

A Comparison of mRNA Sequencing with Random Primed and 3'-Directed Libraries

Yuguang Xiong¹, Magali Soumillon^{2,10}, Jie Wu^{3,4}, Jens Hansen¹, Bin Hu¹, Johan G.C. van Hasselt¹, Gomathi Jayaraman¹, Ryan Lim^{3,4}, Mehdi Bouhaddou¹, Loren Ornelas^{5,6}, Jim Bochicchio², Lindsay Lenaeus^{5,6}, Jennifer Stocksdale⁴, Jaehee Shim¹, Emilda Gomez^{5,6}, Dhruv Sareen^{5,6,7}, Clive Svendsen^{5,6,7}, Leslie M. Thompson^{3,4,8}, Milind Mahajan⁹, Ravi Iyengar¹, Eric A. Sobie¹, Evren U. Azeloglu^{1,#}, Marc R. Birtwistle^{1,11,#}

¹Department of Pharmacological Sciences and DToxS LINCS Center, Icahn School of Medicine at Mount Sinai, New York, NY USA

²Broad Institute of MIT and Harvard, Cambridge, MA USA

³Department of Biological Chemistry, University of California, Irvine, CA USA

⁴UCI MIND, University of California, Irvine, CA USA

⁵Board of Governors-Regenerative Medicine Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA

⁶iPSC Core, The David and Janet Polak Foundation Stem Cell Core Laboratory, Los Angeles, CA, USA

⁷Department of Biomedical Sciences, Cedars-Sinai Medical Center, Los Angeles, CA, USA

⁸Department of Psychiatry and Human Behavior, Neurobiology and Behavior, University of California, Irvine, CA USA

⁹Department of Genetics, Icahn School of Medicine at Mount Sinai, New York, NY

¹⁰Present Address: Berkeley Lights, Inc. 5858 Horton St., Emeryville, CA 94608

¹¹Present Address: Department of Chemical and Biomolecular Engineering, Clemson University, SC USA

[#]To whom all correspondence should be addressed: evren.azeloglu@mssm.edu or marc.birtwistle@mssm.edu

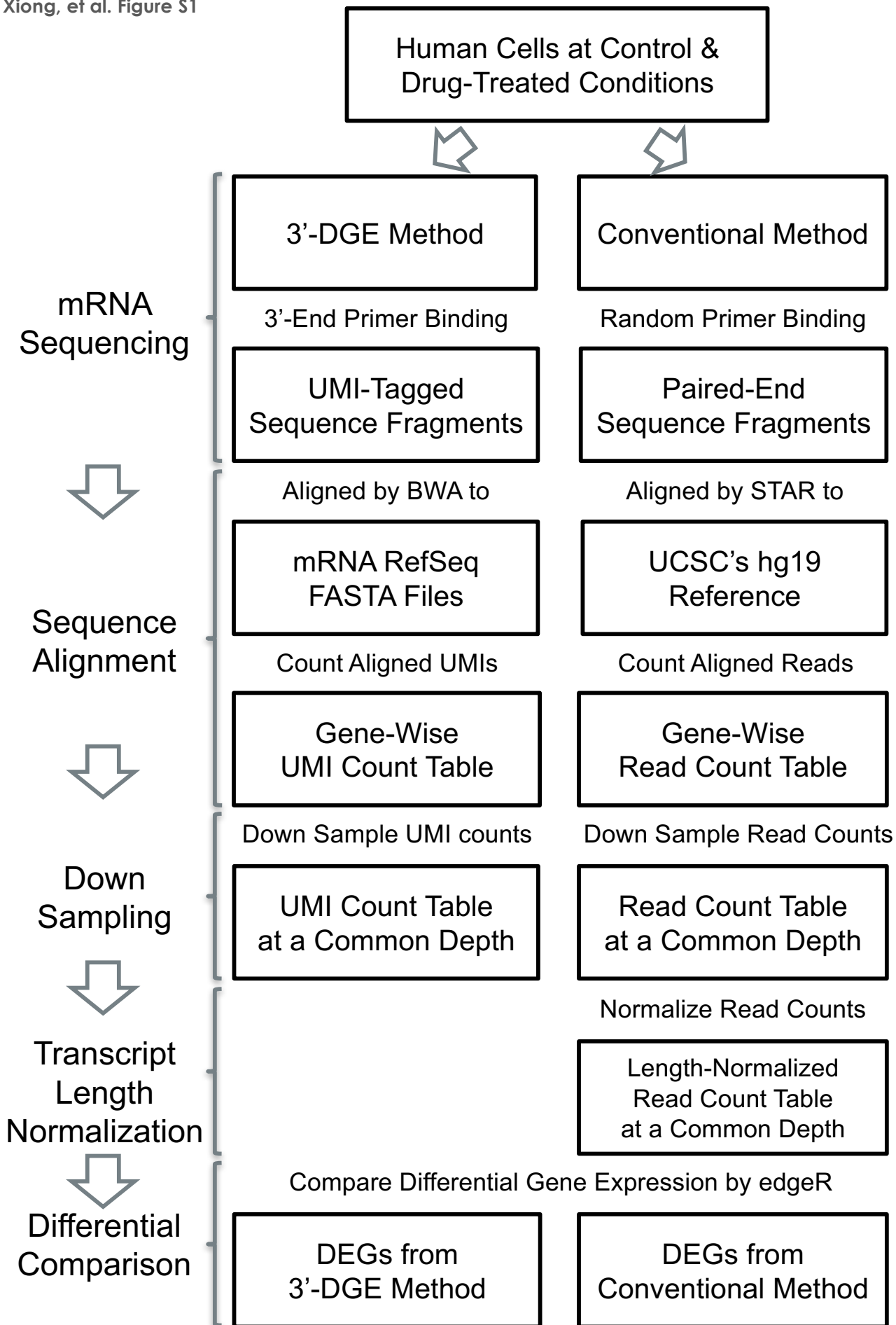


Figure S1. The Workflow of Comparing Two mRNA Sequencing Methods: the 3'-end Digital Gene Expression (3'-DGE) Method and the Conventional Random Primer-binding Method.

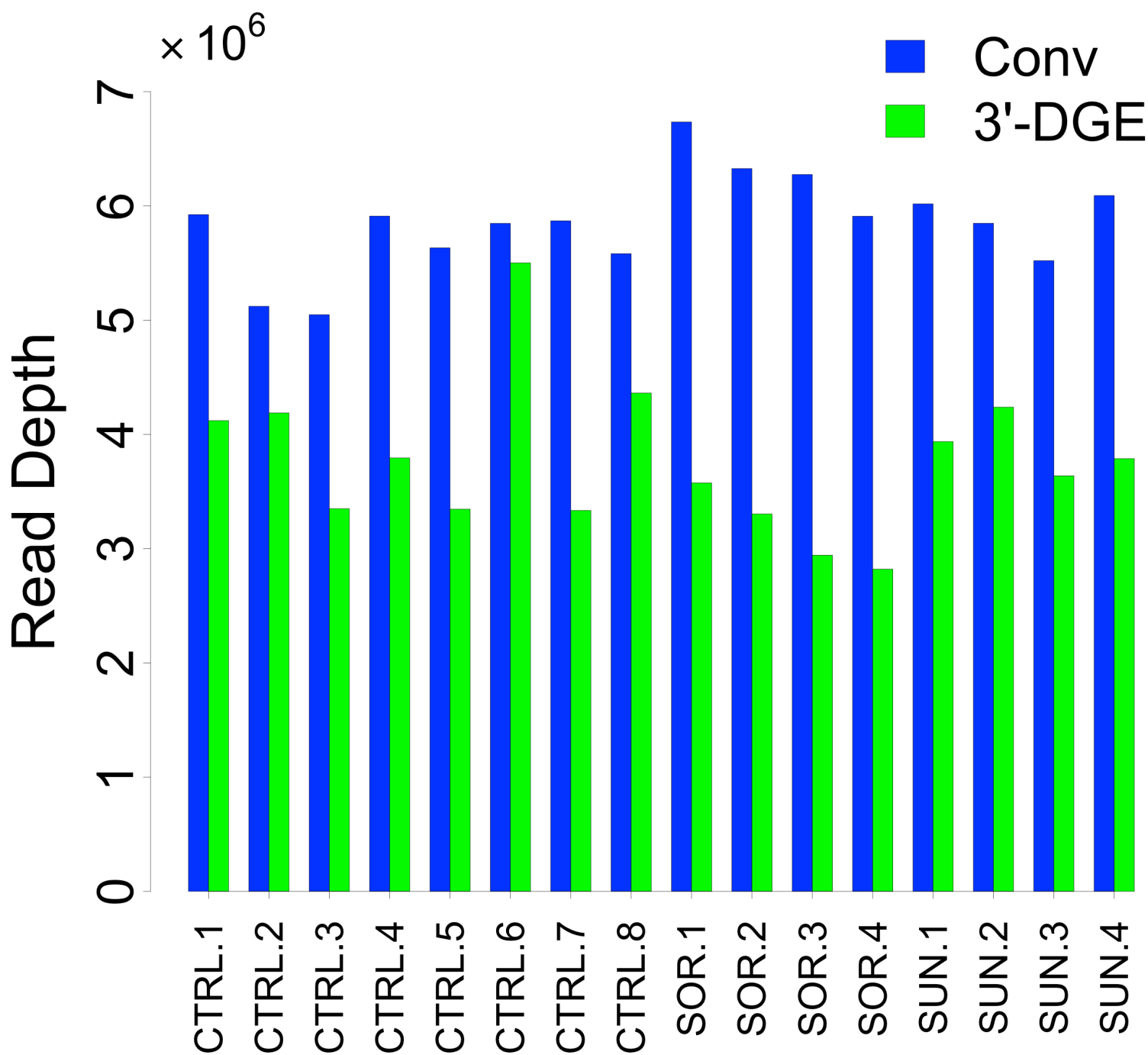


Figure S2. Read Depth Across Multiple Samples for Conventional (Conv) and 3'-end Digital Gene Expression (3'-DGE) Methods. The total number of uniquely aligned reads is plotted for each sample across the three treatment conditions: control (DMSO-CTRL), Sorafenib (SOR), and Sunitinib (SUN). All samples show consistent read depth.

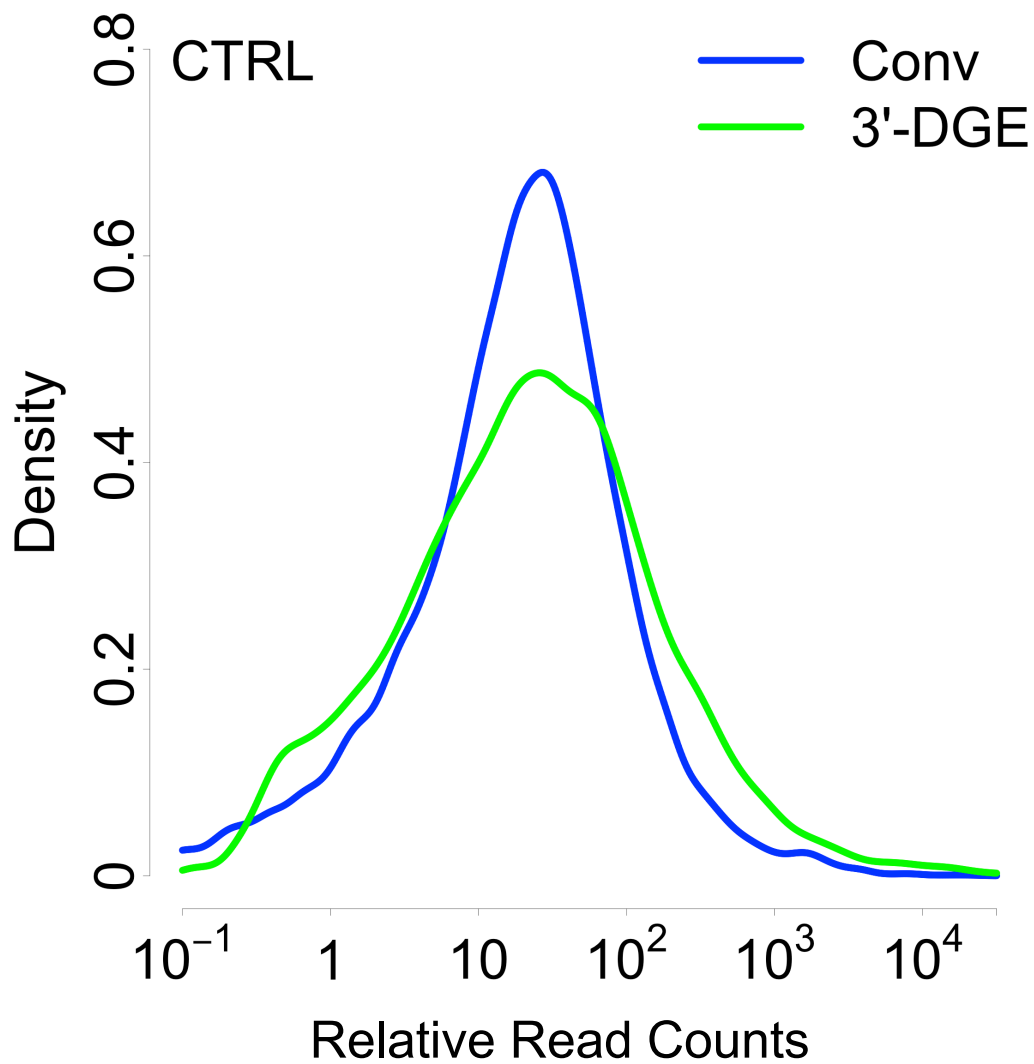


Figure S3. Relative Read Count Distributions for Conventional and 3'-DGE Methods. The mean read counts for each gene across the eight control samples, downsampled to a common read depth (2.8 million per sample) was calculated, and for conventional, this value was divided by transcript length. The probability density was estimated by the *density* function of R package *stats*.

Down Sampling of 3'-DGE/Conv Read Counts of n Genes in a Sample

Update the dependency:

- Update read depth by summing up all read counts.
- Update the number of identified genes by counting the genes with non-zero read counts

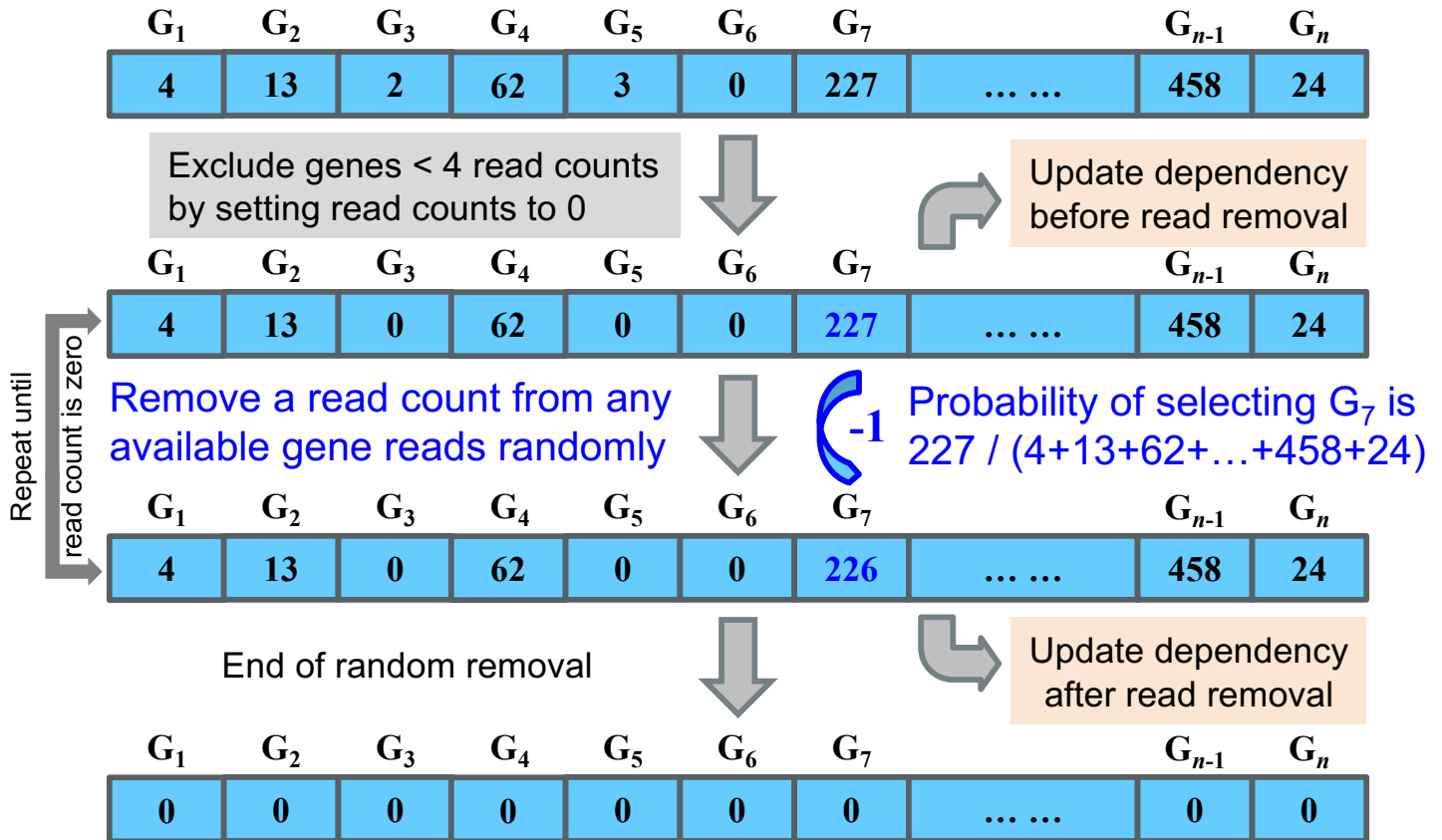


Figure S4. Down-sampling Read Counts in 3'-end Digital Gene Expression (3'-DGE) and Conventional (Conv) Sequencing Methods.

Down Sampling of 3'-DGE Read Counts & UMI Counts of n Genes in a Sample

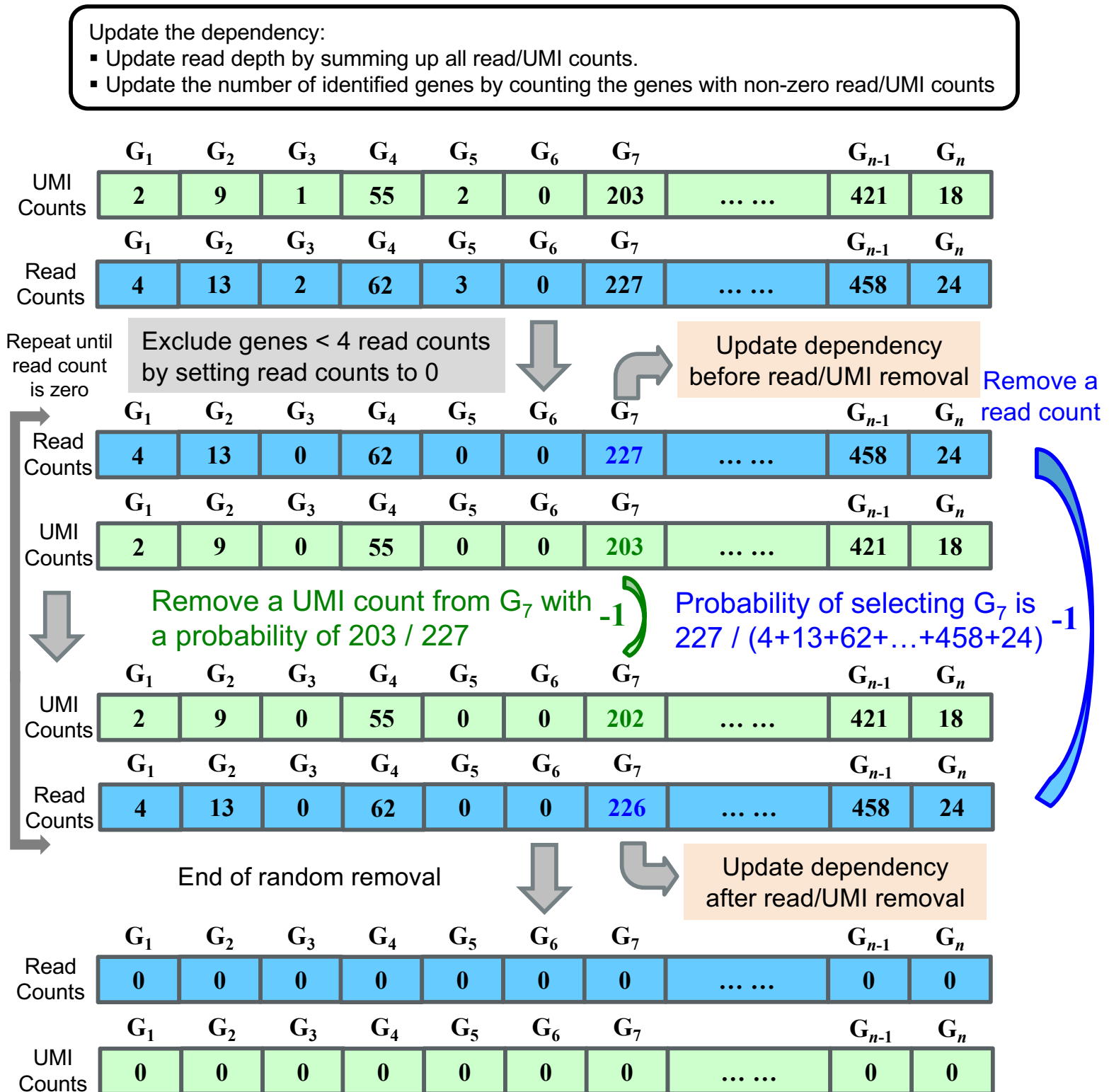


Figure S5. Down-sampling Read Counts and Unique Molecular Identifier Counts for the 3'-Digital Gene Expression (3'-DGE) Method.

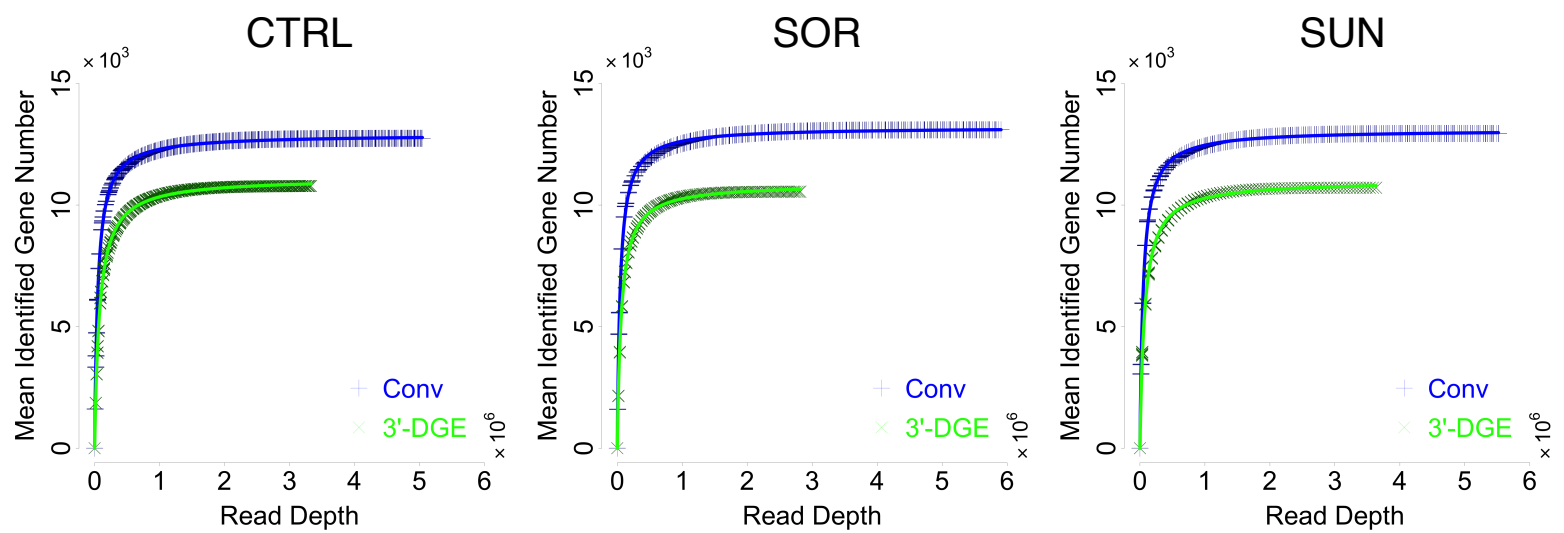
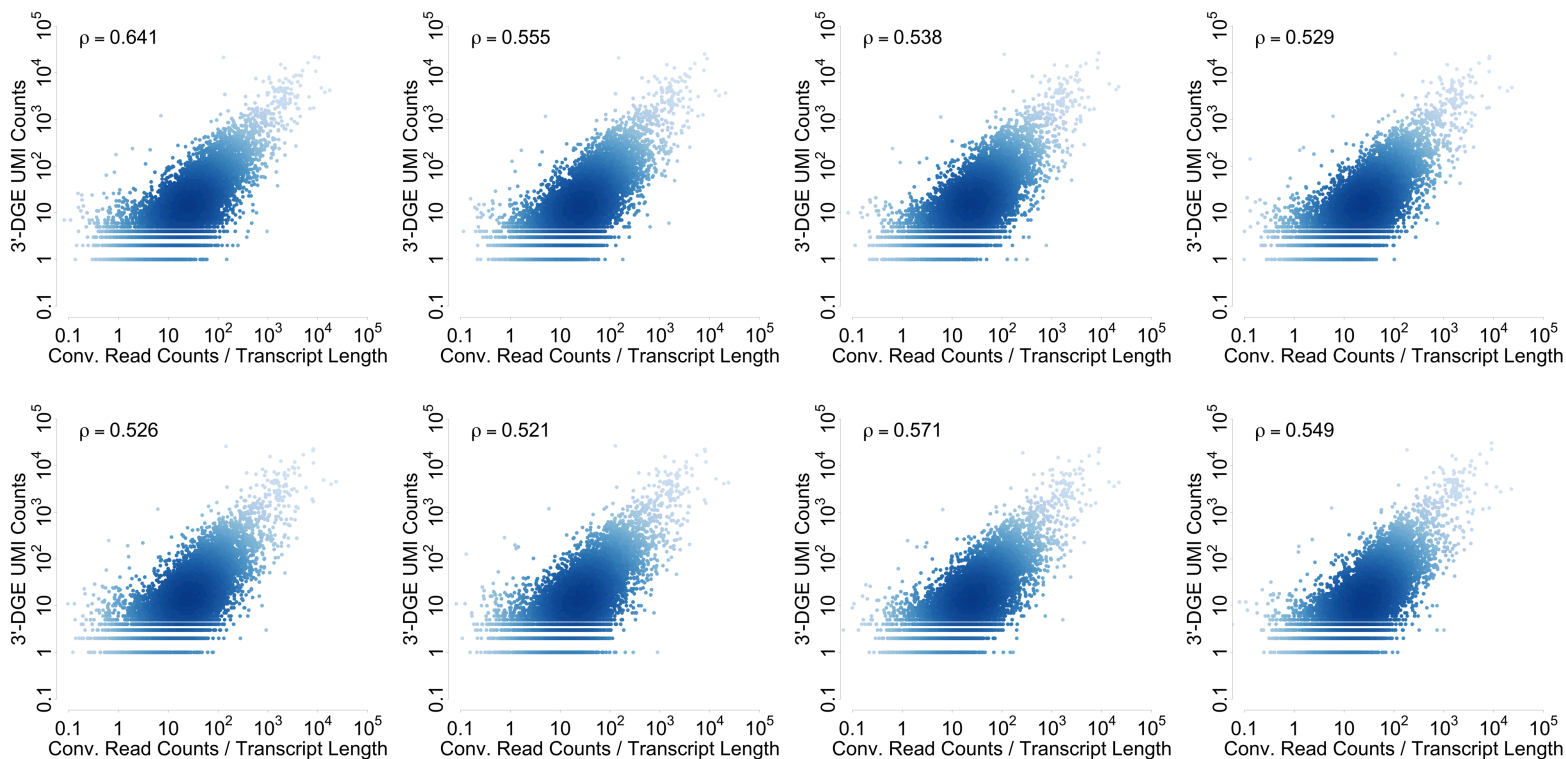


Figure S6. Variability in the Read Removal Process. Random read removal was performed 16 independent times, and the range of variability across those runs is not visible on this chart despite 16 different runs being plotted, indicating a highly reproducible simulation algorithm for read removal on the level of identified genes.

CTRL



SOR

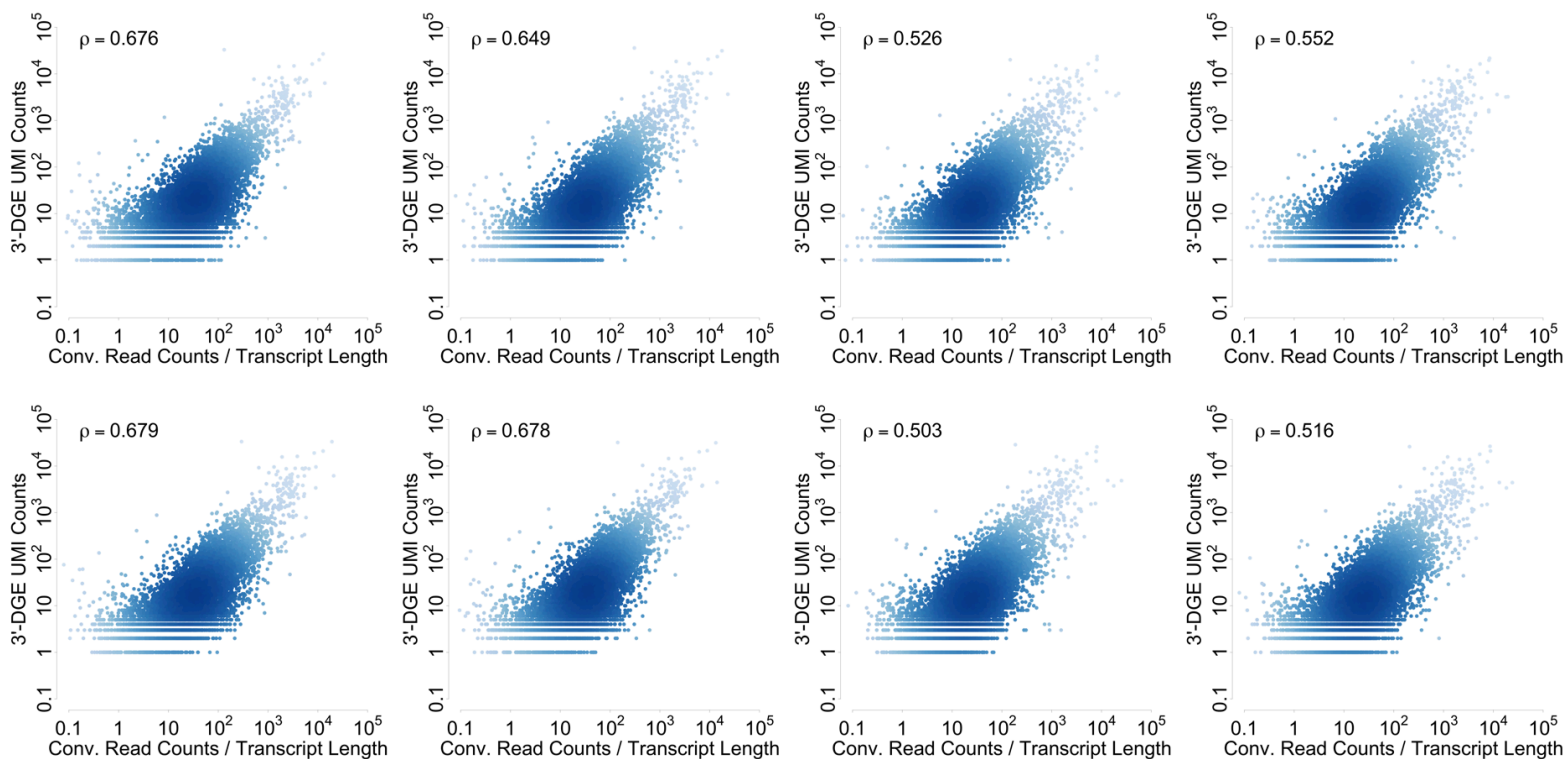
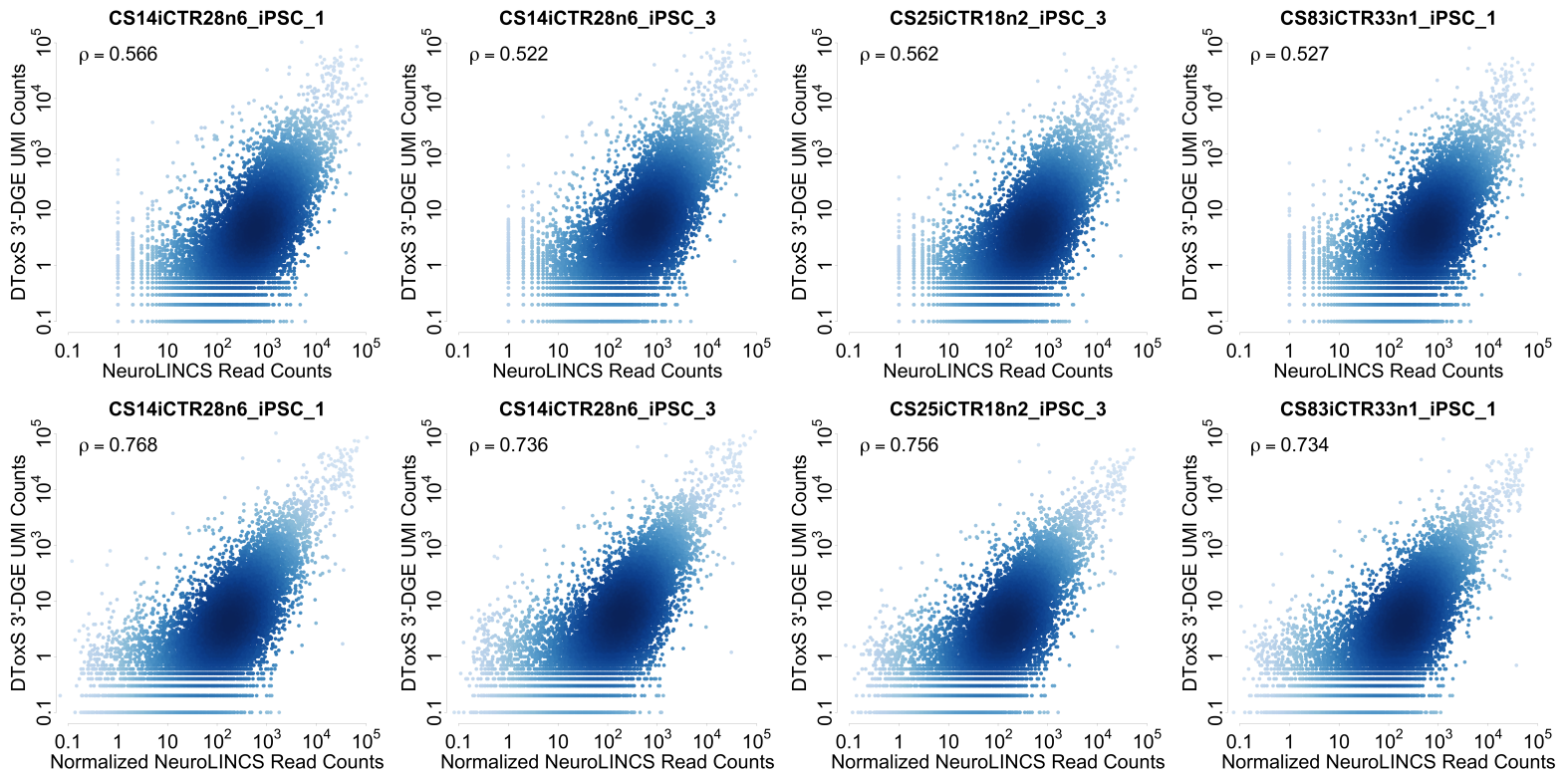


Figure S7. Sample-to-Sample Comparison between Two Techniques. Datasets are down-sampled to a common read depth of 2.8 million reads, and then gene-by-gene comparisons are made via scatter plots. To generate a reduced UMI count dataset, upon removal of a read count, UMI counts were removed with probability proportional to the ratio between UMI counts and read counts for that gene (accounting for PCR bias). Density of points in scatter plots is indicated by depth of color. Inset text box shows Pearson correlation. In all plots, data are scaled so units are comparable. There are eight CTRL samples, four SOR samples, and four SUN samples. All are biological replicates.

CTRL



SMA

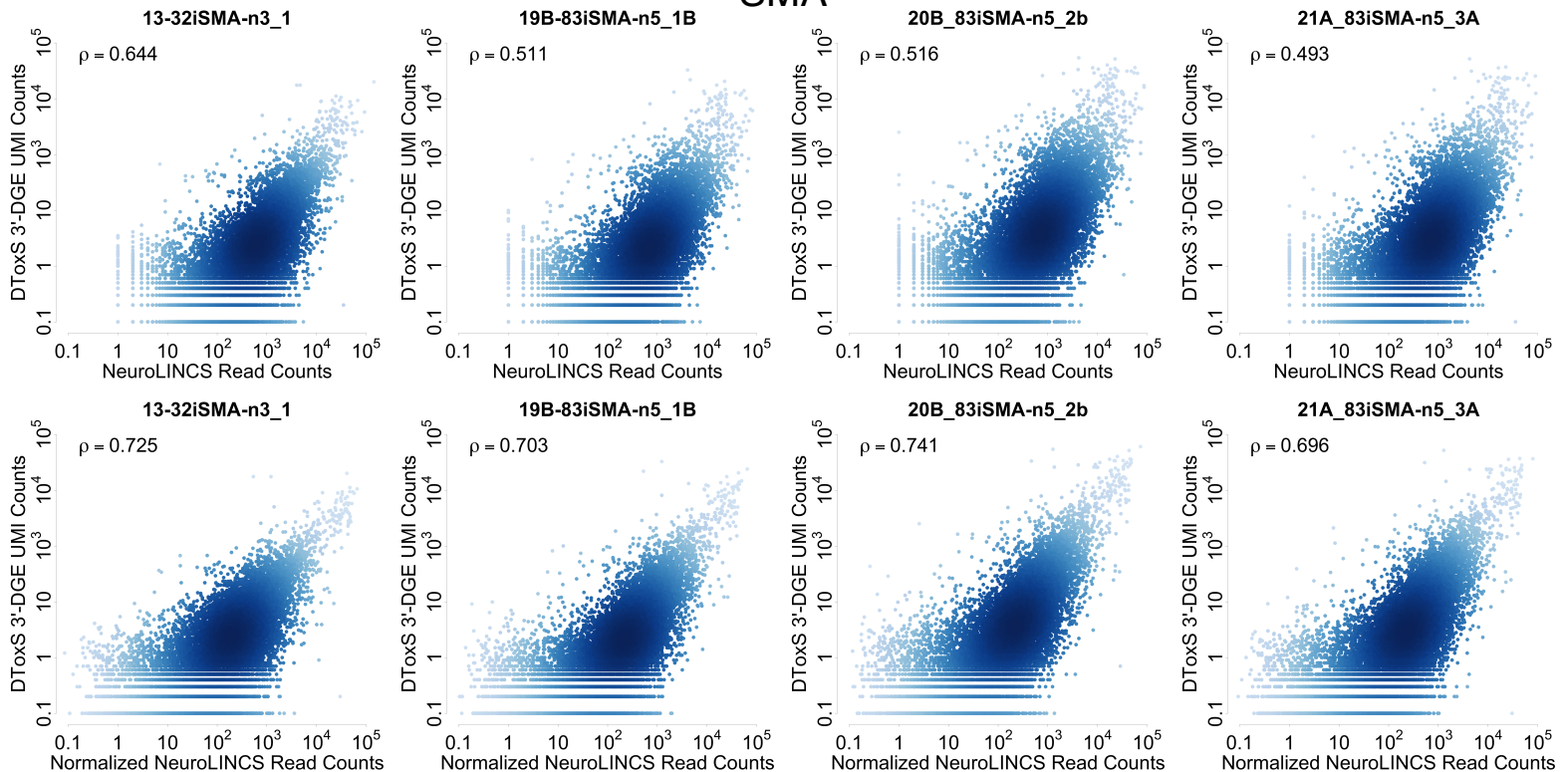


Figure S8. Sample-to-Sample Comparison between 3'-end Digital Gene Expression (3'-DGE) and Independent Conventional Techniques. Density of points in scatter plots is indicated by depth of color. Inset text box shows Pearson correlation. In all plots, data are scaled so units are comparable. There are four CTRL samples and four SMA samples. All are biological replicates.

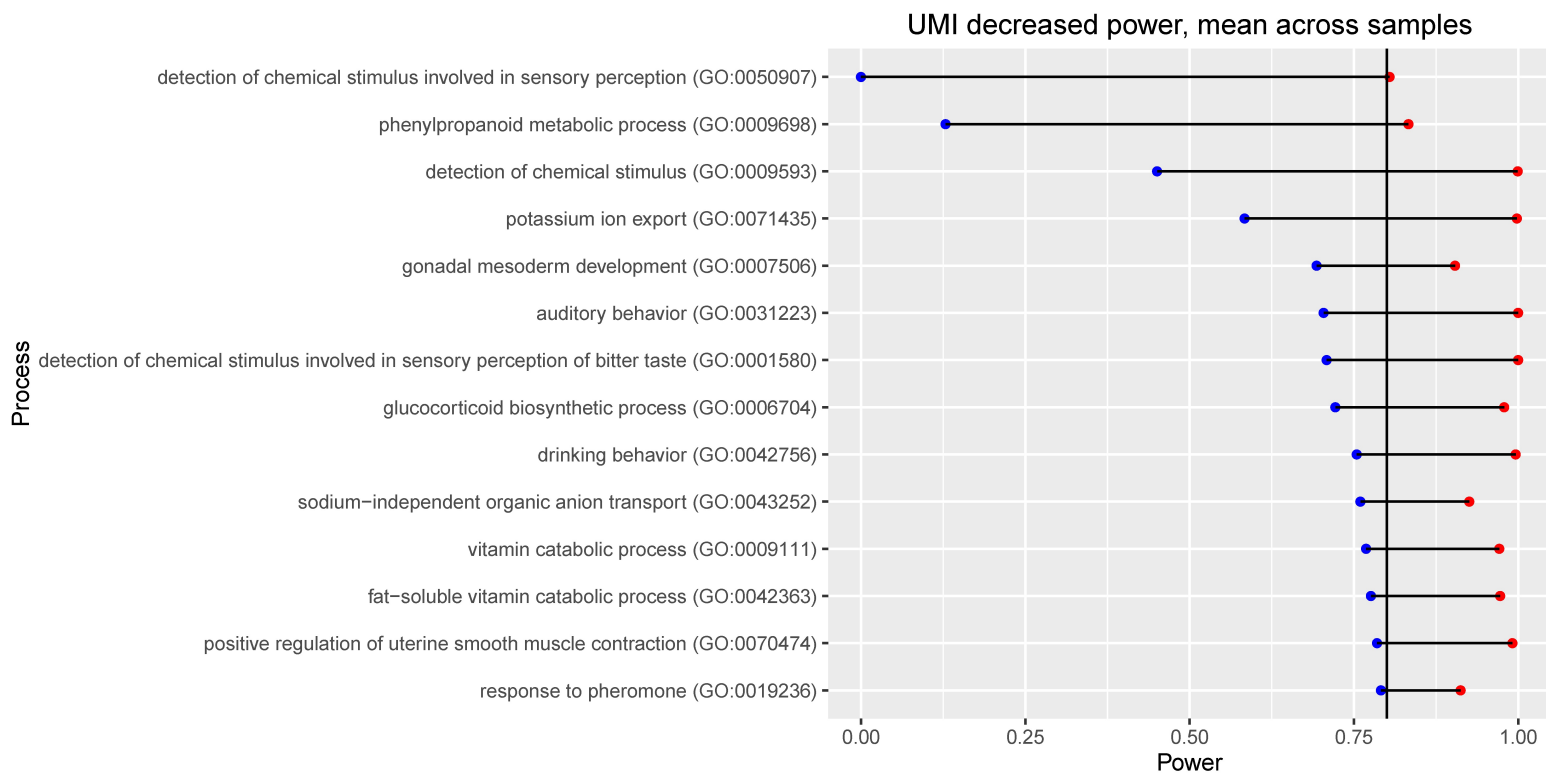


Figure S9. Power of 3'-DGE and Conventional to Detect Biological Processes Based on Genes Detected in the PromoCell Datasets. Mean statistical power across available samples to detect a biological process present in the GO Biological Process ontology using the 3'-DGE method (blue circles) and the conventional mRNA sequencing method (red circles). Shown are only biological processes where the mean power for the 3'-DGE method is <0.80 while the conventional method has a power >0.80 at a significance criterion of 0.05.

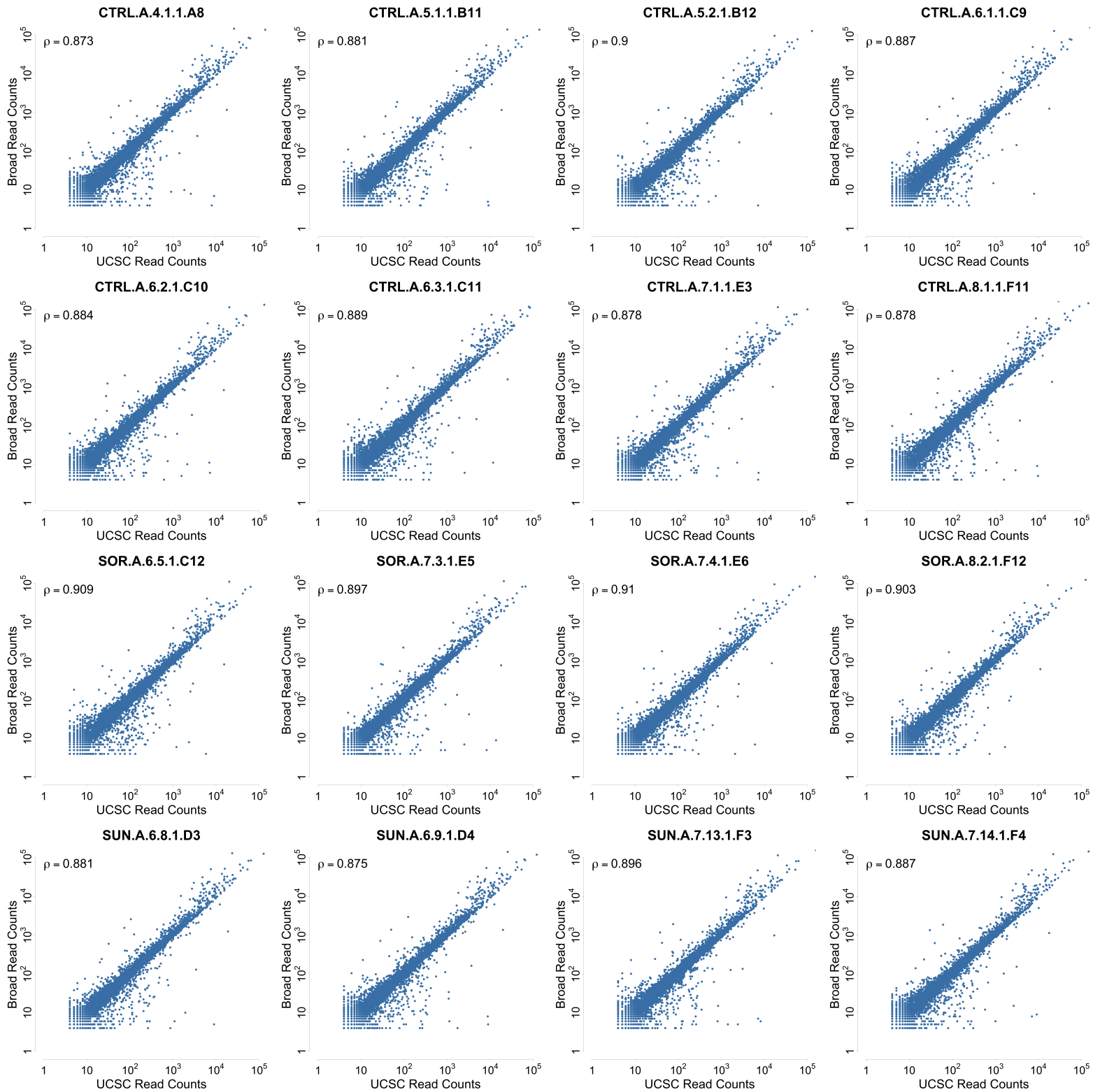


Figure S10. Comparison of read counts for 3'-DGE dataset after alignment to the UCSC RefSeq annotation used for conventional dataset, and to the mRNA RefSeq annotation provided by the Broad.

Sample Read Counts

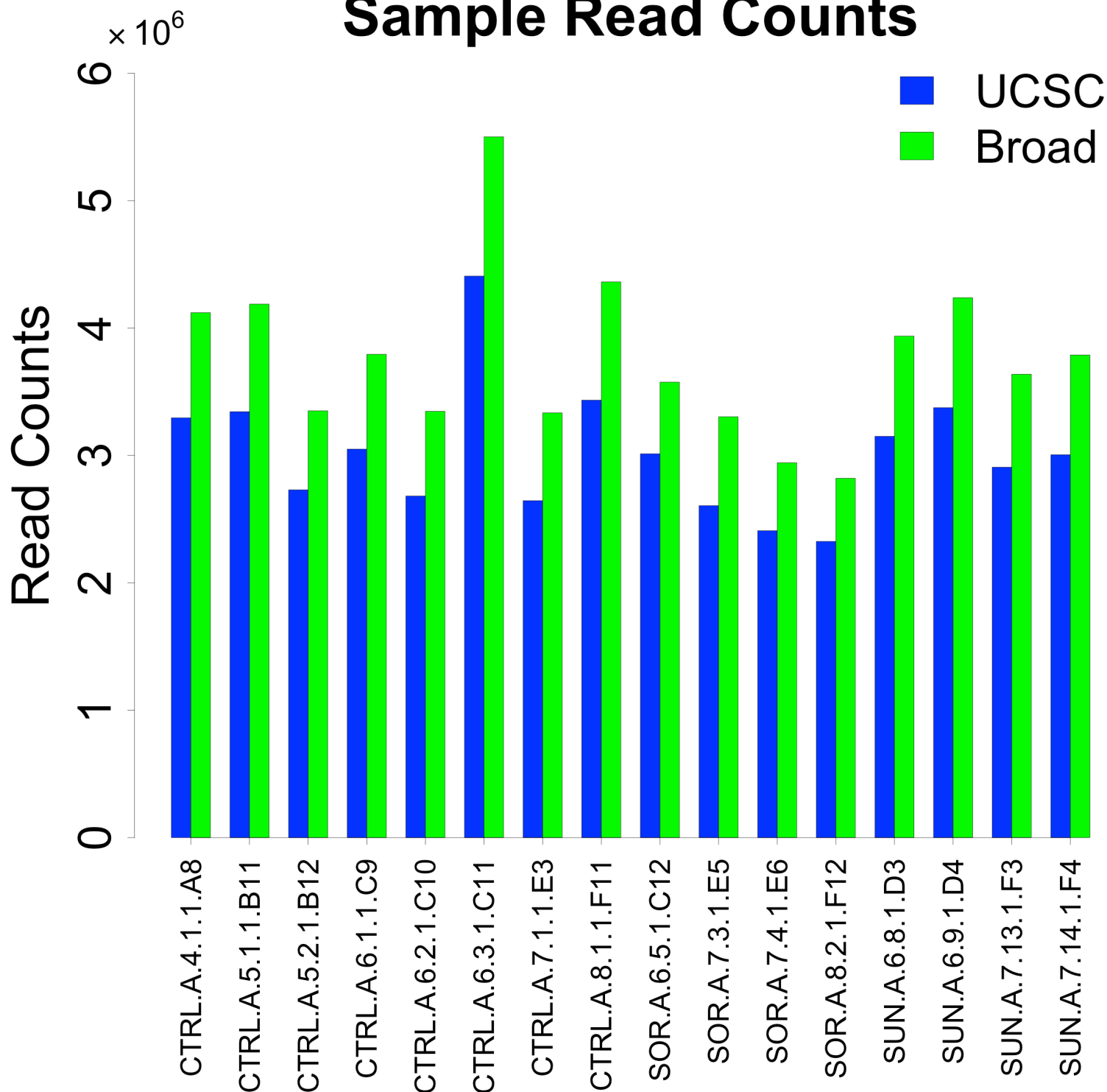


Figure S11. Sample-level read counts of 3'-DGE dataset aligned to the UCSC RefSeq (hg19) annotation versus the read counts when aligned to the mRNA RefSeq FASTA provided by the Broad. Quantitatively, the results are very similar and uniform across samples, with the Broad-provided reference generally providing more aligned reads. These increased read counts are primarily attributable to a group of 67 genes (Table S8). See Methods for details of differences.