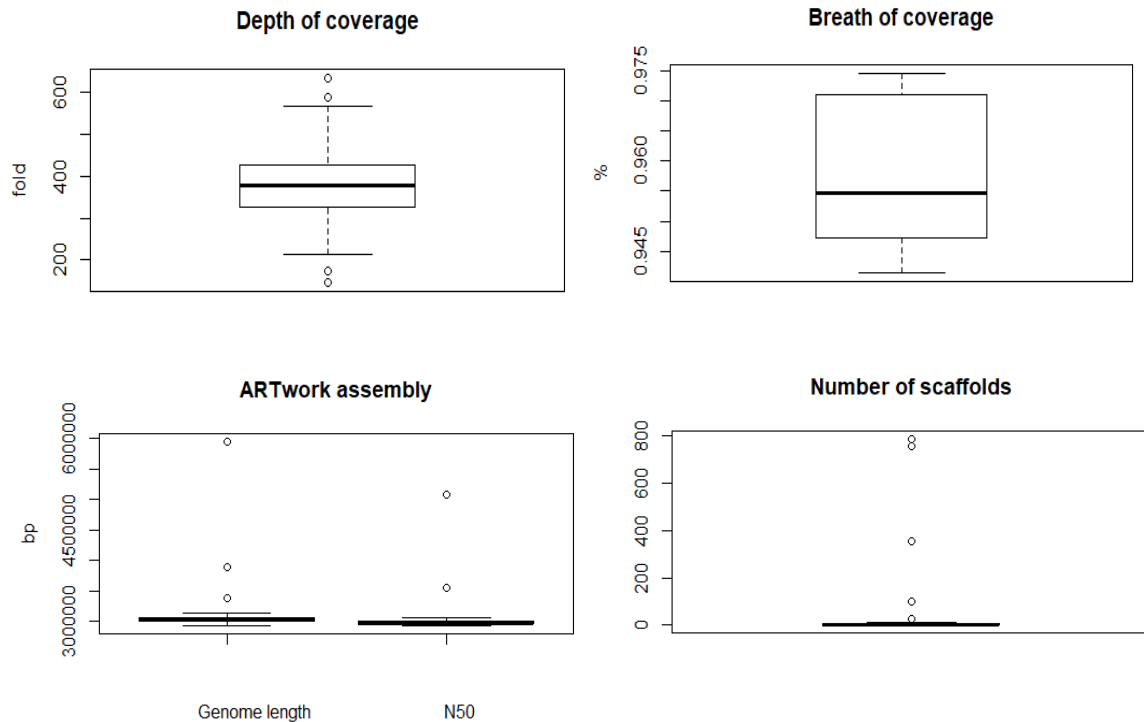1  **Additional file 4. Quality metrics of mapping (i.e. iVARcall2) and *de novo* assembly (i.e. ARTwork)**

2  **from the studied *Listeria monocytogenes* genomes**

3  An essential step in comparative genomics based studies is the Quality Control (QC) of WGS data in

4  order to guarantee the accuracy of sequencing results obtained by *in silico* genome-wide analysis. Poor

5  quality of read sequences as well as contamination of DNA can lead to significant errors in variant

6  calling (low-depth of sequencing impact on false-positive rate) and gene prediction analyses [1]. Even

7  though the harmonization and standardization of WGS data analysis is still an ongoing process [2], a

8  number of metrics for QC of *de novo* draft genome (i.e. contiguity of assemblies) and genome coverage

9  (i.e. number of reads mapped to a specific position within the reference genome, so called

10  "mappability") is currently available [3–5]. In this study, standard quality metrics of reads mapping (i.e.

11  iVARcall2) and *de novo* assembly (i.e. ARTwork) obtained from Illumina paired end reads of 96 *Listeria*

12  *monocytogenes* genomic DNA were assessed and reported in the boxplots. In particular, the quality of

13  reads mapping onto the reference genome was evaluated based on the depth of coverage (average

14  number of times that a base of a genome is sequenced) and breadth of coverage (percentage of bases

15  of a reference genome that are covered with a certain depth). Moreover, contiguity measures, such as

16  the size of assembled genomes (express in total number of bases and representing an indicator of

17  exogenous DNA contamination), and the total number of contigs/scaffolds along with the N50 value

18  (size of the largest contig, or scaffold, for which half the total size is contained in that contigs and those

19  larger), were calculated. Overall, high values of depth and breadth of coverage (1st and 3rd quartile =

20  145-426X and 94-96%, respectively) have been estimated confirming the high quality of the Illumina

21  short reads for further analyses. Accordingly, genome sizes as well as the number of scaffolds and the

22  N50 (median values of 3,021,803 bp, 4 and 2,969,304 bp, respectively) demonstrated high

23  performance of ARTwork pipeline in term of high contiguity of *de novo* assemblies for almost all *L.*

24  *monocytogenes* samples (~98%). However, two out of 96 assemblies were discarded from further

25  analyses since suspected of contamination (total genome length higher than 3.8 Mbp and number of

26    scaffolds > 354).



**Depth of coverage**

**Breath of coverage**

**ARTwork assembly**

**Number of scaffolds**

Genome length    N50

27

28    **References**

29    1. Sims D, Sudbery I, Ilott NE, Heger A, Ponting CP. Sequencing depth and coverage: key
30    considerations in genomic analyses. Nature Reviews Genetics. 2014;15:121–32.

31    2. Revez J, Espinosa L, Albiger B, Leitmeyer KC, Struelens MJ, ECDC National Microbiology Focal Points
32    and Experts Group. Survey on the Use of Whole-Genome Sequencing for Infectious Diseases
33    Surveillance: Rapid Expansion of European National Capacities, 2015-2016. Front Public Health.
34    2017;5:347.

35    3. Ellington MJ, Ekelund O, Aarestrup FM, Canton R, Doumith M, Giske C, et al. The role of whole
36    genome sequencing in antimicrobial susceptibility testing of bacteria: report from the EUCAST
37    Subcommittee. Clinical Microbiology and Infection. 2017;23:2–22.

38    4. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome assemblies.
39    Bioinformatics. 2013;29:1072–5.

40    5. Lüth S, Kleta S, Al Dahouk S. Whole genome sequencing as a typing tool for foodborne pathogens
41    like Listeria monocytogenes – The way towards global harmonisation and data exchange. Trends in
42    Food Science & Technology. 2018;73:67–75.