

Supplementary Information for

deltaRpkM: an R package for a rapid detection of differential gene presence between related genomes

Hatice Akarsu^{1,4#}, Lisandra Aguilar-Bultet^{2,3,4#}, Laurent Falquet^{1,4*}

¹Department of Biology, University of Fribourg, Switzerland,

²Institute of Veterinary Bacteriology, Vetsuisse Faculty, University of Bern, Bern, Switzerland,

³Graduate School for Cellular and Biomedical Sciences, University of Bern, Bern Switzerland,

⁴Swiss Institute of Bioinformatics, BUGFri group, Fribourg, Switzerland.

*To whom correspondence should be addressed.

These authors contributed equally to this work.

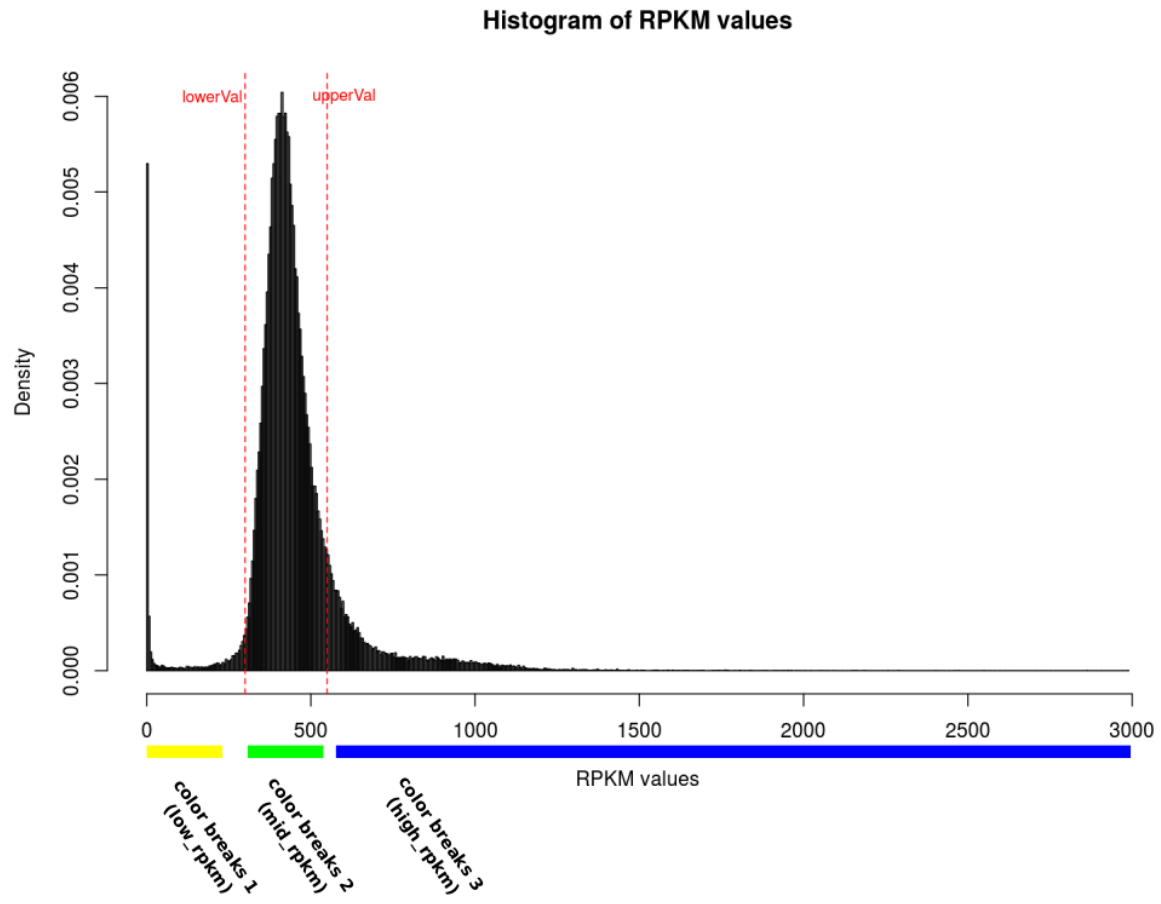


Fig. S1. RPKM values distribution of all genes in the dataset. This can be used to fine tune the heatmap colors breaks parameters.

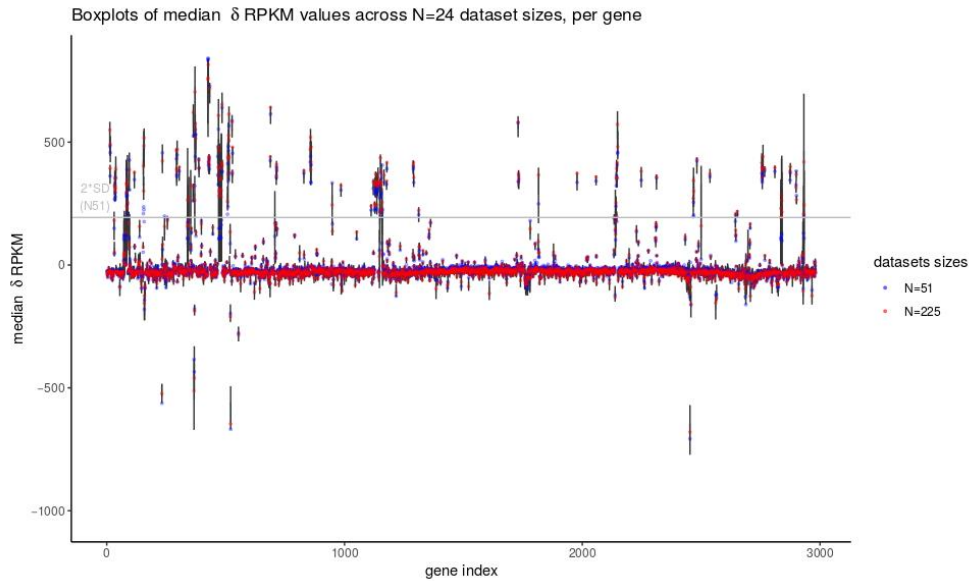
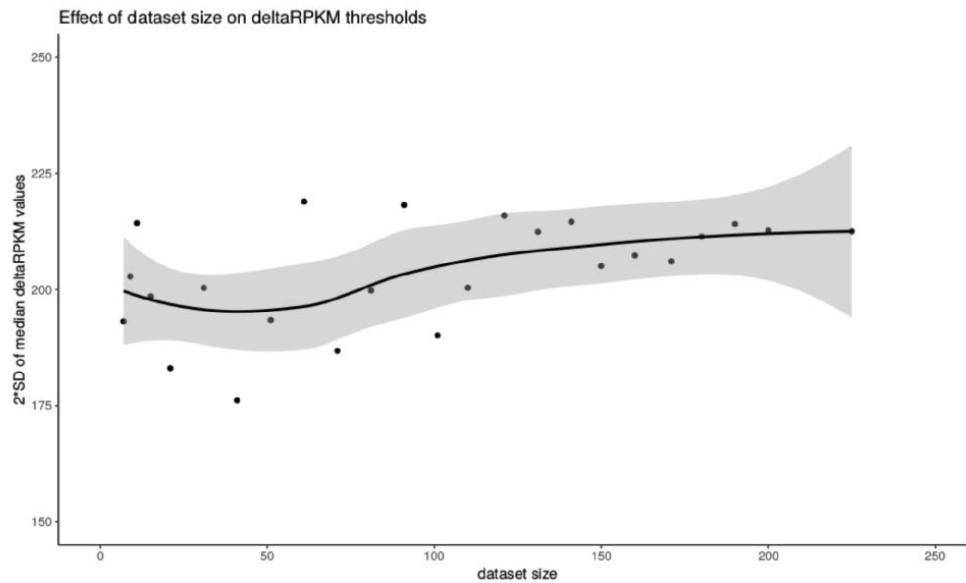
A**B**

Fig. S2. Dataset size effect on δ RPKM values distribution. A. Boxplots for datasets from N=7 to N=225 samples. The dataset size does not influence the median δ RPKM values that are used when computing the differentially present gene selection based on the 2*standard deviation of median δ RPKM values. Two datasets are highlighted for illustration, N=51 samples and N=225 samples. **B. Dataset size effect on threshold value (2*standard deviation) of median δ RPKM.**

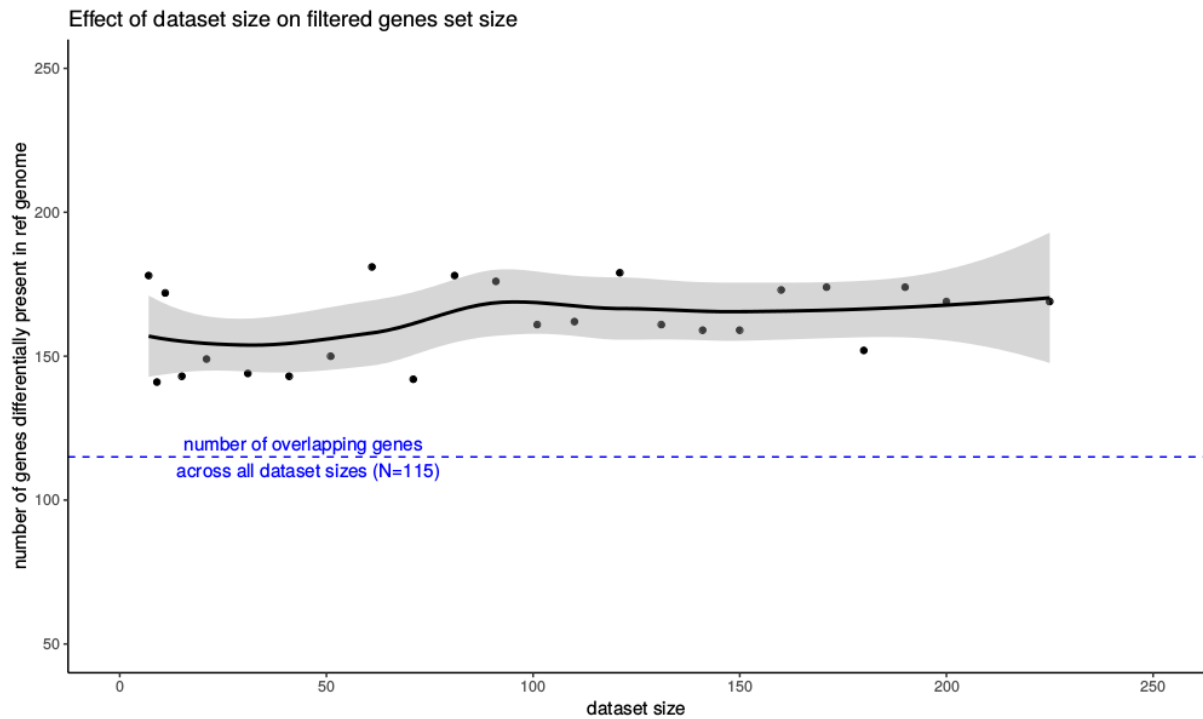


Fig. S3. The selected differentially present gene set is robust. Downsampling shows that even with small size dataset, the identified genes overlap (N=115) highly with the datasets of greater size.

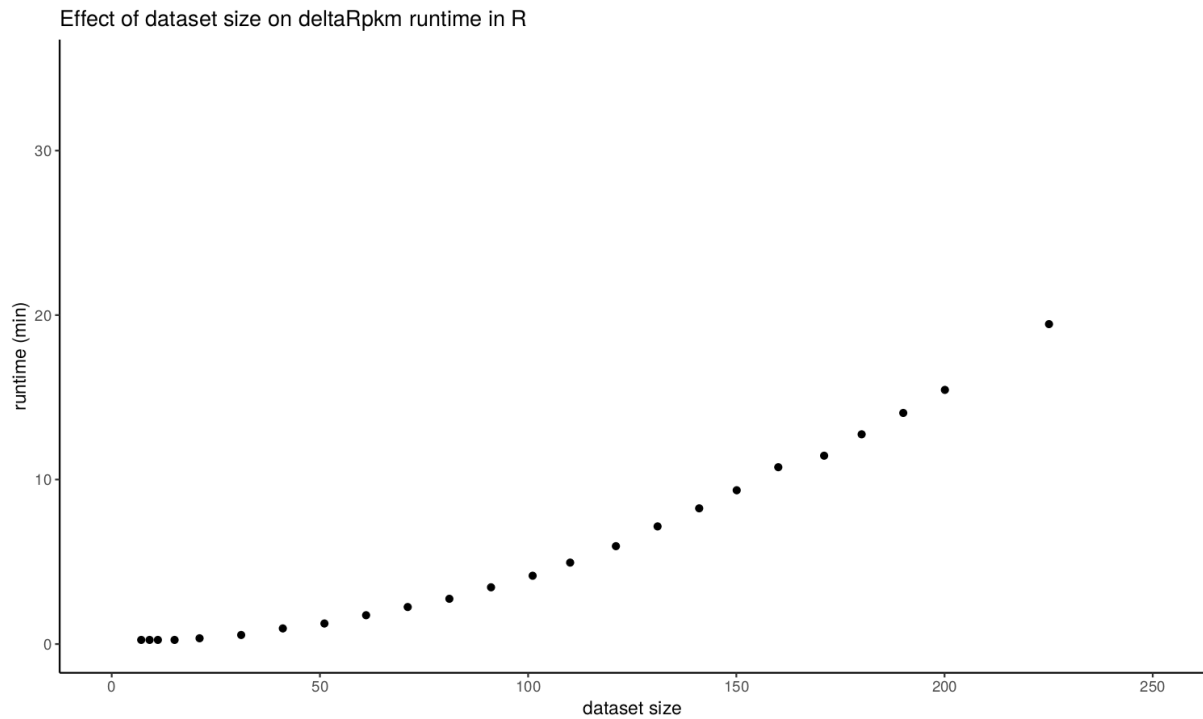


Fig. S4. deltaRpkM performance: dataset size effect on runtime. The whole analysis pipeline with deltaRpkM can be run in less than 20min in R for a dataset with N=225 samples of *Listeria monocytogenes* (~3Mb, ~3K genes). Ubuntu 14.04, R 3.4.4, Intel Core i-4790 CPU @3.60Gzx8.

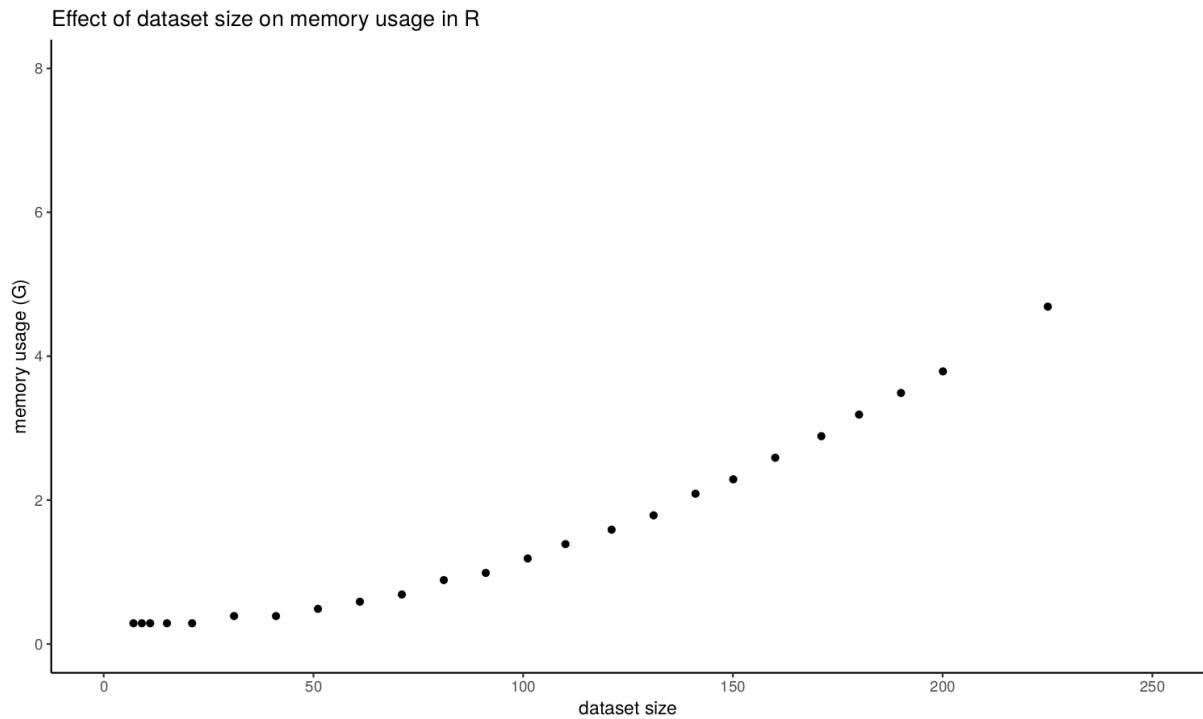


Fig. S5. deltaRpk performance: dataset size effect on memory usage. The whole analysis pipeline with deltaRpk uses less than 4G of memory in R for a dataset with N=225 samples of *Listeria monocytogenes* (~3Mb, ~3K genes). Ubuntu 14.04, R 3.4.4, Intel Core i-4790 CPU @3.60Gzx8.

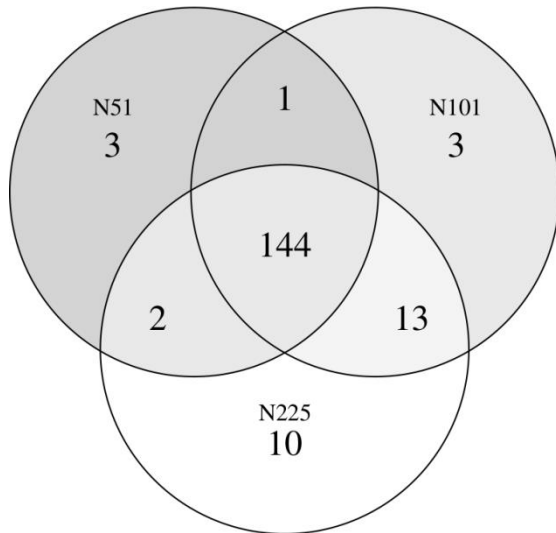
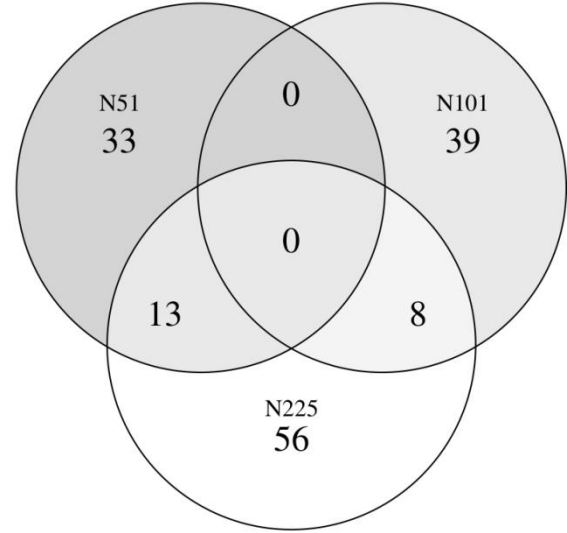
A**B**

Fig. S6. deltaRpk performance: real (A) versus randomized datasets (B). The genes differential presence gives shorter and non-robust list of genes when using randomized datasets of different sizes. Corrected RPKM.