

1. **PURPOSE:** All genomic sequences generated by PulseNet participating laboratories are uploaded in real-time to the sequence read archive (SRA) located at NCBI and the analyzed data are uploaded to the organism-specific surveillance databases located on servers at CDC. In order to ensure the integrity and comparability of the data across the network, PulseNet has set minimum quality requirements for sequences to be uploaded to the national databases and SRA. The purpose of this document is to describe a standardized procedure for Illumina sequence data quality control (QC) prior to upload to the national database and SRA.
2. **SCOPE:** This procedure applies to all whole genome sequence data generated by PulseNet participating laboratories.
3. **DEFINITIONS/ACRONYMS:**
  - 3.1. **Allele:** One of two or more alternative forms of a gene that arise by mutation and are found at the same place on a chromosome.
  - 3.2. **Analysis Certified:** An individual who is certified for checking the quality, performing analysis, and uploading WGS data to the PulseNet National Database and NCBI using BioNumerics.
  - 3.3. **ANI:** Average Nucleotide Intity.
  - 3.4. **BaseSpace:** Illumina cloud-based computing environment for next generation sequencing data analysis, management, and storage, including data sharing.
  - 3.5. **BioNumerics:** Analysis software used by PulseNet, developed by Applied Maths (Sint-Martens-Latem, Belgium).
  - 3.6. **Bp:** Base Pair.
  - 3.7. **CDC:** Centers for Disease Control and Prevention.
  - 3.8. **Core genome:** Genes shared by all strains of the same species.
  - 3.9. **Coverage:** The average number of reads that include a given nucleotide in the reconstructed sequence.
  - 3.10. **Critical Quality Metrics:** Average denovo coverage, average quality (Q score), assembly length, secondary species abundance (contamination detected by MIDAS) and percent core present. Failing to meet the minimum thresholds/acceptable range for any one of these metrics will result in rejection of the sequence.
  - 3.11. **De Novo Assembly:** A sequence assembly generated from the short raw reads without the use of a reference genome.
  - 3.12. **FASTQ:** A text-based file format for storing both sequence and its corresponding quality scores.
  - 3.13. **FastQC:** Raw sequence quality control tool available free-of-charge at [www.bioinformatics.babraham.ac.uk/projects/fastqc](http://www.bioinformatics.babraham.ac.uk/projects/fastqc).
  - 3.14. **Insert:** The stretch of sequence in fragmented target DNA that is between the paired-end adapters. Includes the sequenced portion (reads), as well as the un-sequenced portion. Insert size is determined bioinformatically by mapping the reads back to the reference sequence.
  - 3.15. **MIDAS:** Metagenomic Intra-species Diversity Analysis System. An integrated computational pipeline for quantifying bacterial species abundance.
  - 3.16. **NCBI:** National Center for Bioinformatics, part of the National Institutes of Health (NIH). NCBI houses several databases relevant to biotechnology, including

GenBank for DNA sequence assemblies and Sequence Read Archive (SRA) for raw reads.

- 3.17. Organism-specific Database:** A BioNumerics database, v 7.6 or higher, used for comparing isolates for surveillance. Part of the standard PulseNet workflow.
- 3.18. PF Reads:** Passing Filter Reads, the number of reads which passed filtering (are useable reads) for a sequencing run. % Reads Identified (PF) represents the percentage of the PF reads that have been assigned to an index pair. This is displayed for the entire run (% of reads assigned to an index pair) as well as for a particular set of indices and will vary for each index pair.
- 3.19. PHL:** Public Health Laboratory
- 3.20. PN:** PulseNet
- 3.21. PulseNet Central:** PulseNet team at CDC comprising of the Database Unit ([PulseNet@cdc.gov](mailto:PulseNet@cdc.gov)) and the Next Generation Subtyping Methods Unit ([PulseNetNGSlab@cdc.gov](mailto:PulseNetNGSlab@cdc.gov)).
- 3.22. QC:** Quality Control
- 3.23. Q score:** The sequence quality score for each individual base position in a sequence, indicating the accuracy of the base call. Phred scores are used, where  $Q = -10\log(\text{Error Probability})$ . The higher the quality score, the more reliable the base call. A Q30 means a 1 in 1000 likelihood of an incorrect base call at that position.
- 3.24. Read:** A unit of continuous DNA sequence (based pairs) derived by sequencing a part of the fragmented target DNA.
- 3.25. RefID Database:** A BioNumerics database, v 7.6 or higher, used for quality control of raw sequence data, assembly of sequences, contamination detection, and species identification. Part of the standard PulseNet workflow.
- 3.26. SAV:** Sequencing Analysis Viewer, an application software that allows real-time viewing of quality metrics generated by the real-time analysis (RTA) software on the Illumina sequencing systems.
- 3.27. Size selection:** A step in the sequencing library preparation which involves the targeted capture of DNA fragments of a specific size or a size range.
- 3.28. SOP:** Standard Operating Procedure.
- 3.29. SRA:** Sequencing Read Archive, database at NCBI which stores raw sequence data and alignment information.
- 3.30. Tagmentation:** The first step of the Illumina Nextera sequencing library preparation where the target DNA is enzymatically cleaved and tagged with universal adapters at both ends of the fragment.
- 3.31. WGS:** Whole Genome Sequencing
- 3.32. wgMLST:** Whole Genome Multi-Locus Sequencing Typing

#### 4. RESPONSIBILITIES:

##### 4.1. Certified PulseNet public health laboratory personnel

4.1.1. Sequence isolates and perform quality check of the sequence data.

4.1.1.1. BioNumerics analysis-certified personnel:

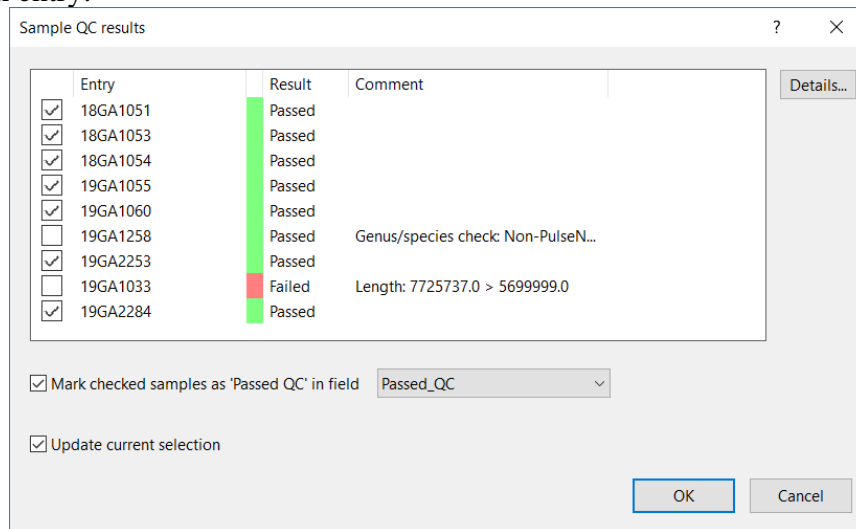
4.1.1.1.1. Use the BioNumerics RefID database workflow to determine average denovo coverage, average read quality, sequence assembly length, secondary species abundance (contamination), species ID, and average read length.

- 4.1.1.1.2. Use the BioNumerics organism-specific databases to obtain percentage (%) values of core genome present.
- 4.1.1.2. Personnel not BioNumerics analysis-certified:
  - 4.1.1.2.1. Coverage can be **estimated** in semi-automated manner using SAV/BaseSpace and the metrics tab in Nextera XT (PNL34.W1) or Nextera DNA Flex (PNL35.W1&2) workbooks or manually using values derived from SAV/BaseSpace or FastQC.
  - 4.1.1.2.2. Quality and average read length can be **roughly** evaluated using FastQC.
- 4.1.2. Re-sequence any isolates which do not meet the minimum thresholds/acceptable ranges for critical quality metrics.
- 4.1.3. Communicate any instrument or sequencing issue with PulseNet Central, as necessary.
- 4.2. **PulseNet Central:**
  - 4.2.1. Perform additional sequence data quality analysis.
  - 4.2.2. Notify public health laboratory if any sequences do not meet quality thresholds.
  - 4.2.3. Assist public health laboratories with troubleshooting, as necessary.

**5. PROCEDURE:**

**5.1. Sequence quality assessment using BioNumerics 7.6 or higher RefID database workflow**

- 5.1.1. Process the fastq-files in the RefID database following the workflow outlined in the SOP PND20 (PulseNet Standard operating procedure for the BioNumerics Reference Identification database).
- 5.1.2. Review sequence quality:
  - 5.1.2.1. The “Sample QC results” window lists the pass/fail results for the critical quality metrics (except % core present) and comments for each selected entry. Select an entry and click on “Details” to see specific metrics that pass/fail for each entry.



- 5.1.2.1.1. If the sequence is flagged green and is a PulseNet organism, the sequence passes quality, remains checked, and the comment is blank.

5.1.2.1.2. If the sequence is flagged green and is not a PulseNet organism, the sequence passes quality, becomes unchecked, and the comment indicates “Genus/species check: Non-PulseNet organism”.

5.1.2.1.2.1. Click on “Details” to pull up a “Detailed QC results” window for all sequences that have comments listed. In the “Detailed QC results” window, the Genus/species check is flagged yellow if it was identified by ANI (e.g. non-PulseNet *Vibrio* or *Campylobacter* species) or red if it was not identified by ANI (e.g. *Morganella*).

Metric	Result	Comment
Genus/species check	Warning	Non-PulseNet organism: <i>Vibrio furnissii</i>
Average quality	Passed	34.6 >= 30.0
Secondary species abundance	Passed	0.0 <= 1.0
Average denovo coverage	Passed	56.0 >= 20.0
Length	Passed	1400000.0 <= 5054375.0 <= 5999999.0

5.1.2.1.3. If the sequence is flagged red, the sequence fails quality, becomes unchecked, and the reason for failed QC is in the “Comment” column.

5.1.2.1.3.1. Failing to meet any of the critical quality metrics will result in the rejection of the sequence. Refer to appendix PNQ07-1 Table 3 for species-specific minimum thresholds/acceptable ranges.

Metric	Result	Comment
Genus/species check	Passed	<i>Salmonella enterica</i>
Average quality	Passed	35.0 >= 30.0
Secondary species abundance	Passed	0.2 <= 1.0
Average denovo coverage	Passed	90.4 >= 30.0
Length	Failed	7725737.0 > 5699999.0

5.1.2.2. Review average read length:

5.1.2.2.1. For 500 cycle chemistry, an average read length  $\geq 225$  bp is acceptable.

5.1.2.2.2. For 300 cycle chemistry, an average read length  $\geq 135$  bp is acceptable.

**NOTE1:** Average read length is a good proxy for estimating the insert size. Shorter than expected read lengths are an indication of sub-optimal library preparation indicating that DNA may be over-tagmented or size selection is not being performed correctly.

**NOTE2:** Average read length is **not** a critical quality metric for routine sequence submission, i.e. shorter than expected read length does not result in an automatic rejection of a sequence. However, shorter than expected read lengths may result in fragmented assemblies and a low percentage of core alleles detected which can result in sequence rejection. Therefore,

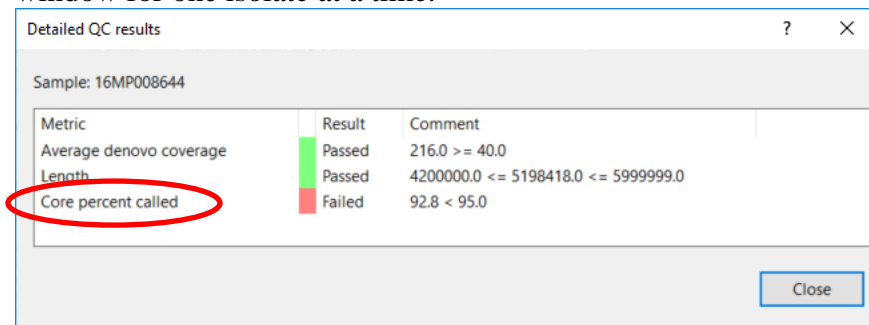
when shorter than expected read lengths are detected, the root cause should always be investigated and addressed.

**NOTE3:** Insert size (and hence the read length) is a critical quality metric for wetlab (fastq) certification only. Sequences with read lengths < 225 bp for 500 cycle chemistry and <135 bp for 300 cycle chemistry will not pass certification. These read lengths correspond to approximately 300 bp insert size which is the PulseNet minimum threshold for insert size. For proficiency tests, shorter than expected read lengths will not result in automatic rejection but will result in points deduction.

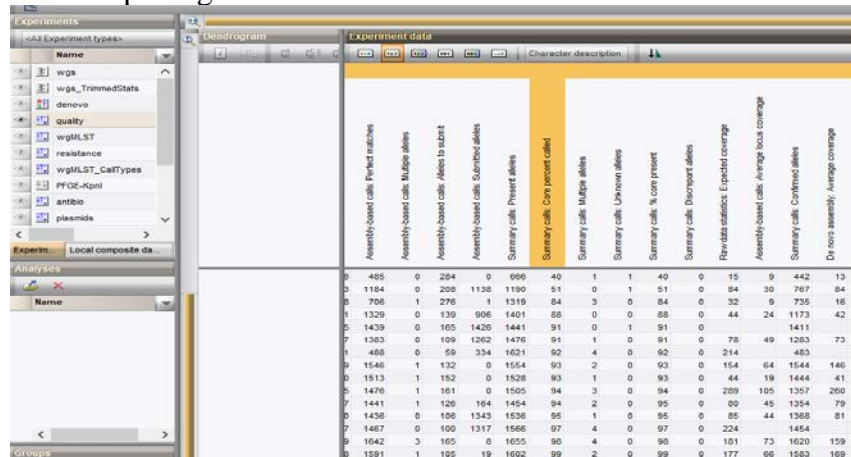
5.1.3. Transfer sequences that belong to PulseNet organisms and pass the critical quality metrics to organism-specific databases.

5.1.3.1. Follow the workflow (SOP under development) in the organism-specific database to obtain the quality metrics for the allele calls, including % of core present (a critical quality metric).

5.1.3.1.1. You can review the percent core present values by selecting the desired entries in the database and opening the “Detailed QC results” window for one isolate at a time.



5.1.3.1.2. Alternatively, you can create a comparison and review all QC metrics for multiple organisms at the same time.



5.1.4. Re-sequence PulseNet organisms and non-PulseNet organisms identified by ANI that did not pass critical quality metrics. For non-PulseNet organisms not identified by ANI, use another validated method (for example phenotypic identification or 16S sequencing) for identification.

5.2. **Estimating coverage without BioNumerics:** Coverage may be estimated manually by using metrics from FastQC (Section 5.2.1.), semi-automatically by using the PN

Workbook and values derived from BaseSpace/SAV (Section 5.2.2.) or manually using the values derived from BaseSpace/SAV (Section 5.2.3.). **NOTE:** This is a mathematical estimate of coverage – actual coverage values may differ from the calculations below due to shorter than expected read lengths, variation of genome size from the estimated value used for coverage determination, etc. True coverage values are obtained using BioNumerics.

5.2.1. **Manually, using FastQC:** Coverage may be calculated using the formula below and values derived from FastQC: **NOTE:** FastQC may be downloaded, free of charge, at: [www.bioinformatics.babraham.ac.uk/projects/fastqc](http://www.bioinformatics.babraham.ac.uk/projects/fastqc).

5.2.1.1. Open FastQC and choose “File” on the main screen.

5.2.1.2. Choose “Open” and then select one read file (.fastq file), either R1 or R2, for the sample to be analysed.

5.2.1.3. Using the **Total Sequences** value and the maximum number depicted for **Sequence Length**, calculate coverage using the formula:

$$(\text{Total Sequences} \times \text{Maximum Sequence Length} \times 2) / \text{Estimated genome length} = \text{Isolate coverage.}$$

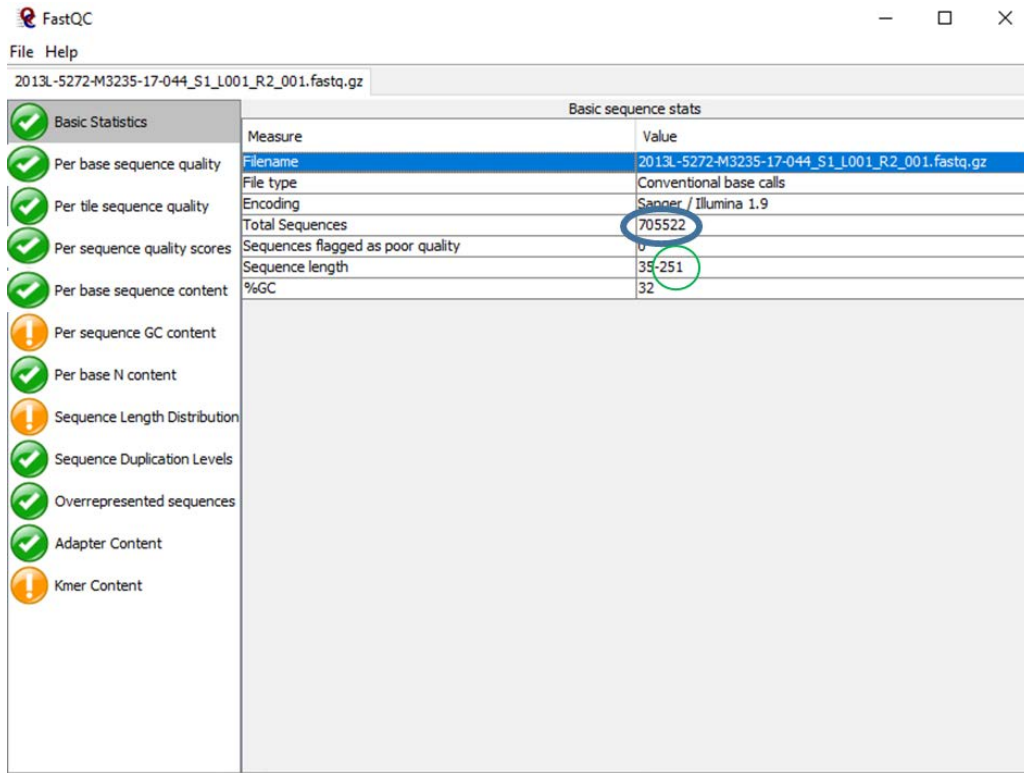
**NOTE:** multiplying by 2 is required for determining coverage for paired-end reads.

5.2.1.3.1. Refer to Table 1 for estimated genome sizes.

<b>Organism</b>	<b>Estimated Genome size (bp)</b>
<i>Listeria monocytogenes</i>	3000000
<i>E. coli/Shigella spp.</i>	5000000
<i>Salmonella spp.</i>	5000000
<i>Campylobacter spp.</i>	1600000
<i>Vibrio spp.</i>	5000000

*Table 1. Estimated genome sizes to be used in the coverage calculations*

5.2.1.4. Refer to Appendix PNQ07-1 Table 3 for minimum coverage thresholds



Example: Coverage calculation for this *Listeria* isolate:

$$(705522 \times 251 \times 2) / 3000000 = 118x$$

5.2.2. **Semi-automatically, using the PN Workbook:** Coverage may be calculated using the “Metrics tab” of the library prep workbooks and values derived from SAV/BaseSpace. **NOTE:** The following fields must be filled in correctly in the workbook in order for the coverage to be accurately calculated:

Workbook	Tab	Field
Nextera XT (PNL 34.W1)	Initial Dilution	Sample ID Genome Size Estimate
	Metrics	Sequencing Kit Type/Chemistry
DNA Flex (PNL 35.W1 & W2)	Library Prep	State Key
		Genome Size Estimate
		Sequencing Kit Type/Chemistry

Table 2. Fields that need to be correctly filled out in PulseNet sequence library workbooks for accurate coverage estimations

5.2.2.1. Select and copy the columns for Index Number, Sample Id, Project, Index 1 (i7), Index 2 (i5) and % Reads Identified (PF) columns **NOTE:** In SAV these values are found in the Indexing tab. On BaseSpace they are found in the Indexing QC tab.

TOTAL READS		PF READS	% READS IDENTIFIED (PF)	CV	MIN	MAX
55,734,818		49,492,984	96.2287	0.4653	0.7027	7.2220

INDEX	SAMPLE ID	PROJECT	INDEX 1 (R)	INDEX 2 (R)	% READS IDENTIFIED (PF)
1	1911500159-CP1b-AK-MD1348-191210	default	TAAGGCGA	GTAAAGAG	2.0937
2	1911500159-EJ1b-AK-MD1348-191210	default	CGTACTAG	GCCTAAGA	1.5678
3	1911500160-ET1b-AK-MD1348-191210	default	AGGCAGAA	ACTGCATA	0.7027
4	1529500176-CP1b-AK-MD1348-191210	default	GGACTCCT	TATCCTCT	3.4454
5	1529500176-EJ1b-AK-MD1348-191210	default	TAGGCATG	CCTAGAGT	3.0909
6	1529500179-CP1b-AK-MD1348-191210	default	CGAGGCTG	TTATGCCA	6.9630
7	1529500179-EJ1b-AK-MD1348-191210	default	AAGAGGCA	TCGACTAG	6.0649
8	1703200181-CP1b-AK-MD1348-191210	default	CGTCTAAT	CGTCTAAT	4.6603
9	1911500161-ET1b-AK-MD1348-191210	default	CGTACTAG	ACTGCATA	1.3619
10	1911500162-ET1b-AK-MD1348-191210	default	AGGCAGAA	GTAAAGAG	2.2979
11	1529500177-CP1b-AK-MD1348-191210	default	TAAGGCGA	GCCTAAGA	7.2220
12	1529500177-EJ1b-AK-MD1348-191210	default	CTCTCTAC	TATCCTCT	3.7932
13	1802400151-ET1b-AK-MD1348-191210	default	GTAGAGGA	CCTAGAGT	4.8883
14	1802400151-ET2b-AK-MD1348-191210	default	AAGAGGCA	CGTCTAAT	4.7928
15	1802400150-EJ1b-AK-MD1348-191210	default	CGTCTAAT	TTATGCCA	3.4456
16	1802400150-EJ1b-AK-MD1348-191210	default	CGAGGCTG	TCGACTAG	4.6690
17	1529500178-ET1b-AK-MD1348-191210	default	AGGCAGAA	GCCTAAGA	6.4289
18	1703200184-ET1b-AK-MD1348-191210	default	TAAGGCGA	ACTGCATA	3.1568
19	1703200186-ET1b-AK-MD1348-191210	default	CGTACTAG	GTAAAGAG	2.1457
20	1802400154-ET1b-AK-MD1348-191210	default	TAGGCATG	TATCCTCT	2.1761
21	1822700122-CP1b-AK-MD1348-191210	default	CGAGGCTG	CCTAGAGT	6.1842
22	1822700122-CP2b-AK-MD1348-191210	default	CGTCTAAT	TCGACTAG	4.5404
23	1822700122-EJ1b-AK-MD1348-191210	default	CGAGGCTG	CGTCTAAT	4.5935
24	1819300049-ET1b-AK-MD1348-191210	default	AAGAGGCA	TTATGCCA	5.9416

- 5.2.2.2. Paste these values into the corresponding columns on the “Metrics” tab of the workbook. **NOTE: Ensure** that the order of the isolates in the “Metrics” tab is the same as the order of the isolates on the previous workbook tabs.
- 5.2.2.3. Enter the “PF Reads” value for the run into column G of the workbook. The total number of reads (bp) for each isolate should now be displayed in Column H of the workbook and the estimated coverage either in Column I or J depending on which MiSeq software version is in use. **NOTE:** Depending on the software version in use, SAV/BaseSpace will either display the **total** number of reads that passed filter for the run, or only **half** of the total reads that passed filter. It is advised to first calculate coverage estimate for one isolate on a run using FastQC to determine which value of PF Reads is displayed. This will allow users to know which column of the “Metrics” tab may be used to automatically generate coverage estimates for the run.
- 5.2.2.4. See Appendix PNQ07-1 Table 3 for passing coverage requirements.



**NOTE2:** Whether or not the **PF Reads** value for the run, displayed in SAV or BaseSpace, is the **total** for the entire run, or only **half** of the total number of reads must be established in order to use the values for manual mathematical estimation of coverage. This may be determined by estimating coverage in FastQC first.

**Maximum read length** = ½ the number of cycles in the run. For a 500 cycle run, the Maximum read length = 250.

**Estimated genome size** = See Table 1

5.2.3.2. See Appendix PNQ07-1 Table 3 for passing coverage requirements.

Example:

Using the values from the SAV example above to determine coverage for 2013L-5272-M3235-17-044, and assuming that this was a 500 cycle run:

**PF Reads (for the run): 13822209**

**% Reads Identified (PF) for 2013L-5272-M3235-17-044: 5.1043 % = 0.051043**

**Maximum read length: ½ of 500 = 250**

**Estimated genome size (*Listeria*): 3000000**

$$((13822209 \times 0.051043) \times 250) / 3000000 = 58.79^{**}$$

**\*\*HOWEVER**, in this example the PF Reads displayed in SAV is only HALF of all of the reads on the run. This is confirmed by noting that  $(13822209 \times 0.051043 = 70552)$ , which corresponds with the Total Sequences value in FastQC for only ONE fastq file. Therefore, the estimated coverage value above (58.79) should be doubled, to find total coverage for the isolate in this instance (117.58). Thus, if the workbook were to be used for semi-automated calculations for this run, the second column (pink header) will yield the most accurate coverage estimate.

### 5.3. Assess sequence data quality and read length using FastQC

**NOTE:** The graphs generated by FastQC are open to subjective interpretation and will not provide actual numeric quality scores and read lengths (which are derived using BioNumerics or command line tools). The graphs depicted below have been correlated with numeric values derived from the CDC in-house command line QC tool

(<https://github.com/lskatz/CG-Pipeline>).

5.3.1. Open FastQC and select “File” from the toolbar on the main screen.

5.3.2. Choose “Open” and select a read file for analysis (.fastq file).

**NOTE1:** More than one sequence read file may be opened at a time.

**NOTE2:** It may be more helpful to analyze R2 data. Generally, R2 will have slightly decreased quality compared to R1. Therefore, if R2 passes initial quality assessment, it may be assumed that R1 will pass as well.

5.3.3. Assess the “Per base sequence quality” graph:

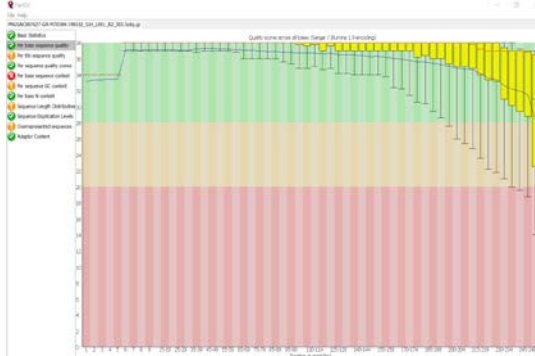
5.3.3.1. The length of the read (in bp) is along the x-axis and the quality score (Q score) is along the y-axis.

5.3.3.2. The yellow box plots indicate the 25th/75th inter-quartile (extremes of the boxes).

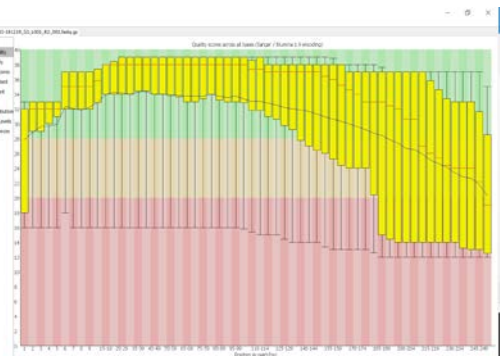
5.3.3.3. The whiskers indicate the 10th and 90th percentage points (ends of whiskers).

5.3.3.4. For this graph, the majority of the length of the reads (i.e. greater than half) should have a quality score  $\geq 30$ . Therefore, most of the yellow box plots should be within the green area of the graph.

R2 Q score 36.18



R2 Q score 30.69



R2 Q score 29.56



R2 Q score 27.08



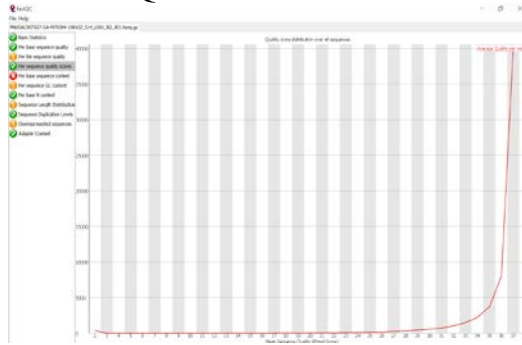
5.3.4. Assess the “Per sequence quality scores” graph:

5.3.4.1. This graph provides a view of the quality score (Q score) along the x-axis by number of reads, along the y-axis.

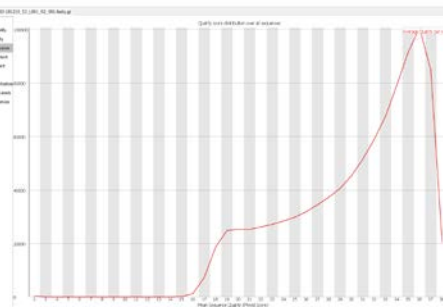
5.3.4.2. The peak should be sharp and  $> 30$ .

5.3.4.3. The peak of the graph being  $< 30$  or having a pronounced shoulder indicates the presence of low quality reads.

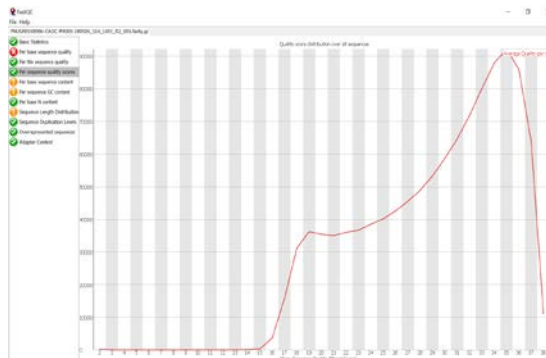
R2 Q score 36.18



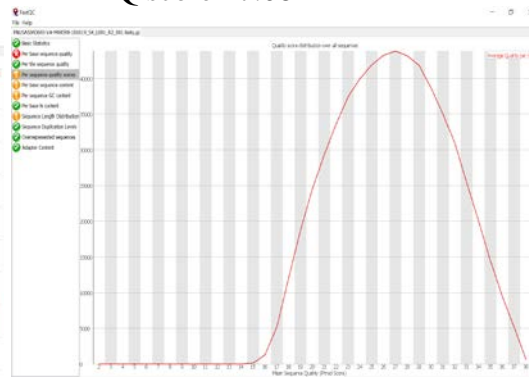
R2 Q score 30.69



R2 Q score 29.56



R2 Q score 27.08



5.3.5. Assess the “Sequence Length Distribution” graph: This graph depicts the number of sequences (y-axis) with varying lengths (in bp, along x-axis).

5.3.5.1. The graph should be flat until reaching the maximum read length, indicating that the reads are of sufficient lengths.

5.3.5.1.1. Acceptable average read length for 500 cycle chemistry is  $\geq 225$  bp.

5.3.5.1.2. Acceptable average read length for 300 cycle chemistry is  $\geq 135$  bp.

5.3.5.2. If the line is above baseline prior to the desired read length, this is indicative of short inserts in the library.

**NOTE1:** Shorter than expected read lengths are an indication of sub-optimal library preparation; indicating that DNA may be over-tagmented or size selection is not being performed correctly.

**NOTE2:** Average read length is not a critical quality metric for routine sequence submission, i.e. shorter than expected read length does not result in an automatic rejection of a sequence. However, shorter than expected read lengths may result in fragmented assemblies and a low percentage of core alleles detected which can result in sequence rejection. Therefore, when shorter than expected read lengths are detected, the root cause should always be investigated and addressed.

**NOTE3:** Inset size (and hence the read length) is a critical quality metric for wetlab (fastq) certification, only. Sequences with read lengths  $< 225$  bp for 500 cycle chemistry and  $< 135$  bp for 300 cycle chemistry will not pass

**PULSENET STANDARD OPERATING PROCEDURE FOR ILLUMINA DATA QUALITY CONTROL**

**Doc. No. PNQ07**

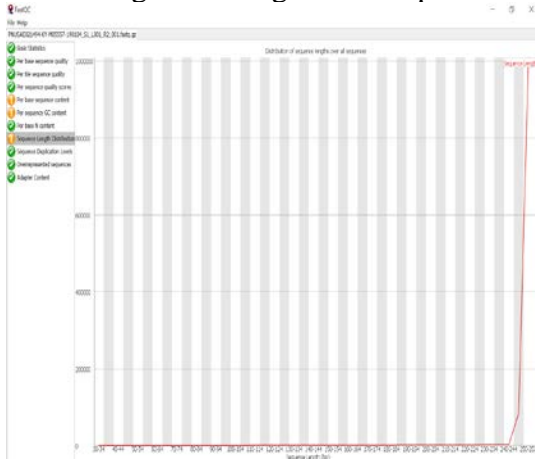
**Ver. No. 08**

**Effective Date:**

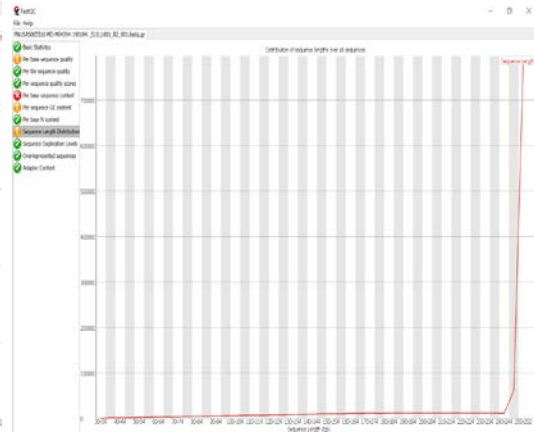
**Page 13 of 18**

certification. These read lengths correspond to approximately 300 bp insert size which is the PulseNet minimum threshold for insert size. For proficiency tests, shorter than expected read lengths will not result in automatic rejection but will result in points deduction.

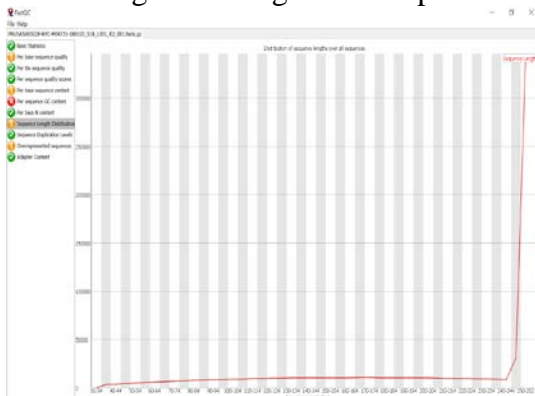
Average read length 246.0 bp



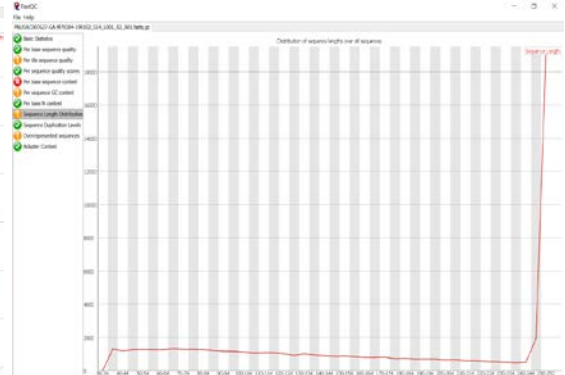
Average read length 225.2 bp



Average read length 200.9 bp

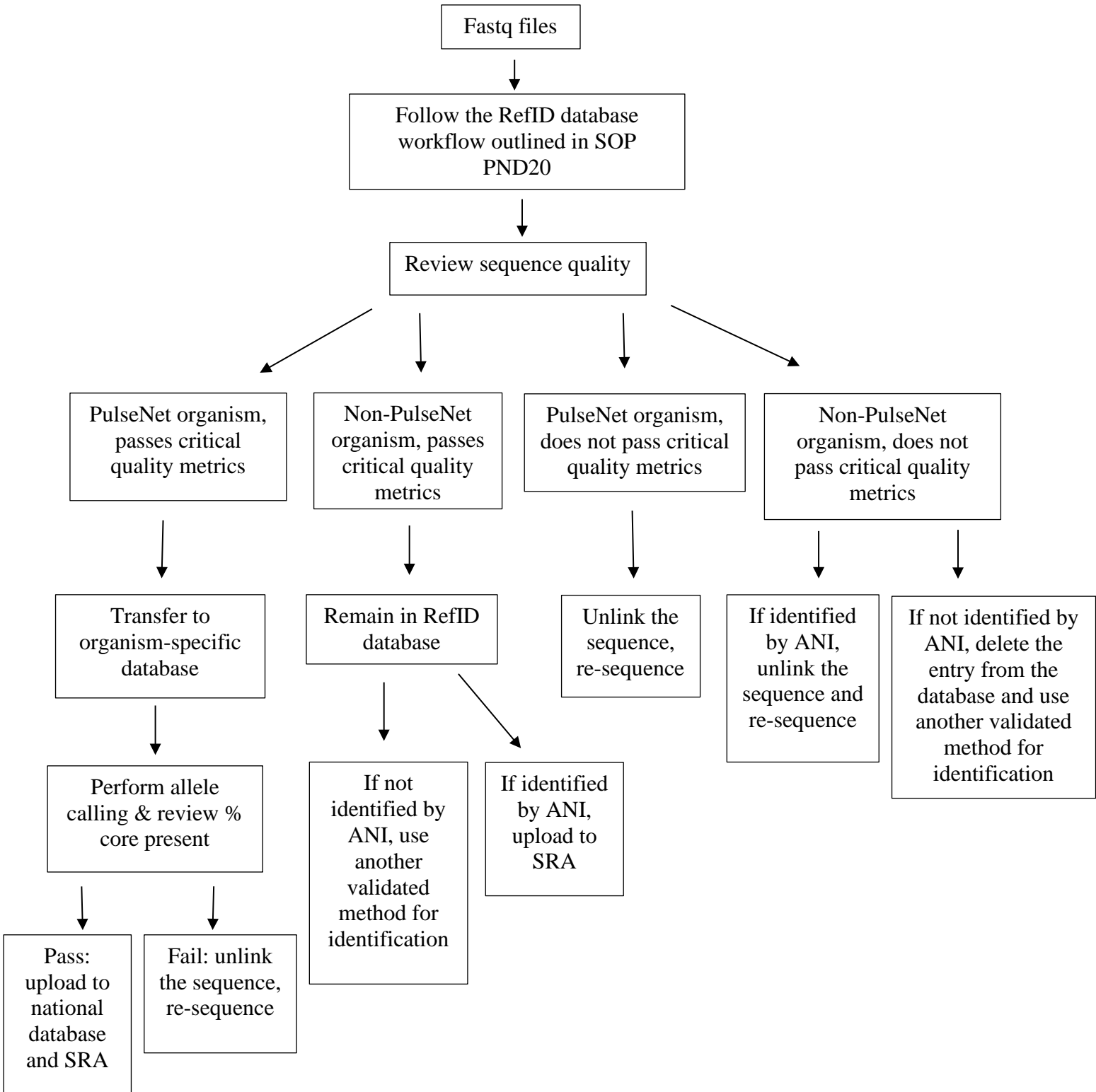


Average read length 170.0 bp

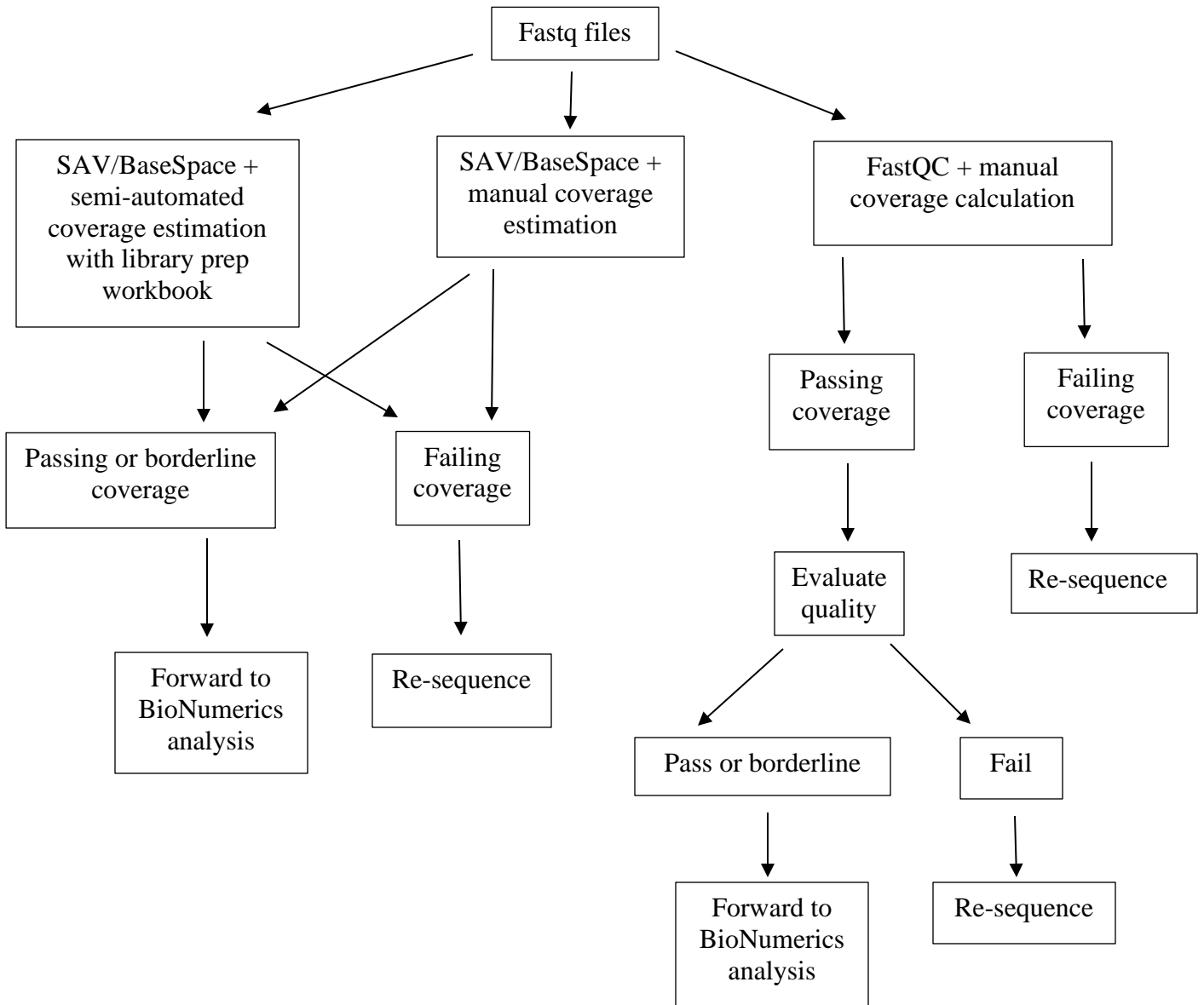


**6. FLOW CHART:**

**Sequence data QC using BioNumerics 7.6 or higher RefID database workflow (analysis certified individuals)**



**Sequence data QC without BioNumerics (individuals not certified for analysis)**



**7. REFERENCES:**

- 7.1. Sequencing Analysis Viewer Software Guide v.2.4. Illumina. 15066069 v03. November 2017
- 7.2. [www.bioinformatics.babraham.ac.uk/projects/fastqc](http://www.bioinformatics.babraham.ac.uk/projects/fastqc) FastQC. Babraham Bioinformatics
- 7.3. [https://support.illumina.com/sequencing/sequencing\\_software/basespace/documentation.html](https://support.illumina.com/sequencing/sequencing_software/basespace/documentation.html) BaseSpace Documentation and Literature. Illumina

**8. CONTACTS:**

- 8.1. CDC PulseNet NGS Laboratory Inbox: [PulsenetNGSlab@cdc.gov](mailto:PulsenetNGSlab@cdc.gov)
- 8.2. CDC PulseNet Database Inbox: [PulseNet@cdc.gov](mailto:PulseNet@cdc.gov)

**9. AMENDMENTS:**

12/22/2015: Added coverage calculation instructions for using the Read Metrics tab in the Nextera XT library prep workbook, and included image in new appendix PNQ07-3.

04/29/2016: Cluster density range corrected from 600-1300 to 600-1200.

06/27/2016:

- Attached updated image for Appendix PNQ07-3

10/13/2016:

- Updated formatting
- Updated information concerning running FastQC within BaseSpace Sequence Hub using iCredits
- Added quality metrics for v3 600 cycle chemistry
- Changed clusters passing filter values so that values were uniform across reagent kits
- Included basic quality guidance and graphical examples for fastq file assessment using FastQC.

01/20/2017:

- Corrected numbering of steps within procedure.
- Reformatted document layout according to new layout (removed footer, updated header, added "Approvals Signatures").
- Added PF to "Definitions".
- Updated formula and wording within step 5.2.2.

04/11/2018:

- Updated Purpose, Responsibilities, Definitions, clarified process for data analysis, added tables and diagrams for assessing quality and references.
- Updated document to include quality metrics for *Vibrio spp.*
- Updated PNQ07-3 to include *Vibrio spp.* Coverage

03/16/2020

- Removed the review of run quality metrics section from this SOP and moved it to PNL38
- Added BioNumerics 7.6 RefID database workflow as the primary option for raw data QC
- Added additional examples on how to interpret FastQC graphs
- Added workflow diagrams for raw data QC using BioNumerics and without BioNumerics
- Average quality score (Q score) minimum threshold increased from 28.0 to 30.0. Sequences with Q scores <30.0 will be rejected regardless of coverage.

**10. APPROVAL SIGNATURES:**

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet QA/QC Personnel

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet Outbreak Detection and Surveillance Unit Chief

Approved By:           N/A           Date:           N/A            
PulseNet PFGE Reference Unit Chief

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet Next Generation Subtyping Methods Unit Chief

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet Reference Outbreak Surveillance Team Lead

**Appendix PNQ07-1. Requirements for Critical Quality Metrics for PulseNet Organisms for Routine Sequence Submissions.**

Organism	Average denovo coverage	Average quality (Q score)	Assembly length (MB)	Secondary species abundance	% core present <sup>1</sup>
<i>Listeria monocytogenes</i>	≥ 20x	≥ 30	2.8-3.2	≤ 1.0	≥ 95
<i>E. coli</i> (most serotypes)	≥ 40x	≥ 30	4.9-6.0	≤ 1.0	≥ 85
<i>Shigella</i> spp./Rare <i>E. coli</i>	≥ 40x	≥ 30	4.2-4.9	≤ 1.0	≥ 85
<i>Salmonella</i> spp.	≥ 30x	≥ 30	4.4-5.7	≤ 1.0	≥ 85
<i>Campylobacter</i> spp.	≥ 20x	≥ 30	1.4-2.2	≤ 1.0	≥ 85 <sup>2</sup>
<i>Vibrio cholerae</i>	≥ 40x	≥ 30	3.8-4.3	≤ 1.0	NA
<i>Vibrio parahaemolyticus</i>	≥ 40x	≥ 30	4.9-5.5	≤ 1.0	NA
<i>Vibrio vulnificus</i>	≥ 40x	≥ 30	4.7-5.3	≤ 1.0	NA

Table 3. Requirements for critical quality metrics for PulseNet organisms

- All critical quality metrics, except % core present can be found in the RefID database (SOP PND20). % core present is determined in the organism-specific databases.
- Rows shaded with the same color represent organisms transferred to the same organism-specific database.

<sup>1</sup> The goal for percent core present is 95% or above. Routine *Escherichia*, *Shigella*, *Salmonella* and *C jejuni* sequences that fall between 85%-94% and pass other quality thresholds can be uploaded to the national database and SRA. The percent core present value has not been evaluated for *Vibrio* spp. because the allele scheme is still under development.

<sup>2</sup> Core scheme applies to *C. jejuni* only. Other *Campylobacter* species may still pass quality with % core present <85%; please contact [PulseNet@cdc.gov](mailto:PulseNet@cdc.gov) before resequencing.