



EURL-CAMPYLOBACTER

REPORT

PROFICIENCY TEST NUMBER 33

WGS and cluster analysis of *Campylobacter*

Publication history

Version	Date
Final version	2022-10-28



**Co-funded by
the European Union**

Funded by the European Union. Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union or the European Health and Digital Executive Agency (HaDEA). Neither the European Union nor HaDEA can be held responsible for them.

Contents

Abbreviations	3
Summary of the proficiency test number 33, 2022	4
Introduction	5
Terms and definitions	5
Outline of the proficiency test	5
Strain selection.....	5
Production and quality control of DNA samples.....	6
Production of reference genomes	6
Distribution of the proficiency test and reporting of results.....	7
Methods for analysis.....	7
Assessing the performance of the NRLs	8
Assessment of sequence quality	8
Assessment of cluster analysis.....	8
Results	8
DNA library preparations and sequencing	9
MLST analyses	9
AMR analyses.....	9
Sequence coverage.....	10
Quantification of high-quality bases.....	10
Sequence contamination	11
Coverage of reference genome <i>k</i> -mers	12
Deviation from expected GC-content.....	12
Assemblies.....	13
Assessment of sequence quality and NRL performance	14
Cluster and phylogenetic analysis	14
Assessment of cluster analysis and NRL performance.....	16
References	18
Appendix A – QC metrics for submitted raw data	19
Appendix B – QC metrics for submitted assemblies	22

Abbreviations

<i>C.</i>	<i>Campylobacter</i>
EU	European Union
EURL	European Union reference laboratory
cgMLST	core genome MLST
wgMLST	whole genome MLST
NGS	next generation sequencing
NRL	national reference laboratory (in this report used for all participating laboratories, also in non-EU Member States)
PT	proficiency test
SNP	single nucleotide polymorphism
ST	sequence type
WGS	whole genome sequencing

Summary of the proficiency test number 33, 2022

The EU reference laboratory for *Campylobacter* organised proficiency test (PT) number 33 on WGS and cluster analysis of *Campylobacter* in March 2022. The PT included WGS and cluster analysis of seven samples of *Campylobacter*. The objective was to assess the quality of whole genome sequence (WGS) data and accuracy of cluster analysis of *Campylobacter* performed by participating laboratories.

Participation in PT 33 was voluntary for all NRLs. Twenty-three NRLs in 18 EU member states (some member states have more than one NRL) and in Norway, Switzerland and United Kingdom received the PT and responses were reported from 20 NRLs.

The individual parts (sequence quality and cluster analysis) have been assessed through different criteria as satisfactory/needs improvement and no overall performance criteria has been applied for this PT.

In summary, the majority of the NRLs met the criteria for satisfactory performance in sequence quality and all NRLs met the criteria for satisfactory performance in cluster analysis.

Introduction

Proficiency test (PT) number 33 on WGS and cluster analysis of *Campylobacter* was organised by the EU reference laboratory (EURL) for *Campylobacter* in March 2022. Participation in the PT was voluntary. Twenty-three national reference laboratories (NRLs) in 18 EU member states (some member states have more than one NRL) and in Norway, Switzerland and United Kingdom registered for the PT. The test results and operational details were reported to the EURL from 20 NRLs in 17 EU Member States (some Member States have more than one NRL) and in Norway and United Kingdom.

The PT included whole genome sequencing (WGS) and cluster analysis of seven samples of DNA from *Campylobacter jejuni*. The main purpose of this PT was to help laboratories in the implementation of WGS and cluster analysis. The test was also designed to test the joint capability of the network to solve a multi-country *Campylobacter* outbreak based on WGS data. The objective was to assess the quality of WGS data and accuracy of cluster analysis of *Campylobacter* performed by participating laboratories.

Terms and definitions

Only some selected terms are defined here. For additional definitions of terms used in this document, please see ISO 23418:2022 [1].

- Assembly: output from a process of aligning and merging sequencing reads into larger contiguous sequences (contigs).
- Coverage: number of times that a given base position is read in a sequencing run.
- Library: collection of genomic DNA fragments from a single isolate intended for determining genome sequence(s).
- N50: length (N) such that sequence contigs of N or longer include half the bases in the assembly.

Outline of the proficiency test

Strain selection

All strains used in PT 33 were *C. jejuni* of sequence type (ST) 19. The strains were selected based on cgMLST analysis using the ‘Oxford scheme’ (PubMLST scheme based on 1343 targets) [2] and SNP analysis using Snippy [3], aiming for a relevant challenge in cluster analysis and confirming similar topology with both comparison approaches.

Sample PT33-1 and PT33-6 were from the same stock of DNA, whereas all other samples were from strains isolated at different timepoints and mostly from different farms (Table 1).

Table 1. Information about the 7 DNA samples distributed to the NRLs in proficiency test No. 33 (2022).

Sample	Strain	Matrix	Location	Sampling
PT33-1, PT33-6	20C120	Chicken caeca	Sweden, farm A	October, 2020
PT33-2	20C028	Chicken caeca	Sweden, farm B	July, 2020
PT33-3	20C126	Chicken caeca	Sweden, farm C	July, 2020
PT33-4	20C060	Chicken caeca	Sweden, farm A	July, 2020
PT33-5	Val_Cj015	Milk filter	Sweden, farm D	2011
PT33-7	20C102	Chicken caeca	Sweden, farm E	October, 2020

Production and quality control of DNA samples

Strains were cultivated on horse blood agar and overnight cultures were prepared using 2-3 colonies inoculated in BHI with 0.6 % (w/v) yeast extract. The cultures were grown until OD600 values reached >0.5. The cultures were collected by centrifugation, washed with PBS and the pellets were frozen at -20 °C. Multiple pellets were prepared for each sample. The overnight cultures were checked for contamination by cultivation on blood agar plates.

DNA was extracted using the Genomic Tip 20/G kit (Qiagen) according to the kit protocol, except that Ready-Lyse reagent (Biosearch Technologies) was used for cell lysis instead of 100 mg/ ml lysozyme. The concentration of the extracted DNA was measured using a Qubit 2.0 with a DNA HS kit (Thermo Fisher Scientific) and quality checked with a Nanopore instrument. Multiple DNA solutions were pooled to generate a homogeneous stock solution for each sample. The stock was then quantified, and quality checked as described above.

For stabilisation, DNA stocks were mixed with GenTegra-DNA in 0.5 ml screw cap tubes (GenTegra). The solution was then further aliquoted to DNase free 1.5 ml screw cap microtubes and dried by leaving the cap off in a biosafety hood for at least 48 h. The tubes were then closed and stored at room temperature. The expected yield for each tube was >1 µg.

The whole test, including DNA quantification and quality check, library preparation, sequencing and cluster analysis of the PT samples, was performed by the EURL before dispatch and one week after final date to report the results to control that the test was stable. The output was satisfactory in terms of both DNA and sequence quality and the cluster topology was the same in both datasets.

Production of reference genomes

Reference genomes of all six strains were generated by sequencing the strains on both Illumina MiSeq and Oxford Nanopore instruments. Hybrid assemblies using both short-read and long-read data were then generated using Tricycler v0.5.3 [4] for PT33-1 – PT33-6 and Unicycler v0.5.0 [5] for PT33-7 due to low depth of long-read data for this sample. Complete (gap-free) genomes were obtained for all the samples. The assembly sizes are listed in Table 2.

Table 2. Statistics for reference assemblies for proficiency test No. 33 (2022).

Sample	No. of contigs	Assembly size (bp)	GC %	Assembly pipeline
PT33-1, PT33-6	1	1711096	30,47	Trycycler
PT33-2	1	1753524	30,40	Trycycler
PT33-3	1	1673246	30,48	Trycycler
PT33-4	1	1753518	30,40	Trycycler
PT33-5	1	1760653	30,34	Trycycler
PT33-7	1	1711097	30,48	Unicycler

Distribution of the proficiency test and reporting of results

The PT samples were distributed from the EURL on the 7th of March 2022 together with PT 31 and PT 32. The samples were placed in foam boxes along with freezing blocks. The foam boxes were packed in cardboard boxes for transport and were sent from the EURL using courier service. One test was distributed by ordinary mail and was sent to an NRL that did not participate in the other two PTs.

Each participant received a plastic bag containing seven numbered tubes, each containing stabilised and dried DNA from *Campylobacter*. A Micro-T-Log was included in each package to record the temperature every second hour during transport.

Fifteen NRLs received the PT within one day after the packages had been dispatched from the EURL, five NRLs within two days and the one sent by ordinary mail within three weeks.

The PT samples were recommended to be stored at room temperature until start of analysis. Instructions for rehydration of each sample were included in the packages and were also sent out by e-mail a few days before the PT distribution.

All results and information about the procedures had to be reported in the Questback Essentials system before 1st of June 2022. Additional data requested were: raw sequence files (i.e., FASTQ files), assembly files in FASTA format (only requested if assembling was part of the analysis), tree used to draw conclusions from the analysis (e.g. phylogenetic tree or minimum spanning tree) and raw clustering data used to create the tree (distance matrix or alignment). Participants were instructed to upload requested files onto a personal OneDrive folder. Each NRL was given a unique LabID number that was used as an identifier for reporting and uploading of sequence data. LabID has been shortened to L# in the text and figures of this report.

Methods for analysis

The test included to perform library preparations, sequencing, identification of Multi Locus Sequence Type (MLST) and voluntarily antibiotic resistance (AMR) genes and/or point mutations potentially causing AMR, and cluster analysis. The NRLs were instructed to use their standard laboratory procedures for all parts of the analysis. Cluster analysis could be performed using SNP methods, gene-by-gene methods (wgMLST or cgMLST) or other types of comparisons. The participants were instructed to use their own interpretation (cut-off value) of a cluster.

Assessing the performance of the NRLs

Different criteria for the individual steps of each part (sequence quality and cluster analysis) were assessed. The results were ‘Satisfactory’ when all criteria were met for all the samples, whereas failure to reach one or more criteria for one or more samples was marked as ‘Needs improvement’. No overall performance grade was applied for this PT. Overall comments on the data and possible focus areas for improving performance were commented further in each laboratory’s individual report.

Assessment of sequence quality

Cut-off values were defined for five different criteria to assess the sequence quality through submitted information (MLST) or in fastq files for each sample (Table 3). The criteria were: definition of ST, percentage of Q30 bases, percentage of contaminating reads, percentage coverage of the corresponding reference genome, and percentage GC-deviation in the sequence reads from the corresponding reference genome.

Table 3. Overview of the criteria and cut-off values used for assessment of sequence quality in proficiency test No. 33 (2022).

Criteria	Cut-off value for satisfactory performance
MLST	Must match ST-19
Q30	>70 %, 75 % or 80 % depending on read length (300, 250, 150-100 bp)
Contamination	<5 % from non-target species
Reference coverage	>98 % of reference genome ^a
GC-deviation	<4 % deviation from reference genomes

^aThe maximum amount of data used for the assessment was 80X coverage for NRLs using Nextera XT and 30X coverage for NRLs using other library preparation kits.

Assessment of cluster analysis

The assessment of cluster analysis was done on the topology of MSTs or phylogenetic trees provided by participants. If no tree was provided, the topology was derived from the distance matrices submitted by the participants. Three statements were used to capture the topology; (i) “PT33-6 and PT33-7 are the two closest samples to PT33-1”, (ii) “PT33-4 is the closest sample to PT33-2”, and (iii) “PT33-5 is most distant to the other samples”.

Results

Proficiency test number 33 was distributed to 23 NRLs and 20 of them reported the results of the analysis. The analysis was started at different timepoints between March and end of May 2022. The sequence quality measures were calculated in raw data submitted by participants. If adapters were left in the data, they were removed before the analysis using Trimmomatic

ILLUMINACLIP. Adapter-free sequence data was used for all proceeding evaluations. A summary of all sequence quality measures in each submitted dataset can be found in Appendix A.

DNA library preparations and sequencing

For library preparations, 11 NRLs used the Illumina DNA Prep kit (previously known as Nextera DNA Flex Library Preparation kit), six NRLs used the Illumina Nextera XT, one NRL used the Illumina TrueSeq DNA Nano/PCR-Free, one NRL used the Invitrogen Collibri ES DNA Library Prep, and one NRL used the NEBNext Ultra II DNA Library Prep. One NRL using Nextera XT used 1/5 of the volume of reagents, and four NRLs using the DNA Prep kit used half the volume or less.

For quantification of the library preparations, 12 NRLs used Qubit, two NRLs used Agilent TapeStation, one NRLs used Agilent BioAnalyzer, three NRLs used quantitative PCR, one NRL used Promega QuantiFluor dsDNA and one NRL used the Quant-IT kit and Promega GloMax fluorometer. Library quality control was performed by 15 NRLs; six using Agilent TapeStation, five using Agilent BioAnalyzer, one using Agilent Fragment Analyzer, one used LabChip, one Qsep100 Fragment Analyzer, and one used capillary electrophoresis.

All the participating NRLs used Illumina technology for sequencing. The majority used Illumina MiSeq for sequencing, except two used NextSeq, one NovaSeq, one MiniSeq and one HiSeq. The read length was: 2x100 (1), 2x150 (7), 2x250 (5) and 2x300 (7). Nineteen NRLs reported targeted theoretical coverage: 40X-60X (7), 80X-100X (7) and >100X (5). NRLs using Nextera XT for library preparation aimed for 40X-60X or 80X-100X in theoretical coverage. All NRLs aiming for >100X used the DNA Prep kit for library preparation.

MLST analyses

Nineteen NRLs correctly identified the STs for all the samples, and one NRL correctly identified all individual alleles of ST-19 for all the samples but did not report the ST.

AMR analyses

Nineteen NRLs performed the optional AMR analyses. The participants were instructed to report any genes or point mutations that could possibly lead to AMR. Thirteen NRLs used ResFinder (version 4.0 or later), two used AMRFinderPlus, three used ABRicate (one in combination with Pointfinder) and one used an in-house developed software.

Sixteen of the 19 NRLs that performed the AMR analyses reported a beta-lactam resistance (*bla*OXA) gene in all the samples and three NRLs reported no AMR gene in any of the samples.

Eighteen of the 19 NRLs that performed the AMR analyses were able to identify the quinolone resistance (*gyrA*. P. T86I) point mutation in samples PT33-1, PT33-3, PT33-6 and PT33-7. Two NRLs identified the quinolone resistance (*gyrA*. P. T86I) point mutation in all the samples whereas 17 NRLs reported no point mutation in the remaining samples.

Sequence coverage

Quantifications and QC measurements were made with the tngs script [6]. A minimum coverage threshold of 30X was applied for unbiased library preps (all except Nextera XT) and 80X for NRLs using the Nextera XT, which has GC-dependent coverage bias and therefore requires higher coverage for similar performance. All NRLs submitted data with coverage over the expected threshold (Figure 1A). The sequence depths in relation to the threshold were also quantified for each dataset, illustrating the excess amount of data produced by each NRL (Figure 1B).

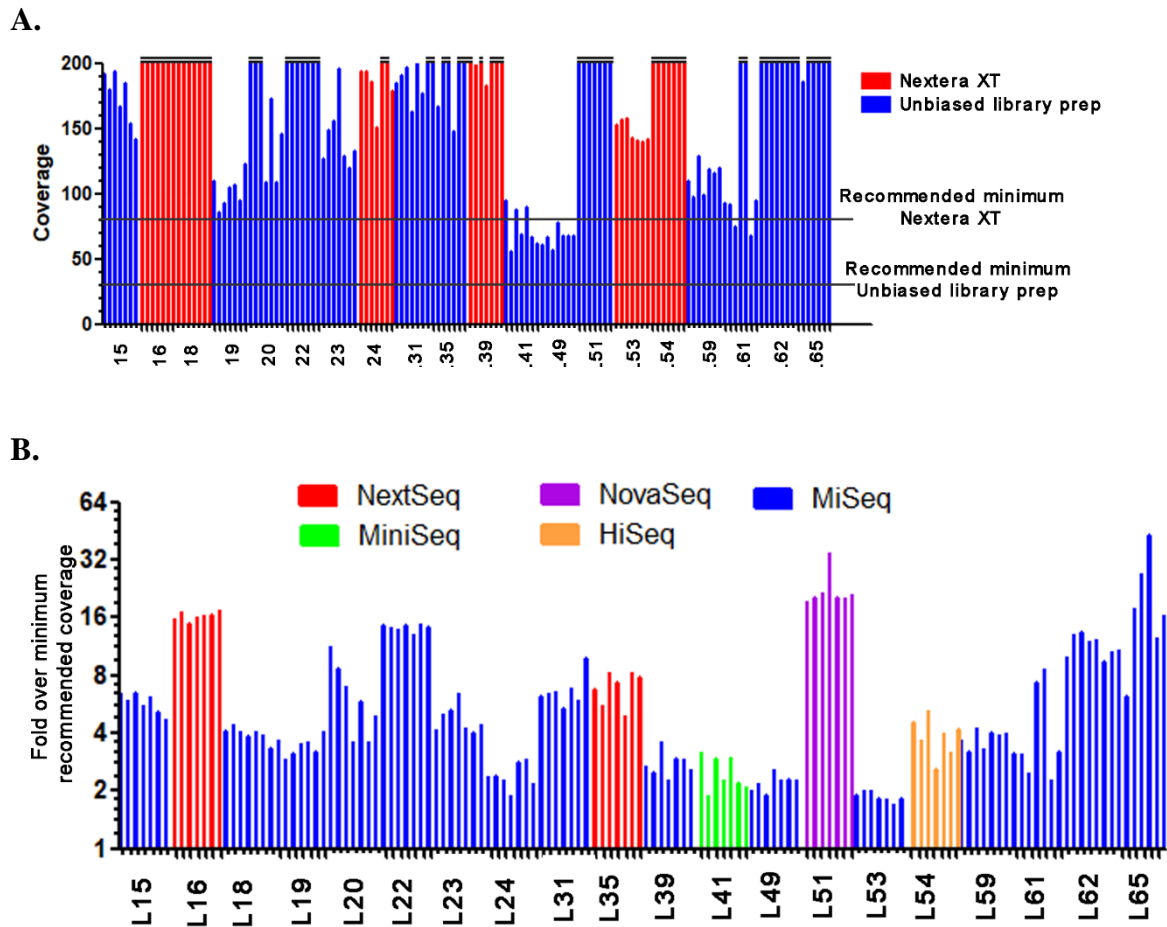


Figure 1. A) Sequence read coverage, i.e., average number of sequenced bases per base position in the reference genome. The recommended coverage is indicated and is higher when Nextera XT has been used. Data submitted by each NRL is colour coded based on whether Nextera XT was used or not.

B) Quantification of how many fold excess data that was produced in relation to the recommended minimum (30X or 80X depending on the type of library prep used).

Quantification of high-quality bases

The quantification of high-quality bases was done using the tngs script [6]. The percentage of bases with at least a quality score Q30 was calculated. The minimum quality threshold was set depending on the number of cycles sequenced (read length). Short read lengths (2x100-2x150 bp) were expected to have at least 80 % Q30 bases, 2x250 bp were expected to have at least 75

% Q30 bases and 2x300 bp at least 70 % Q30 bases [1]. All NRLs except one produced data of quality above the threshold (Figure 2).

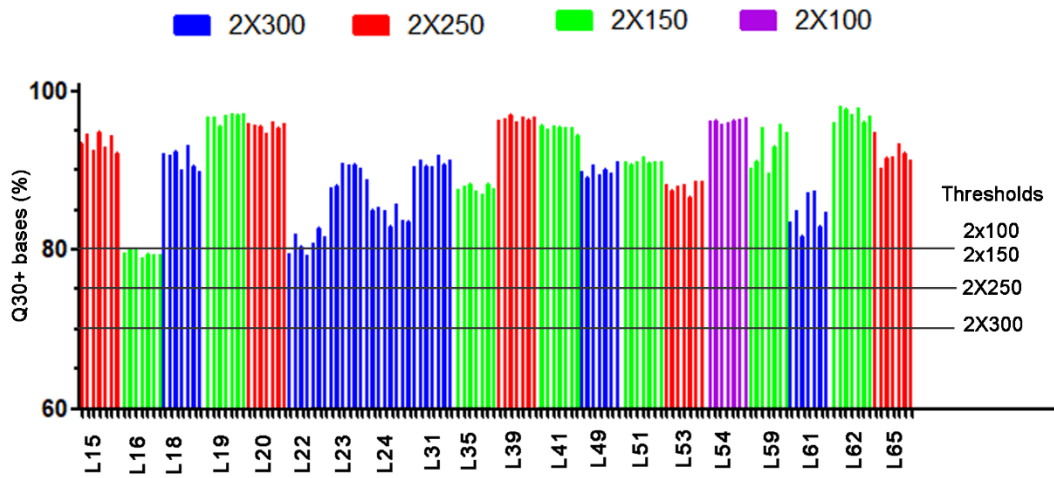


Figure 2. Quantification of the percentage bases having quality score Q30 or higher, coloured based on the number of cycles used. The threshold for satisfactory quality is dependent on the number of sequencing cycles used.

Sequence contamination

Contamination levels were estimated using the Kraken2 [7] software to obtain metagenomic information about the sequencing datasets. Kraken2 classifies reads as belonging to different phylogenetic taxa and this indicates if the correct species was sequenced and if the samples contained contaminating reads from a different organism. The threshold was set to 5 % [1]. The contamination levels were low (below 1 %) for all samples from all NRLs apart from sample PT33-4. Almost all NRLs had the highest levels of contamination in sample PT33-4, and it was especially high in one case where it exceeded the threshold of 5 % (Figure 3). The dominating contaminating species was the same, which indicates that the contamination occurred before the samples were dispatched from the EURL. Therefore, sample PT33-4 has been excluded from the evaluation of sequence quality.

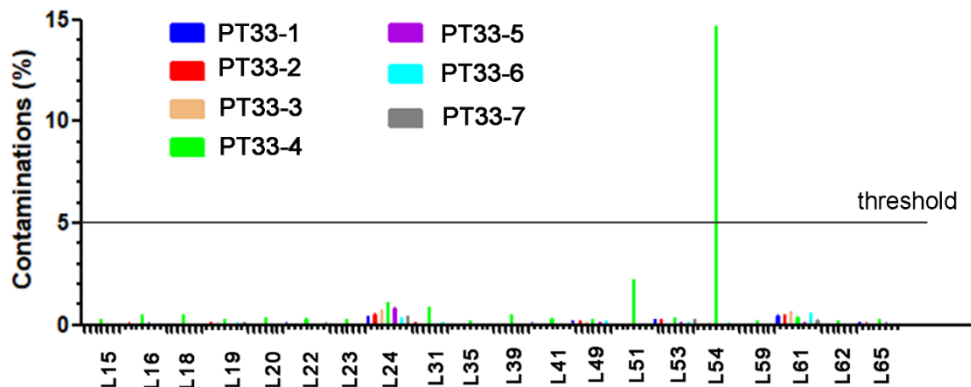


Figure 3. Quantification of the percentage reads originating from a contaminating species. The threshold for satisfactory contamination level is indicated (5 %). Sample PT33-4 was not used in the evaluation because the contamination may have arisen before the samples were distributed.

Coverage of reference genome k -mers

The percentage of k -mers present in the reference genomes that was covered by the read data submitted by participants was quantified using the tngs script [6]. The percentage of k -mers present in the reference genomes that was found in the raw data was quantified using 30X of raw data (Figure 4A) and 80X of raw data (Figure 4B). The threshold was set to 98 % reference coverage, but the NRLs using Nextera XT were evaluated at 80X and the NRLs using unbiased library prep at 30X of data. All NRLs reached the expected coverage of reference genome k -mers.

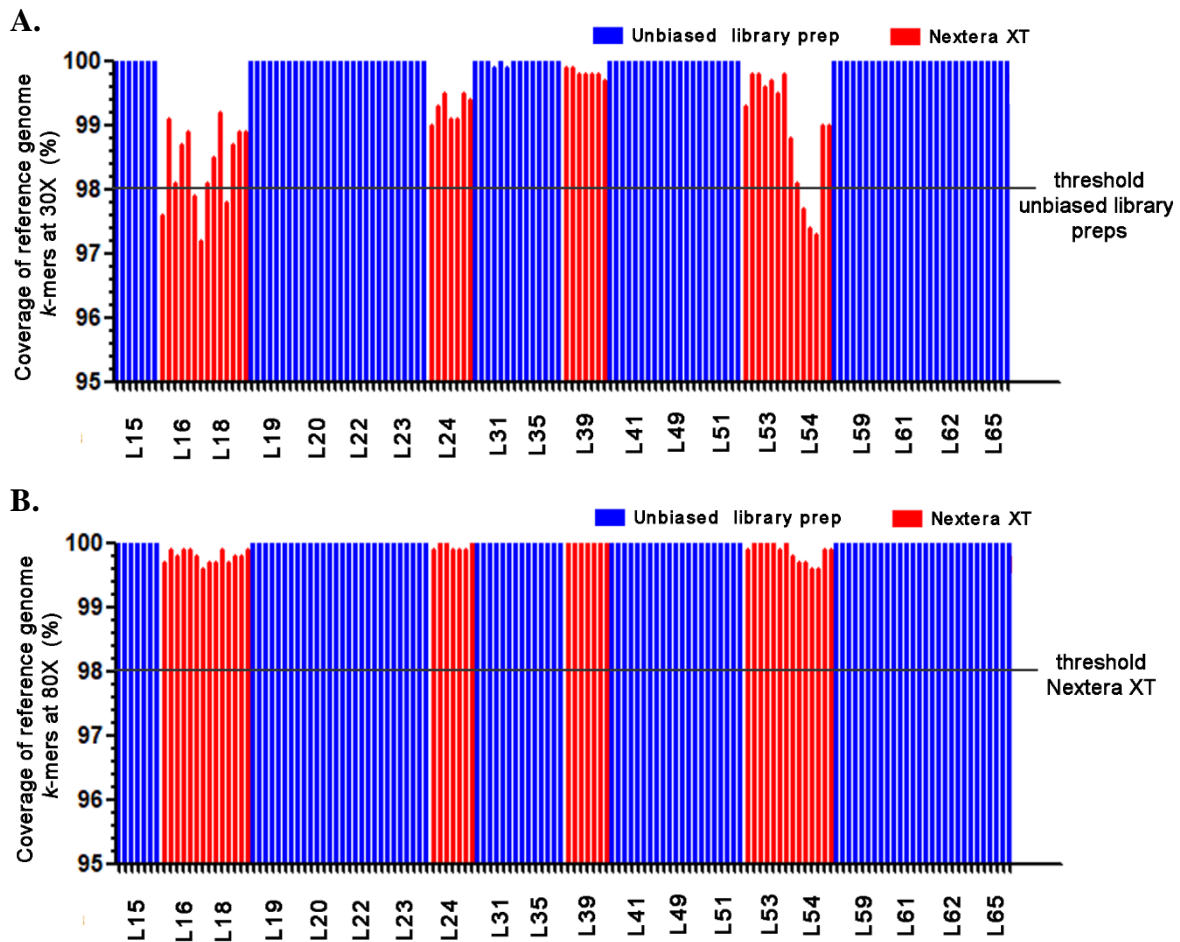


Figure 4. Coverage of the reference genome, i.e., the percentage of the k -mers present in the reference genome sequence that were also present in the 30X read data submitted by NRLs (A) or 80X read data (B). The threshold for satisfactory reference coverage was set at 98 %, but was evaluated at 80X read data for NRLs that used Nextera XT and at 30X read data for NRLs that used other library preps. Thus, the calculations were done at the recommended minimum read depth coverage.

Deviation from expected GC-content

Quantification of the GC-content was made using the tngs script [6]. The GC-content of respective reference genome sequence was used as an expected value. Deviation from this value is seen in Nextera XT data because of GC-dependent coverage bias but it can also arise by large number of contaminating reads from a species with different GC-content. The cut-off value

applied for this criterion was a deviation larger than 4 % [1]. Two NRLs using Nextera XT had samples above this threshold (Figure 5).

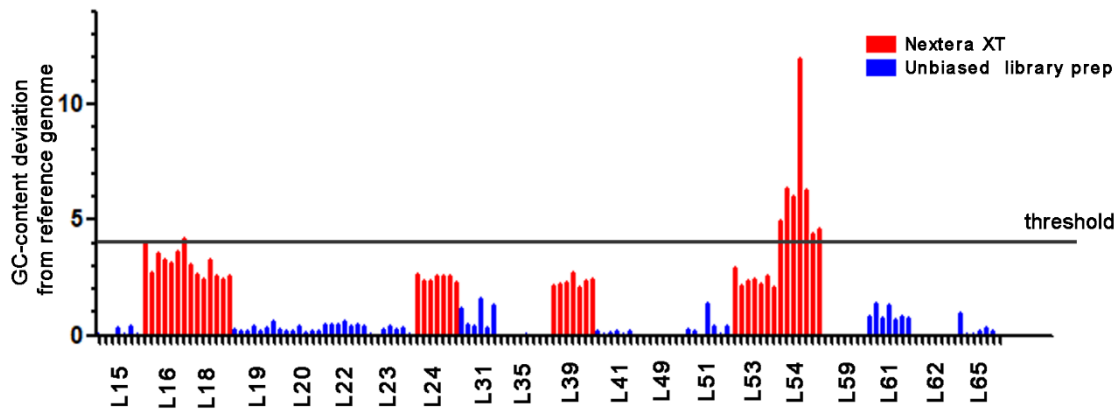


Figure 5. Deviation of the GC-content in the read data submitted by each NRLs from the GC-content of the reference genome, colour coded based on whether Nextera XT was used or not in order to illustrate that this is a main factor affecting the results in this measurement. Threshold for satisfactory performance is indicated (4 %).

Assemblies

Participants were asked to submit assemblies if this was part of their analysis. Nineteen of 20 NRLs reported to have generated assemblies. One NRL used Velvet while the remaining NRLs used SPAdes, three as integrated in Shovill and two as integrated in Unicycler.

Assemblies were submitted by 17 participants. Four NRLs had used coverage filtering, thus removed contigs before submitting the assembly files.

The following QC metrics were calculated for each assembly submitted by the participants:

- Total size of assembly (bp)
- k -mer coverage over the reference genome (%)
- Total number of contigs
- Total number of contigs > 1kb
- Longest contig
- N50 length

The QC metrics for each assembly is summarised in Appendix B.

For some NRLs, the contamination present in sample PT33-4 became part of the assembly, affecting size and number of contigs negatively.

The sizes of the assemblies never deviated more than 2 % from the size of the reference genomes (except for sample PT33-4). One assembly, PT33-1-L61, had a slightly larger assembly than the reference, which might be explained by low levels of contamination included in the assembly.

The assemblies k -mer coverage over the reference genomes was never lower than 99 % with most of the assemblies covering 99.9 % of the respective reference genomes. With sample PT33-4 excluded, the total number of contigs ranged from 1 to 68 with the median value around

30 contigs for each sample. L20 used Shovill for their assembly and managed to generate assemblies containing only 1-6 contigs for all the samples, including sample PT33-4. With L20 excluded, the N50 size varied from 79,979 bp to 216,952 bp for the assemblies and the longest contig varied from 188,052 bp to 720,947 bp in size.

Assessment of sequence quality and NRL performance

The results using the defined criteria for assessment of sequence quality of each NRL is summarised in Table 4. According to the assessment, 18 NRLs fulfilled the criteria for satisfactory performance on all samples and 2 NRLs scored below the criteria for one or more samples.

Table 4. Overview of assessment of the sequence quality of each NRL in proficiency test No. 33 (2022). Number of samples out of six included in the evaluation reaching the criteria cut-offs. PT33-4 was excluded from the evaluation.

Lab code	MLST	Q30	Contamination	Reference coverage	GC deviation	Overall evaluation sequence quality
L15	6/6	6/6	6/6	6/6	6/6	Satisfactory
L16	6/6	2/6	6/6	6/6	5/6	Needs improvement
L18	6/6	6/6	6/6	6/6	6/6	Satisfactory
L19	6/6	6/6	6/6	6/6	6/6	Satisfactory
L20	6/6	6/6	6/6	6/6	6/6	Satisfactory
L22	6/6	6/6	6/6	6/6	6/6	Satisfactory
L23	6/6	6/6	6/6	6/6	6/6	Satisfactory
L24	6/6	6/6	6/6	6/6	6/6	Satisfactory
L31	6/6	6/6	6/6	6/6	6/6	Satisfactory
L35	6/6	6/6	6/6	6/6	6/6	Satisfactory
L39	6/6	6/6	6/6	6/6	6/6	Satisfactory
L41	6/6	6/6	6/6	6/6	6/6	Satisfactory
L49	6/6	6/6	6/6	6/6	6/6	Satisfactory
L51	6/6	6/6	6/6	6/6	6/6	Satisfactory
L53	6/6	6/6	6/6	6/6	6/6	Satisfactory
L54	6/6	6/6	6/6	6/6	0/6	Needs improvement
L59	6/6	6/6	6/6	6/6	6/6	Satisfactory
L61	6/6	6/6	6/6	6/6	6/6	Satisfactory
L62	6/6	6/6	6/6	6/6	6/6	Satisfactory
L65	6/6	6/6	6/6	6/6	6/6	Satisfactory

Cluster and phylogenetic analysis

For cluster analysis, 16 NRLs used gene-by-gene based comparisons (cgMLST/wgMLST) and nine NRLs used SNP-based comparisons. Five of the 20 NRLs used both gene-by-gene and SNP-based comparisons. Among NRLs using commercially available software for cluster analysis, six used Ridom SeqSphere+ and two BioNumerics. Of those using online tools, one used CGE cgMLSTfinder and one used NDtree. Of those using “Open-source / in-house developed pipeline”, two used Snippy, four used chewBBACA, one used PyMLST, one used BCFtools, one used CSI phylogeny and two used SAM tools to call SNPs. Of the 19 NRLs responding to the question, four NRLs used read-mapping for allele or variant calling, and 15 used assemblies. Of the four NRLs using read-mapping, one used it for cgMLST analysis and three for SNP analysis.

For gene-by-gene based comparisons, six used the scheme in Ridom SeqSphere+ (two used only 637 core targets, and four used core targets plus 958 accessory targets), 10 used the ‘Oxford scheme’ (PubMLST scheme based on 1,343 targets) [2], and one used the INNUENDO wgMLST scheme with 2,795 targets [8].

Thirteen NRLs created Minimum Spanning Trees (MST) for the analysis using: Ridom SeqSphere+ (6), GrapeTree (5) and BioNumerics (2). Thirteen NRLs performed a phylogenetic analysis using: NDtree (2), BioNumerics (1), RAxML (2), Phylml (1), GrapeTree (2), IQTree (2), CGE tools (2) and mash dist (1).

For gene-by-gene based comparisons, cluster cut off values varied between 3 and 13. The NRLs using cut-off values below 10; 3, 4, 5 or 7 allelic differences, used either the Oxford PubMLST scheme or the INNUENDO wgMLST scheme. Where a cut-off value had been defined for SNP analysis, it varied from 4 to 25, with four NRLs using a cut-off value of 10.

The distances of selected PT samples to their closest neighbours were compared between the data submitted by NRLs for sample PT33-5 (most distantly related to all other samples), sample PT33-3 (related to other samples but not part of a cluster) and sample PT33-7 (part of a cluster with sample PT33-1 and PT33-6). The distances were taken from the MSTs or the distance matrices, depending on what was available. One NRL did not submit data that could be used to determine allele or SNP distances (L62). The distances and cut-off values for clusters reported by the NRLs are depicted in Figure 6. Some NRLs did not report a cut-off value.

Most NRLs divided the samples into the same cluster structure despite that different cut-off values were used. All NRLs placed PT33-1, PT33-6 and PT33-7 into the same cluster, and PT33-2 and PT33-4 were also placed together. However, there were three exceptions, and this was L51, L53 and L61, which included PT33-3 into the cluster with PT33-1, 6 and 7. One NRL (L61) placed all samples except PT33-5 into a single cluster.

Several participants mentioned that since all samples were closely related, except the more distant PT33-5, further epidemiological investigation would be needed to evaluate if they could be part of the same outbreak.

The distance measurements displayed most variability for sample PT33-5. Half of the NRLs performing SNP analysis reported distances up to ten-fold higher than was typically seen by cg/wgMLST results, whereas the other half reported distances more similar to the cg/wgMLST results. One NRL performing cgMLST analysis (L59) reported a unproportionally large distance to the PT33-5 sample compared to the other NRLs using cg/wgMLST. The reason for this was that the distances were based on mutations rather than alleles. Two NRLs determined there was one or two allele/SNP differences between the two identical DNA samples (PT33-1 and PT33-6). Overall, the same topology was seen in trees from all NRLs with the only difference that the closest neighbour to the distantly related sample PT33-5 differed for some SNP performing NRLs.

To summarise, the interpretations and division into clusters were the same for the majority of the NRLs, but deviations were present (3 NRLs) and sample PT33-3 was very close to the cluster cut-off value for some NRLs while deemed more distinctly unrelated for other NRLs. Thus, comparability in a multi-country outbreak situation could benefit from usage of a joint database/common analysis of data.

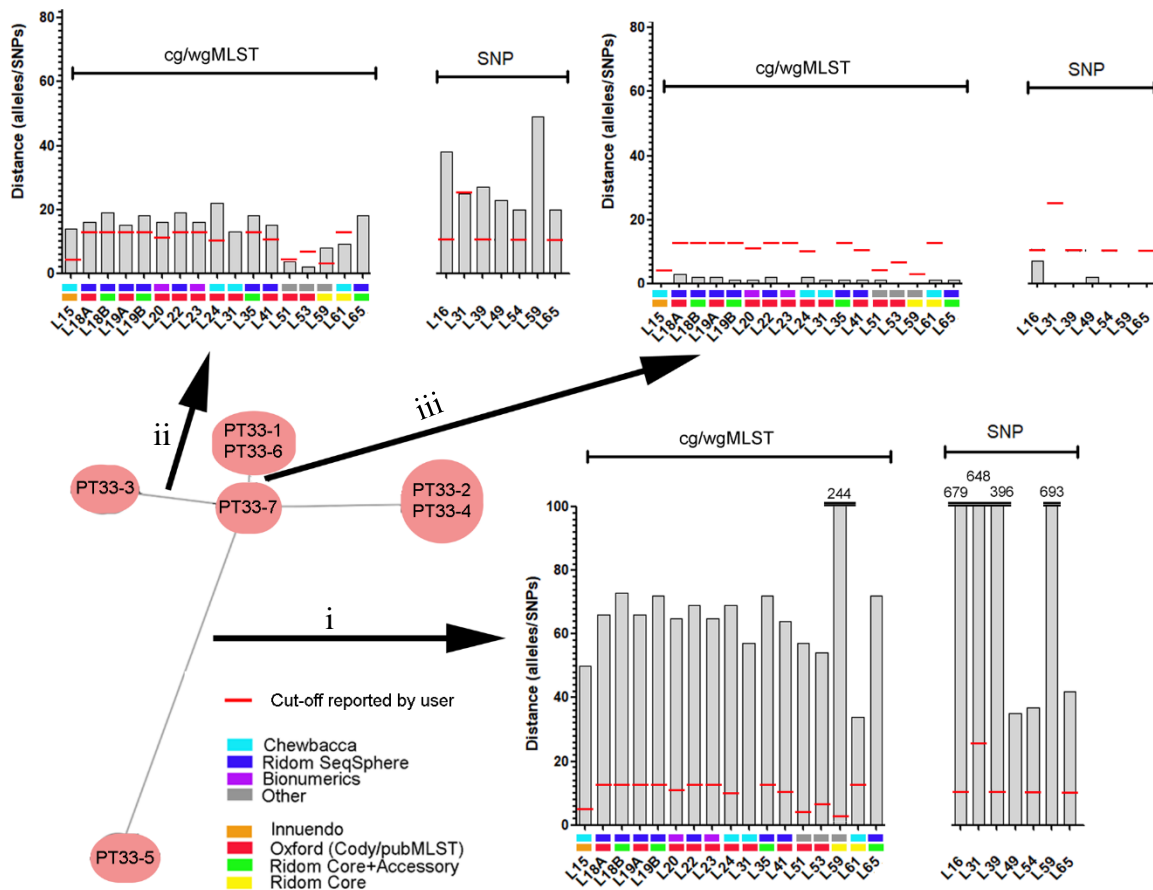


Figure 6. A summarisation of the clustering results. Results for three genetic distances reported by the NRLs are compared. These were distances between the closest neighbour and (i) the sample PT33-5 which represents a relatively distantly related isolate, (ii) the sample PT33-3 which represents what the EURL have judged as a close, but still distinct isolate and (iii) PT33-7 which is part of a cluster together with PT33-1 and PT33-6. The distances in alleles or SNPs are plotted and the cut-off value reported by each NRL is also indicated. Colour codes show which software and cg/wgMLST schema was used for the NRLs using cg/wgMLST.

Assessment of cluster analysis and NRL performance

The results using the defined criteria for assessment of cluster analysis of each NRL is summarised in Table 5. According to the assessment, all 20 NRLs fulfilled the criteria for satisfactory performance on all samples.

Table 5. Overview of assessment of the cluster analysis by each NRL in proficiency test No. 33 (2022).

Lab code	PT33-6 and PT33-7 are the two closest samples to PT33-1	PT33-4 is the closest sample to PT33-2	PT33-5 is most distant to other samples	Overall evaluation sequence quality
15	+	+	+	Satisfactory
16	+	+	+	Satisfactory
18	+	+	+	Satisfactory
19	+	+	+	Satisfactory
20	+	+	+	Satisfactory
22	+	+	+	Satisfactory
23	+	+	+	Satisfactory
24	+	+	+	Satisfactory
31	+	+	+	Satisfactory
35	+	+	+	Satisfactory
39	+	+	+	Satisfactory
41	+	+	+	Satisfactory
49	+	+	+	Satisfactory
51	+	+	+	Satisfactory
53	+	+	+	Satisfactory
54	+	+	+	Satisfactory
59	+	+	+	Satisfactory
61	+	+	+	Satisfactory
62	+	+	+	Satisfactory
65	+	+	+	Satisfactory

References

- [1] “ISO 23418:2022 - Microbiology of the food chain — Whole genome sequencing for typing and genomic characterization of bacteria — General requirements and guidance.”
- [2] A. J. Cody, J. E. Bray, K. A. Jolley, N. D. McCarthy, and M. C. J. Maiden, “Core Genome Multilocus Sequence Typing Scheme for Stable, Comparative Analyses of *Campylobacter jejuni* and *C. coli* Human Disease Isolates,” *Journal of Clinical Microbiology*, vol. 55, no. 7, pp. 2086–2097, Jul. 2017, doi: 10.1128/JCM.00080-17.
- [3] T. Seemann, “Snippy.” Sep. 04, 2022. Accessed: Sep. 18, 2022. [Online]. Available: <https://github.com/tseemann/snippy>
- [4] R. R. Wick *et al.*, “Trycycler: consensus long-read assemblies for bacterial genomes,” *Genome Biology*, vol. 22, no. 1, p. 266, Sep. 2021, doi: 10.1186/s13059-021-02483-z.
- [5] R. R. Wick, L. M. Judd, C. L. Gorrie, and K. E. Holt, “Unicycler: Resolving bacterial genome assemblies from short and long sequencing reads,” *PLOS Computational Biology*, vol. 13, no. 6, p. e1005595, Jun. 2017, doi: 10.1371/journal.pcbi.1005595.
- [6] B. Segerman, “tngs.” Jun. 01, 2022. Accessed: Sep. 18, 2022. [Online]. Available: <https://github.com/BoSegerman/tngs>
- [7] D. E. Wood, J. Lu, and B. Langmead, “Improved metagenomic analysis with Kraken 2,” *Genome Biology*, vol. 20, no. 1, p. 257, Nov. 2019, doi: 10.1186/s13059-019-1891-0.
- [8] M. Rossi *et al.*, “INNUENDO whole genome and core genome MLST schemas and datasets for *Campylobacter jejuni*.” Zenodo, Jul. 27, 2018. doi: 10.5281/zenodo.1322564.

Appendix A – QC metrics for submitted raw data

Sample ID	Q30 (bp)	Contamination (%)	Reference coverage (%)	GC deviation (%)
PT33-1-L15	93.38	0.04	100.0	0.06
PT33-1-L16	79.65	0.05	99.7	3.95
PT33-1-L18	92.08	0.05	99.7	3.08
PT33-1-L19	96.74	0.07	100.0	0.26
PT33-1-L20	95.99	0.06	100.0	0.28
PT33-1-L22	79.53	0.06	100.0	0.48
PT33-1-L23	87.85	0.05	100.0	0.02
PT33-1-L24	85.00	0.45	99.9	2.67
PT33-1-L31	90.45	0.07	100.0	1.15
PT33-1-L35	87.69	0.01	100.0	-0.09
PT33-1-L39	96.36	0.04	100.0	2.16
PT33-1-L41	95.62	0.02	100.0	0.18
PT33-1-L49	89.93	0.19	100.0	-0.14
PT33-1-L51	91.15	0	100.0	0.28
PT33-1-L53	88.24	0.25	99.9	2.91
PT33-1-L54	96.27	0.03	99.8	4.93
PT33-1-L59	90.34	0.03	100.0	-0.01
PT33-1-L61	83.52	0.4	100.0	0.79
PT33-1-L62	96.03	0.01	100.0	-0.17
PT33-1-L65	94.82	0.11	100.0	0.99
Median value	90.80	0.05	99.99904	
PT33-2-L15	94.61	0.03	100.0	-0.01
PT33-2-L16	80.00	0.05	99.9	2.73
PT33-2-L18	91.81	0.06	99.7	2.65
PT33-2-L19	96.78	0.09	100.0	0.21
PT33-2-L20	95.84	0.06	100.0	0.16
PT33-2-L22	82.06	0.04	100.0	0.46
PT33-2-L23	87.95	0.04	100.0	0.01
PT33-2-L24	85.45	0.47	100.0	2.38
PT33-2-L31	91.21	0.14	100.0	0.45
PT33-2-L35	87.92	0.01	100.0	-0.28
PT33-2-L39	96.63	0.03	100.0	2.20
PT33-2-L41	95.28	0.03	100.0	0.08
PT33-2-L49	89.13	0.2	100.0	-0.29
PT33-2-L51	90.75	0.01	100.0	0.19
PT33-2-L53	87.48	0.31	100.0	2.14
PT33-2-L54	96.16	0.08	99.7	6.33
PT33-2-L59	91.14	0.01	100.0	-0.08
PT33-2-L61	84.91	0.48	100.0	1.37
PT33-2-L62	98.05	0.03	100.0	-0.26
PT33-2-L65	90.35	0.03	100.0	0.02
Median value	90.94	0.045	99.99795	
PT33-3-L15	92.46	0.04	100.0	0.00
PT33-3-L16	80.20	0.04	99.8	3.56
PT33-3-L18	92.29	0.05	99.9	2.46
PT33-3-L19	95.67	0.08	100.0	0.20
PT33-3-L20	95.58	0.05	100.0	0.20
PT33-3-L22	80.39	0.07	100.0	0.45
PT33-3-L23	90.95	0.07	100.0	0.25
PT33-3-L24	84.90	0.69	100.0	2.35
PT33-3-L31	90.52	0.08	100.0	0.39
PT33-3-L35	88.20	0.01	100.0	-0.03
PT33-3-L39	97.05	0.02	100.0	2.27
PT33-3-L41	95.64	0.02	100.0	0.10
PT33-3-L49	90.63	0.15	100.0	-0.29
PT33-3-L51	91.03	0.01	100.0	-0.35

PT33-3-L53	88.04	0.07	100.0	2.35
PT33-3-L54	95.85	0.02	99.7	5.96
PT33-3-L59	95.45	0.02	100.0	-0.06
PT33-3-L61	81.59	0.63	100.0	0.72
PT33-3-L62	97.73	0.02	100.0	-0.37
PT33-3-L65	91.44	0.03	100.0	0.03
Median value	91.23	0.045	99.99895	
PT33-4-L15	94.81	0.3	100.0	0.32
PT33-4-L16	79.08	0.51	99.9	3.23
PT33-4-L18	90.15	0.51	99.7	3.26
PT33-4-L19	97.00	0.25	100.0	0.38
PT33-4-L20	94.81	0.38	100.0	0.39
PT33-4-L22	79.50	0.26	100.0	0.58
PT33-4-L23	90.63	0.28	100.0	0.43
PT33-4-L24	82.95	1.07	99.9	2.55
PT33-4-L31	90.56	0.88	99.9	1.57
PT33-4-L35	87.32	0.2	100.0	0.03
PT33-4-L39	96.25	0.51	100.0	2.72
PT33-4-L41	95.44	0.29	100.0	0.21
PT33-4-L49	89.50	0.25	100.0	-0.12
PT33-4-L51	91.72	2.26	100.0	1.35
PT33-4-L53	88.31	0.34	100.0	2.40
PT33-4-L54	95.93	14.67	99.6	11.98
PT33-4-L59	89.70	0.2	100.0	-0.03
PT33-4-L61	87.16	0.38	100.0	1.34
PT33-4-L62	97.12	0.23	100.0	0.00
PT33-4-L65	91.64	0.25	100.0	0.15
Median value	90.60	0.32	99.997348	
PT33-5-L15	92.96	0.03	100.0	0.02
PT33-5-L16	79.32	0.03	99.9	3.13
PT33-5-L18	93.19	0.06	99.8	2.59
PT33-5-L19	97.27	0.05	100.0	0.21
PT33-5-L20	96.22	0.03	100.0	0.14
PT33-5-L22	80.82	0.04	100.0	0.38
PT33-5-L23	90.60	0.06	100.0	0.28
PT33-5-L24	85.75	0.8	99.9	2.59
PT33-5-L31	91.81	0.07	100.0	0.34
PT33-5-L35	87.08	0.02	100.0	-0.14
PT33-5-L39	96.80	0.03	100.0	2.08
PT33-5-L41	95.42	0.01	100.0	0.02
PT33-5-L49	90.13	0.15	100.0	-0.27
PT33-5-L51	90.83	0	100.0	0.42
PT33-5-L53	86.59	0.13	100.0	2.20
PT33-5-L54	96.30	0.03	99.6	6.24
PT33-5-L59	92.94	0.02	100.0	-0.10
PT33-5-L61	87.50	0.1	100.0	0.69
PT33-5-L62	97.89	0.01	100.0	-0.57
PT33-5-L65	93.44	0.05	100.0	0.33
Median value	92.37	0.035	99.99955	
PT33-6-L15	94.29	0.04	100.0	0.42
PT33-6-L16	79.39	0.05	99.8	3.60
PT33-6-L18	90.53	0.05	99.8	2.45
PT33-6-L19	97.02	0.07	100.0	0.30
PT33-6-L20	95.38	0.05	100.0	0.15
PT33-6-L22	82.75	0.05	100.0	0.44
PT33-6-L23	90.20	0.05	100.0	0.30
PT33-6-L24	83.69	0.33	99.9	2.53
PT33-6-L31	90.69	0.1	99.9	1.28
PT33-6-L35	88.23	0.01	100.0	-0.02
PT33-6-L39	96.52	0.03	100.0	2.36

PT33-6-L41	95.43	0.03	100.0	0.18
PT33-6-L49	89.76	0.2	100.0	-0.14
PT33-6-L51	91.17	0.01	100.0	0.05
PT33-6-L53	88.70	0.14	99.9	2.58
PT33-6-L54	96.46	0.04	99.9	4.42
PT33-6-L59	95.85	0.01	100.0	-0.02
PT33-6-L61	82.90	0.57	100.0	0.79
PT33-6-L62	96.04	0	100.0	-0.05
PT33-6-L65	92.05	0.08	100.0	0.17
Median value	90.93	0.05	99.9889	
PT33-7-L15	92.11	0.03	100.0	0.03
PT33-7-L16	79.44	0.07	99.6	4.21
PT33-7-L18	89.82	0.06	99.9	2.55
PT33-7-L19	97.23	0.09	100.0	0.57
PT33-7-L20	95.99	0.05	100.0	0.20
PT33-7-L22	81.67	0.05	100.0	0.38
PT33-7-L23	88.78	0.04	100.0	0.07
PT33-7-L24	83.47	0.4	100.0	2.31
PT33-7-L31	91.31	0.05	100.0	-0.02
PT33-7-L35	87.84	0.01	100.0	-0.08
PT33-7-L39	96.75	0.02	100.0	2.46
PT33-7-L41	94.45	0.02	100.0	-0.03
PT33-7-L49	91.05	0.05	100.0	-0.12
PT33-7-L51	91.03	0	100.0	0.42
PT33-7-L53	88.56	0.25	100.0	2.06
PT33-7-L54	96.55	0.03	99.9	4.61
PT33-7-L59	94.83	0.01	100.0	-0.06
PT33-7-L61	84.79	0.18	100.0	0.76
PT33-7-L62	96.81	0.01	100.0	-0.13
PT33-7-L65	91.26	0.03	100.0	-0.01
Median value	91.16	0.045	99.99983	

Appendix B – QC metrics for submitted assemblies

Sample ID	Total size of assembly (bp)	k-mer coverage over reference genome (%)	Total number of contigs	Total number of contigs > 1kb	Longest contig (bp)	N50 length (bp)
PT33-1-L15	1,695,784	99.999	20	15	440 801	154 508
PT33-1-L18	1,692,848	99.845	53	30	232 124	96 853
PT33-1-L19	1,695,063	99.997	37	15	440 729	154 472
PT33-1-L20	1,696,633	99.994	2	2	1 695 600	1 695 600
PT33-1-L22	1,696,675	99.996	27	16	631 792	154 047
PT33-1-L23	1,694,663	99.956	25	16	456 064	154 047
PT33-1-L24	1,706,687	99.991	52	18	265 330	159 972
PT33-1-L31	1,696,943	99.999	30	17	407 642	154 508
PT33-1-L35	1,694,550	99.985	23	13	635 255	154 458
PT33-1-L41	1,695,252	99.987	36	19	351 055	216 952
PT33-1-L49	1,686,403	99.698	28	16	631 273	153 822
PT33-1-L51	1,694,175	99.965	26	16	631 503	154 015
PT33-1-L53	1,690,901	99.771	21	21	283 045	175 515
PT33-1-L54	1,691,396	99.908	37	30	379 586	153 979
PT33-1-L59	1,695,089	99.986	34	15	632 656	153 947
PT33-1-L61	1,712,826	99.999	65	21	201 903	154 047
PT33-1-L65	1,697,186	99.972	41	21	431 828	121 505
Median value	1,695,089	99.986	30	16	440 729	154 047
PT33-2-L15	1,736,665	99.991	30	22	404,359	154,127
PT33-2-L18	1,736,750	99.851	55	32	273,880	107,896
PT33-2-L19	1,738,029	99.995	38	20	284,315	121,483
PT33-2-L20	1,739,226	99.991	2	2	1,738,088	1,738,088
PT33-2-L22	1,738,702	99.993	31	17	435,470	154,127
PT33-2-L23	1,737,007	99.973	27	18	260,323	154,127
PT33-2-L24	1,749,183	99.993	63	22	260,323	154,127
PT33-2-L31	1,739,797	99.994	37	20	404,359	154,127
PT33-2-L35	1,735,563	99.986	25	17	656,174	153,948
PT33-2-L41	1,735,674	99.968	35	21	345,584	153,934
PT33-2-L49	1,727,887	99.671	28	18	596,688	153,873
PT33-2-L51	1,736,480	99.973	28	16	596,945	154,095
PT33-2-L53	1,730,415	99.610	22	22	227,165	121,432
PT33-2-L54	1,732,859	99.813	42	30	201,109	79,979
PT33-2-L59	1,738,002	99.983	51	18	649,010	153,948
PT33-2-L61	1,742,457	99.958	39	16	674,163	154,127
PT33-2-L65	1,737,762	99.976	31	17	421,415	154,127
Median value	1,737,007	99.976	31	18	404,359	154,095
PT33-3-L15	1,657,147	99.984	36	25	371,435	120,512
PT33-3-L18	1,659,993	99.975	54	23	343,149	154,010
PT33-3-L19	1,658,763	99.995	60	26	371,363	121,430
PT33-3-L20	1,660,072	99.992	6	2	1,257,493	1,257,493
PT33-3-L22	1,659,128	99.993	42	21	371,435	127,924
PT33-3-L23	1,656,540	99.961	26	18	593,589	154,049
PT33-3-L24	1,665,892	99.978	68	24	194,360	121,533
PT33-3-L31	1,661,465	99.991	57	22	561,543	154,049
PT33-3-L35	1,655,007	99.978	37	17	570,290	153,949
PT33-3-L41	1,656,334	99.975	30	20	312,710	189,644
PT33-3-L49	1,647,726	99.661	28	17	593,327	153,795
PT33-3-L51	1,656,249	99.964	28	16	593,557	154,017
PT33-3-L53	1,645,926	99.363	24	24	195,964	127,923
PT33-3-L54	1,651,692	99.879	22	21	367,432	153,981
PT33-3-L59	1,655,952	99.981	45	18	575,824	153,949
PT33-3-L61	1,663,657	99.993	48	22	227,403	120,513
PT33-3-L65	1,657,147	99.970	29	18	418,119	154,049
Median value	1,657,147	99.978	36	21	371,435	153,949

PT33-4-L15	1,737,296	99.992	22	16	596,978	154,127
PT33-4-L18	1,847,256	99.820	243	38	292,540	87,973
PT33-4-L19	1,737,461	99.996	35	18	421,393	154,055
PT33-4-L20	1,739,255	99.993	3	1	1,738,144	1,738,144
PT33-4-L22	1,753,644	99.995	61	19	545,080	154,127
PT33-4-L23	1,737,016	99.976	28	18	545,080	154,127
PT33-4-L24	1,785,750	99.977	126	25	260,323	106,649
PT33-4-L31	2,000,956	99.995	512	28	435,470	154,127
PT33-4-L35	1,927,703	99.977	655	21	544,776	121,253
PT33-4-L41	1,736,886	99.981	35	22	480,944	189,429
PT33-4-L49	1,727,906	99.674	29	19	596,689	153,873
PT33-4-L51	2,026,082	99.974	137	93	596,946	121,400
PT33-4-L53	1,731,749	99.714	17	17	227,165	175,515
PT33-4-L54	3,899,870	99.759	102	80	414,161	114,420
PT33-4-L59	1,736,860	99.984	36	16	673,603	153,948
PT33-4-L61	1,816,137	99.995	170	18	596,978	154,127
PT33-4-L65	1,803,357	99.979	173	18	545,080	121,432
Median value	1,753,644	99.979	61	19	545,080	154,055
PT33-5-L15	1,743,883	99.990	27	17	441,117	154,127
PT33-5-L18	1,745,235	99.887	56	29	194,235	99,355
PT33-5-L19	1,744,072	99.997	40	19	441,369	154,055
PT33-5-L20	1,745,233	99.995	3	1	1,744,726	1,744,726
PT33-5-L22	1,745,125	99.997	33	15	681,361	154,127
PT33-5-L23	1,744,443	99.974	27	16	681,104	154,127
PT33-5-L24	1,753,755	99.996	54	17	227,398	154,127
PT33-5-L31	1,745,073	99.997	30	14	681,104	154,127
PT33-5-L35	1,742,013	99.983	30	16	655,809	154,027
PT33-5-L41	1,744,415	99.988	34	20	488,315	189,528
PT33-5-L49	1,736,015	99.715	28	17	680,842	189,523
PT33-5-L51	1,744,033	99.973	26	15	681,072	154,095
PT33-5-L53	1,741,613	99.893	25	25	269,925	120,427
PT33-5-L54	1,740,242	99.871	35	28	188,052	106,567
PT33-5-L59	1,743,856	99.980	46	17	680,578	153,948
PT33-5-L61	1,749,419	99.997	41	14	681,104	154,127
PT33-5-L65	1,745,026	99.977	32	17	656,792	154,127
Median value	1,744,415	99.983	32	17	656,792	154,127
PT33-6-L15	1,694,157	99.951	27	16	631,798	154,508
PT33-6-L18	1,697,082	99.978	40	21	350,053	107,979
PT33-6-L19	1,694,273	99.996	28	17	440,729	153,975
PT33-6-L20	1,696,332	99.994	2	1	1,696,087	1,696,087
PT33-6-L22	1,695,924	99.996	29	19	631,792	154,047
PT33-6-L23	1,694,726	99.957	27	17	440,801	154,047
PT33-6-L24	1,700,896	99.995	38	20	265,330	121,966
PT33-6-L31	1,696,143	99.996	30	17	407,642	154,508
PT33-6-L35	1,694,061	99.986	21	13	635,255	154,458
PT33-6-L41	1,693,973	99.973	33	21	543,003	189,641
PT33-6-L49	1,686,414	99.698	28	16	631,273	153,823
PT33-6-L51	1,694,163	99.965	26	16	631,503	154,015
PT33-6-L53	1,691,588	99.826	23	23	232,592	175,515
PT33-6-L54	1,694,576	99.877	32	24	300,444	150,125
PT33-6-L59	1,695,318	99.987	36	15	632,656	153,947
PT33-6-L61	1,710,137	99.997	65	25	235,336	111,789
PT33-6-L65	1,695,812	99.970	32	17	631,535	154,047
Median value	1,694,726	99.978	29	17	543,003	154,047
PT33-7-L15	1,694,982	99.993	25	19	440,477	154,047
PT33-7-L18	1,697,411	99.906	42	24	274,102	175,515
PT33-7-L19	1,694,927	99.998	47	18	431,213	153,963
PT33-7-L20	1,695,845	99.995	1	1	1,695,845	
PT33-7-L22	1,695,447	99.997	29	17	631,792	154,047
PT33-7-L23	1,694,842	99.961	26	18	265,330	154,047

PT33-7-L24	1,698,209	99.996	32	19	265,006	154,047
PT33-7-L31	1,695,823	99.998	22	14	642,857	154,047
PT33-7-L35	1,693,880	99.988	21	14	720,947	189,647
PT33-7-L41	1,695,037	99.992	44	22	561,458	153,933
PT33-7-L49	1,686,417	99.699	28	16	631,273	153,822
PT33-7-L51	1,694,178	99.967	26	16	631,503	154,015
PT33-7-L53	1,692,382	99.833	21	21	411,358	154,047
PT33-7-L54	1,691,519	99.919	26	19	407,017	153,968
PT33-7-L59	1,695,157	99.988	34	15	632,655	153,947
PT33-7-L61	1,699,088	100.000	46	22	440,477	154,047
PT33-7-L65	1,694,870	99.971	27	17	631,535	154,047
Median value	1,694,982	99.988	27	18	561,458	154,047