

SUPPLEMENTAL FIGURES

Rapid and accurate demultiplexing of direct RNA nanopore sequencing datasets with SeqTagger

Leszek P Pryszcz^{1,#,*}, Gregor Diensthuber^{1,2,#}, Laia Llovera¹, Rebeca Medina¹, Anna Delgado-Tejedor^{1,2}, Luca Cozzuto¹, Julia Ponomarenko¹ and
Eva Maria Novoa^{1,2,3,*}

¹*Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Dr. Aiguader
88, Barcelona 08003, Spain*

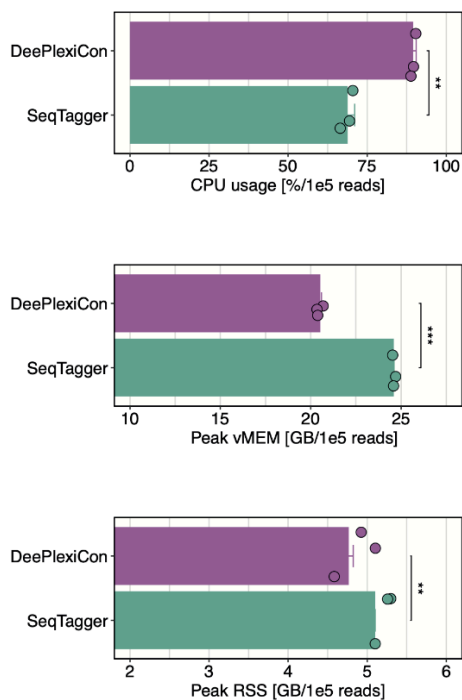
²*Universitat Pompeu Fabra (UPF), Barcelona, Spain*

³*ICREA, Pg. Lluís Companys 23, Barcelona, España*

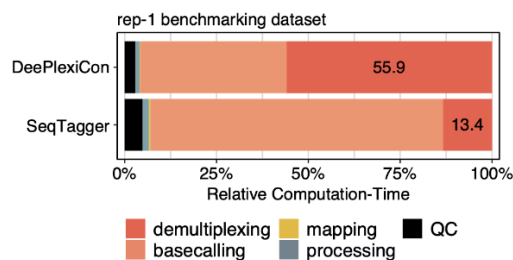
[#] *These authors contributed equally*

* Correspondence to: Leszek P Pryszcz (leszek.pryszcz@crg.eu) and Eva Maria Novoa
(eva.novoa@crg.eu)

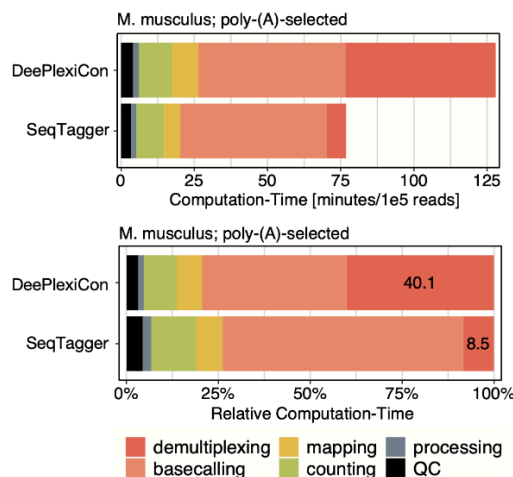
A



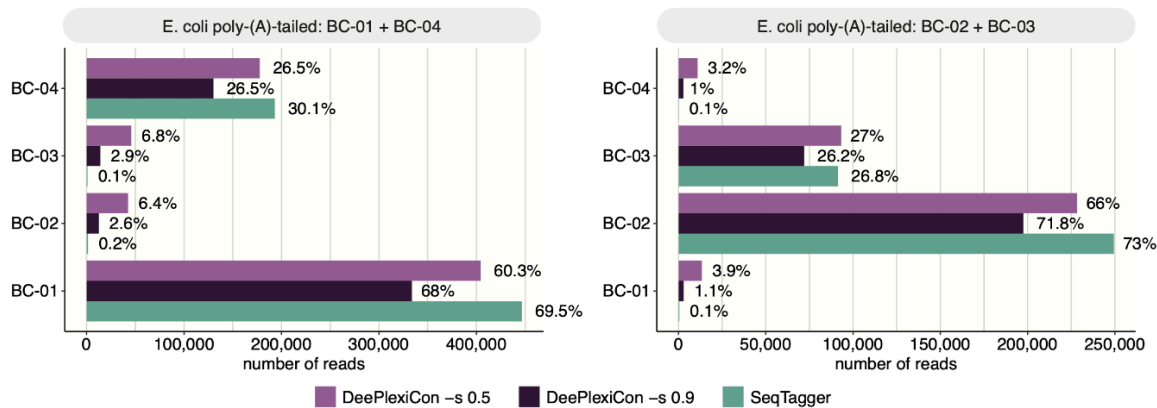
B



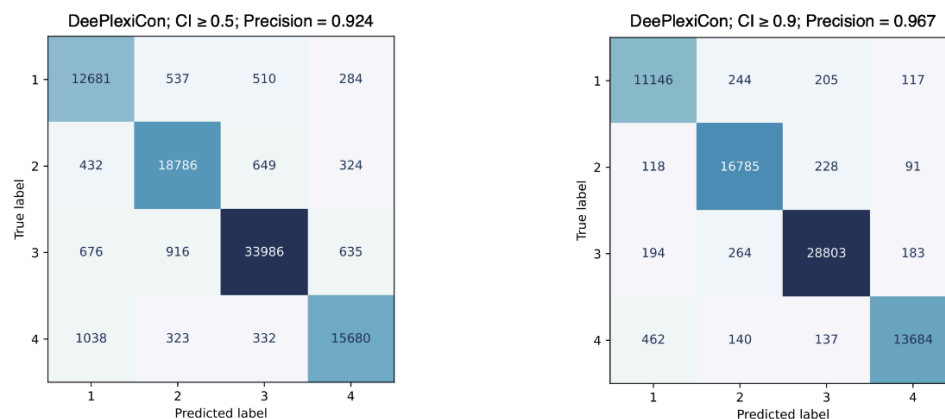
C



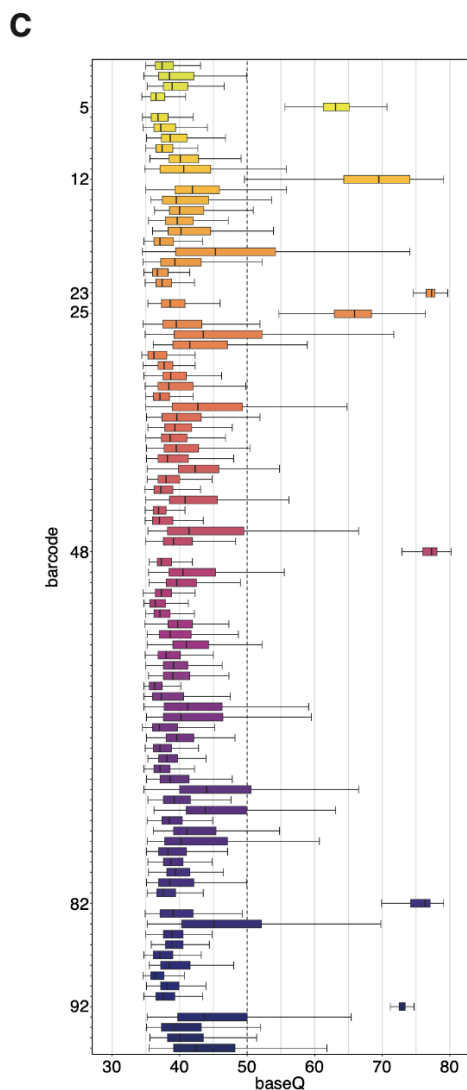
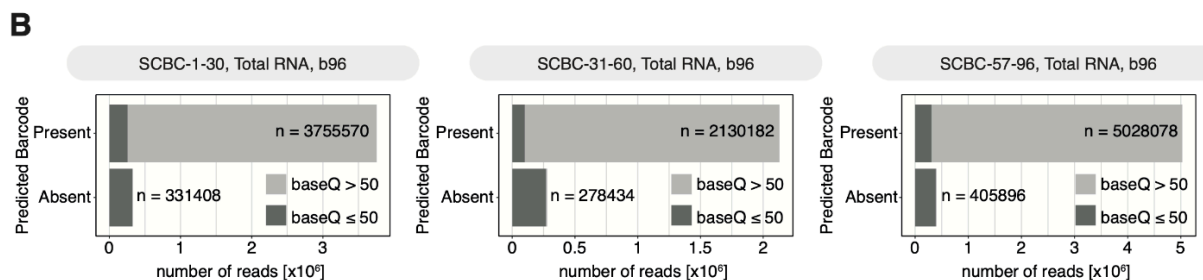
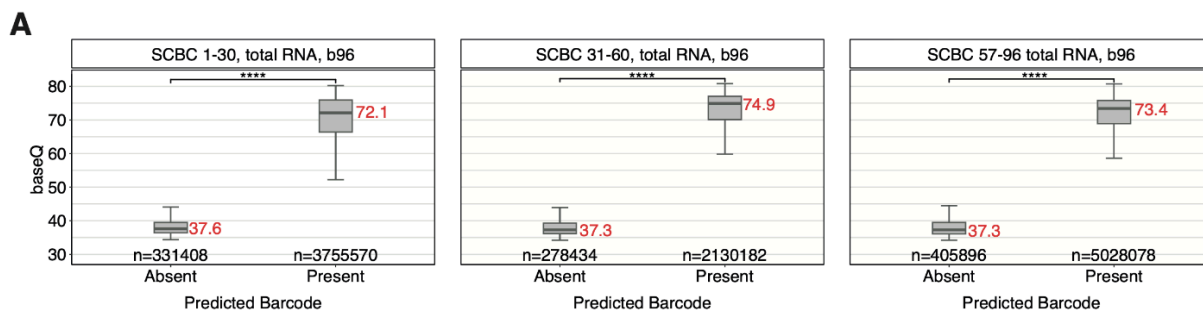
D



E



Supplemental Figure S1. Benchmarking results for computational requirements and model performance (A) Barplots depicting system requirements (CPU usage; peak vMEM; peak RSS) of DeePlexiCon and SeqTagger on the benchmarking dataset. Dots represent individual replicates with bars representing the mean value and error bars depicting ± 1 standard deviation. To determine statistical significance, a two-sided t-test was performed and results were corrected for multiple hypothesis testing using the Benjamini-Hochberg procedure. **(B)** Barplot depicting the relative contributions of each preprocessing workflow to the overall computation time on rep-1 of the benchmarking dataset. **(C)** Barplots depicting the absolute and relative contributions to the overall computation time of 100,000 reads sampled from a mouse poly(A)-selected sample aligned to the mm39 genome. **(D)** Barplots representing the percentage of reads assigned to each barcode for two runs of total RNA from *E.coli* (poly(A)-tailed). The first run (left) contained BC-01 and BC-04 while the second run (right) contained barcodes BC-02 and BC-03. Runs were demultiplexed with either SeqTagger (b04_RNA002) or DeePlexiCon (resnet20-final.h5) with high recovery (-s 0.5) or high accuracy (-s 0.9) settings (see *Methods*). **(E)** Confusion matrices corresponding to DeePlexiCon results for high recall (-s 0.5) and high precision (-s 0.9) on rep-1 of the benchmarking dataset.



Supplemental Figure S2. Performance of SeqTagger's 96 barcode model on independent test data. (A) Boxplots showing the baseQ distribution for barcodes present and absent in the three independent test runs. The Number of reads is indicated by n with the median value shown in red. Statistical analysis was performed using a two-sided non-parametric Wilcoxon test. Results were corrected for multiple-hypothesis testing using the Bonferroni procedure to obtain adjusted p-values (ns: $p > 0.05$, *: $p \leq 0.05$, **: $p \leq 0.01$, ***: $p \leq 0.001$, ****: $p \leq 0.0001$). **(B)** Barplots representing the total number of reads (n) for three independent test runs demultiplexed with SeqTagger's 96 barcode model (b96_RNA002). Colors indicate different baseQ thresholds. **(C)** Boxplots of base quality (baseQ) per barcode for an additional independent test run containing SCBC-05, SCBC-12, SCBC-23, SCBC-25, SCBC-48, SCBC-82 and SCBC-96. For Figures S2A and S2C the box is limited by the lower quartile Q1 (bottom) and upper quartile Q3 (top). Whiskers are defined as $1.5 * IQR$ with outliers not shown for visualization purposes.
