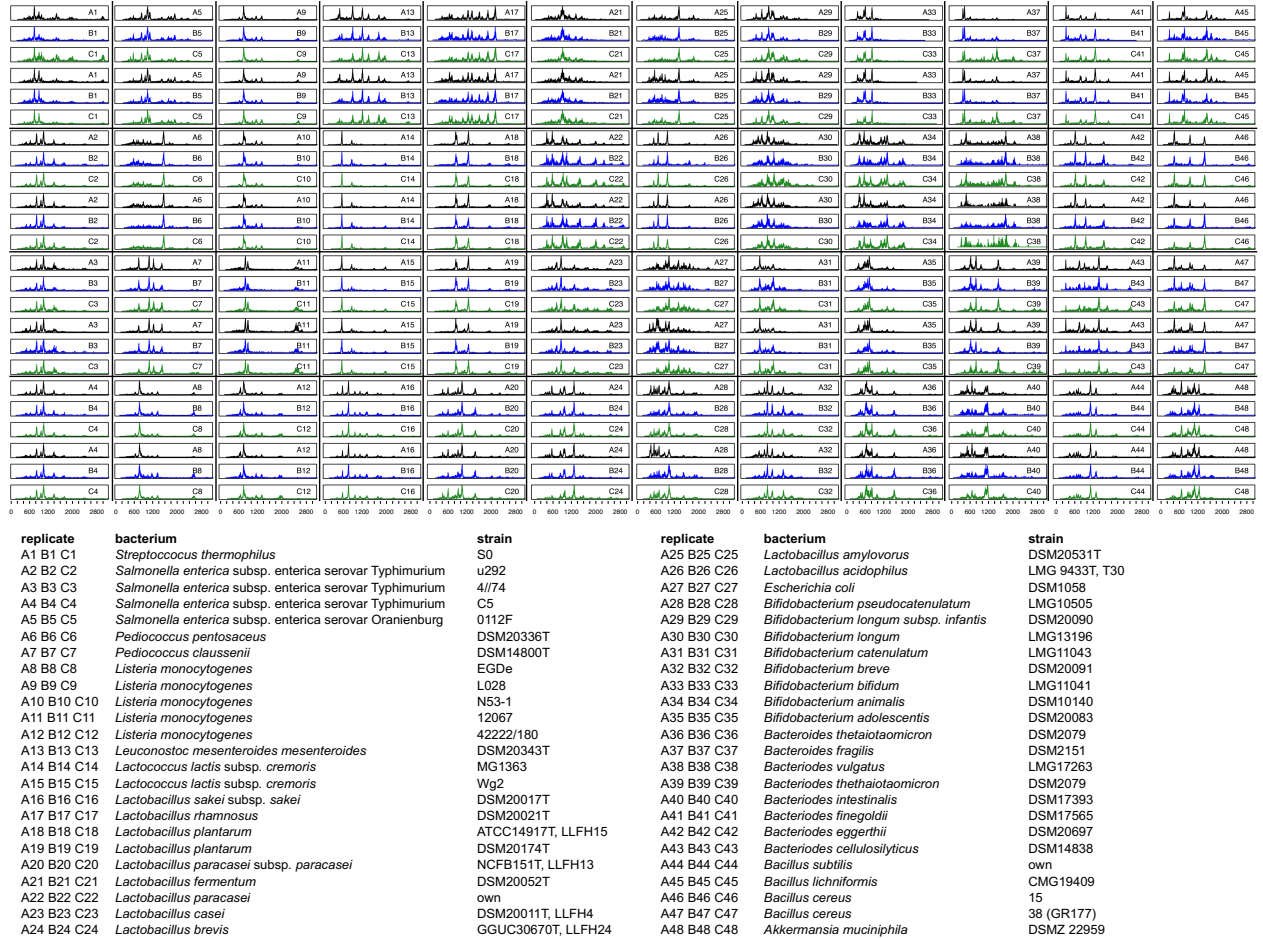
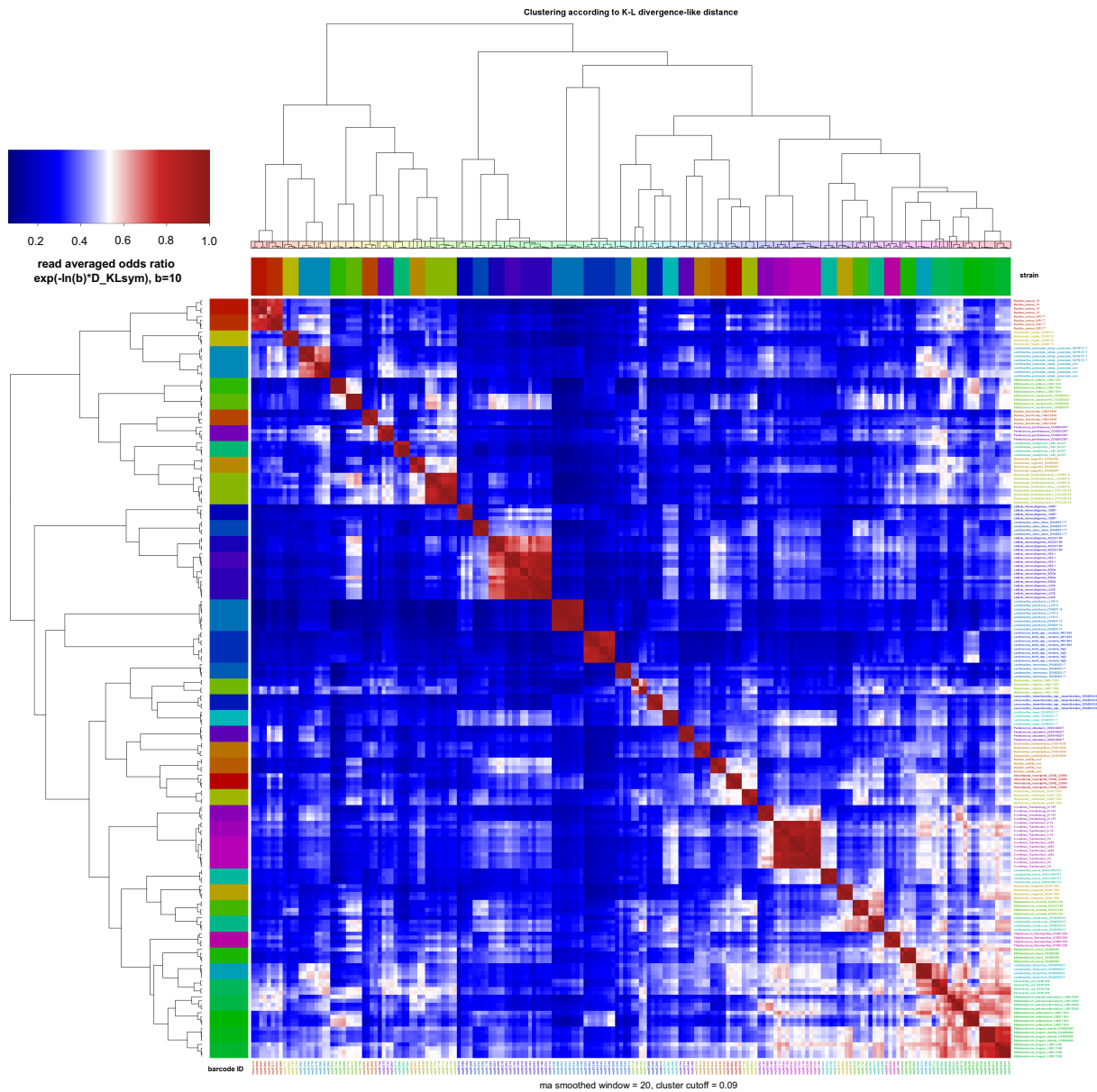


Supplementary Figures

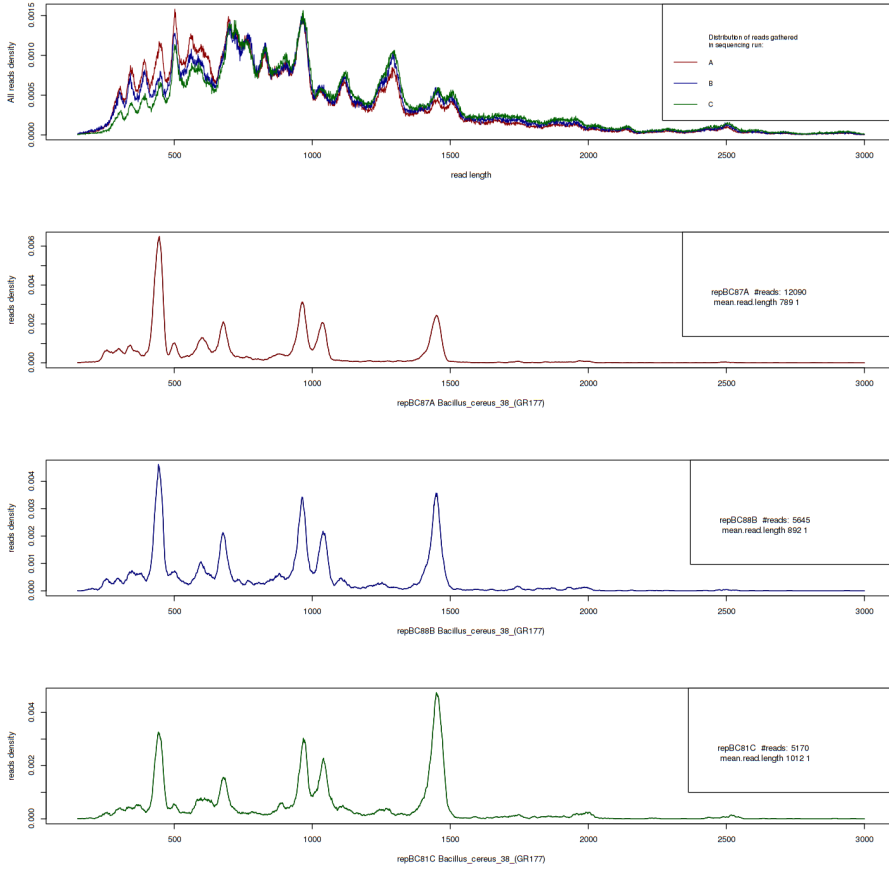


Supplementary Figure 1. LCp generated for 48 bacterial isolates using Oxford Nanopore Technology based rep-PCR amplicon sequencing (ON-rep-seq). The black, blue and green profiles indicate data collected during run A, B and C respectively for which each technical replicate received different barcode. All isolates were analysed in duplicates within each run. The list of bacterial taxa matching given LCp is given in the table.



Supplementary Figure 2. Row/Column clustering according to "Ward.D2" hierarchical clustering on D_{KLsym} distance of all 48 isolates. Heatmap showing similarity ($\exp(-\ln(b) \cdot D_{KLsym}), b=10$), and clustering according to cutoff=0.09. The detailed analysis using varying cutoff value (no single cutoff achieves exact separation between all and only different LCp, see Supplementary Figure 4 C, D ROC curves) and LCp visual inspection allowed for accurate differentiation between all except two pairs of bacterial strains described thoroughly in the results section (see Figure 3 and Supplementary figure 2 for details). Technical replicates from the third run "repC" were removed from the analysis due to higher short/long reads imbalance.

A

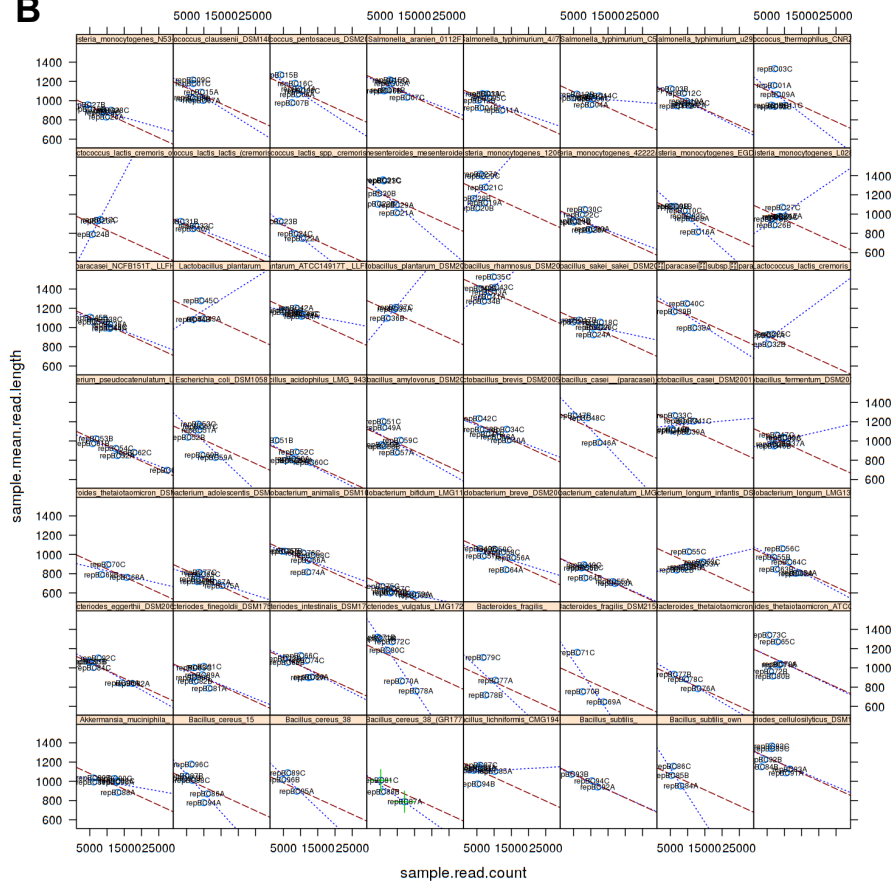


Supplementary Figure 3

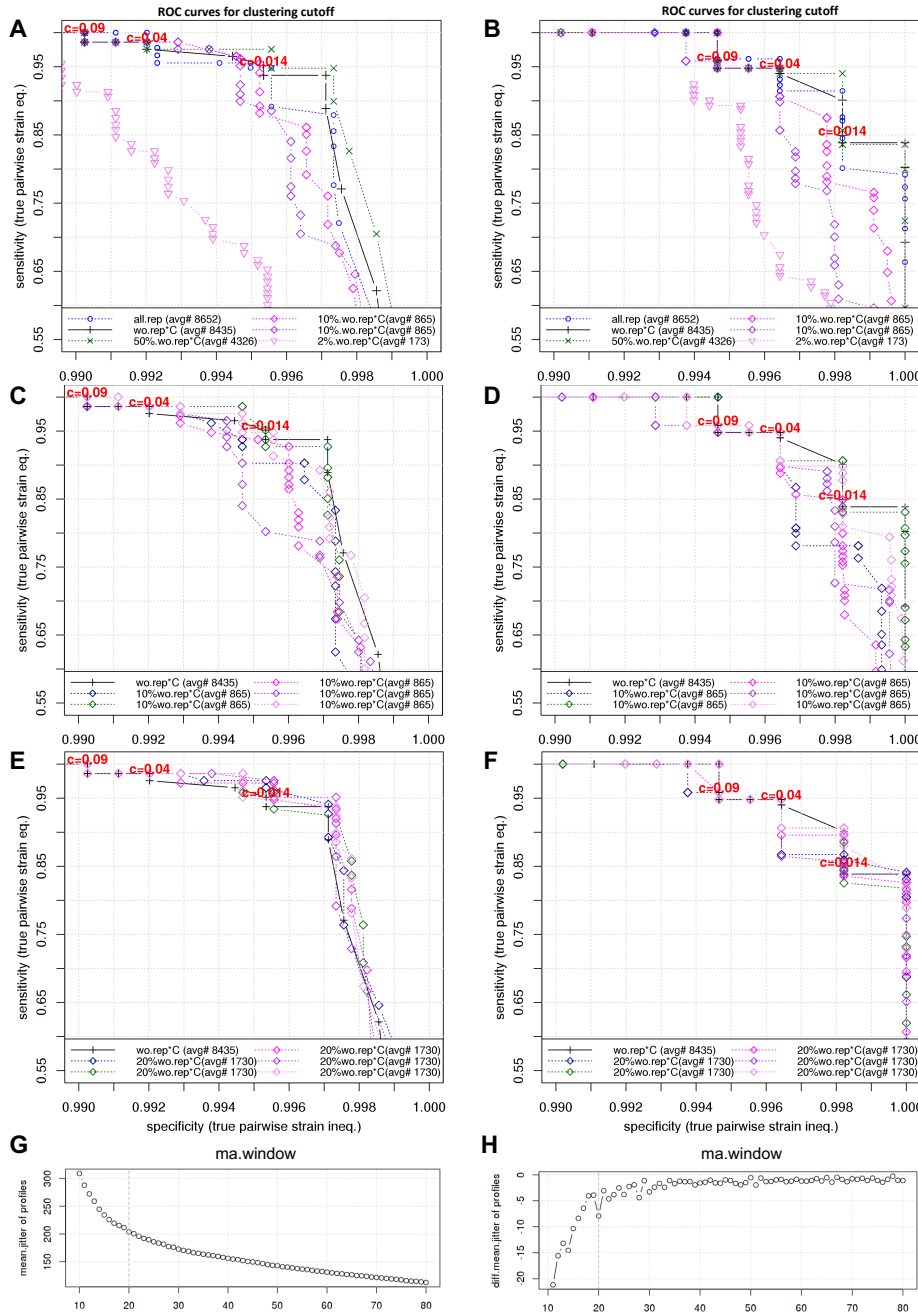
A) Top panel presents distribution of lengths of reads obtained in 3 separate consecutive sequencing runs A, B, C on the same flow cell. Third run C obtained less short reads, some differences are also visible in second run B, compared to the first run A. Bottom 3 panels show LCps of *Bacillus_cereus_38* (GR177) strain obtained from runs A,B,C.

B) Regression analysis of mean read length from LCp vs read count in LCp, data shown in separate panels for each strain replicates. Red dashed line is regression line obtained in all samples analysis, blue lines are regression lines for each strain only. Green markers mark runs A, C for *Bacillus_cereus_38*(GR177) (panels 2 and 4 in A).

B



Supplementary Figure 4. Peaks profiles comparison



A-F) Receiver operating characteristic (ROC) curves of pairwise "same/not-the-same" strain discrimination in various cutoffs c (diff. step=0.005), for various subsets of data: "all", "wo.rep*C" dataset without the third sequencing run "C" on twice used flow cells, "50%.wo.rep*C" subsample half the size of original, "20%.wo.rep*C" five subsamples $1/20^{\text{th}}$ of reads, "10%.wo.rep*C" seven subsamples $1/10^{\text{th}}$ of reads and "2%.wo.rep*C" $1/50^{\text{th}}$ of reads. On x-axis specificity, the percentage of correctly identified "not-the-same strain" pairs out of all such pairs (36096 for wo.rep*C), on y-axis sensitivity, the percentage of correctly identified "same strain" pairs, out of all such pairs (768 for wo.rep*C). A) Clustering according to sample strain label, viewed as a whole method performance, in contrast to B. B) Clustering according to sample strain similarity derived from visual inspection of profiles, thus these curves correspond more to D_KLsym-based profile comparison performance, than to the whole method. Values on the plot: $c=0.09$ (sp 0.9947, se 0.9583), $c=0.014$ (sp 0.9982, se 0.8490). All cutoffs "c" values marked for "wo.rep*C". C) Clustering according to sample strain label using 5 iterations of 10% subsets. D) Clustering according to sample strain similarity derived from visual inspection of profiles using 5 iterations of 20% subsets. The

analysis shows that 20% subsets perform similarly to the whole dataset what indicates the theoretical throughput of ON-rep-seq to range from 960 (for ~1.5M reads) to 1440 (for ~2.5M reads) isolates per flow cell. The analysis on panels E, F) shows that 20% subsets perform similarly to the whole dataset, what indicates the theoretical throughput of ON-rep-seq to range from 960 (for ~1.5M reads) to 1440 (for ~2.5M reads) isolates per flow cell. G) mean jitter of all profiles dependence on smoothing moving average "ma" window size. Jitter was defined as an average number of times when profile's discrete derivative changes sign (change to 0 was counted as 0.5). H) discrete derivative (diff lag=1) of the (top) mean jitter. Sizes of ma.window > 20 change mean jitter slowly and steadily suggestive of stabilization (noise decoupling) of information content in higher smoothing window results.

Supplementary Tables

Supplementary Table 1. Identification of MLST genes alleles among selected strains of *Salmonella enterica* subsp. *enterica* serovar Typhimurium and *Listeria monocytogenes*

Strain	Gene	Allele type	Strain	Gene	Allele type
<i>Salmonella enterica</i> U292	<i>aroC</i>	10	<i>Listeria monocytogenes</i> EDGe	<i>abcZ</i>	6
	<i>dnaN</i>	7		<i>bgIA</i>	5
	<i>hemD</i>	12		<i>cat</i>	6
	<i>hisD</i>	9		<i>dapE</i>	20
	<i>purE</i>	5		<i>dat</i>	176
	<i>sucA</i>	9		<i>ldh</i>	4
	<i>thrA</i>	2		<i>lhk</i>	1
<i>Salmonella enterica</i> C5	<i>aroC</i>	10	<i>Listeria monocytogenes</i> L028	<i>abcZ</i>	6
	<i>dnaN</i>	7		<i>bgIA</i>	5
	<i>hemD</i>	12		<i>cat</i>	6
	<i>hisD</i>	9		<i>dapE</i>	51
	<i>purE</i>	5		<i>dat</i>	176
	<i>sucA</i>	9		<i>ldh</i>	4
	<i>thrA</i>	2		<i>lhk</i>	1
<i>Salmonella enterica</i> 4/74	<i>aroC</i>	10			
	<i>dnaN</i>	7			
	<i>hemD</i>	12			
	<i>hisD</i>	9			
	<i>purE</i>	5			
	<i>thrA</i>	2			

Supplementary Table 2. Details regarding benchmarking of two R9.4.1 flow cells.

Run ID	Flow cell 1				Flow cell 2		
	A	B	C	D	A	B	C
Run Time (h)	4	4	4	4	4	4	12
Break between the next run (day)	1	4	3	7	1	1	1
Active pores at start	1347	1324	1098	925	1034	779	615
Voltage at start (mV)	-180	-180	-190	-195	-180	-180	-190
Initial sequences in strand	~300	~200	~150	~50	~200	~120	~70
Total number of high quality reads collected	9.4×10 ⁵	7.9×10 ⁵	5.7×10 ⁵	2.2×10 ⁵	10.5×10 ⁵	5.7×10 ⁵	8.7×10 ⁵
Library concentration loaded in 12 µl (ng/ µl)	2.5	1.8	3.0	1.6	3.2	2.1	2.4

Both flow cells generated in total similar amount of data, although flow cell 1 was in much better condition and had more active pores at start what allowed to perform four consecutive runs. Flow cell 2 had lower number of active pores at arrival and seemed to deteriorate faster therefore only three runs were conducted. Last run was elongated to 12 h in order to collect maximum amount of data from declining flow cell. The data from the first benchmarked flow cell were used solely to test the optimal concentration of DNA needed and viability of the flow cell while data from the second flow cell are presented herein