence types (STs) using the seven-locus scheme [65], capsular polysaccharide (K) and lipopolysaccharide O locus types, and serotype predictions, acquired virulence traits including the siderophores aerobactin (*iuc*), yersiniabactin (*ybt*) and salmochelin (*iro*), the genotoxin colibactin (*clb*) and the hypermucoidy locus (*rmpADC*). Pathogenwatch also deploys Kleborate to identify established AMR determinants (acquired genes and chromosomal mutations) [47] for the following antimicrobial classes: aminoglycosides, carbapenems, third-generation cephalosporins, third-generation cephalosporins plus β -lactamase inhibitors, colistin, fluoroquinolones, fosfomycin, penicillins, penicillins + β -lactamase inhibitors, phenicols, sulfonamides, tetracyclines, tigecycline and trimethoprim.

Clustering

We downloaded the Pathogenwatch pairwise distance matrix and corresponding neighbour-joining tree for the full set of assemblies. The distance matrix is available at <https://figshare.com/s/6026855223031e769d8a> (DOI: 10.6084/m9.figshare.19745608), and the tree for interactive viewing at Microreact (<https://microreact.org/project/sUrpBsvXi1aiKD7ssPv9pu-nanopore-only-assemblies-for-genomic-surveillance-of-klebsiella-pneumoniae>). Pathogenwatch calculates pairwise SNP distances between genomes based on a concatenated alignment of 1972 genes (2172367 bp) that make up the core gene library for *K. pneumoniae* in Pathogenwatch and infers a neighbour-joining tree from the resulting pairwise distance matrix [46]. Here, we assessed the feasibility of identifying potential nosocomial transmission clusters using these distance matrices. Several studies have proposed thresholds in the range of 21–25 genome-wide SNPs for identifying nosocomial transmission clusters of *K. pneumoniae* [66–68]. However, as Pathogenwatch calls SNPs only in 1972 core genes and not genome-wide, we compared the SNP distances calculated by Pathogenwatch with genome-wide SNP counts obtained by mapping short reads to a reference genome to determine the equivalent cut-off for clustering analysis using Pathogenwatch distances. To do this, we used the genome-wide SNP alignment generated previously for $n=270$ *K. pneumoniae* isolated at Alfred Health, based on mapping of Illumina reads to the *K. pneumoniae* NTUH-K2044 reference genome using the RedDog pipeline [69] (see full details in [43]). Pairwise SNP counts were extracted using snp-dist [70]. Assemblies for these 270 genomes (assembled from Illumina reads *de novo* using SPAdes optimised with Unicycler v0.4.74, see full details in [43]) were uploaded to Pathogenwatch, and the pairwise distance matrix was downloaded and compared against that generated from RedDog. We then used R to fit a linear regression model for Pathogenwatch distances as a function of genome-wide mapping-based SNP distances (see Fig. S1). This indicated that a Pathogenwatch distance threshold of 10 SNPs would be approximately equivalent to the established genome-wide distance threshold of 25. These thresholds assume accurate basecalling from Illumina data. To ascertain a corresponding threshold distance using ONT-only data, we compared pairwise Pathogenwatch distances calculated using ONT-only (SUP+Medaka) assemblies vs Illumina assemblies, for pairs of strains linked via probable transmission clusters. Using R to fit a linear regression model indicated that an ONT-only Pathogenwatch distance of 50 SNPs would approximate the Illumina-based Pathogenwatch distance of 10 SNPs or genome-wide distance of 25 SNPs (see Fig. 2).

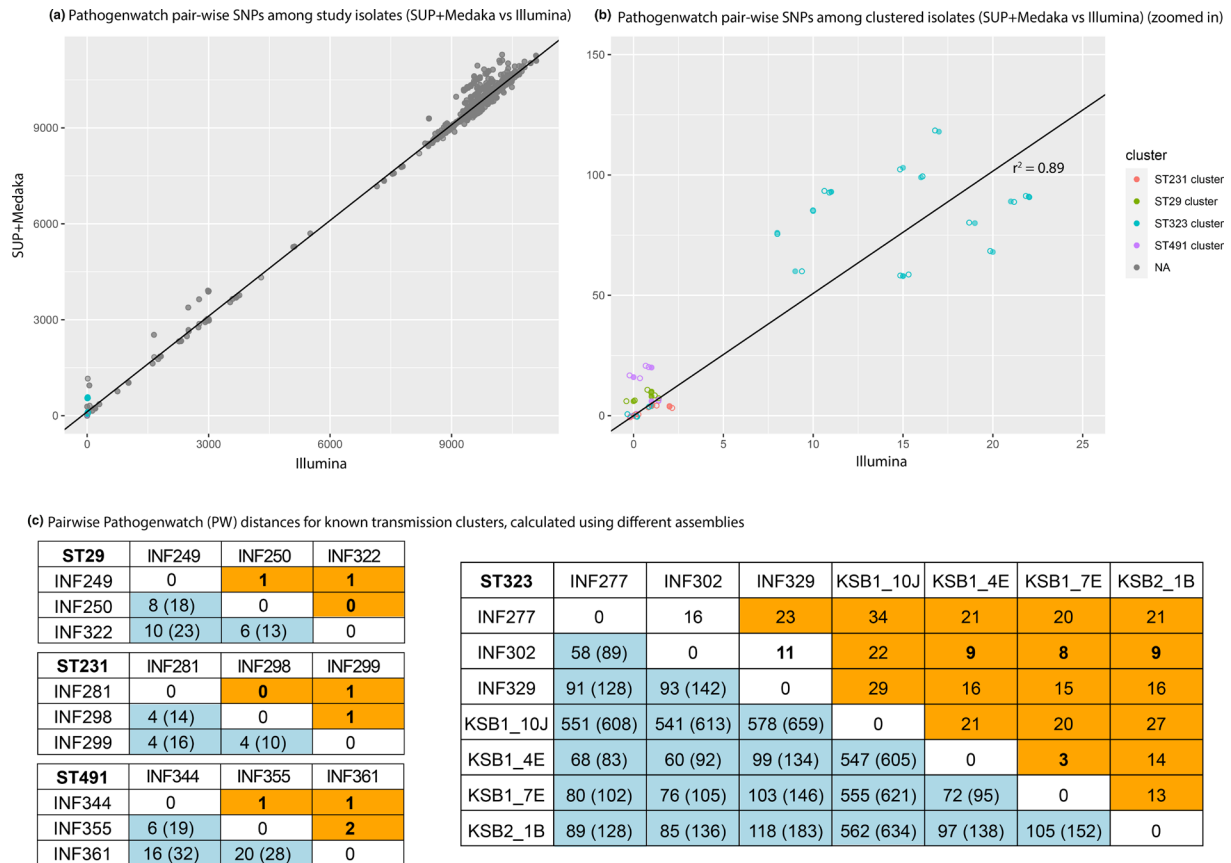


Fig. 2. Comparison of pairwise Pathogenwatch distances between $n=54$ isolates, calculated from ONT-only vs Illumina assemblies. The y -axis (a and b) shows pairwise distances between isolates, calculated using SUP+Medaka assemblies for all isolates; the x -axis shows pairwise SNP distances between isolates, calculated using Illumina assemblies for all isolates (based on the core gene set of 1972 genes within the Pathogenwatch *Klebsiella pneumoniae* core genome scheme). Each data point represents a pair of isolates, coloured to indicate isolate pairs belonging to the same transmission cluster (as per the inset legend). (a) All isolate pairs; line shows $y=x$. (b) Zoom in to transmission cluster isolates with ONT-based pairwise distances ≤ 150 ; linear regression line is shown (solid line), adjusted $R^2=0.8923$ indicating very good model fit, slope=5.1. (c) Pairwise Pathogenwatch (PW) distances for known transmission clusters, calculated using different assemblies. PW distances calculated from reference assemblies are highlighted in orange; those calculated from ONT-only with SUP basecalling are highlighted in blue (numbers in parentheses are those without Medaka polishing). Distances below the clustering threshold for PW distances, i.e. ≤ 10 , are in bold type.

We compared the topologies of neighbour-joining trees generated from Pathogenwatch distance matrices calculated using SUP+Medaka, HAC+Medaka or Fast+Medaka assemblies against the reference tree (calculated from hybrid SUP+Medaka+pilon assemblies), using the tanglegram function in the R package dendextend v1.15.2 to generate comparative tree plots and calculate entanglement coefficients. We also used the phytools package v1.0–3 in R to compute the Robinson–Foulds distance [71, 72] between tree topologies, which represents a sum of the number of partitions inferred by the first tree but not the second tree and that inferred by the second tree but not the first tree.

Data analysis and visualisation

Data were analysed and visualised using R v4.1.0. Linear regression was done using the lm function in base R, forcing the line through the origin. Genotyping and clustering results were visualised and compared using R packages ape v5.5 [73], cowplot v1.1.1 [74], data.table v1.14.2 [75], DECIPHER v2.20.0 [76], dendextend v1.15.2 [77], dplyr v1.0.7 [78], ggtree v3.0.4 [79], ggpubr v0.4.0 [80], gridExtra v2.3 [81], janitor 2.1.0 [82], kableExtra v1.3.4 [83], phytools v1.0.3 [72], tidyverse v1.3.1 [84], treeio v1.16.2 [85], treemap v2.4–3 [86] and stringr v1.4.0 [87].

RESULTS AND DISCUSSION

We investigated the accuracy of ONT-only assemblies for identifying ST, K and O antigen loci, AMR determinants, plasmid- and ICEKp-borne virulence factors and phylogenetic clusters, compared with assemblies generated via Illumina or hybrid methods that

combine Illumina and ONT reads (Fig. 1). We used a diverse set of 54 unique *K. pneumoniae* isolates (see Methods), spanning 30 STs and displaying a wide range of K types, O types, AMR and virulence profiles (see Table S1, available at <https://figshare.com/s/6026855223031e769d8a>, DOI: 10.6084/m9.figshare.19745608). The genome sizes were in the range 5075945 bp – 6163371 bp, and spanned G+C content of 49.6–57.7% (Fig. S2). Notably, the choice of ONT basecalling model did not appear to have a significant impact on the N50 or G+C content. Full details of sequence quality and assembly metrics are presented in Table S1 and Supplementary File 1.

MLST and clustering analysis

Perfect ST calls require exact sequence matches at all seven MLST loci; however, ST is often treated as a proxy for lineage, which can sometimes be correctly identified based on single-locus or double-locus variants of the correct ST (i.e. where 6/7 or 5/7 of the MLST loci have exact matches, these are reported as STx-1LV or STx-2LV respectively by Kleborate). The Illumina-only and SUP+Illumina hybrid assemblies gave identical MLST results for all genomes. Compared to this gold-standard result, for ONT-only assemblies, the proportion of correct ST calls was 87.3% for SUP basecalled assemblies with or without polishing, 78.2% for HAC+Medaka and 32.7% for non-polished HAC (see Fig. 3a). Fast basecalled assemblies showed poor results (20% correct with polishing, 0 without). ONT-only assemblies based on SUP basecalled reads reliably identified lineage in all cases when allowing for single- or double-locus variant calls. Assemblies based on HAC basecalling also performed well, with 96.3% (HAC+Medaka) and 90.7% (HAC) matching the true ST within 0–2 locus variants, respectively. Interestingly, for isolate KSB1_10J, the SUP+Medaka assembly yielded two loci mismatches with the expected ST323, while the SUP assembly without polishing yielded a single mismatch (i.e. a more accurate result).

Pathogenwatch calculates pairwise SNP distances (hereafter referred to as PW distances) across a set of 1972 *K. pneumoniae* core genes and constructs a neighbour-joining tree based on these distances. Fig. 4 shows a PW-distance tree for all 540 assemblies analysed here. PW distances between ONT-only assemblies vs their corresponding hybrid reference ranged from 2–882 SNPs (median 17) for SUP+Medaka to 198–3390 SNPs (median 374) for unpolished Fast basecalled assemblies (Fig. 5). Accordingly, the tree shows clear clustering of alternative assemblies for the same isolate, with SUP and HAC assemblies clustering closely with their corresponding hybrid references and Fast basecalled assemblies being distant relatives (see Figs 4 and 5, and interactive version of the tree in Microreact; <https://microreact.org/project/5chGLxaT1eVHKrThyc4b4J-nanopore-only-assemblies-for-genomic-surveillance-of-the-global-priority-drug-resistant-pathogen-klebsiella-pneumoniae>).

We then constructed assembly method-specific neighbour-joining trees for the set of 54 isolates, using PW distance matrices constructed from either reference assemblies, or SUP, HAC or Fast basecalled and polished ONT-only assemblies. The trees based on SUP or HAC ONT-only assemblies yielded generally similar topologies to the reference tree (entanglement coefficients 0.31 and 0.35, Robinson–Foulds distances 20 and 22, respectively), but the Fast basecalled ONT-only tree was more divergent (entanglement coefficient 0.96, Robinson–Foulds distance 102; see Figs 6, S3 and S4). However, clustering by ST was evident in all trees regardless of basecalling and assembly method.

Pairwise SNP distances can also be used directly, in the absence of trees, as a rule-in/rule-out indicator for potential transmission clusters. A consensus has emerged recently that a pairwise distance exceeding a threshold of 21–25 genome-wide SNPs is a reasonable basis for ruling out recent transmission of *K. pneumoniae* [66–68]. Based on comparisons using Illumina assemblies for $n=270$ unique *K. pneumoniae* clinical isolates, we estimate this equates to a PW distance of 10 SNPs (see Methods and Fig. S1). Fig. 2 shows PW distances calculated using reference assemblies, or SUP+Medaka ONT-only assemblies, for four groups of isolates in our data set that have previously been identified as part of nosocomial transmission clusters (based on a combined analysis of patient movement data and Illumina sequence, using genome-wide SNPs called from read mapping to an in-cluster reference sequence) [48, 87]. Clusters of ST29, ST231 and ST491 had pairwise PW distances of 0–2 using reference genomes, but much higher PW distances using ONT-only assemblies (ranges 6–10, 4–4 and 6–20, respectively using SUP+Medaka; 13–23, 10–16, 19–28, respectively using SUP without polishing). Cluster ST323 was more diverse, with PW distances of 3–34 between reference genomes and 58–659 using ONT-only (excluding one outlier strain, KSB1_10J, the distances were 3–23 for reference assemblies, 58–118 for SUP+Medaka and 83–183 for SUP). Linear regression of polished ONT-only PW distances vs Illumina PW distances yielded a slope of ~ 5 (adjusted $r^2=0.89$, see Fig. 2), suggesting a relaxed PW distance threshold of ~ 50 could be suitable for calling transmission clusters using ONT-only (SUP+Medaka) data, although the relationship does seem to vary by ST (see Fig. 2). We therefore conclude that Pathogenwatch analyses of ONT-only assemblies are quite reliable for identification of lineage or ST (as described above), and thus could be reliably used to ‘rule out’ isolates of distinct lineages from suspected transmission clusters (i.e. unrelated genomes/isolates). However, more work is needed to establish suitable analytical methods for calculating and interpreting smaller distances within lineages, as is required to differentiate recent transmission chains from coincidental but independent infections with the same lineage.

K and O locus typing

In *K. pneumoniae*, K and O locus types vary in their gene content, and typing via the Kaptive software is achieved by searching genomes for full-length reference K or O loci and confirming the presence of the expected genes within the best-matching locus

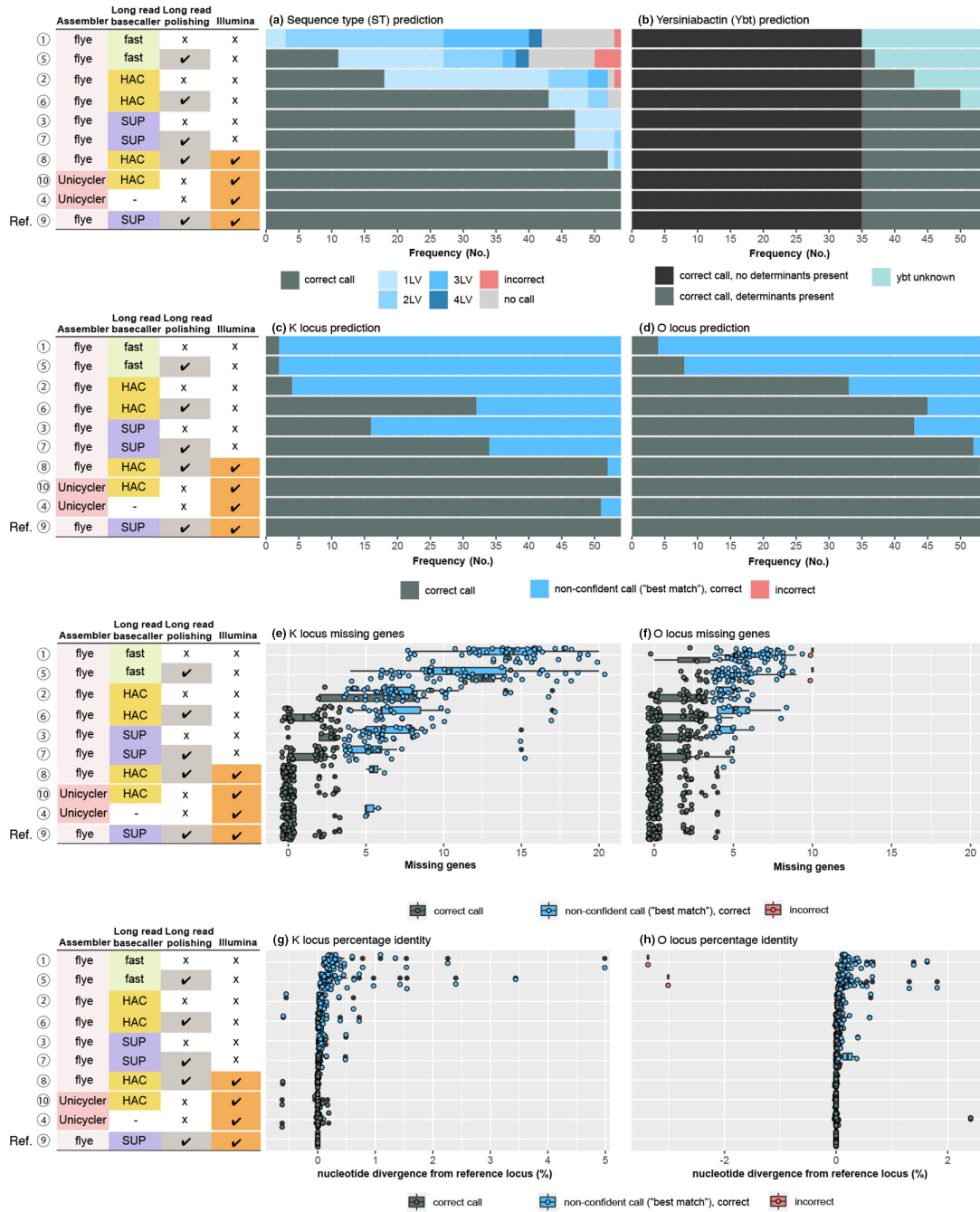


Fig. 3. Genotyping accuracy for MLST, yersiniabactin and K/O loci using different assemblies. Each panel summarises the accuracy of genotyping results across $n=54$ isolates, for each type of assembly (one per row as labelled on the left; numbered 1–10 as in Fig. 1, and ordered based on overall performance), compared to the reference assembly (i.e. long-read-first hybrid assembly using Flye and Medaka to assemble and polish with super-accurate basecalled reads, then polished with Illumina; last row per panel). (a) Multi-locus sequence type (ST) calling; 1LV, 2LV, 3LV and 4LV indicate the correct ST was predicted but the allele call was incorrect for one, two, three or four alleles. (b) Detection and typing of yersiniabactin; ‘correct call, no determinants present’ indicates the *ybt* locus was absent in the reference assembly and was (correctly) not detected in the test assembly; ‘correct call, determinants present’ indicates the *ybt* locus was present in the reference assembly and this was correctly identified and subtyped (to lineage level) in the test assembly; ‘determinants present, called as unknown’ indicates the *ybt* locus was present in the reference assembly, and this was correctly identified as present but could not be subtyped in the test assembly. (c) K locus and (d) O locus typing; ‘correct call’ indicates the correct locus was identified with a confidence ranking of good or better; ‘non-confident call (‘best-match’), correct’ indicates that the correct locus was identified but with low or no confidence. (e, f) The number of ORFs that are encoded in the reference K and O loci but were not detected in the test assemblies (labelled as ‘missing genes’ by the Kaptive genotyper, although in this case the nucleotide sequences are present but basecall errors result in disruption of the ORFs). (g, h) The nucleotide divergence between each assembly and the reference K and O loci, across the length of the detected loci.

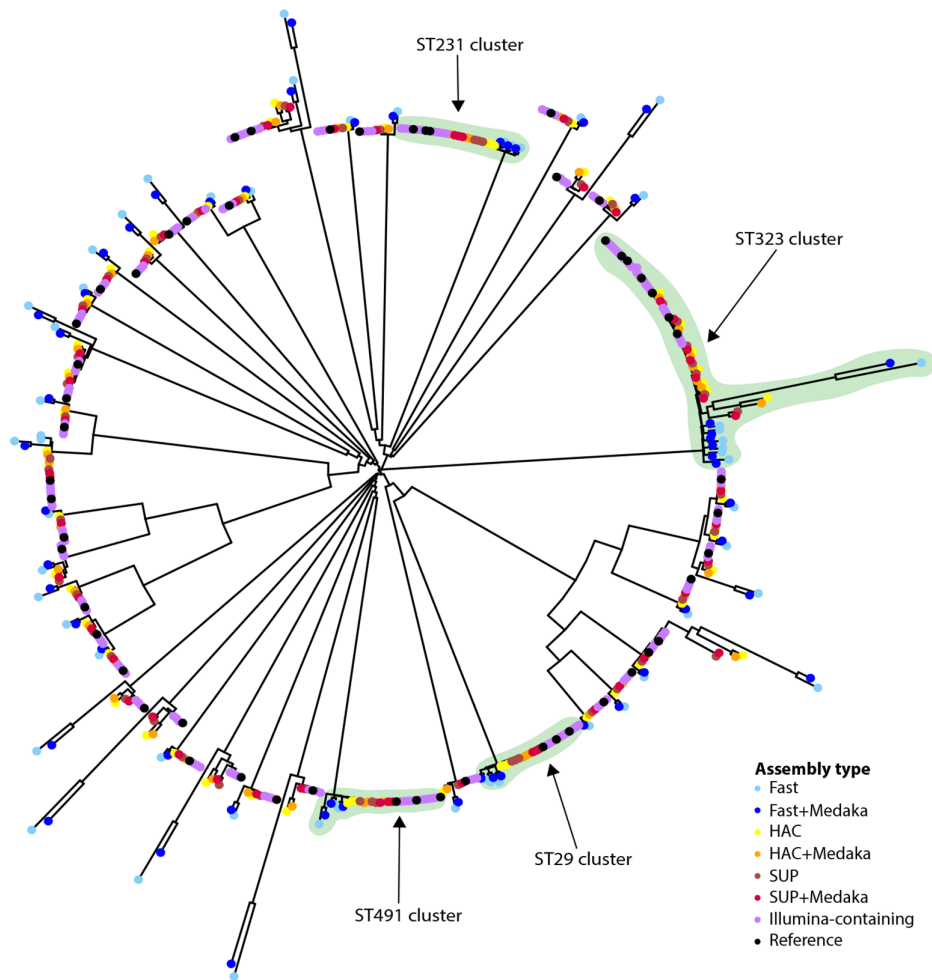


Fig. 4. Neighbour-joining tree representing the phylogenetic relationships among the 540 genome assemblies, derived from $n=54$ *K. pneumoniae* isolates as outlined in Fig. 1. The tree was constructed from the Pathogenwatch distance matrix (i.e. based on SNPs called in 1972 core genes) and is midpoint-rooted. Tips are coloured by assembly type, as indicated by the inset legend; note the reference assembly (i.e. long-read first hybrid assembly using Flye and Medaka to assemble and polish with super-accurate basecalled reads, then polished with Illumina) is coloured black and all other Illumina-containing genomes are grouped and coloured purple. Alternative assemblies of the same genome tended to cluster together but are non-zero branches because the sequences are non-identical due to basecalling errors. ONT/Illumina hybrid assemblies of the same isolates clustered more closely than ONT-only assemblies. Assemblies derived from the Fast model (\pm Medaka, blue colours) consistently formed outliers separated by long branches, due to high rates of basecalling errors. Most of the study isolates belong to distinct *K. pneumoniae* lineages and are unrelated, except for 16 isolates belonging to four transmission clusters (see Fig. 2); these clades are highlighted in green and labelled by multi-locus sequence type (ST).

[17, 18, 88]. K locus (KL) and O locus (OL) calls were in perfect agreement between Illumina and hybrid assemblies. Compared to these results as reference, all ONT-only assemblies correctly identified the best-matching KL type. However, the confidence reported by the Kaptive typing tool – which depends on the detection of intact genes (i.e. ORFs) in the K locus and is thus expected to be impacted by basecall accuracy – varied widely for the different assemblies (see Fig. 3c). The distribution of missing gene count, stratified by confidence, is shown in Fig. 3e; this supports the expectation that the lack of confidence in KL calls is due to disruption of ORFs in the locus, presumably due to basecall errors that result in frameshifts or stop codons in the expected coding genes. Fig. 3g shows that nucleotide similarity (compared to KL reference sequences) along the K locus was quite high for all assemblies, with nucleotide divergence well below the Kaptive threshold for a confident call ($\leq 5\%$ divergence). This confirms that even small numbers of basecall errors are enough to disrupt ORFs and thus reduce confidence in genotype calls.

The results were similar for O typing, although some incorrect OL calls were made with the Fast assemblies (1.85% each with Fast and Fast+Medaka; red colours in Fig. 3d). A correct O1 prediction requires the detection of an extra gene (*wbbY*) outside the O locus [18, 89]; if this ORF is disrupted, O1 may be miscalled as O2. Unsurprisingly, the discrepancies with the Fast OL calls arose from O1 being erroneously reported as O2, due to frameshift mutations arising from basecalling/assembly errors affecting the unlinked *wbbY* gene.

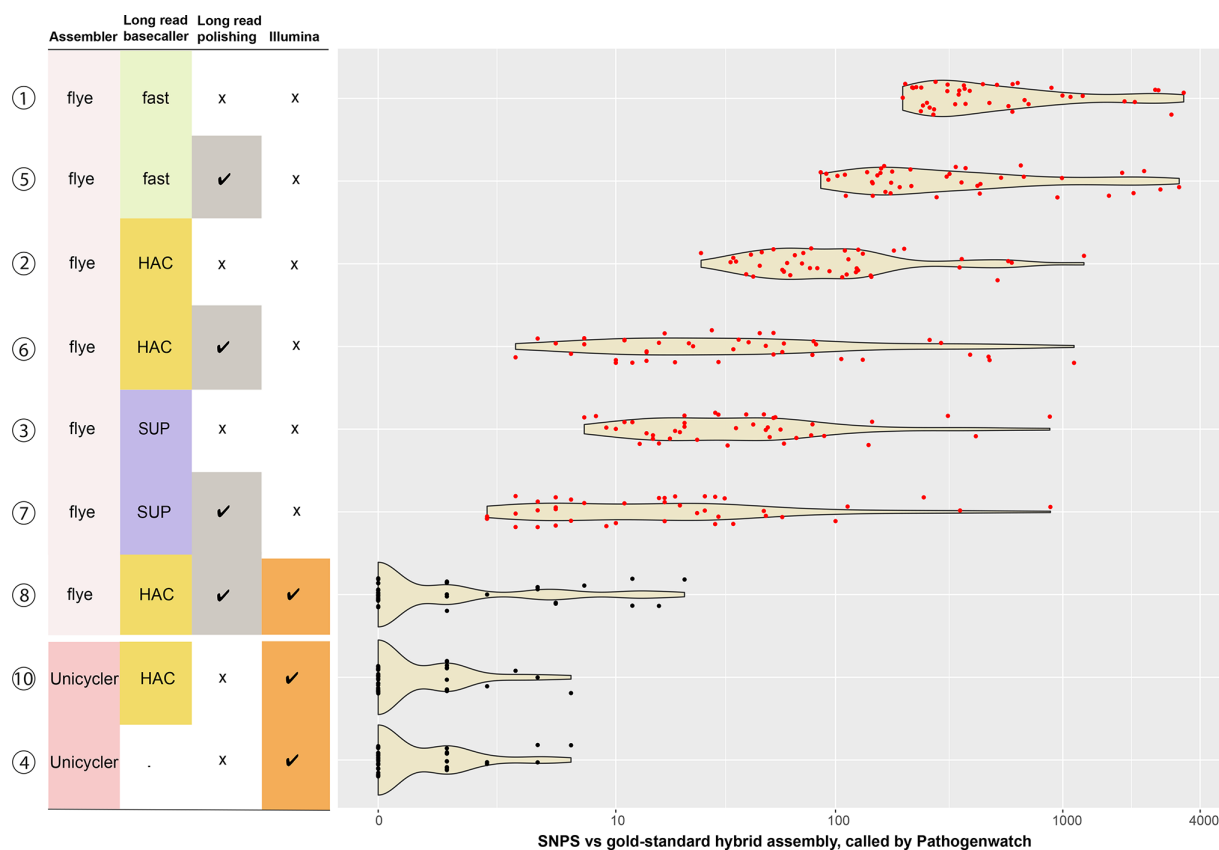


Fig. 5. SNPs called vs reference assembly. A violin plot depicting the distribution of SNPs between each assembly vs the gold-standard hybrid reference assembly for the 54 study genomes based on a concatenated alignment of 1972 genes (2172367 bp) that make up the core gene library for *K. pneumoniae* in Pathogenwatch. Assembly types are numbered 1–10 as in Fig. 1. Data points corresponding to ONT-only assemblies are coloured red. Within the cluster of ten assemblies for each sample, the SUP and HAC were consistently separated by a reasonable SNP threshold (≤ 10 SNPs); however, the Fast assemblies tended to form outliers (>100 SNPs). For clustering, only the SUP+Medaka assemblies reliably identified outbreak clusters (≤ 25 genome-wide SNPs threshold, which corresponds to ≤ 10 SNPs within the Pathogenwatch *K. pneumoniae* core genome typing scheme). HAC+Medaka polishing performed similarly to SUP without polishing.

AMR determinants

Kleborate screens for acquired AMR genes by interrogating each input assembly against a species-aware modified version of the Comprehensive Antibiotic Resistance Database [47]. A total of 270 acquired genes were identified across the 54 reference genome assemblies (median nine per genome, range 0–18). These acquired genes were mostly identified correctly from SUP+Medaka (95.8%), SUP (84.8%), HAC+Medaka (82.7%) and HAC (60.3%) assemblies, with lower recovery rates from Fast+Medaka (53.2%) and Fast (25.7%) assemblies. Kleborate also screens for mutations in chromosomal core genes known to be associated with AMR [47]. Here, reportable mutations were identified from reference assemblies in the *gyrA* and *parC* genes (substitutions associated with fluoroquinolone resistance, $n=9$ isolates), in *blaSHV* (substitutions associated with extended-spectrum β -lactamase activity, $n=27$ genomes), and in *ompK35* and *ompK36* (truncations associated with reduced susceptibility to β -lactams, $n=2$ genomes). Truncations in *pmrB* and *mgrB* (associated with colistin resistance) are screened by Kleborate, but these genes were intact in all reference assemblies. The substitutions in *gyrA*, *parC* and *blaSHV* were recovered quite accurately from the SUP+Medaka (91.7%) and SUP (88.9%) assemblies, but less so using HAC+Medaka (83.3%), HAC (75%), Fast+Medaka (55.6%) and Fast (44.4%) assemblies (see Fig. 7a, b). The two *omp* gene truncations were identified by all assembly methods, however, truncations of *omp*, *pmrB* and *mgrB* were also falsely identified in many cases (7.4% with SUP+Medaka, 7.4% SUP, 14.8% HAC+Medaka, 9.3% HAC, 77.8% Fast+Medaka, 27.8% Fast; see Fig. 7c, d).

As the presence of AMR determinants in genomes is often used to predict or explain resistance to clinically relevant drugs, we also considered the impact of ONT basecalling and assembly method on the task of identifying which isolates are likely to be non-susceptible to specific drug classes. Table 1 shows the proportion of genomes that would be correctly classified as non-susceptible due to the detection of AMR determinants known to be present in the reference genome. SUP+Medaka assemblies performed well, with 88–100% of genomes with AMR determinants correctly identified as such across the various drug classes.

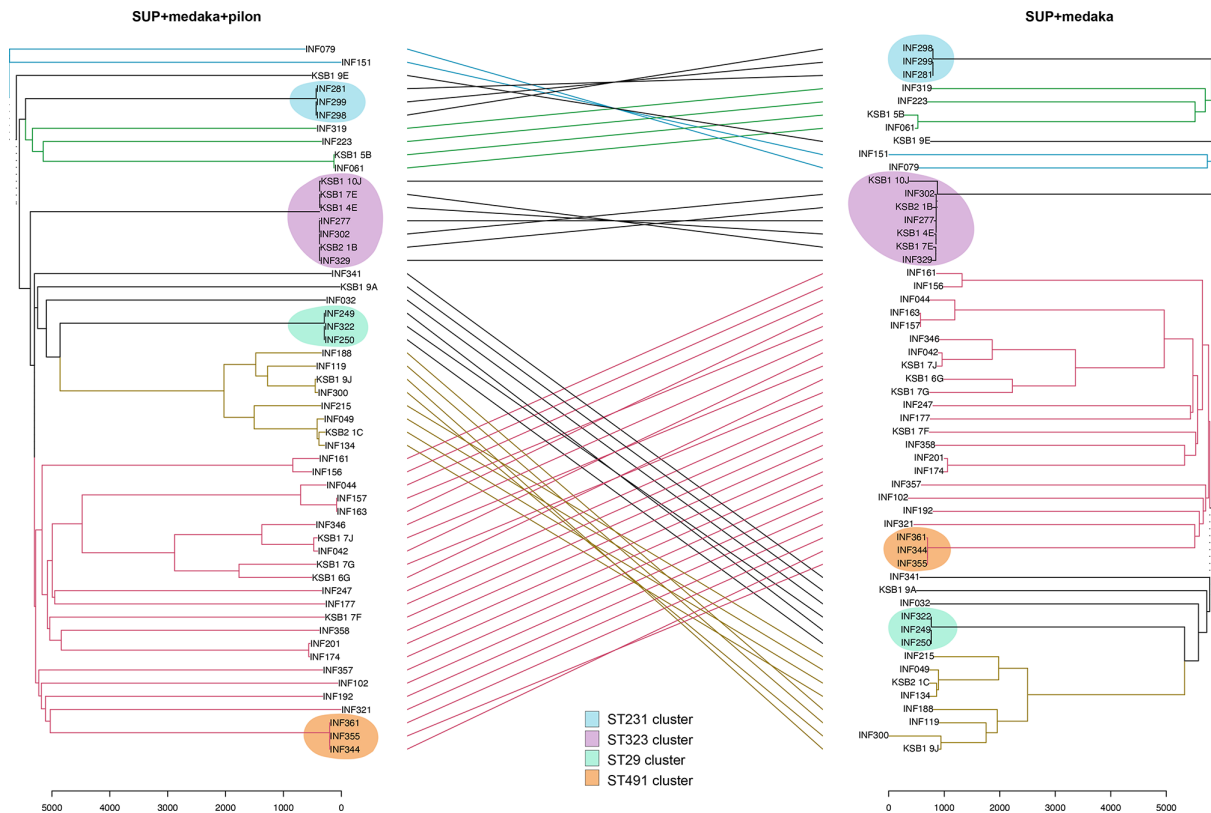


Fig. 6. Comparison of trees constructed from ONT-only vs reference assemblies. Trees are midpoint-rooted neighbour-joining trees inferred from Pathogen pairwise distance matrices calculated from ONT-only (SUP+Medaka) assemblies or ONT/Illumina long-read-first hybrid reference assemblies. Trees are shown as a tanglegram, which attempts to render the tree visualisations to maximise the alignment of tips, and links tips with lines (coloured to highlight shared subtrees), to facilitate comparison of topology between the two trees. The tanglegram yields an entanglement coefficient of 0.33 and a Robinson–Foulds distance of 20, indicating a good alignment. Coloured lines join matching tip labels for common subtrees, while black lines indicate subtrees that are not common between the two trees. Discordant branch lengths on the ONT-only tree arise from basecalling errors, despite a good alignment.

HAC+Medaka also showed good detection rates (60–100%), while SUP detection rates were more variable, ranging from 42 to 100% (see Table 1).

Acquired virulence

Kleborate reports the presence of five major acquired virulence loci, including the ICEKp-encoded siderophore yersiniabactin (*ybt*) and genotoxin colibactin (*clb*), plus three loci that are typically plasmid-borne and commonly linked with invasive infections caused by hypervirulent *K. pneumoniae*, namely the siderophores aerobactin (*iuc*) and salmochelin (*iro*), and the hypermucoidy locus *rmpADC*. Each of these loci comprises several genes, which are included in locus-specific MLST schemes for each [47, 90, 91]. Each locus will be reported, and an ST (or the closest match) assigned following the same logic to assign STs as for seven-gene MLST, if >50% of the genes contained in the locus are detected. Lineages of each virulence locus are inferred based on virulence STs, with each lineage pre-assigned to a set of virulence STs [90, 91]. Acquired virulence genes detected in our reference genomes include *ybt*, *clb* and *iuc*. The genes *iuc* and *clb* were present in only two isolates each (*iuc* in INF079 and INF151, *clb* in INF079 and INF341), while *ybt* was detected in 20 isolates (36.4%). When present, the Illumina and hybrid assemblies correctly detected both *iuc* and *clb*, including correct lineage assignment. For the ONT-only assemblies, SUP (\pm Medaka) and HAC+Medaka correctly reported the *iuc* and *clb* STs and lineages, while the HAC and Fast (\pm Medaka) detected the loci as present but failed to identify an ST match and thus could not identify a lineage. Similarly, for *ybt*, the Illumina and hybrid assemblies correctly reported the presence or absence of the *ybt* locus, and accurately predicted the lineage when the locus was present (100% each). However, for ONT-only assemblies, only the SUP (\pm Medaka) matched the Illumina and hybrid in uniformly detecting *ybt* and assigning the correct lineage (Fig. 3b). The other ONT-only assemblies all reported the presence or absence of *ybt* but failed to call the ST and thus to report a lineage (HAC+Medaka, 7.4%; HAC, 20.4%; Fast+Medaka, 31.5%, Fast, 35.2%). Notably, there were no false positives (i.e., reporting the locus present when absent), nor any false negatives (i.e., failing to detect presence of the locus when present).

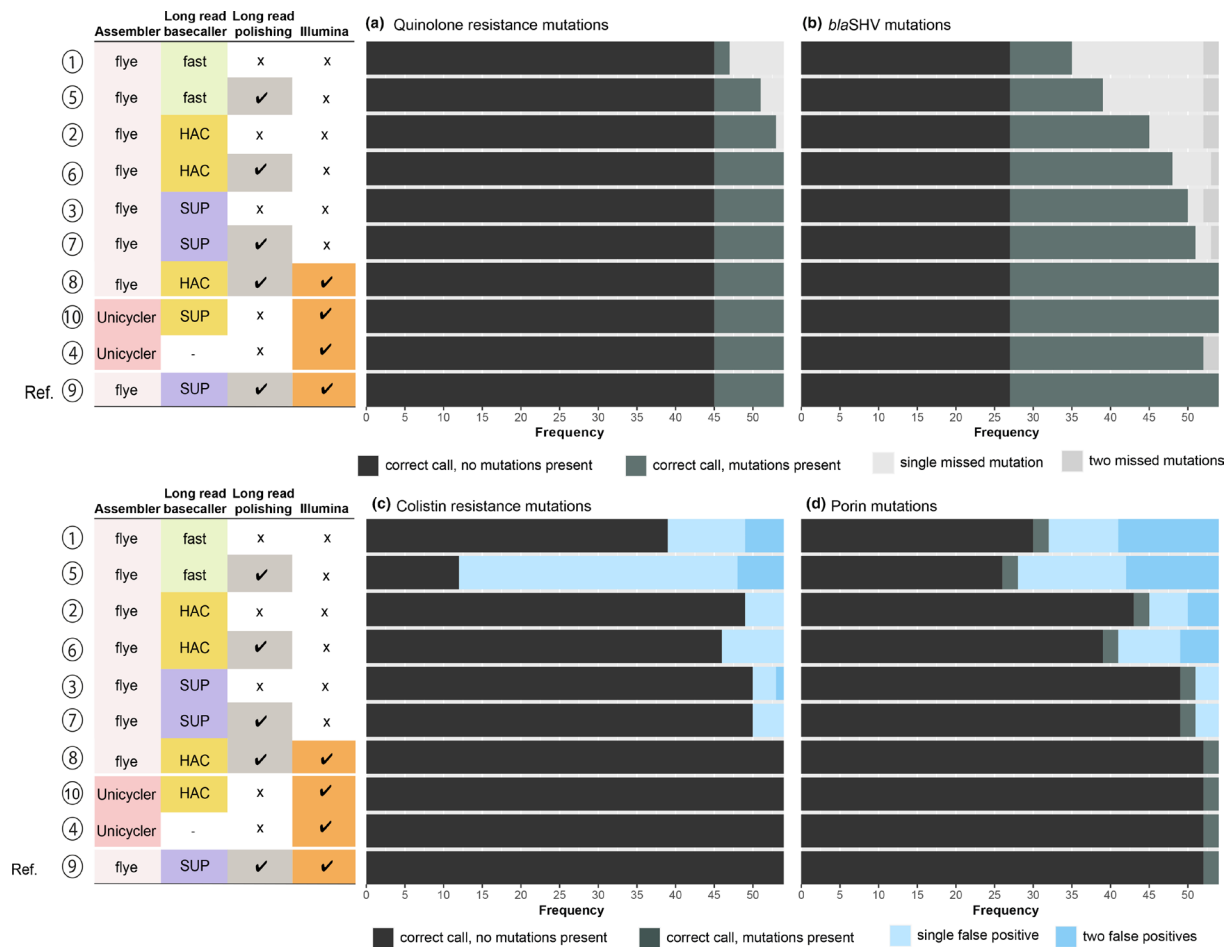


Fig. 7. Detection of AMR-associated mutations in core chromosomal genes. Each panel summarises the accuracy of genotyping results across $n=54$ isolates, for each type of assembly (one per row as labelled on the left; numbered 1–10 as in Fig. 1, and ordered based on overall performance), compared to the reference assembly (i.e. long-read-first hybrid assembly using Flye and Medaka to assemble and polish with super-accurate basecalled reads, then polished with Illumina; last row per panel). (a) Detection of fluoroquinolone resistance-associated mutations (GyrA-83I, ParC-80I, GyrA-83Y, GyrA-87N and GyrA-87Y). (b) Detection of mutations in *bla*SHV that are associated with a change in enzyme activity (35Q, 146V, 238S and 156D). (c) Detection of colistin resistance-associated mutations (disruption of MgrB or PmrB). (d) Detection of mutations in outer membrane porins associated with a change in carbapenem susceptibility (truncation of OmpK35 or OmpK36, GD or TD insertions in OmpK36 loop 3).

Limitations and future directions

An artefact of our sample collection is that there were few virulence plasmid-positive strains, so we could not assess other virulence genes apart from yersiniabactin. However, it is clear from the available results that genotyping based on detection of gene presence/absence is generally quite reliable using assemblies inferred from SUP-basecalled ONT reads, hence virulence gene detection is expected to be well tolerated generally. There were also very few strains in our collection with genuine AMR-associated truncations, so the call rate for these could not be accurately established. However, clearly, basecall errors frequently disrupt ORFs, such that AMR-related truncations are subject to relatively high false-positive call rates. The same issue also affects the confidence of K and O locus typing: whilst the correct locus is almost always detected, the reported confidence of the call depends on the detection of ORFs for all genes in the locus, which is rarely achieved using ONT data alone.

Our study considered *K. pneumoniae* only. However, the basecalling, assembly and genotyping tasks explored here are commonly applied across other bacterial pathogens. In particular, since *Enterobacterales* share common methylation patterns (which is a major determinant of basecalling accuracy) and AMR mechanisms, the results presented here can be considered informative as to the general accuracy of ONT-only analysis of the family more broadly.

Notably, laboratory and informatics methods for ONT sequencing are under constant development—for example, new flowcells have recently been released with new (R10) pores; basecalling models can be trained on a larger or species-specific dataset, and assembly and polishing methods are under active development [92]. These developments are expected to lead to improvements in

Table 1. Accuracy of identifying non-susceptibility to specific drugs/classes based on detection of known AMR determinants

	Aminoglycosides	Carbapenems	Third-generation cephalosporins	Fluoroquinolones	Phenicol	Sulfamethoxazole	Trimethoprim	Tetracycline
Reference genomes with AMR determinants, N (% of total)	26 (48.1%)	6* (11.1%)	25 (46.3%)	27† (50.0%)	17 (31.5%)	26 (48.1%)	24 (44.4%)	20 (37.0%)
ONT-only assemblies with AMR determinants, N (% of reference genome positives)								
SUP+Medaka	24 (92.3%)	6 (100%)	22 (88%)	26 (96.3%)	16 (94.1%)	26 (100%)	24 (100%)	20 (100%)
SUP	11 (42.3%)	6 (100%)	21 (84%)	27 (100%)	14 (82.4%)	25 (96.2.3%)	24 (100%)	11 (55%)
HAC+Medaka	20 (76.9%)	6 (100%)	15 (60%)	27 (100%)	14 (82.4%)	24 (92.3%)	21 (87.5%)	15 (75%)
HAC	2 (7.7%)	6 (100%)	19 (76.0%)	27 (100%)	14 (82.4%)	8 (30.8%)	14 (58.3%)	2 (10%)
Fast+Medaka	3 (11.5%)	6 (100%)	8 (32%)	17 (63.0%)	13 (76.5%)	21 (80.9%)	19 (79.2%)	4 (20%)
Fast	1 (3.8)	5 (71%)	6 (24%)	13 (48.1%)	12 (70.6%)	4 (15.4)	8 (33.3%)	1 (5%)

Reference assemblies are long-read-first ONT/Illumina hybrid assemblies generated using SUP-based called ONT data.

*Carbapenem resistance determinants detected were mainly acquired genes ($n=4$ genomes) rather than porin mutations ($n=2$ genomes).

†Fluoroquinolone resistance mutations detected were mainly acquired genes ($n=18$ genomes) rather than point mutations ($n=9$ genomes). Details of resistance determinants detected in each reference genome are given in Table S1.

the accuracy of ONT-only assemblies and therefore the accuracy of genotyping from such assemblies [37]. Therefore, our results should be considered a baseline from which improvements are expected to accrue. In addition, it is likely that the task of pairwise SNP calling could be more accurately performed using read-based variant calling approaches, rather than the assembly-based methods used here. However, there is currently a lack of SNP-calling tools designed specifically for bacterial ONT data (notably those trained on human data do not perform well on bacterial data due to differences in methylation [38]), so exploring this was beyond the scope of the current study, which focused on the use of a consistent assembly-based approach to all analysis tasks (implemented in the freely accessible online Pathogenwatch platform).

Despite these limitations, our study provides a timely update on the level of accuracy that can be achieved with currently well-established and widely available protocols, including (i) R9.4.1 pores (which have been widely and stably used for several years now); (ii) out-of-the-box basecalling with Guppy, which can be done in real-time, and on-board with devices such as Mk1C or GridION; (iii) a simple, rapid low-resource assembly with Flye (assessed in several benchmarking studies as the best singular assembler for ONT-only); and (iv) the well-established and stable Medaka polisher.

Conclusions

Overall, our results show that MLST, K/O locus type, virulence and AMR determinants can be reliably identified from ONT-only genome assemblies. However, pairwise SNP distance estimation was less reliable and thus we propose that ONT-only analysis should be considered reliable only to rule-out potential clusters using MLST-based lineage assignments, rather than being the sole trigger for specific actions that are usually based on a high suspicion of transmission (e.g., enhanced containment procedures in hospital settings), as it currently lacks the sensitivity needed for public health or infection control investigations. Where compute resources and/or time are limiting, our data indicate that compute resources are best directed towards basecalling (SUP model), noting that polishing is also worthwhile for improved performance.

Funding information

This work was supported by the Bill and Melinda Gates Foundation, Seattle (OPP1175797, KlebNet Project to D.M.A., K.E.H., S.B.). The funding bodies had no role in study design or in data collection, analysis and interpretation.

Author contribution

Conceptualisation: K.E.H. Methodology and Software: K.E.H., K.L.W., M.M.C.L., R.R.W. Formal Analysis and Visualisation: E.F.N., H.C., K.L.W., M.M.C.L., K.E.H. Data curation: K.E.H., H.C. Supervision: K.E.H. Funding: K.E.H. Writing – Original Draft Preparation: E.F.N., K.E.H. All authors read and approved the final manuscript.

Conflicts of interest

The authors declare they have no conflicts of interest.

Ethical statement

Ethical approval for the collection and sequencing of clinical isolates was granted by the Alfred Hospital Ethics Committee, Melbourne, Australia (Project numbers #550/12 and #526/13).

References

- Foster-Nyarko E, Holt KE, Cottingham H, Wick R, Judd LM, et al. Nanopore-only assemblies for genomic surveillance of the global priority drug-resistant pathogen, *Klebsiella pneumoniae*. *Figshare* 2023.
- Meatherall BL, Gregson D, Ross T, Pitout JDD, Laupland KB. Incidence, risk factors, and outcomes of *Klebsiella pneumoniae* bacteraemia. *Am J Med* 2009;122:866–873.
- Anderson DJ, Moehring RW, Sloane R, Schmader KE, Weber DJ, et al. Bloodstream infections in community hospitals in the 21st century: a multicenter cohort study. *PLoS One* 2014;9:e91713.
- Vading M, Naucleur P, Kalin M, Giske CG. Invasive infection caused by *Klebsiella pneumoniae* is a disease affecting patients with high comorbidity and associated with high long-term mortality. *PLoS One* 2018;13:e0195258.
- Collaborators AR. Global burden of bacterial antimicrobial resistance in 2019: a systematic analysis. *Lancet* 2022;399:629–655.
- Jung Y, Lee MJ, Sin H-Y, Kim N-H, Hwang J-H, et al. Differences in characteristics between healthcare-associated and community-acquired infection in community-onset *Klebsiella pneumoniae* bloodstream infection in Korea. *BMC Infect Dis* 2012;12:239.
- Giske CG, Monnet DL, Cars O, Carmeli Y, ReAct-Action on Antibiotic Resistance. Clinical and economic impact of common multidrug-resistant gram-negative bacilli. *Antimicrob Agents Chemother* 2008;52:813–821.
- World Health Organisation. *Prioritization of Pathogens to Guide Discovery, Research and Development of New Antibiotics for Drug-Resistant Bacterial Infections, Including Tuberculosis*. 2017.
- Navon-Venezia S, Kondratyeva K, Carattoli A. *Klebsiella pneumoniae*: a major worldwide source and shuttle for antibiotic resistance. *FEMS Microbiol Rev* 2017;41:252–275.
- Temkin E, Fallach N, Almagor J, Gladstone BP, Tacconelli E, et al. Estimating the number of infections caused by antibiotic-resistant *Escherichia coli* and *Klebsiella pneumoniae* in 2014: a modelling study. *Lancet Glob Health* 2018;6:e969–e979.
- European Centre for Disease Prevention and Control, Antimicrobial resistance (EARS-Net). *ECDC. Annual epidemiological report for 2014*. 2018.
- World Health Organisation. *Antimicrobial resistance. Draft global action plan on antimicrobial resistance*. Geneva; 2015. <https://apps.who.int/iris/handle/10665/193736> [accessed 10 November 2021].
- Cassini A, Högberg LD, Plachouras D, Quattrocchi A, Hoxha A, et al. Attributable deaths and disability-adjusted life-years caused by infections with antibiotic-resistant bacteria in the EU and the European Economic Area in 2015: a population-level modelling analysis. *Lancet Infect Dis* 2019;19:56–66.
- Nagaraj G, Shamanna V, Govindan V, Rose S, Sravani D, et al. High-resolution genomic profiling of carbapenem-resistant *Klebsiella pneumoniae* isolates: a multicentric retrospective Indian study. *Clin Infect Dis* 2021;73:S300–S307.

15. Saavedra SY, Bernal JF, Montilla-Escudero E, Arévalo SA, Prada DA, et al. Complexity of genomic epidemiology of carbapenem-resistant *Klebsiella pneumoniae* isolates in Colombia urges the reinforcement of whole genome sequencing-based surveillance programs. *Clin Infect Dis* 2021;73:S290–S299.
16. Aanensen DM, Carlos CC, Donado-Godoy P, Okeke IN, Ravikumar KL, et al. Implementing whole-genome sequencing for ongoing surveillance of antimicrobial resistance: exemplifying insights into *Klebsiella pneumoniae*. *Clin Infect Dis* 2021;73:S255–S257.
17. Lam MMC et al. Kaptive 2.0: updated capsule and lipopolysaccharide locus typing for the. *Microb Genom* 2022;8.
18. 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.
19. de Sousa JAM, Buffet A, Haudiquet M, Rocha EPC, Rendueles O. Modular prophage interactions driven by capsule serotype select for capsule loss under phage predation. *ISME J* 2020;14:2980–2996.
20. Feldman MF et al. A promising bioconjugate vaccine against hyper-virulent. *Proc Natl Acad Sci* 2019;116:18655–18663.
21. Ravinder M, Liao K-S, Cheng Y-Y, Pawar S, Lin T-L, et al. A synthetic carbohydrate-protein conjugate vaccine candidate against *Klebsiella pneumoniae* serotype K2. *J Org Chem* 2020;85:15964–15997.
22. Campbell WN, Hendrix E, Cryz S Jr, Cross AS. Immunogenicity of a 24-valent *Klebsiella* capsular polysaccharide vaccine and an eight-valent *Pseudomonas* O-polysaccharide conjugate vaccine administered to victims of acute trauma. *Clin Infect Dis* 1996;23:179–181.
23. Wyres KL, Nguyen TNT, Lam MMC, Judd LM, van Vinh Chau N, et al. Genomic surveillance for hypervirulence and multi-drug resistance in invasive *Klebsiella pneumoniae* from South and Southeast Asia. *Genome Med* 2020;12:11.
24. Ashton PM, Nair S, Dallman T, Rubino S, Rabsch W, et al. MinION nanopore sequencing identifies the position and structure of a bacterial antibiotic resistance island. *Nat Biotechnol* 2015;33:296–300.
25. Loman NJ, Pallen MJ. Twenty years of bacterial genome sequencing. *Nat Rev Microbiol* 2015;13:787–794.
26. Arredondo-Alonso S, Willems RJ, van Schaik W, Schürch AC. On the (im)possibility of reconstructing plasmids from whole-genome short-read sequencing data. *Microb Genom* 2017;3:e000128.
27. Wick RR, Judd LM, Gorrie CL, Holt KE. Completing bacterial genome assemblies with multiplex MinION sequencing. *Microb Genom* 2017;3:e000132.
28. De Maio N, Shaw LP, Hubbard A, George S, Sanderson ND, et al. Comparison of long-read sequencing technologies in the hybrid assembly of complex bacterial genomes. *Microb Genom* 2019;5:e000294.
29. Ben Khedher M, Ghedira K, Rolain J-M, Ruimy R, Croce O. Application and challenge of 3rd generation sequencing for clinical bacterial studies. *Int J Mol Sci* 2022;23:1395.
30. Murgia M. Pandemic puts oxford nanopore “on the map.” *Financial Times* 2021 2022.
31. Oxford Nanopore Technologies. Medaka; 2022 Jun 8. <https://github.com/nanoporetech/medaka>
32. Lee JY, Kong M, Oh J, Lim J, Chung SH, et al. Comparative evaluation of Nanopore polishing tools for microbial genome assembly and polishing strategies for downstream analysis. *Sci Rep* 2021;11:20740.
33. Wan YK, Hendra C, Pratanwanich PN, Göke J. Beyond sequencing: machine learning algorithms extract biology hidden in Nanopore signal data. *Trends Genet* 2022;38:246–257.
34. Vasiljevic N, Lim M, Humble E, Seah A, Kratzer A, et al. Developmental validation of Oxford Nanopore Technology MinION sequence data and the NGSpeciesID bioinformatic pipeline for forensic genetic species identification. *Forensic Sci Int Genet* 2021;53:102493.
35. Chen Z, Erickson DL, Meng J. Polishing the Oxford Nanopore long-read assemblies of bacterial pathogens with Illumina short reads to improve genomic analyses. *Genomics* 2021;113:1366–1377.
36. Delahaye C, Nicolas J. Sequencing DNA with nanopores: troubles and biases. *PLoS One* 2021;16:e0257521.
37. Sanderson N, Kapel N, Rodger G, Webster H, Lipworth S, et al. Comparison of R9.4.1/Kit10 and R10/Kit12 Oxford Nanopore flowcells and chemistries in bacterial genome reconstruction. *Genomics* 2022. DOI: 10.1101/2022.04.29.490057.
38. Wick RR, Judd LM, Holt KE. Performance of neural network basecalling tools for Oxford Nanopore sequencing. *Genome Biol* 2019;20:129.
39. Steinig E, Duchêne S, Aglua I, Greenhill A, Ford R, et al. Phylodynamic inference of bacterial outbreak parameters using Nanopore sequencing. *Mol Biol Evol* 2022;39:msac040.
40. Khezri A, Avershina E, Ahmad R. Hybrid assembly provides improved resolution of plasmids, antimicrobial resistance genes, and virulence factors in *Escherichia coli* and *Klebsiella pneumoniae* clinical isolates. *Microorganisms* 2021;9:12.
41. Boostrom I, Portal EAR, Spiller OB, Walsh TR, Sands K. Comparing long-read assemblers to explore the potential of a sustainable low-cost, low-infrastructure approach to sequence antimicrobial resistant bacteria with oxford Nanopore sequencing. *Front Microbiol* 2022;13:796465.
42. Craddock HA, Motro Y, Zilberman B, Khalfin B, Bardenstein S, et al. Long-read sequencing and hybrid assembly for genomic analysis of clinical *Brucella melitensis* isolates. *Microorganisms* 2022;10:619.
43. Gorrie CL, Mirčeta M, Wick RR, Judd LM, Lam MMC, et al. Genomic dissection of *Klebsiella pneumoniae* infections in hospital patients reveals insights into an opportunistic pathogen. *Nat Commun* 2022;13:3017.
44. Hunt M, Mather AE, Sánchez-Busó L, Page AJ, Parkhill J, et al. ARIBA: rapid antimicrobial resistance genotyping directly from sequencing reads. *Microb Genom* 2017;3:e000131.
45. Inouye M, Dashnow H, Raven L-A, Schultz MB, Pope BJ, et al. SRST2: Rapid genomic surveillance for public health and hospital microbiology labs. *Genome Med* 2014;6:90.
46. Argimón S, David S, Underwood A, Abrudan M, Wheeler NE, et al. Rapid genomic characterization and global surveillance of *Klebsiella* using pathogenwatch. *Clin Infect Dis* 2021;73:S325–S335.
47. Lam MMC, Wick RR, Watts SC, Cerdeira LT, Wyres KL, et al. A genomic surveillance framework and genotyping tool for *Klebsiella pneumoniae* and its related species complex. *Nat Commun* 2021;12:4188.
48. Gorrie CL, Mirceta M, Wick RR, Judd LM, Wyres KL, et al. Antimicrobial-resistant *Klebsiella pneumoniae* carriage and infection in specialized geriatric care wards linked to acquisition in the referring hospital. *Clin Infect Dis* 2018;67:161–170.
49. Hawkey J, Wyres KL, Judd LM, Harshegyi T, Blakeway L, et al. ESBL plasmids in *Klebsiella pneumoniae*: diversity, transmission, and contribution to infection burden in the hospital setting. *Genetic and Genomic Medicine* 2021.
50. Wyres KL, Hawkey J, Mirčeta M, Judd LM, Wick RR, et al. Genomic surveillance of antimicrobial resistant bacterial colonisation and infection in intensive care patients. *BMC Infect Dis* 2021;21:683.
51. Oxford Nanopore Technologies. Guppy v4.0.14; 2022. <https://github.com/nanoporetech/pyguppyclient> [accessed 8 June 2022].
52. Oxford Nanopore Technologies. Guppy v5.0.7 release note (21st May 2021); 2021 [accessed 1 June 2021].
53. Oxford Nanopore Technologies. qcat v1.1.0; 2022. <https://github.com/nanoporetech/qcat> [accessed 10 March 2022].
54. Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* 2018;34:i884–i890.
55. 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.

56. 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.
57. 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.
58. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs. *Nat Biotechnol* 2019;37:540–546.
59. Wick RR, Holt KE. Benchmarking of long-read assemblers for prokaryote whole genome sequencing. *F1000Res* 2019;8:2138.
60. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One* 2014;9:e112963.
61. Zhang P, Jiang D, Wang Y, Yao X, Luo Y, et al. Comparison of de novo assembly strategies for bacterial genomes. *Int J Mol Sci* 2021;22:14.
62. Wick RR, Judd LM, Holt KE. Assembling the perfect bacterial genome using oxford nanopore and illumina sequencing. *In SciELO Preprints* 2022. <https://preprints.scielo.org/index.php/scielo/preprint/view/5053/version/5357>
63. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* 2013;29:1072–1075.
64. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res* 2015;25:1043–1055.
65. Gorrie CL, Mirceta M, Wick RR, Edwards DJ, Thomson NR, et al. Gastrointestinal carriage is a major reservoir of *Klebsiella pneumoniae* infection in intensive care patients. *Clin Infect Dis* 2017;65:208–215.
66. Sherry NL, Lane CR, Kwong JC, Schultz M, Sait M, et al. Genomics for molecular epidemiology and detecting transmission of carbapenemase-producing *Enterobacteriales* in Victoria, Australia, 2012 to 2016. *J Clin Microbiol* 2019;57:e00573-19.
67. David S, Reuter S, Harris SR, Glasner C, Feltwell T, et al. Epidemic of carbapenem-resistant *Klebsiella pneumoniae* in Europe is driven by nosocomial spread. *Nat Microbiol* 2019;4:1919–1929.
68. Holt KE. RedDog; 2022. <https://github.com/katholt/RedDog> [accessed 10 March 2022].
69. Seemann T. snp-dist; 2022. <https://github.com/tseemann/snp-dists> [accessed 12 April 2022].
70. Robinson DF, Foulds LR. Comparison of phylogenetic trees. *Mathematical Biosciences* 1981;53:131–147.
71. Revell LJ. phytools: an R package for phylogenetic comparative biology (and other things). *Methods Ecol Evol* 2012;3:217–223.
72. Paradis E, Schliep K. ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. *Bioinformatics* 2019;35:526–528.
73. Wilke CO. cowplot: streamlined plot theme and plot annotations for “ggplot2”; 2022. <https://wilkelab.org/cowplot> [accessed 10 January 2022].
74. Dowle M. data.table: Extension of “data.frame”; 2022. <https://rdatatable.gitlab.io/data.table> [accessed 3 April 2022].
75. Wright E. Using DECIPHER v2.0 to analyze big biological sequence data in R. *The R Journal* 2016;8:352.
76. Galili T. dendextend: an R package for visualizing, adjusting and comparing trees of hierarchical clustering. *Bioinformatics* 2015;31:3718–3720.
77. Wickham H. dplyr: a grammar of data manipulation; 2022. <https://dplyr.tidyverse.org>, <https://github.com/tidyverse/dplyr> [accessed 7 June 2022].
78. Yu G. Using ggtree to visualize data on tree-like structures. *Curr Protoc Bioinformatics* 2020;69:e96.
79. Kassambara A. ggpubr: ‘ggplot2’ based publication ready plots; 2022. <https://rpkgs.datanovia.com/ggpubr> [accessed 10 February 2022].
80. Auguie B. gridExtra: miscellaneous functions for “Grid” graphics; 2022. <https://cran.r-project.org/web/packages/gridExtra/gridExtra.pdf> [accessed 6 June 2022].
81. Firke S. janitor: simple tools for examining and cleaning dirty data; 2022. <https://cran.r-project.org/web/packages/janitor/index.html> [accessed 3 June 2022].
82. Zhu H. kableExtra: construct complex table with “kable” and pipe syntax; 2022. <http://haozhu233.github.io/kableExtra/>, <https://github.com/haozhu233/kableExtra> [accessed 8 June 2022].
83. Wickham H, Grolemund G. Tidyverse: R packages for data science; 2022. <https://www.tidyverse.org/> [accessed 8 June 2022].
84. Wang L-G, Lam T-Y, Xu S, Dai Z, Zhou L, et al. Treeio: an R package for phylogenetic tree input and output with richly annotated and associated data. *Mol Biol Evol* 2020;37:599–603.
85. Tennekes M, Ellis P. treemap: treemap visualization; 2022. <https://cran.r-project.org/web/packages/treemap/treemap.pdf> [accessed 4 January 2022].
86. Wickham H, Grolemund G. stringr: Simple, consistent wrappers for common string operations; 2022. <http://stringr.tidyverse.org>, <https://github.com/tidyverse/stringr> [accessed 5 June 2022].
87. Gorrie CL, Mirceta M, Wick RR, Judd LM, Lam MMC, et al. Genomic dissection of *Klebsiella pneumoniae* infections in hospital patients reveals insights into an opportunistic pathogen. *Nat Commun* 2022;13:3017.
88. Follador R, Heinz E, Wyres KL, Ellington MJ, Kowarik M, et al. The diversity of *Klebsiella pneumoniae* surface polysaccharides. *Microb Genom* 2016;2:e000073.
89. Lam MMC, Wick RR, Wyres KL, Gorrie CL, Judd LM, et al. Genetic diversity, mobilisation and spread of the yersiniabactin-encoding mobile element ICEKp in *Klebsiella pneumoniae* populations. *Microb Genom* 2018;4:e000196.
90. Lam MMC, Wyres KL, Judd LM, Wick RR, Jenney A, et al. Tracking key virulence loci encoding aerobactin and salmochelin siderophore synthesis in *Klebsiella pneumoniae*. *Genome Med* 2018;10:77.
91. Wick RR, Judd LM, Cerdeira LT, Hawkey J, Méric G, et al. Trycycler: consensus long-read assemblies for bacterial genomes. *Genome Biol* 2021;22:266.
92. Sereika M, Kirkegaard RH, Karst SM, Michaelsen TY, Sørensen EA, et al. Oxford Nanopore R10.4 long-read sequencing enables near-perfect bacterial genomes from pure cultures and metagenomes without short-read or reference polishing. *Microbiology* 2021. DOI: 10.1101/2021.10.27.466057.