

# A multicenter study on accuracy and reproducibility of nanopore sequencing-based genotyping of bacterial pathogens

Johanna Dabernig-Heinz,<sup>1</sup> Mara Lohde,<sup>2</sup> Martin Hölzer,<sup>3</sup> Adriana Cabal,<sup>4</sup> Rick Conzemius,<sup>5</sup> Christian Brandt,<sup>2</sup> Matthias Kohl,<sup>6</sup> Sven Halbedel,<sup>7,8</sup> Patrick Hyden,<sup>4</sup> Martin A. Fischer,<sup>9</sup> Ariane Pietzka,<sup>10</sup> Beatriz Daza,<sup>4</sup> Evgeny A. Idelevich,<sup>11</sup> Anna Stöger,<sup>4</sup> Karsten Becker,<sup>11</sup> Stephan Fuchs,<sup>3</sup> Werner Ruppitsch,<sup>4</sup> Ivo Steinmetz,<sup>1</sup> Christian Kohler,<sup>11</sup> Gabriel E. Wagner<sup>1</sup>

**AUTHOR AFFILIATIONS** See affiliation list on p. 14.

**ABSTRACT** Nanopore sequencing has shown the potential to democratize genomic pathogen surveillance due to its ease of use and low entry cost. However, recent genotyping studies showed discrepant results compared to gold-standard short-read sequencing. Furthermore, although essential for widespread application, the reproducibility of nanopore-only genotyping remains largely unresolved. In our multicenter performance study involving five laboratories, four public health-relevant bacterial species were sequenced with the latest R10.4.1 flow cells and V14 chemistry. Core genome MLST analysis of over 500 data sets revealed highly strain-specific typing errors in all species in each laboratory. Investigation of the methylation-related errors revealed consistent DNA motifs at error-prone sites across participants at read level. Depending on the frequency of incorrect target reads, this either leads to correct or incorrect typing, whereby only minimal frequency deviations can randomly determine the final result. PCR preamplification, recent basecalling model updates and an optimized polishing strategy notably diminished the non-reproducible typing. Our study highlights the potential for new errors to appear with each newly sequenced strain and lays the foundation for computational approaches to reduce such typing errors. In conclusion, our multicenter study shows the necessity for a new validation concept for nanopore sequencing-based, standardized bacterial typing, where single nucleotide accuracy is critical.

**KEYWORDS** nanopore sequencing, multicenter performance study, bacterial typing, genomic surveillance, cgMLST, molecular surveillance

Genomic surveillance of microbial pathogens is crucial for modern infection control and hence for our health systems addressing the emergence (1), spread, and transmission of (new) pathogens (2, 3), drug-resistant strains (4), and vaccine-evading variants (5). As such, it plays a key role in data-driven decision-making and the implementation of countermeasures across clinical, animal health, and food safety sectors (4, 6–9). Ideally, surveillance involves global, long-term monitoring for population changes in circulating pathogens, coupled with local, (real-time) analysis to handle outbreaks swiftly, trace infection sources, and expedite patient management (6, 7). Whole-genome sequencing of pathogens by short-read (SR) next-generation sequencing has revolutionized the field (10, 11). This approach not only enabled investigations of complex relationships between different isolates at an unprecedented level (3, 6, 10), but it also tackled issues regarding standardization and reproducibility of existing methods (6). Additionally, the reconstructed genomes enable a reference-free, in-depth analysis of the genetic profiles of individual isolates, including factors associated with resistance, virulence, and pathogenicity (4, 6, 12).

Despite these advantages and their immense importance from a One Health viewpoint, widespread implementation and regular application of genomic surveillance

**Editor** Daniel D. Rhoads, Cleveland Clinic, Cleveland, Ohio, USA

Address correspondence to Ivo Steinmetz, ivo.steinmetz@medunigraz.at, Christian Kohler, christian.kohler@med.uni-greifswald.de, or Gabriel E. Wagner, gabriel.wagner-lichtenegger@medunigraz.at.

Johanna Dabernig-Heinz, Mara Lohde, Martin Hölzer, Adriana Cabal, Rick Conzemius, Christian Brandt, Stephan Fuchs, Werner Ruppitsch, Christian Kohler, and Gabriel E. Wagner contributed equally to this article. Author order was determined in order of increasing seniority per affiliation.

R.C. was an employee of the company Ares Genetics. This does not affect the authors' adherence to all the journal's policies on sharing data and materials. Twenty flow cells were provided free of charge by Oxford Nanopore Technologies. However, the manufacturer did not participate in the study's design, data collection, interpretation, or any other aspects of the research.

See the funding table on p. 15.

**Received** 26 April 2024

**Accepted** 25 July 2024

**Published** 19 August 2024

Copyright © 2024 Dabernig-Heinz et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

via SR sequencing face hurdles such as high equipment cost, required laboratory space, low cost-efficiency with low sample numbers, complex bioinformatics workflows, and the need for trained personnel (8, 12).

The severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) pandemic, a prime example of the positive impact of extensive genomic surveillance and variant analysis (9), has shown how these limitations can be overcome through the use of straightforward protocols, the implementation of nanopore sequencing and automated, user-friendly analysis pipelines (13–16). Nanopore sequencing addressed the sequencing capacity bottleneck with its affordability, ease of implementation, and suitability for decentralized applications (3, 13, 14). Given the positive user experiences, the competitiveness with SR methods in SARS-CoV-2 surveillance (17) and the above-mentioned advantages, especially in combination with the real-time data availability (18) and selective sequencing (19, 20), its expansion into additional domains of routine surveillance is a logical next step for many.

This is enhanced by the fact that the increased error rates previously associated with nanopore sequencing compared to SR methods have been significantly reduced by recent advances in DNA library preparation chemistry (Q20+), flow cells (R10), and bioinformatics (basecalling models), resulting in a read accuracy above Q20 (21). In public health and clinical microbiology, routine nanopore sequencing holds promise in near-patient/real-time metagenomics for culture-free pathogen identification (22), plasmid characterization (20), and AMR marker screening (23). In the case of bacterial WGS, current studies show a very positive development of nanopore sequencing. Sereika et al. showed that recent nanopore advancements allow assemblies of a quality similar to those of Illumina data but with significantly less fragmentation. Lermineaux et al. evaluated R10.4.1 flow cells and the V14 chemistry and showed that they yield high-quality and contiguous assemblies for Gram-negative bacteria, while finer-scale analyses on single nucleotide levels would still benefit from SR data (24). Another recent study with 14 different bacterial species has even shown that nanopore data using deep learning-based variant callers match or exceed Illumina data and are therefore excellently suited for variant calling approaches (25). However, during our study, there were also emerging reports of incorrect base calls attributed to methylation when sequencing native DNA from bacterial microorganisms (26, 27).

Classical genetic profiling/genetic typing seems particularly plausible as a near-term application of nanopore sequencing with considerable potential for the surveillance of bacterial pathogens (18). For in-depth genetic profiling and high-resolution typing of bacterial pathogens, e.g., for outbreak analysis, gene-by-gene approaches like core genome multilocus sequence typing (cgMLST) have become widely established offering high resolution in combination with standardization, ease of data exchange, and minimal hardware and bioinformatics knowledge requirements (28–31). Due to these properties, they have been extensively used in the field of genomic surveillance, especially in cross-laboratory/cross-national activities, e.g., the multicentric initiatives “PulseNet” and “Pathogenwatch.” In contrast to SNP analyses, alleles of a defined set of genes (e.g., the core genome) are called and the isolates are then characterized or compared on the basis of this genetic profile (28). Despite the continuously increasing raw read accuracy of nanopore sequencing and our promising results in the case of *Bordetella pertussis* typing (18), the reported results of other bacterial typing studies were discrepant (26, 32). Linde et al. reported good performance in plasmid assembly (32), an important advantage for analyzing AMR transmission (33, 34). They also report substantial problems for high-resolution genotyping in *Brucella suis* when comparing Illumina data to Oxford Nanopore Technologies (ONT) data from flow cell versions R9 and R10.4 (32). Similarly, Lohde et al. report unexpected sequencing errors of problematic magnitude that obscure typing results in *Klebsiella pneumoniae* and other species, which they attribute to nearby methylation sites causing systematic basecalling errors (26). Such errors can be highly problematic for allele-calling approaches, as by principle it is irrelevant whether one or one hundred bases change within a gene, both lead to the same result, namely to a

different allele designation (28). Even a small number of erroneous SNPs when spread across random genes in the genome can accumulate incorrect allele calls and ultimately hinder analyses that demand high resolution, like outbreak analysis or source tracing investigations (26, 32).

These inconsistent results in the typing performance, particularly with regard to different species (18, 26, 32), underscore the necessity for systematic validation of nanopore sequencing in the domain of bacterial typing. Moreover, a crucial question, that largely lacks a definitive answer, pertains to the accuracy against gold-standard SR and especially the reproducibility of nanopore sequencing-based bacterial typing across diverse laboratories. This assessment is essential prior to implementing and establishing this innovative method in delicate domains like genomic surveillance and diagnostics.

A multicenter approach makes it possible to investigate whether the performance of the method and the reproducibility of the results are also guaranteed across various settings and thus enables a realistic assessment of the effectiveness and practical utility under real-world conditions.

In this study, we present the results of an international, multi-laboratory performance assessment aiming to thoroughly validate the use of nanopore sequencing for comprehensive genomic surveillance of bacterial pathogens in terms of typing accuracy and comparability. Using four healthcare-relevant species with 18–20 isolates each, we investigated the performance of nanopore R10.4.1 flow cells and V14 sequencing chemistry for routine high-resolution typing by cgMLST using data sets from five participants, including recognized public health institutions and established clinical microbiology laboratories.

## MATERIALS AND METHODS

### Strains and DNA isolation

Strains from four bacterial species comprising 19 *Enterococcus faecium* (EF), 20 *K. pneumoniae* (KP), 20 *Listeria monocytogenes* (LM), and 18 *Staphylococcus aureus* (SA) isolates were selected for this performance test. Of note, three additional isolates, EF21, SA64, and SA70, were only used as controls and were not analyzed in this study. Starting from a single colony and subsequent propagation, each laboratory (LAB1-5) was in a blind-coded manner provided with identical pure cultures of each strain in stitch-agar. Details on the cultivation of the strains and DNA preparation of the individual participants can be found in Table S1.

### Nanopore library preparation

The library preparation in all labs was carried out using the manufacturer's ligation sequencing gDNA protocol using the Native Barcoding Kit 24 V14 SQK-NBD114.24 (Oxford Nanopore, UK). To investigate basecalling errors potentially resulting from DNA methylation 11 suspicious strains were resequenced in LAB1 according to the manufacturer's ligation sequencing V14 — PCR barcoding protocol using the ligation kit SQK-LSK114 (Oxford Nanopore, UK) and its PCR-barcoding expansion EXP-PBC001 (Oxford Nanopore, UK).

### Nanopore sequencing

The initial sequencing procedure was performed on R10.4.1 flow cells run at 260 bp/s (4 kHz) using MinKNOW in all labs. Per run 10–20 strains were sequenced and the data were basecalled in super-accurate (SUP) mode. Lab-specific details including the number of multiplexed samples, DNA concentrations, and software versions used are denoted in Table S1. The PCR library (LAB1) was sequenced on R10.4.1 flow cells run at 400 bp/s (5 kHz) and basecalled in SUP mode using MinKNOW.

## Assembly pipeline

The nanopore reads were assembled with Flye (35) and assemblies were polished with Medaka (ONT) and its respective models: “r1041\_e82\_260bps\_sup\_g632” for the native barcoding kit data and “r1041\_e82\_400bps\_sup\_variant\_v4.2.0” for the PCR barcoding data. Program versions are stated in Table S1.

## Performance study under improved sequencing conditions using a bacterial methylation basecalling model

For evaluation of recent sequencing advances made by the manufacturer during our study, LAB2 resequenced the same DNA preps again using the Native Barcoding kit as described above. Sequencing, however, was performed with the latest sequencing conditions using 400 bp/s and 5 kHz instead of the discontinued “accuracy” mode and its 260 bp/s and 4 kHz. Raw data were basecalled with Dorado version 0.4.0 utilizing ONT’s specifically trained research model for bacterial methylation (res\_dna\_r10.4.1\_e8.2\_400bps\_sup@2023-09-22\_bacterial-methylation). Subsequently reads were assembled with Flye (35), and optional Racon (36) polishing was included to assess its impact on typing performance. Afterward, the Flye-only assembly as well as the Racon polished assembly were subjected to one round of Medaka polishing either with the recommended consensus model (-m r1041\_e82\_400bps\_sup\_v4.2.0) or with the variant model (-m r1041\_e82\_400bps\_sup\_variant\_v4.2.0). cgMLST analysis was used to evaluate which of these four polishing variants delivered the best typing results. A subset of 12 strains (three for each species) was also resequenced in LAB1 and LAB3 to evaluate these updates in terms of reproducibility of typing results. Program versions are stated in Table S1.

## Short-read sequencing

Two independent short-read data sets were acquired for comparison with the long-read data sets. For the first short-read whole-genome sequencing data set, the same extracted DNA already used for long-read sequencing in Lab 2 was utilized for library preparation in the case of *S. aureus*, *K. pneumoniae*, and *E. faecium*. In the case of *L. monocytogenes*, an additional DNA preparation was used, but coming from a culture grown under the same culture conditions and from the same initial stock as described for the long-read DNA extraction. DNA libraries were generated using the Illumina DNA Prep Kit (Illumina, USA) combined with the Illumina DNA/RNA UD Indexes Set B and Tagmentation Kit (Illumina, USA). Five hundred nanograms of DNA were employed for the normalization process. Libraries were amplified using the Illumina DNA Prep PCR Kit (Illumina, USA) following the manufacturer’s protocol. All cleanup and size selection procedures were conducted using Illumina beads (Illumina, USA) according to the manufacturer’s guidelines. Library concentration was quantified using the Qubit 1X dsDNA Assay Kit (Thermo Fisher Scientific, USA). Sequencing was performed on the Illumina MiSeq DX with the MiSeq reagent Kit v3 2 × 300 Cycles (Illumina, USA) using 14 pM of prepared DNA library.

For the second independent short-read whole-genome sequencing data set, the starting material was the same culture as used for the DNA extraction for long-read sequencing in Lab 5. DNA was extracted with MagMAX Viral/Pathogen Ultra Nucleic Acid Isolation Kit (ThermoFisher) on a KingFisher Apex robot. DNA libraries were generated with Illumina DNA prep Kit using the Illumina DNA/RNA UD Indexes (Sets A to D) and Tagmentation Kit (Illumina, USA). Between 100 to 500 ng of input DNA were used. Library concentration was quantified using the Qubit 1X dsDNA Assay Kit (Thermo Fisher Scientific, USA), and the fragment length of the library pool was measured in Qsep100 (Nippon Genetics). Sequencing was performed on an Illumina NextSeq 2000 device with the NextSeq 1000/2000 P1 Reagent cartridge of 300 Cycles (Illumina, USA) using 750 pM of prepared DNA library.

Subsequently, short-read genome assembly for Illumina data was carried out using the SKESA assembler version 2.4.0 (37), integrated into the Ridom SeqSphere+ version 9.0.8 (Ridom, Germany) (38).

### Typing with respective cgMLST schemes

All assemblies were analyzed with the cgMLST scheme for the respective species (28, 29, 39–41) in SeqSphere+. Typing results of the respective Illumina data served as a reference for allele calls of nanopore data. The exact number of mismatched cgMLST loci was calculated by cross-comparing all strains in distance matrices. Minimum spanning trees (MSTs) were built using default parameters to assess and visualize relations between assemblies of different laboratories, both between long-read (LR) data sets and between LR and SR. The parameter “pairwise ignore missing values” was selected in SeqSphere for all analyses.

### Exploration of the methylation error

A subset of suspicious strains with a high number of typing errors compared to the Illumina data were selected to thoroughly analyze reads and assemblies. Potential patterns of incorrect basecalls in reads covering erroneous loci were investigated via mapping of reads to reference alleles using minimap version 2.26 and visual inspection. To investigate patterns adjacent to ambiguous positions, and for the generation of the sequence logos the MPOA pipeline version 1.4.1 was applied as described previously (26). Seqtk version 1.3 (<https://github.com/lh3/seqtk>) with seed 99 and seed 100 was used for the downsampling of the data.

### Statistical testing

The assembly and polishing pipelines were compared by the Quade test followed by pairwise Wilcoxon signed rank tests as post hoc tests. The *P*-values of the post hoc tests were adjusted by Holm’s method. In all cases, results with an (adj.) *P*-values < 0.05 were considered significant.

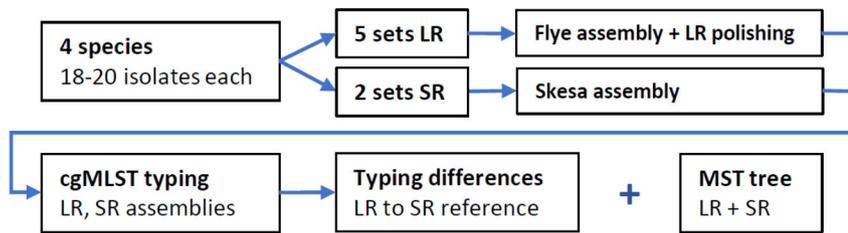
## RESULTS

### Design and implementation of our nanopore validation study

We selected the species *E. faecium*, *K. pneumoniae*, *L. monocytogenes*, and *S. aureus* for our extensive validation due to their significant importance in the medical, public health, and food safety sectors. Moreover, whole-genome sequencing and genomic surveillance by cgMLST analysis are well-established and extensively used for these species (28, 29, 39–41). Each participating institution generated long-read sequencing data based on strictly defined conditions: Native barcoding kit V14 library prep, R10.4.1 flow cells, and basecalling in super accurate mode (SUP) for the highest data accuracy. Data sets from isolates with a sequencing depth below the recommended 40× (21) were excluded from further analysis (see Table S2). In accordance with ONT guidelines, assemblies were generated using Flye and polished with Medaka (42, 43). In total, we generated and analyzed 526 assemblies in our performance study (Table S2). Assembly metrics can be found in Table S3.

### cgMLST-based genotyping of short-read reference data sets shows consistent results

The assemblies underwent high-resolution genotyping utilizing the respective cgMLST species schemes (28, 29, 39–41), followed by a comparative analysis. We specifically focused on a comparison of cgMLST results because they are pivotal in outbreak investigations, transmission analysis, and pinpointing infection sources. Two independent SR data sets facilitated a direct comparison with the established gold standard. The summarized methodological workflow is depicted in Fig. 1.



**FIG 1** The methodological workflow of the study consists of LR and SR sequencing of 77 isolates. The assemblies of both methods were compared using cgMLST schemes for the respective species to assess differences in the typing results for the respective strain. In addition, the results were visualized in the form of MSTs containing assemblies of both methods.

In accordance with previous studies (44), the two SR data sets yielded consistent cgMLST results. In an MST, the sequencing replicates of distinct strains matched perfectly, as exemplified by *L. monocytogenes* (Fig. 2a). The same applies to the other species, as illustrated in Fig. S1a through c.

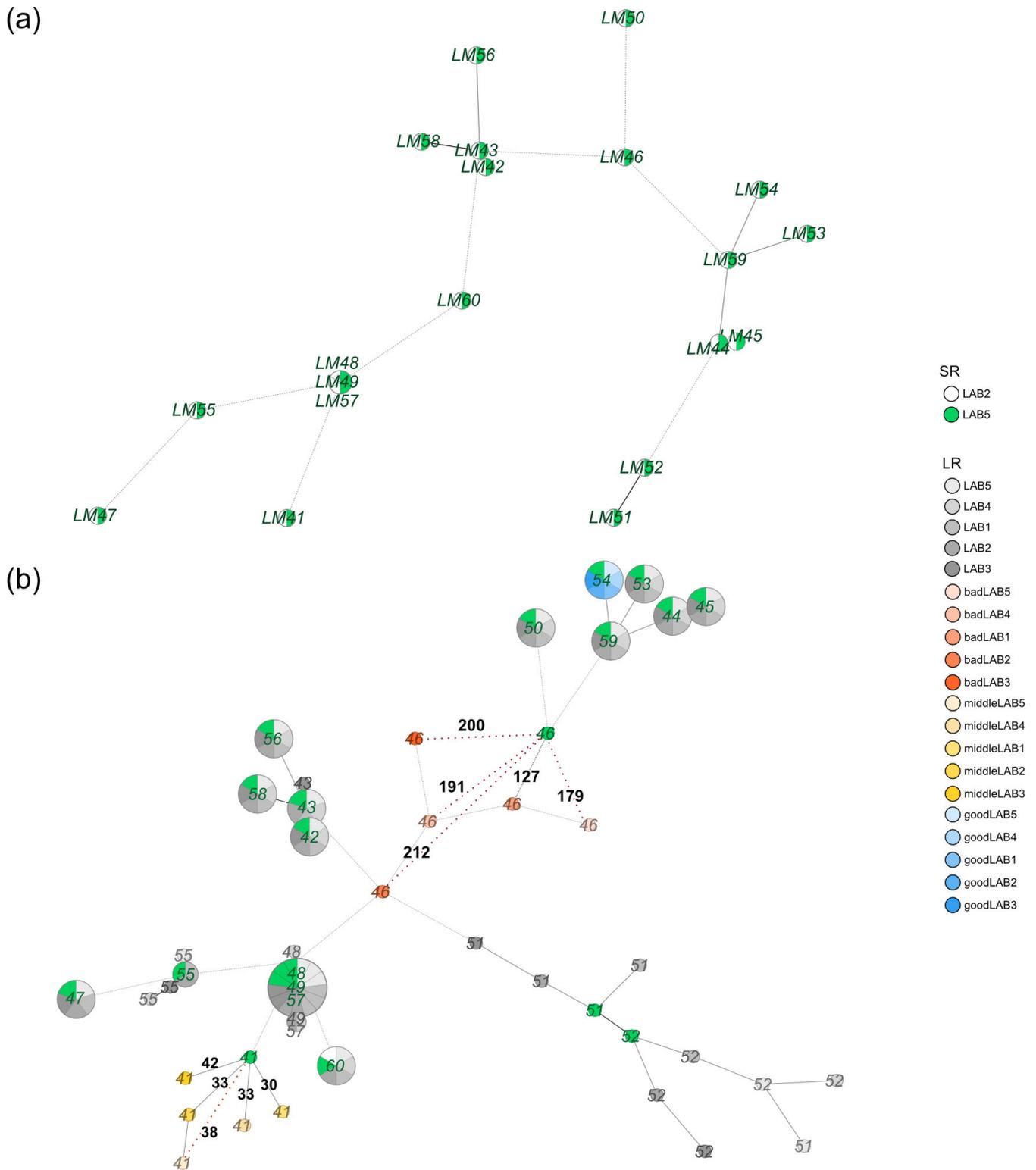
### Strain-specific inconsistencies revealed in cgMLST-based genotyping of long-read data sets, impeding typing correctness and reproducibility

The observations with LR data sets differed substantially from those with SR data. Although there were strains in which there were no/minimal differences in the typing results between the LR data of all participants and the SR reference, there were also strains where the typing of the LR data of different participants not only diverged from the reference but also from each other (Fig. 2b; Fig. S2 to S4). The magnitude of the observed typing errors varied depending on the strain (exemplary strains are shown in blue, yellow, and red shades in the figures). For each single isolate, the number of incorrect alleles compared to the reference was in a similar range for all participants, but the affected targets varied. This led to non-reproducible typing, clearly visible in the minimum spanning trees where it prevents the assemblies from coinciding in a single node. We could exclude the possibility that these observations arose from inappropriate extraction/handling of the DNA before its application to sequencing, as such issues would have similarly impacted the SR data and all isolates alike. Additionally, data from different labs still led to consistent LR typing results for several strains. Notably, LAB2's SR and LR data sets originated from identical DNA preparation for *S. aureus*, *K. pneumoniae*, and *E. faecium*. Yet, discrepancies in typing were observed between the two datasets, comparable to those observed in other participants. This clearly shows that the typing errors are related to LR sequencing.

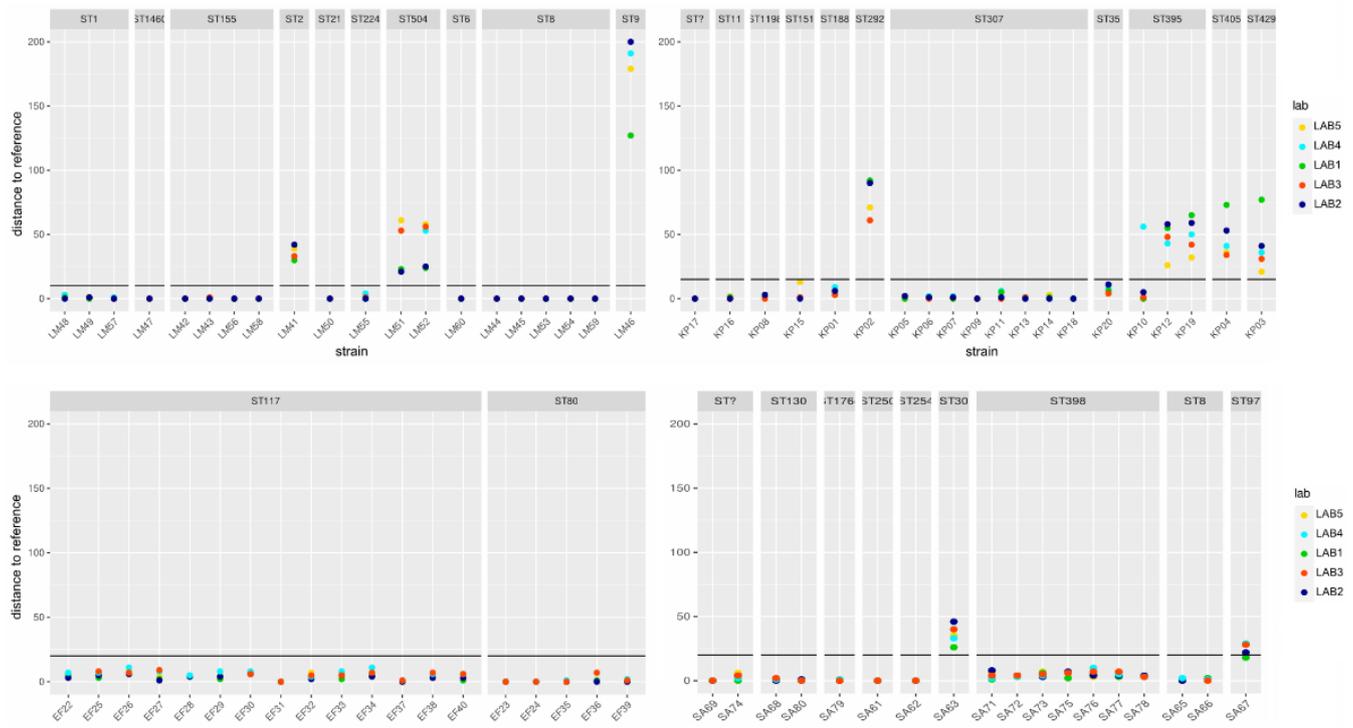
### Consistency in data quality from sequencing genetically similar strains

To assess the magnitude and variation of errors between participants for each isolate, we quantified the count of erroneous alleles across all assemblies in relation to their corresponding SR references. Based on the minimum spanning trees, we hypothesized that clusters of genetically similar isolates produced similar good/bad sequencing results and hence grouped the strains according to their MLST sequence type (ST). The analysis showed (Fig. 3) that there were differences between the species in (i) the number of strains affected and (ii) the extent of the observed error but also that (iii) the latter varied enormously between different isolates within a species.

Among all species, *L. monocytogenes*, e.g., had the highest number of isolates where the assemblies of all participants matched the reference perfectly but also included the strain with the highest number of differences to the short-read reference. It is also striking that strains with the same ST typically showed deviations from the reference of a similar magnitude between the participants. This indicates that the ST is indeed a good, first predictor of nanopore sequencing quality if strains with a certain ST have been previously sequenced and analyzed. In substantiating our hypothesis, we utilized



**FIG 2** (a) cgMLST-based MST of *L. monocytogenes* isolates using short read data. Sequencing replicates of identical strains exhibit consistent cgMLST profiles, leading to direct clustering irrespectively of the executing laboratory. (b) MST of the same *L. monocytogenes* isolates, augmented with LR assemblies from participating laboratories (in grayscale). Only one SR data set is shown (in green). Depending on the strain under investigation, LR assemblies of different participants showed inconsistent typing results. There are isolates where the typing of the LR matches that of the SR (an exemplary one in blue shades), but also others with differences not only to the SR but also between the LR assemblies of the participants. Furthermore, the magnitude of the observed differences varied between isolates (an exemplary one in yellow and one in red shades). For a clearer presentation, we only show the differences for selected strains; the differences at the isolate level are detailed in Fig. 3.



**FIG 3** Allelic differences between the assemblies of the different participants compared to the short-read reference of the respective strains, showing species and isolate-specific differences regarding the number of affected strains and the magnitude of the differences. Isolates with the same ST show a similar error range compared to the assemblies of the individual participants. The cluster threshold for the respective cgMLST is added as a black line. Typing errors in assemblies that are close to this threshold are highly problematic in, e.g., outbreak investigations.

our comprehensive collection of *Listeria* strains. We selected, sequenced, and typed three new isolates with previously good (Mismatches to the reference  $<3$ ; ST1, ST8, ST155) and three isolates with previously poorly performing STs (Mismatches to the reference  $\geq 3$ ; ST2, ST9, ST504). As predicted, the former showed no (ST8, ST155) or minimal difference to the reference (ST1—1 difference), but the others showed 31 (ST2), 211 (ST9), and 62 (ST504) allelic differences.

Based on the criterion of  $\geq 3$  cgMLST mismatches (MM) compared to the SR reference in the majority of participants, we identified seven highly problematic isolates for *K. pneumoniae* (KP01, 02, 03, 04, 12, 19, 20), 12 for *E. faecium* (EF22, 25, 26, 27, 28, 29, 30, 32, 33, 34, 38, 40), four for *L. monocytogenes* (LM41, 46, 51, 52), and 10 for *S. aureus* (SA63, 67, 71, 72, 73, 74, 75, 76, 77, 78)—summarized in Table S4. For 10 additional isolates, at least one assembly of one participant was above the threshold, increasing the overall number of strains with at least one erroneous assembly to 43, which is more than half of the isolates. The criterion of  $\geq 3$  MM is based on the consideration that in investigations/studies where two problematic strains are involved, the genetic distance between them is artificially increased by at least six mismatches (at least three incorrect alleles in each isolate). Such an increase approaches the cluster threshold of cgMLST schemes in magnitude.

### Read-level analysis reveals base ambiguities due to wrong basecalls in the case of DNA methylation as a major source of error

To understand how the errors in the polished consensus assemblies arise, we examined the data at the read level by performing a mapping and subsequent visual inspection of the problematic sites in an initial screening of randomly selected targets in all species. At the position of the incorrect base that was responsible for the erroneous allele call, the mapping showed a base ambiguity that is not present in the short reads. One of the two

ambiguous bases was predominantly found in forward mapping reads, while the other was primarily present in reverse mapping reads (Fig. S5a). This suggests errors associated with strand-specific methylation and is consistent with reports from other groups that emerged during our study (26, 27).

Interestingly, minimal differences in the data frequency of incorrect target reads between participants could influence the final typing result, even though the error position was conserved at the read level (Fig. S5). As an example given for the respective *E. faecium* reads of isolate EF26 at target “EF01658,” base frequencies—focusing only on the two prominent, ambiguous bases C and T that show a strand-bias—differed among participants: LAB1 (48% C—50%T), LAB2 (47% C—53%T), LAB3 (47% C—52%T), LAB4 (62% C—38%T), LAB5 (53% C—44%T). This influences which base appears at this position in an assembly’s sequence and thus finally determines the (non-reproducible) allele call.

Conserved sequence motifs also emerged when analyzing the surrounding positions around ambiguous bases using the MPOA pipeline (26). For this purpose, the reads of all participants of the most problematic strains, based on the cgMLST distance to the SR reference, were mapped to the respective assembly, ambiguous sites were identified and sequence logos were generated for the areas around ambiguous bases. As can be seen in Fig. 4, the patterns of the affected sites in the respective species were reproducible, again indicating a common cause of error.

The minimal differences in read counts between correct and incorrect bases suggest that random subsampling might slightly shift these ratios, affecting typing results. This was confirmed by subsampling into two datasets, showing that assemblies of well-performing strains remained consistent while problematic strains exhibited errors and typing discrepancies, highlighting data fragility (see Supplementary material for details).

Of note, we also investigated how the typing results are affected by the re-sequencing of the same DNA in the same lab. As expected due to the subsampling results, this did not have any significant influence on the typing results of good strains. In the case of problematic strains, however, also re-sequencing led to deviating and incorrect results (see Supplementary material and Tables S5 and S6 for details).

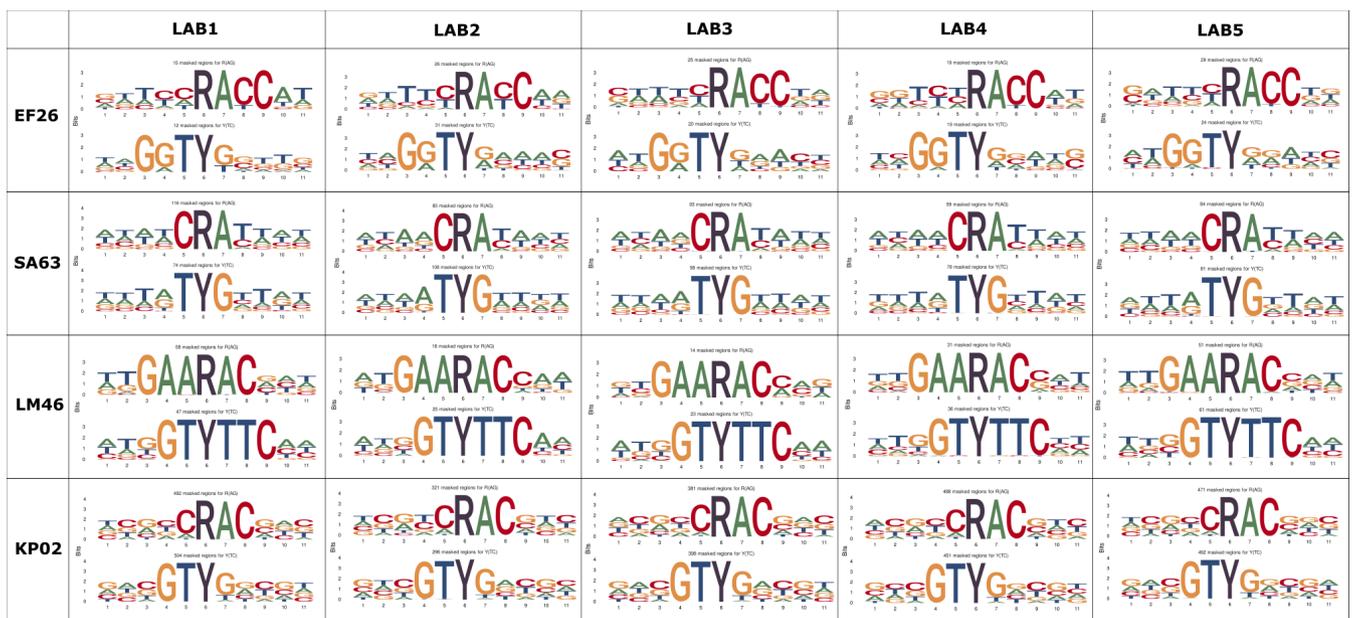


FIG 4 Sequence logos based on ambiguous positions and surrounding bases in a genome reveal conserved sequence patterns and high strain-level agreement between participants. Ambiguous sites identified were purine (R; A or G) or pyrimidine (Y; C or T) discrepancies.

## PCR whole-genome amplification validates methylation as a predominant issue and resolves typing issues

As the observed errors strongly suggested issues arising from DNA methylation, we performed whole-genome amplification via PCR on some of the most problematic strains to remove DNA methylation. In line with our hypothesis, this approach significantly mitigated the observed issues in typing (Table 1). As previously observed (26, 27), however, in *K. pneumoniae*, it resulted in substantial assembly fragmentation (44–110 contigs), and consequently, a problematic increase in the number of missing targets in typing. Of 92 mismatching loci in the original KP02 assembly, 54 were correct in the PCR assembly and the remaining 38 were not found due to missing targets. Similarly, the mismatches of KP03 and KP12 were corrected in the PCR assemblies, with only one of those loci missing in the KP12 PCR assembly. In *E. faecium*, differences compared to the reference persisted in the low single digits, indicating additional sequencing errors not resolvable by PCR. However, as shown below, these errors could be corrected by updates in the basecalling model (Table S7), which is why we did not examine it in detail.

## Sequencing updates improve typing performance, but challenges with problematic strains persist

Throughout our investigation, the manufacturer implemented various enhancements, including a more precise data sampling process at an increased sequencing speed (5 kHz, 400 bp/s instead of 4 kHz, 260 bp/s) and the introduction of a dedicated research basecalling model tailored for native bacterial DNA and methylation (`res_dna_r10.4.1_e8.2_400bps_sup@2023-09-22_bacterial-methylation`). While a complete rerun of the performance test wasn't feasible, we aimed to evaluate these improvements by resequencing all isolates by one participant. Additionally, we explored different polishing strategies to identify the optimal tool combination for this updated sequencing data. Surprisingly, despite contrary official recommendations for bacterial assemblies (42, 43), the combination of the polishing tools Racon followed Medaka using its variant (!) polishing model (`r1041_e82_400bps_sup_variant_v4.2.0`) yielded the most accurate cgMLST typing results of the assemblies in our study. This tool combination performed significantly better according to the Quade test followed by pairwise Wilcoxon signed rank tests as post hoc test (Holm adj. *P*-value < 0.05) than pipelines with the non-variant Medaka mode or Flye-only assemblies without polishing (adj. *P*-values < 0.001) and had better but not significantly different results than Medaka variant without Racon (adj. *P*-value = 0.124) (Table S7). By applying the best pipeline only 4 of the 77

**TABLE 1** Mismatches of problematic isolates compared to the short-read reference with and without PCR whole-genome-amplification reveals methylation as the most common source of sequencing errors affecting typing results<sup>a</sup>

	Native DNA Mismatches to reference	PCR amplified DNA Mismatches to reference
EF26	6	1
EF30	7	2
EF34	5	3
KP02	92	0 (994 missing)
KP03	77	1 (128 missing)
KP12	55	0 (89 missing)
LM41	30	0
LM46	127	0
LM51	23	0
SA67	18	0
SA73	6	0

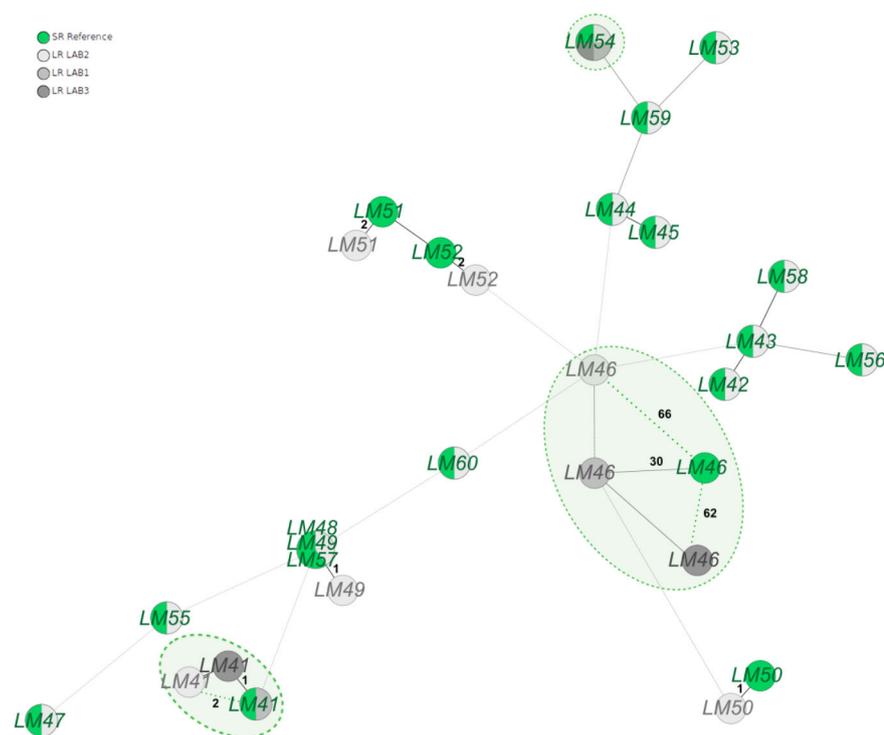
<sup>a</sup>Pronounced fragmentation of the PCR assembly and an elevated number of missing cgMLST targets was evident in *K. pneumoniae* (KP). Still, most of the mismatching targets from the original assembly (LAB1) were found and corrected in the respective PCR assembly (KP02 54/92, KP03 77/77, KP12 54/55—the other ones were missing in the PCR data sets).

isolates have  $\geq 3$  mismatches compared to the SR reference, hence belonging to the problematic isolate category (Table S7). This is a decrease of 30 isolates compared to the initial results in the same laboratory. Notably, most LR assemblies now align seamlessly with the corresponding SR assemblies, as depicted in the MST in Fig. 5; Fig. S6 to S8.

Expanding upon our encouraging findings, we conducted a final multicenter assessment involving 12 selected isolates, comprising one good and two challenging isolates for each species. The enhancements were also clearly recognizable in this evaluation (refer to Fig. 5; Fig. S6 to S8; Table 2). However, it is noteworthy that the previously mentioned pattern of non-reproducibility of LR assemblies of problematic isolates, while significantly diminished, is still recognizable (LM46, KP04).

## DISCUSSION

Nanopore-based whole-genome sequencing of bacteria increasingly finds its way into genomic surveillance of pathogens and has the potential to revolutionize the field due to its low cost and simple application. However, our study demonstrated that, in the case of high-resolution typing of bacteria nanopore sequencing does not yet achieve uniform sequencing quality across strains and reproducibility between labs compared to gold-standard SR methods. For the latter, in agreement with previous studies (44), all replicates of individual strains in different laboratories provided virtually identical results in the case of standardized and high-resolution cgMLST analysis (compare Fig. 2 or Fig. S1 to S4). Remarkably, the discrepancies between long-read typing results from different participants of the same isolate were in some cases near or even surpassing the cgMLST



**FIG 5** The MST of *L. monocytogenes* (LM), constructed using both SR (green) reference data and LR (in grayscale) data under the latest optimal conditions (sampling at 5 kHz, 400 bp/s, bacterial methylation model basecalling, polishing with Racon and Medaka using its variant model), vividly illustrates the significant enhancements achieved in nanopore sequencing. To further assess this progress, LR data from two additional participants were generated and analyzed for three selected isolates per species, encompassing one good and two challenging isolates (marked in circles). While this analysis highlights improvements, it also underscores the presence of errors that still pose challenges to the reproducibility of typing, as exemplified by strains like LM46. All distances from the re-sequenced LR assemblies to the reference SR assemblies not being zero are denoted in black numbers.

**TABLE 2** Disparities in allelic variants within the LR assemblies of 12 selected and re-sequenced samples (sampling at 5 kHz, 400 bp/s, bacterial methylation model basecalling, polishing with Racon, and Medaka using its variant model) compared to the SR reference, with the improvements relative to the initial LR typing results noted in parentheses

	LAB1	LAB2	LAB3
EF22	1 (-4)	1 (-4)	2 (-8)
EF26	2 (-4)	0 (-4)	0 (-7)
EF35	1 (+1)	0 (-0)	0 (-0)
KP02	0 (-92)	1 (-89)	1 (-60)
KP04	37 (-36)	27 (-26)	17 (-17)
KP13	0 (-1)	0 (-0)	0 (-1)
LM41	0 (-30)	2 (-40)	1 (-32)
LM46	30 (-97)	66 (-134)	62 (-150)
LM54	0 (-0)	0 (-0)	0 (-0)
SA62	0 (-0)	0 (-0)	0 (-0)
SA63	2 (-24)	1 (-45)	1 (-39)
SA67	0 (-18)	1 (-21)	0 (-28)

scheme's cluster threshold for assigning isolates to outbreak clusters (see Fig. 3). Since it is not guaranteed that related isolates or even replicates of an individual isolate can be assigned to the same outbreak cluster, the method presently might not meet the requirements for high-resolution typing in such a critical setting. Furthermore, the issue of reproducibility in typing results for problematic isolates would also have an impact on cross-laboratory and cross-national surveillance.

The problem was significant, as all four studied species with public health relevance were affected to a certain degree - with the number of affected isolates varying by the species. Even within a single species, the frequency of errors varied between isolates, as the range spanned from none to several hundred incorrect alleles compared to the reference, indicating the origin of the problem at the strain level. The issue did not occur randomly. Our multicenter analyses showed that genetically similar strains, sharing the sequence type (ST), typically provided similar typing results and error ranges (Fig. 3). The ST might therefore be used for an initial estimate of data quality if comparative data is available. Particular caution should be taken when interpreting data from new or problematic STs.

Our results also clarify the discrepant results in the literature concerning the accuracy of the recent Nanopore Q20+ chemistry, when applied to a given bacterial species (28, 29, 39–41). It is important to realize that not only the species but also the strain(s) under investigation are relevant. Although results may appear contradictory at first, they might result from this (technical) limitation. Contrary to SR technologies, high typing quality of one strain (and species) did not guarantee success with other isolates from the same species in LR sequencing (see Fig. 2b and 3; Fig. S2 to S4). Consequently, general conclusions should be drawn cautiously, especially when studies and observations are based on only one or a few strains per species. In contrast to SR sequencing, nanopore sequencing would require a fundamentally different approach to performance studies of high-resolution genotyping: the crucial factor lies not solely in a high number of strains but primarily in achieving a high genetic diversity in strains and species to ensure sampling of problematic strains. This is evident in our *L. monocytogenes* strain collection. The 15 strains, in which the typing of all participants showed negligible differences to the reference, were contrasted with the most problematic isolate of the entire performance study, LM46. Conversely, our results of *E. faecium* isolates were generally good, but diversity was low (only covering two STs). It is conceivable that more severely affected strains may also be present within this species.

Upon analyzing the reads of problematic strains, consistency emerged, revealing similar error-prone sequence patterns (Fig. 4) and strand-specific basecalling errors (Fig. S5) in the sequencing data of the participants. The plausible explanation that these problems arise from DNA methylation, as reported in other recent studies (26, 27), could

be confirmed by whole-genome PCR amplification, which largely resolved the typing problems (see Table 1). Apart from the fact, that this approach requires additional wet lab work, PCR also introduces new challenges (26, 27, 45), such as increased cost, issues in amplifying GC-rich species, the introduction of new errors due to amplification and shorter reads resulting in less contiguous *de novo* assemblies and, subsequently, an increased number of missing targets in cgMLST analysis (Table 1).

A problem with far-reaching implications only became apparent through this study's multicenter approach. Although the identified methylation issue in cgMLST typing resulted in a similar order of magnitude of incorrect alleles in individual isolates between participants (Fig. 3), the affected targets varied. The difference between the correctly or the incorrectly called alleles of a target was typically reflected in a single different (incorrect) base. Even minimal differences in the data sets (Fig. S5), either a few additional reads with the correct or the incorrect base—could ultimately have significant implications on the allelic typing being correct or incorrect in the final assemblies. If the error related to methylation accumulates over several targets, it leads to completely non-reproducible allele profiles of individual isolates, which has a detrimental effect on the downstream analysis (Fig. 2b; Fig. S2 to S4). To exclude a contribution from other sources of error, we randomly subsampled data sets, and even these subsamples resulted in different allelic profiles compared to each other (Table S5), illustrating the inherent fragility of the read data. Consequently, our investigation highlights the potential for detecting problematic isolates by sequencing biological replicates or merely subsampling the data, provided that the sequencing depth permits.

The consistent errors in reads found among participants suggest the possibility and need for technical and software interventions to address the issue. Indeed, technical improvements during our study and our optimization of the polishing strategy improved typing, nevertheless, there were still problematic results for four isolates (KP03, KP04, KP12, and LM46). While this number may seem low, it is reasonable to assume additional cases, as the strains examined in our study represent only a minimal fraction of the population of a given species. This is complicated by the fact that in, e.g., outbreak scenarios, one is dealing with genetically similar strains. Should the error manifest in one strain, it is likely to recur in others.

Thus, the application of nanopore sequencing in high-resolution typing, e.g., for outbreak analysis or source tracings, must be evaluated critically currently. The potential for inaccurate data and flawed genotyping, which could have significant implications for decision-making, is deemed too high, even when dealing with a limited number of affected strains.

Nevertheless, the potential of nanopore sequencing is evident in our study. When the DNA modifications did not interfere, sequencing and typing of unproblematic strains yielded consistently good results across all participants. These results were robust despite minimal differences in laboratory workflows, a characteristic crucial for widespread applicability.

Although our study still shows a clear need for action, the substantial enhancements made by the manufacturer recently instill optimism that further developments can resolve the issue. The main challenge will be to ensure robust performance even for emerging strains with novel patterns of DNA modifications that are not represented in the models. Furthermore, it would be beneficial to conduct a detailed investigation of the specific characteristics that differentiate the “problematic” strains from the others in a future study. This should entail examining whether there are differences in the frequency of methylation and whether genetic markers can be used to identify such strains.

Since major improvements are based on models for basecalling and polishing, trained using machine learning, open communication of the training data set is required for a systematic evaluation to exclude a training bias, as observed for many A.I. models. Ideally, the improvements should apply to the entire species population and not just to a subgroup (newly integrated into the training set). Furthermore, to ensure improved but consistent results with previous iterations and future applications in routine genomic

surveillance, quality control parameters need to be developed and published together with details about the new models. Besides the species and strain composition used for the model training, this should also include how and to what extent the models were validated.

## Conclusion

As impressive as nanopore sequencing performs in many areas, it is not yet generally applicable in the clinical/public health sector for high-resolution bacterial genotyping, where even a few incorrect bases infer with the analysis, especially given the massive consequences this can ultimately have. Because the observed problem of incorrect basecalls is highly isolate-dependent even within a given species, we argue that future improvements must be evaluated using strain collections of a given species with an emphasis on a high genetic diversity. Based on our observations, we also recommend including parameters such as read ambiguity at specific positions, agreement between forward and reverse reads, and analysis of sequencing replicates for comprehensive quality assessment of nanopore data. Furthermore, for the basecalling and polishing models, the strains of the training data set as well as their quality control parameters should be disclosed to ensure a systematic evaluation and subsequently a safe use in genomic pathogen monitoring.

## ACKNOWLEDGMENTS

We thank Oxford Nanopore Technologies for supporting our study with 20 flow cells.

We would like to thank the Sequencing Core Facility of the Genome Competence Center of the Robert Koch Institute (RKI) for providing excellent sequencing services. In particular, we would like to thank Tanja Pilz from the Sequencing Core Facility, and Anne Kohler and Claudia Wiede from the Friedrich Loeffler Institute for Medical Microbiology, for their excellent technical support in preparing the Nanopore libraries and sequencing. Furthermore, we would like to thank Simone Dumschat from Unit 11 (Enteropathogenic Bacteria and Legionella) of the RKI for her excellent technical support in the extraction of high molecular weight DNA.

M.L. and C.B. received financial support from the Ministry for Economics, Sciences and Digital Society of Thuringia (TMWWDG) under the framework of the Landesprogramm ProDigital (DigLeben-5575/10-9).

During the preparation of this publication, the authors used DeepL and ChatGPT for the sole purpose of improving readability and language. After using these tools/services, the authors reviewed and edited the content as needed and took full responsibility for the content of the publication.

## AUTHOR AFFILIATIONS

<sup>1</sup>Diagnostic and Research Institute of Hygiene, Microbiology and Environmental Medicine, Medical University of Graz, Graz, Austria

<sup>2</sup>Institute for Infectious Diseases and Infection Control, Jena University Hospital, Jena, Germany

<sup>3</sup>Genome Competence Center (MF1), Robert Koch Institute, Berlin, Germany

<sup>4</sup>Austrian Agency for Health and Food Safety, Vienna, Austria

<sup>5</sup>Ares Genetics GmbH, Vienna, Austria

<sup>6</sup>Medical and Life Sciences Faculty, Furtwangen University, Villingen-Schwenningen, Germany

<sup>7</sup>Nosocomial Pathogens and Antibiotic Resistances (FG13), Robert Koch Institute, Wernigerode, Germany

<sup>8</sup>Institute for Medical Microbiology and Hospital Hygiene, Otto von Guericke University Magdeburg, Magdeburg, Germany

<sup>9</sup>Enteropathogenic bacteria and Legionella (FG11), Consultant Laboratory for Listeria, Robert Koch Institute, Wernigerode, Germany

<sup>10</sup>Austrian Agency for Health and Food Safety, Graz, Austria

<sup>11</sup>Friedrich Loeffler Institute for Medical Microbiology, F-Sauerbruch-Str., Greifswald, Germany

### AUTHOR ORCID*s*

Johanna Dabernig-Heinz  <http://orcid.org/0009-0007-2800-3659>

Sven Halbedel  <http://orcid.org/0000-0002-5575-8973>

Ariane Pietzka  <http://orcid.org/0000-0003-2987-2905>

Beatriz Daza  <http://orcid.org/0000-0002-7138-5441>

Evgeny A. Idelevich  <http://orcid.org/0009-0009-4207-5290>

Karsten Becker  <https://orcid.org/0000-0002-6391-1341>

Werner Ruppitsch  <http://orcid.org/0000-0001-9940-3333>

Ivo Steinmetz  <http://orcid.org/0000-0003-0510-7336>

Christian Kohler  <http://orcid.org/0000-0003-3921-6776>

Gabriel E. Wagner  <http://orcid.org/0000-0002-5704-3955>

### FUNDING

Funder	Grant(s)	Author(s)
Ministry for Economics, Sciences and Digital Society of Thuringia (TMWWDG)	DigLeben-5575/10-9	Mara Lohde Christian Brandt

### AUTHOR CONTRIBUTIONS

Johanna Dabernig-Heinz, Data curation, Formal analysis, Investigation, Visualization, Writing – original draft, Writing – review and editing | Mara Lohde, Data curation, Formal analysis, Investigation, Visualization, Writing – review and editing | Martin Hölzer, Data curation, Formal analysis, Investigation, Writing – review and editing | Adriana Cabal, Data curation, Formal analysis, Investigation, Writing – original draft, Writing – review and editing | Rick Conzemius, Data curation, Formal analysis, Investigation, Resources, Software, Writing – original draft, Writing – review and editing | Christian Brandt, Data curation, Formal analysis, Investigation, Methodology, Visualization, Writing – review and editing | Matthias Kohl, Data curation, Formal analysis, Investigation, Writing – review and editing | Sven Halbedel, Methodology | Patrick Hyden, Data curation, Formal analysis, Investigation | Martin A. Fischer, Methodology | Ariane Pietzka, Data curation, Resources | Beatriz Daza, Methodology | Evgeny A. Idelevich, Methodology, Resources | Anna Stöger, Methodology | Karsten Becker, Resources | Stephan Fuchs, Formal analysis, Investigation | Werner Ruppitsch, Resources, Formal analysis, Investigation | Ivo Steinmetz, Conceptualization, Methodology, Project administration, Resources, Supervision, Writing – review and editing | Christian Kohler, Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Project administration, Supervision, Writing – original draft, Writing – review and editing | Gabriel E. Wagner, Conceptualization, Data curation, Investigation, Methodology, Project administration, Supervision, Validation, Visualization, Writing – original draft, Writing – review and editing

### DATA AVAILABILITY

Nanopore and Illumina read data sets of all participants have been deposited under BioProject accession no. [PRJNA1091452](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA1091452) in the National Center for Biotechnology Information Sequence Read Archive repository.

### ADDITIONAL FILES

The following material is available [online](#).

## Supplemental Material

**Supplemental material (JCM00628-24-s0001.pdf).** Figures S1 to S8; Tables S5 and S6.

**Table S1 (JCM00628-24-s0002.xlsx).** Experimental details.

**Table S2 (JCM00628-24-s0003.xlsx).** Information about obtained sequencing data.

**Table S3 (JCM00628-24-s0004.xlsx).** Assembly metrics.

**Table S4 (JCM00628-24-s0005.xlsx).** Typing differences for problematic strains.

**Table S7 (JCM00628-24-s0006.xlsx).** Distance to the reference for best pipeline for updated assemblies.

## REFERENCES

- Wu F, Zhao S, Yu B, Chen Y-M, Wang W, Song Z-G, Hu Y, Tao Z-W, Tian J-H, Pei Y-Y, Yuan M-L, Zhang Y-L, Dai F-H, Liu Y, Wang Q-M, Zheng J-J, Xu L, Holmes EC, Zhang Y-Z. 2020. A new coronavirus associated with human respiratory disease in China. *Nature* 579:265–269. <https://doi.org/10.1038/s41586-020-2008-3>
- Meehan CJ, Goig GA, Kohl TA, Verboven L, Dippenaar A, Ezewudo M, Farhat MR, Guthrie JL, Laukens K, Miotto P, et al. 2019. Whole genome sequencing of *Mycobacterium tuberculosis*: current standards and open issues. *Nat Rev Microbiol* 17:533–545. <https://doi.org/10.1038/s41579-019-0214-5>
- Quick J, Grubaugh ND, Pullan ST, Claro IM, Smith AD, Gangavarapu K, Oliveira G, Robles-Sikisaka R, Rogers TF, Beutler NA, et al. 2017. Multiplex PCR method for MinION and Illumina sequencing of Zika and other virus genomes directly from clinical samples. *Nat Protoc* 12:1261–1276. <https://doi.org/10.1038/nprot.2017.066>
- Djordjevic SP, Jarocki VM, Seemann T, Cummins ML, Watt AE, Drigo B, Wyrsh ER, Reid CJ, Donner E, Howden BP. 2024. Genomic surveillance for antimicrobial resistance - a One Health perspective. *Nat Rev Genet* 25:142–157. <https://doi.org/10.1038/s41576-023-00649-y>
- Snell LB, Cliff PR, Charalampous T, Alcolea-Medina A, Ebie S, Sehmi JK, Flaviani F, Batra R, Douthwaite ST, Edgeworth JD, Nebbia G. 2021. Rapid genome sequencing in hospitals to identify potential vaccine-escape SARS-CoV-2 variants. *Lancet Infect Dis* 21:1351–1352. [https://doi.org/10.1016/S1473-3099\(21\)00482-5](https://doi.org/10.1016/S1473-3099(21)00482-5)
- Armstrong GL, MacCannell DR, Taylor J, Carleton HA, Neuhaus EB, Bradbury RS, Posey JE, Gwinn M. 2019. Pathogen genomics in public health. *N Engl J Med* 381:2569–2580. <https://doi.org/10.1056/NEJMs1813907>
- Gardy JL, Loman NJ. 2018. Towards a genomics-informed, real-time, global pathogen surveillance system. *Nat Rev Genet* 19:9–20. <https://doi.org/10.1038/nrg.2017.88>
- Inzaule SC, Tessema SK, Kebede Y, Ogwel Ouma AE, Nkengasong JN. 2021. Genomic-informed pathogen surveillance in Africa: opportunities and challenges. *Lancet Infect Dis* 21:e281–e289. [https://doi.org/10.1016/S1473-3099\(20\)30939-7](https://doi.org/10.1016/S1473-3099(20)30939-7)
- Oude Munnink BB, Worp N, Nieuwenhuijse DF, Sikkema RS, Haagmans B, Fouchier RAM, Koopmans M. 2021. The next phase of SARS-CoV-2 surveillance: real-time molecular epidemiology. *Nat Med* 27:1518–1524. <https://doi.org/10.1038/s41591-021-01472-w>
- Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW. 2012. Transforming clinical microbiology with bacterial genome sequencing. *Nat Rev Genet* 13:601–612. <https://doi.org/10.1038/nrg3226>
- Priest NK, Rudkin JK, Feil EJ, van den Elsen JMH, Cheung A, Peacock SJ, Laabei M, Lucks DA, Recker M, Massey RC. 2012. From genotype to phenotype: can systems biology be used to predict *Staphylococcus aureus* virulence? *Nat Rev Microbiol* 10:791–797. <https://doi.org/10.1038/nrmicro2880>
- Maljkovic Berry I, Melendrez MC, Bishop-Lilly KA, Rutvisuttinunt W, Pollett S, Talundzic E, Morton L, Jarman RG. 2020. Next generation sequencing and bioinformatics methodologies for infectious disease research and public health: approaches, applications, and considerations for development of laboratory capacity. *J Infect Dis* 221:S292–S307. <https://doi.org/10.1093/infdis/jiz286>
- Brandt C, Krautwurst S, Spott R, Lohde M, Jundzill M, Marquet M, Hölzer M. 2021. poreCov-an easy to use, fast, and robust workflow for SARS-CoV-2 genome reconstruction via nanopore sequencing. *Front Genet* 12:711437. <https://doi.org/10.3389/fgene.2021.711437>
- Tyson JR, James P, Stoddart D, Sparks N, Wickenhagen A, Hall G, Choi JH, Lapointe H, Kamelian K, Smith AD, Prystajecy N, Goodfellow I, Wilson SJ, Harrigan R, Snutch TP, Loman NJ, Quick J. 2020. Improvements to the ARTIC multiplex PCR method for SARS-CoV-2 genome sequencing using nanopore. *bioRxiv:2020.09.04.283077*. <https://doi.org/10.1101/2020.09.04.283077>
- Rambaut A, Holmes EC, O’Toole Á, Hill V, McCrone JT, Ruis C, du Plessis L, Pybus OG. 2020. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* 5:1403–1407. <https://doi.org/10.1038/s41564-020-0770-5>
- Wagner GE, Totaro MG, Volland A, Lipp M, Saiger S, Lichtenegger S, Forstner P, von Laer D, Oberdorfer G, Steinmetz I. 2021. A novel high-throughput nanopore-sequencing-based strategy for rapid and automated S-protein typing of SARS-CoV-2 variants. *Viruses* 13:2548. <https://doi.org/10.3390/v13122548>
- Bull RA, Adikari TN, Ferguson JM, Hammond JM, Stevanovski I, Beukers AG, Naing Z, Yeang M, Verich A, Gamaarachchi H, Kim KW, Luciani F, Stelzer-Braid S, Eden J-S, Rawlinson WD, van Hal SJ, Deveson IW. 2020. Analytical validity of nanopore sequencing for rapid SARS-CoV-2 genome analysis. *Nat Commun* 11:6272. <https://doi.org/10.1038/s41467-020-20075-6>
- Wagner GE, Dabernig-Heinz J, Lipp M, Cabal A, Simantzik J, Kohl M, Scheiber M, Lichtenegger S, Ehrlich R, Leitner E, Ruppsch W, Steinmetz I. 2023. Real-time nanopore Q20+ sequencing enables extremely fast and accurate core genome MLST typing and democratizes access to high-resolution bacterial pathogen surveillance. *J Clin Microbiol* 61:e0163122. <https://doi.org/10.1128/jcm.01631-22>
- Weilguny L, De Maio N, Munro R, Manser C, Birney E, Loose M, Goldman N. 2023. Dynamic, adaptive sampling during nanopore sequencing using Bayesian experimental design. *Nat Biotechnol* 41:1018–1025. <https://doi.org/10.1038/s41587-022-01580-z>
- Ulrich J-U, Epping L, Pilz T, Walther B, Stingl K, Semmler T, Renard BY. 2024. Nanopore adaptive sampling effectively enriches bacterial plasmids. *mSystems* 9:e0094523. <https://doi.org/10.1128/msystems.00945-23>
- Sereika M, Kirkegaard RH, Karst SM, Michaelsen TY, Sørensen EA, Wollenberg RD, Albertsen M. 2022. Oxford Nanopore R10.4 long-read sequencing enables the generation of near-finished bacterial genomes from pure cultures and metagenomes without short-read or reference polishing. *Nat Methods* 19:823–826. <https://doi.org/10.1038/s41592-022-01539-7>
- Ko KKK, Chng KR, Nagarajan N. 2022. Metagenomics-enabled microbial surveillance. *Nat Microbiol* 7:486–496. <https://doi.org/10.1038/s41564-022-01089-w>
- Viehweger A, Marquet M, Hölzer M, Dietze N, Pletz MW, Brandt C. 2023. Nanopore-based enrichment of antimicrobial resistance genes - a case-based study. *GigaByte* 2023:gigabyte75. <https://doi.org/10.46471/gigabyte75>
- Lerminiaux N, Fakhruddin K, Mulvey MR, Mataseje L. 2024. Do we still need Illumina sequencing data? Evaluating Oxford Nanopore technologies R10.4.1 flow cells and the rapid v14 library prep kit for Gram negative bacteria whole genome assemblies. *Can J Microbiol* 70:178–189. <https://doi.org/10.1139/cjm-2023-0175>
- Hall MB, Wick RR, Judd LM, Nguyen ANT, Steinig EJ, Xie O, Davies MR, Seemann T, Stinear TP, Coin LJM. 2024. Benchmarking reveals superiority of deep learning variant callers on bacterial nanopore sequence data. *bioRxiv*. <https://doi.org/10.7554/eLife.98300.1>

26. Lohde M, Wagner GE, Dabernig-Heinz J, Viehweger A, Braun SD, Monecke S, et al. 2023. Nanopore sequencing for accurate bacterial outbreak tracing. *bioRxiv*. <https://doi.org/10.1101/2023.09.15.556300>
27. Chiou CS, Chen BH, Wang YW, Kuo NT, Chang CH, Huang YT. 2023. Correcting modification-mediated errors in nanopore sequencing by nucleotide demodification and reference-based correction. *Commun Biol* 6:1215. <https://doi.org/10.1038/s42003-023-05605-4>
28. Maiden MCJ, Jansen van Rensburg MJ, Bray JE, Earle SG, Ford SA, Jolley KA, McCarthy ND. 2013. MLST revisited: the gene-by-gene approach to bacterial genomics. *Nat Rev Microbiol* 11:728–736. <https://doi.org/10.1038/nrmicro3093>
29. Ruppitsch W, Pietzka A, Prior K, Bletz S, Fernandez HL, Allerberger F, Harmsen D, Mellmann A. 2015. Defining and evaluating a core genome multilocus sequence typing scheme for whole-genome sequence-based typing of *Listeria monocytogenes*. *J Clin Microbiol* 53:2869–2876. <https://doi.org/10.1128/JCM.01193-15>
30. Deneke C, Uelze L, Brendebach H, Tausch SH, Malorny B. 2021. Decentralized investigation of bacterial outbreaks based on hashed cgMLST. *Front Microbiol* 12:649517. <https://doi.org/10.3389/fmicb.2021.649517>
31. Lichtenegger S, Trinh TT, Assig K, Prior K, Harmsen D, Pesl J, Zauner A, Lipp M, Que TA, Mutsam B, Kleinhappl B, Steinmetz I, Wagner GE. 2021. Development and validation of a *Burkholderia pseudomallei* core genome multilocus sequence typing scheme to facilitate molecular surveillance. *J Clin Microbiol* 59:e0009321. <https://doi.org/10.1128/JCM.00093-21>
32. Linde J, Brangsch H, Hölzer M, Thomas C, Elschner MC, Melzer F, Tomaso H. 2023. Comparison of Illumina and Oxford Nanopore technology for genome analysis of *Francisella tularensis*, *Bacillus anthracis*, and *Brucella suis*. *BMC Genomics* 24:258. <https://doi.org/10.1186/s12864-023-09343-z>
33. George S, Pankhurst L, Hubbard A, Votintseva A, Stoesser N, Sheppard AE, Mathers A, Norris R, Navickaite I, Eaton C, Iqbal Z, Crook DW, Phan HTT. 2017. Resolving plasmid structures in *Enterobacteriaceae* using the MinION nanopore sequencer: assessment of MinION and MinION/Illumina hybrid data assembly approaches. *Microb Genom* 3:e000118. <https://doi.org/10.1099/mgen.0.000118>
34. Peter S, Bosio M, Gross C, Bezdán D, Gutierrez J, Oberhettinger P, Liese J, Vogel W, Dörfel D, Berger L, Marschal M, Willmann M, Gut I, Gut M, Autenrieth I, Ossowski S. 2020. Tracking of antibiotic resistance transfer and rapid plasmid evolution in a hospital setting by Nanopore sequencing. *mSphere* 5:e00525-20. <https://doi.org/10.1128/mSphere.00525-20>
35. Kolmogorov M, Armstrong J, Raney BJ, Stretter I, Dunn M, Yang F, Odum D, Flicek P, Keane TM, Thybert D, Paten B, Pham S. 2018. Chromosome assembly of large and complex genomes using multiple references. *Genome Res* 28:1720–1732. <https://doi.org/10.1101/gr.236273.118>
36. Vaser R, Sović I, Nagarajan N, Šikić M. 2017. Fast and accurate *de novo* genome assembly from long uncorrected reads. *Genome Res* 27:737–746. <https://doi.org/10.1101/gr.214270.116>
37. Souvorov A, Agarwala R, Lipman DJ. 2018. SKESA: strategic k-mer extension for scrupulous assemblies. *Genome Biol* 19:153. <https://doi.org/10.1186/s13059-018-1540-z>
38. Jünemann S, Sedlazeck FJ, Prior K, Albersmeier A, John U, Kalinowski J, Mellmann A, Goesmann A, von Haeseler A, Stoye J, Harmsen D. 2013. Updating benchtop sequencing performance comparison. *Nat Biotechnol* 31:294–296. <https://doi.org/10.1038/nbt.2522>
39. Leopold SR, Goering RV, Witten A, Harmsen D, Mellmann A. 2014. Bacterial whole-genome sequencing revisited: portable, scalable, and standardized analysis for typing and detection of virulence and antibiotic resistance genes. *J Clin Microbiol* 52:2365–2370. <https://doi.org/10.1128/JCM.00262-14>
40. de Been M, Pinholt M, Top J, Bletz S, Mellmann A, van Schaik W, Brouwer E, Rogers M, Kraat Y, Bonten M, Corander J, Westh H, Harmsen D, Willems RJL. 2015. Core genome multilocus sequence typing scheme for high-resolution typing of *Enterococcus faecium*. *J Clin Microbiol* 53:3788–3797. <https://doi.org/10.1128/JCM.01946-15>
41. Weber RE, Pietsch M, Frühauf A, Pfeifer Y, Martin M, Luft D, Gatermann S, Pfennigwerth N, Kaase M, Werner G, Fuchs S. 2019. IS26-mediated transfer of *bla*<sub>NDM-1</sub> as the main route of resistance transmission during a polyclonal, multispecies outbreak in a German hospital. *Front Microbiol* 10:2817. <https://doi.org/10.3389/fmicb.2019.02817>
42. Oxford\_Nanopore\_Technologies. 2023. Bacterial assembly and annotation workflow. Available from: <https://labs.epi2me.io/workflows/wf-bacterial-genomes/>
43. Oxford\_Nanopore\_Technologies. Assembling bacterial genomes using long nanopore sequencing reads 2022. Retrieved 25 Mar 2024. Accessed March 25, 2024
44. Mellmann A, Andersen PS, Bletz S, Friedrich AW, Kohl TA, Lilje B, Niemann S, Prior K, Rossen JW, Harmsen D. 2017. High interlaboratory reproducibility and accuracy of next-generation-sequencing-based bacterial genotyping in a ring trial. *J Clin Microbiol* 55:908–913. <https://doi.org/10.1128/JCM.02242-16>
45. Ordóñez CD, Redrejo-Rodríguez M. 2023. DNA polymerases for whole genome amplification: considerations and future directions. *Int J Mol Sci* 24:9331. <https://doi.org/10.3390/ijms24119331>