

MCB 516a Final Project Report

Alex Borowiec, Bryan Jacobs, Kate Johnson

Reference Paper: Habowski, A. N., Flesher, J. L., Bates, J. M., Tsai, C.-F., Martin, K., Zhao, R., Ganesan, A. K., Edwards, R. A., Shi, T., Wiley, H. S., Shi, Y., Hertel, K. J., & Waterman, M. L. (2020). Transcriptomic and proteomic signatures of stemness and differentiation in the colon crypt. *Communications Biology*, 3(1). <https://doi.org/10.1038/s42003-020-01181-z>

Introduction:

The colonic epithelium's remarkable ability to regenerate itself every 3 – 5 days is maintained by balancing stem cell renewal with terminal differentiation.⁶ Stem cells, and their associated supporting cells, are located in stem cell niches at the base of folds in the tissue, known as crypts. As stem cells differentiate into mature cell populations, they move up and out of the stem cell niche into the transient amplifying (TA) compartment of the colonic crypt where they adopt a progenitor cell fate. These progenitor cells within the TA compartment will then further differentiate into final cell fates within the colon, namely enterocytes, tuft cells, goblet cells, and enteroendocrine cells. Transition from a pluripotent state within the stem cell niche to terminally differentiated cells outside of the niche coincides with dramatic changes to mRNA expression in each cell population. However, it was unknown if changes in mRNA expression are an early event in stem cell differentiation (i.e. when stem cells first differentiate into progenitors within the TA compartment) or if other processes, such as mRNA processing, alternative polyadenylation, or protein expression have a more significant role.

This paper sought to understand these questions through a multiomics approach that looked at mRNA expression levels, mRNA splicing, polyadenylation, post-translational modifications, and proteomics across all cell populations within the colon. To achieve this multiplexing the authors optimized a Fluorescent Activated Cell Sorting (FACS) protocol using well established markers of each cell population on mouse colons to enrich each cell type found within the colon. In total, six populations of cells were sorted for downstream analysis; three of which were terminally differentiated mature cells (enterocytes, enteroendocrine cells, and tuft cells), two progenitor populations that give rise to mature differentiated cells (secretory progenitors and absorptive progenitors), as well as the colonic stem cells that give rise to the progenitor populations. With this optimized FACS protocol the authors then comprehensively analyzed each cell population using their multiomics approach. Interestingly, this paper did not use single cell RNA

sequencing, however their FACS approach allows for a pseudo-single cell analysis since each population that was enriched is one cell type.

Through these investigations the authors found that the initial transition from stem cell into progenitor cell is not driven by large changes in mRNA levels. Rather, alternative splicing, alternative polyadenylation, and protein changes dominate this early transition. For example, the authors showed that the transition from stem cell to absorptive progenitors coincides with over 900 skipped exon alternative splicing events. In contrast, their data showed that only 301 mRNAs showed significant changes in levels comparing stem and absorptive progenitor cells; 3-fold lower than alternative splicing events. Additionally, the authors were able to use their multiomics approach to tease out fine details about previously well studied signaling pathways in the intestines such as Notch signaling. It has been well established that the lateral inhibition from Notch signaling in stem cells influences whether they differentiate down a secretory or absorptive lineage.⁵ Typically, secretory cells express Notch ligands whereas absorptive cells express Notch receptors.⁷ The authors showed in this study that early on in commitment to these fates there is alternative splicing of the Notch regulator *Spen* between stem cells and their progenitors. Only 50% of stem cells maintain full length *Spen*, the other 50% are missing one RRM RNA binding domain. In contrast, 100% of progenitor cells show full length *Spen* – indicating a switch in splicing early on in the transition from stem to progenitor with a subsequent shift in mRNA levels as maturation of each cell type occurs.

The authors of this paper were also able to gain new insight into how the regenerative capacity of the colon is maintained during injury. One fascinating observation in this area is that stem cells are not strictly required to regenerate the intestines after injury.² Current models propose that some level of de-differentiation occurs in the colonic epithelium when the tissue is damaged, where differentiated cells revert into a plastic state and can replace stem cells that were lost during injury.² The authors found significant support for this model in their data stating that, “*...stem cells are not so much defined by what they express, but by what they do not express.*” They argue, and their data supports, the view that classic markers of stem cells and pluripotency, such as *Lgr5*, *Smoc2*, and *Cd44*, are also expressed in differentiated cells. The difference between the populations is that mature differentiated cells express more mRNA compared to stem cells, and that this final difference in expression is driven by large changes in early alternative splicing.

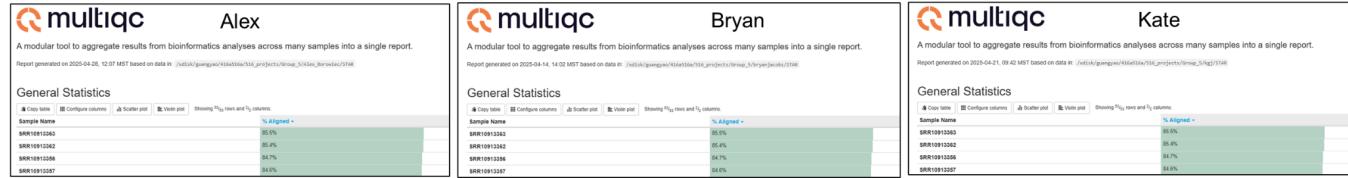
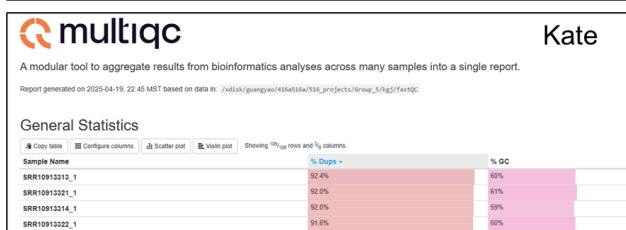
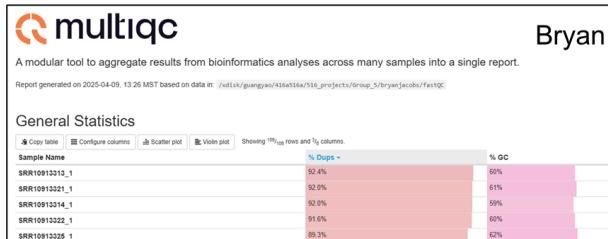
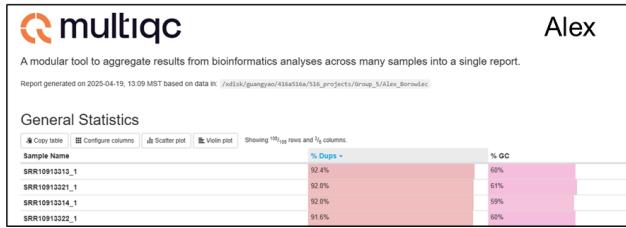
While the authors of this paper were focused on the initial transition of a stem cell into a progenitor cell, our group focused our analysis on comparing stem cells to fully differentiated enterocytes – a cell population that is classically thought of as the furthest from stemness. We chose this comparison since mRNA levels do not change dramatically when a stem cell becomes a progenitor cell. We reasoned that comparing stem cells and enterocytes would give us large changes in mRNA expression that we could run our analyses on. Our group was able to successfully reproduce figures and data from the reference paper with good agreement. We also found the data generated in this investigation to be robust to pre-processing and alignment – i.e., even though we changed parameters during RNAseq data processing our downstream analysis still somewhat recapitulated the results of the paper.

Results and Discussion:

HPC Analysis:

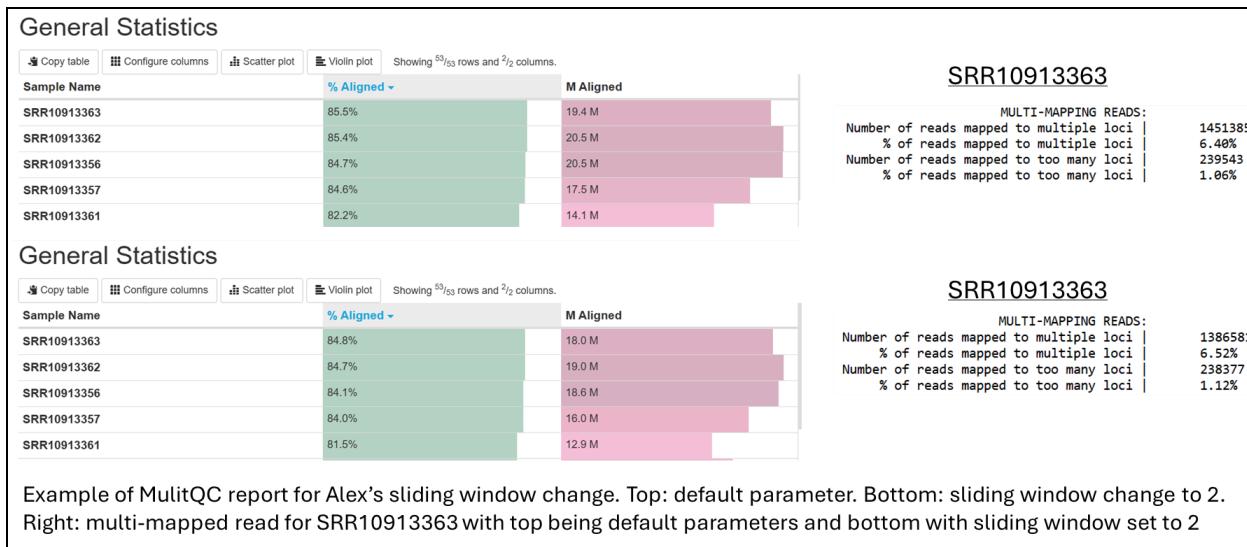
FastQC and STAR MultiQC

Each group member independently performed MultiQC on fastQC and STAR using default parameters. MultiQC reports for each group member are identical for both fastQC and STAR. For fastQC, percent duplicate reads ranged from 39.2% - 92.4% across all samples analyzed and total sequences ranged from 15.0 - 34.1 million. Overall, our fastQC results were normal for basic statistics, per base sequence quality, per sequence quality scores, per base N content, sequence length distribution, and adapter content. However, indication of abnormalities in the data was detected for per base sequence content, per sequence GC content, sequence duplication levels, and overrepresented sequences. For multiQC on STAR alignment using default parameters, the percentage of uniquely mapped reads ranged from 45.9 - 85.5% across all samples and uniquely mapped reads ranged from 7.5 - 21.9 million. The high percentage of duplicate reads from fastQC results could potentially be a result of how the RNA was prepared for sequencing in these experiments. We investigated the kits used to isolate the RNA as described in the paper's methods (Direct-zol RNA Micro-Prep kit Zymo #11-330 M & Clontech Low Input Pico Kit Takara #634940). As far as we can tell the authors did not specifically enrich mRNA or deplete other RNAs from the samples such as rRNA. Excessive rRNA in the sample preparation could lead to high percent duplicate reads.

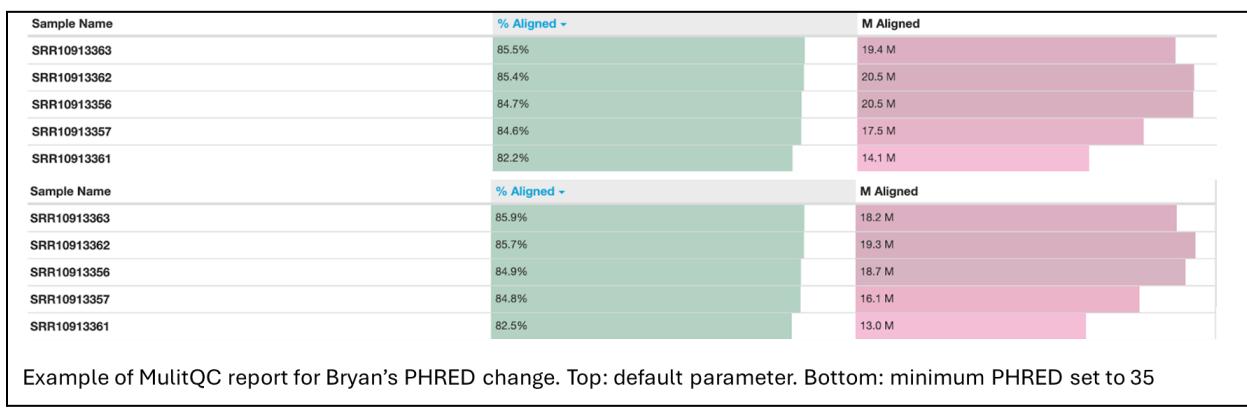


Example MultiQC reports for all group members for STAR default parameter results

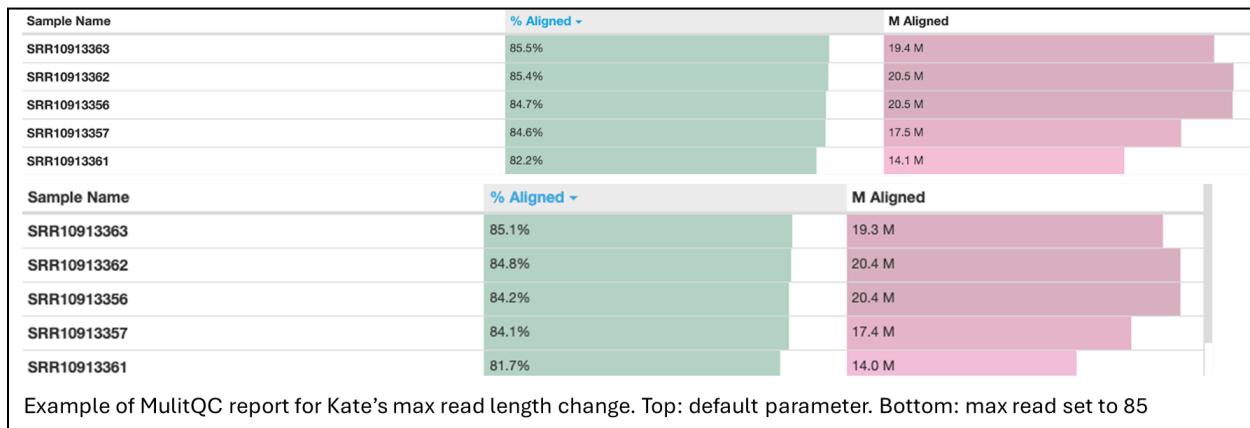
Sliding Window Parameter Change. For the sliding window parameter in fastp, we chose to increase the window size from 1 to 2. The analysis should then average the quality of 2 consecutive bases instead of evaluating each base individually. This is a modest increase in size and from it we expected slightly less trimming of the reads. With this parameter change we did see modest changes to STAR alignment. For default parameters we saw 85.5% uniquely mapped reads for SRR10913363, with the increased window size we observed a decrease in that number to 84.4%. This decrease in uniquely mapped reads with the increased window coincides with an increase in the number of reads mapped to multiple loci: 6.40% in default versus 6.52% with parameter change. These results are consistent with how the sliding window affects trimming of reads. Increasing the window to 2 averages two bases which could lead to the retention of repetitive or ambiguous bases that would normally be trimmed. This could lead to the mapping of reads to multiple loci in downstream analysis.



Min PHRED Score Parameter Change. We chose to increase the minimum PHRED score parameter in fastp to 35, from its default value of 20. While PHRED scores of 28 and above are generally considered high quality, we chose a much stricter threshold due to the majority of our reads falling above the threshold of 35. In general, with a threshold this high, we would expect to see very aggressive filtering and therefore a large decrease in the number of mapped reads in STAR alignment. However, due to the high quality of the data, we saw modest changes across samples. With default parameters, SRR10913363 had 19.4 million aligned reads, and 85.5% of reads aligned. With the PHRED score increased to 35, SRR10913363 had 18.2 million reads aligned and 85.9% of reads aligned. This trend was generally seen throughout the samples, and can be explained by the stricter PHRED threshold filtering out more reads, increasing the overall quality of the remaining reads. It is important to note that while the remaining data leads to a slight increase in alignment percentage, a PHRED threshold of this strictness has potential to harm the overall coverage of the data.



Max Read Length Parameter Change. The third parameter we changed for fastp was maximum read length using the command `--max_len1 # --max_len2 #`. Setting a limit for maximum read length means trimming all longer reads on the 3' end to the specified length, which removes any low quality bases for long reads. Reads shorter than the specified length remain unchanged. When we initially capped at 100 base pairs, there was no change in the fastp results or PCA or heatmap plots, suggesting all read lengths were less than 100 bp for the original PCA. We then tried capping at 10 bp, which made the reads too short for alignment using STAR. Next, we attempted capping max read length at 40 bp, which was also too short; we then realized reads need to be a minimum of 50 bp for STAR, and ideally between 75-100 bp. With this information, we finally settled on a maximum read length of 85 bp which subtly altered our fastp results and subsequent plots.

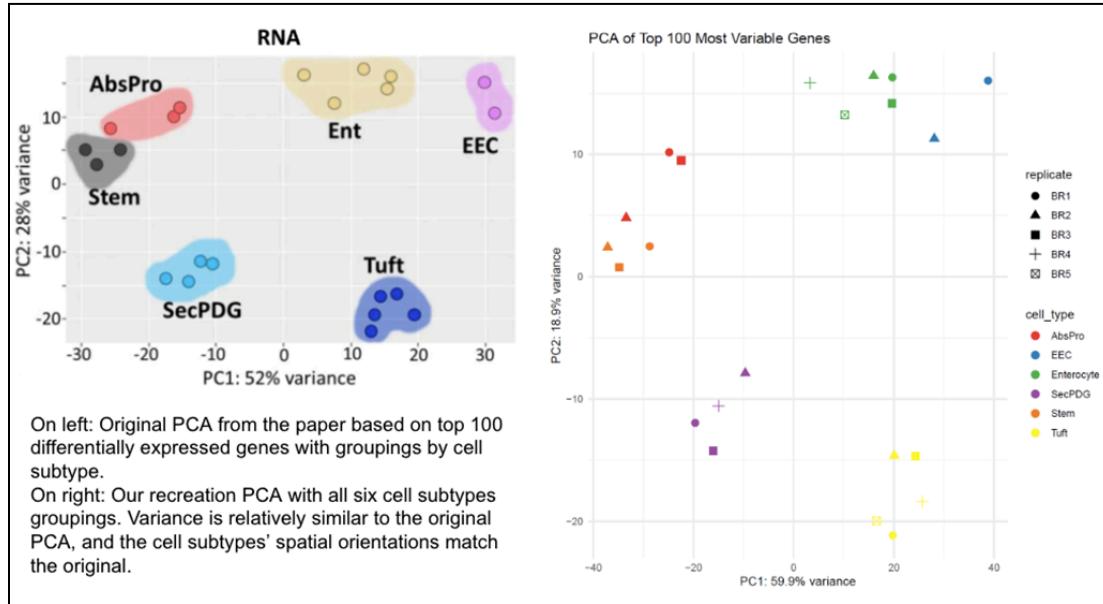


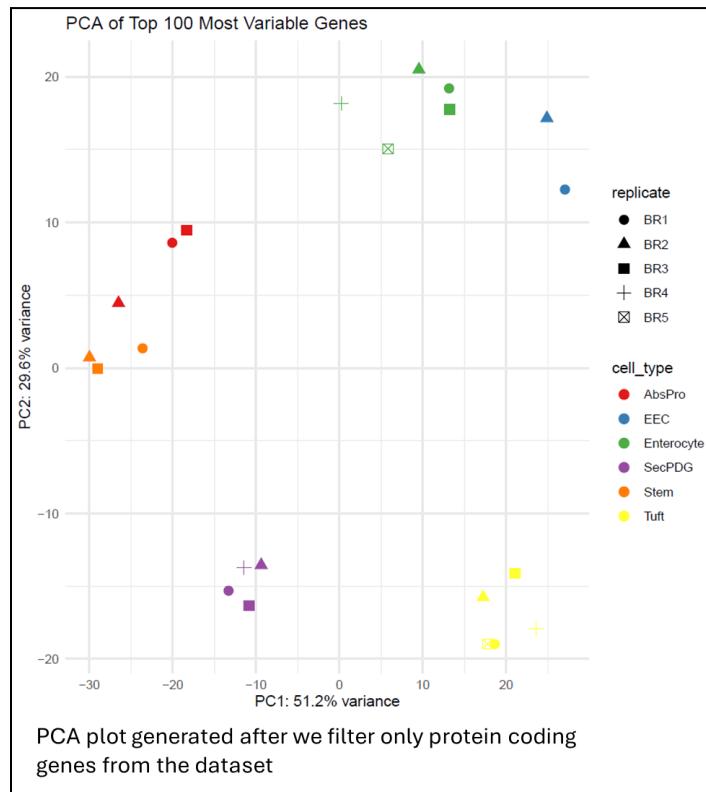
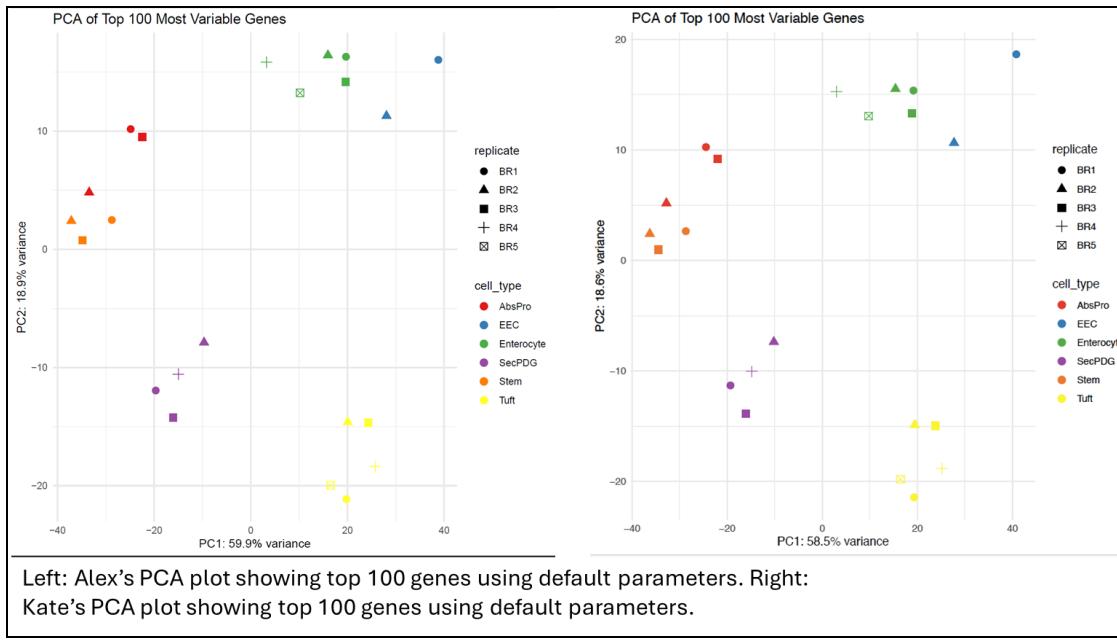
Differential Expression Analysis:

PCA

PCA maps generated using default parameters matched closely with what the authors published in the reference paper. There were changes with variance explained percentages and some of the grouping was slightly different. We attribute these differences to different versions of DESeq2 used in the paper versus what we used in our analysis (reference paper DESeq2 = 1.16.1. Our version = 1.48.0). The other area that can explain the difference we see comes from how the authors defined biological replicates as well as how the authors filtered genes after STAR. In total there were 53 files for this paper, however only 22 biologic replicates were used in the paper (stem = 3, AbsPro = 3, SecPDG = 4, tuft = 5, Ent = 5, and EEC = 2). Eventually, we

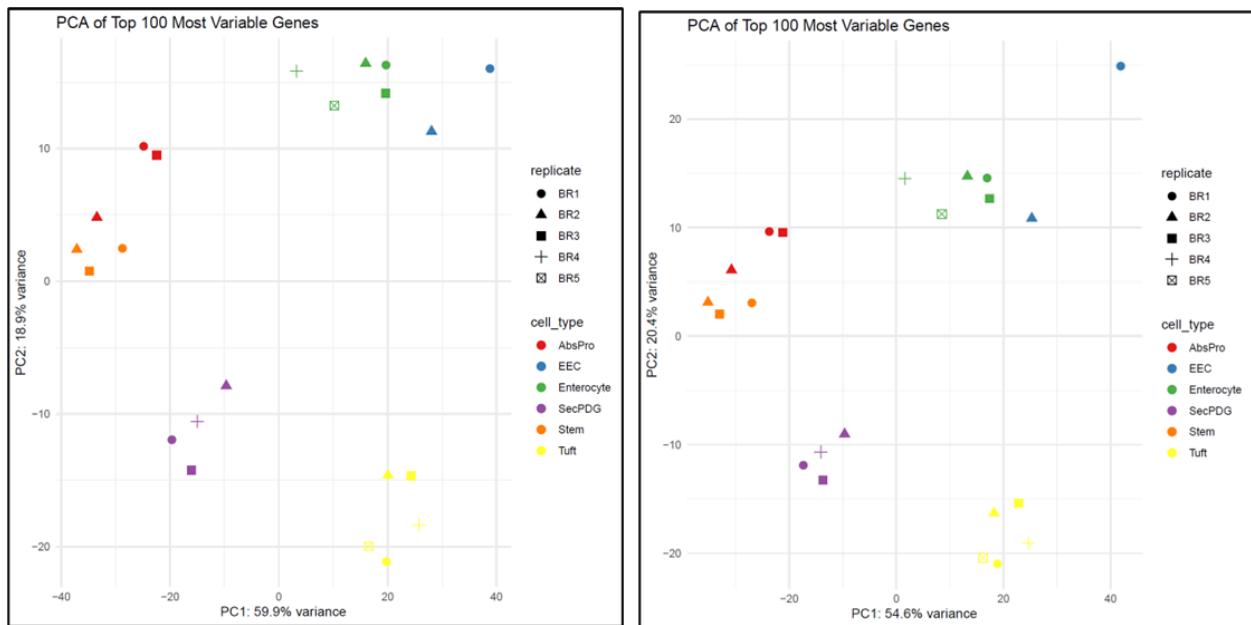
found that each of the 22 biologic replicates had differing numbers of technical replicates within each biological replicate. The authors did not directly address how technical replicates were handled in their work. For our analysis, we kept one of the technical replicates that had the highest counts for each biological replicate. If the authors used all 53 samples in their analysis versus our 22 then this could explain the slight difference we see in the PCA. In addition to replicates, we also found that at some point during the author's pipeline they specifically filtered out protein coding genes for DESeq2 analysis. The raw counts from their dataset maps to over 70,000 genes, however their final published DESeq2 results show a little over 20,000 genes - indicating that the genes were parsed at some point that was not addressed in the text. We made this discovery after initializing our individual data analysis. To confirm that our results using the full geneset still matched what the authors published, we filtered protein coding genes from our dataset and made the PCA plot again. After filtering, the number of genes in our dataset was reduced down to 21,748 genes which matched closely with the authors 24,421. Our results from this filtered set were similar to our original PCA plots using all values. This gave us confidence to keep and move forward with our analysis and we chose to keep the full gene list for our analysis. Comparing PCA plots generated by Alex and Kate using default parameters show similar plots with slight differences to the variance explained percentages on the x- and y-axis as well as grouping for EEC. We were careful to run our analysis in a similar manner, however the difference in grouping on the PCA could potentially come from the slight differences in variance explained that was calculated in R.





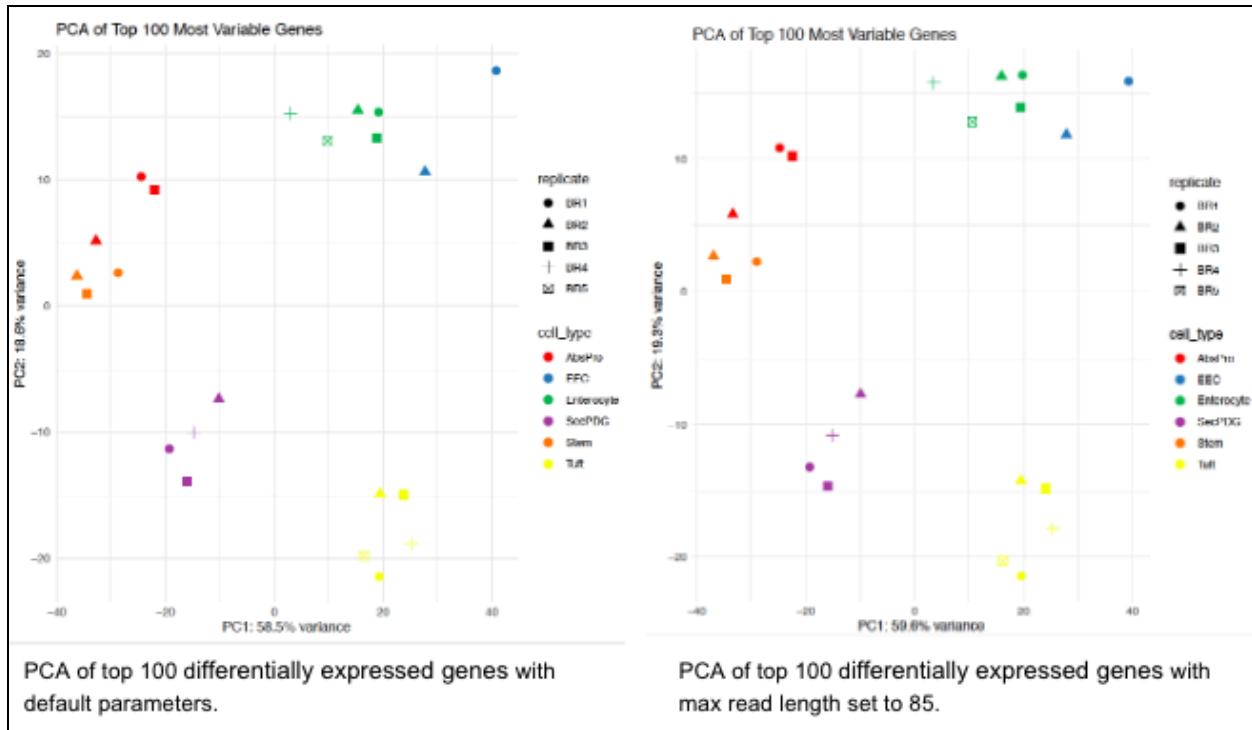
Sliding Window Parameter Change. Changing the sliding window from 1 to 2 showed modest changes to the grouping of the samples in the PCA plot. For secPDG cells the clusters became

tighter. Indicating that for secPDG the larger window improved the biologic signal potentially due to a reduction in noise when averaging bases. On the other hand, for the EEC population, one of the biologic replicates no longer clusters with the other biologic replicate. This could be due to the averaging losing single bases resolution that is important for this cell type. The paper discussed how the EEC population was the most difficult to sort and enrich for this analysis due to low cell numbers. They had to use multiple mice to get sufficient cells for RNAseq. It could be that the increased window is exaggerating inherent sample to sample variability due to the difficulty in enriching this population.



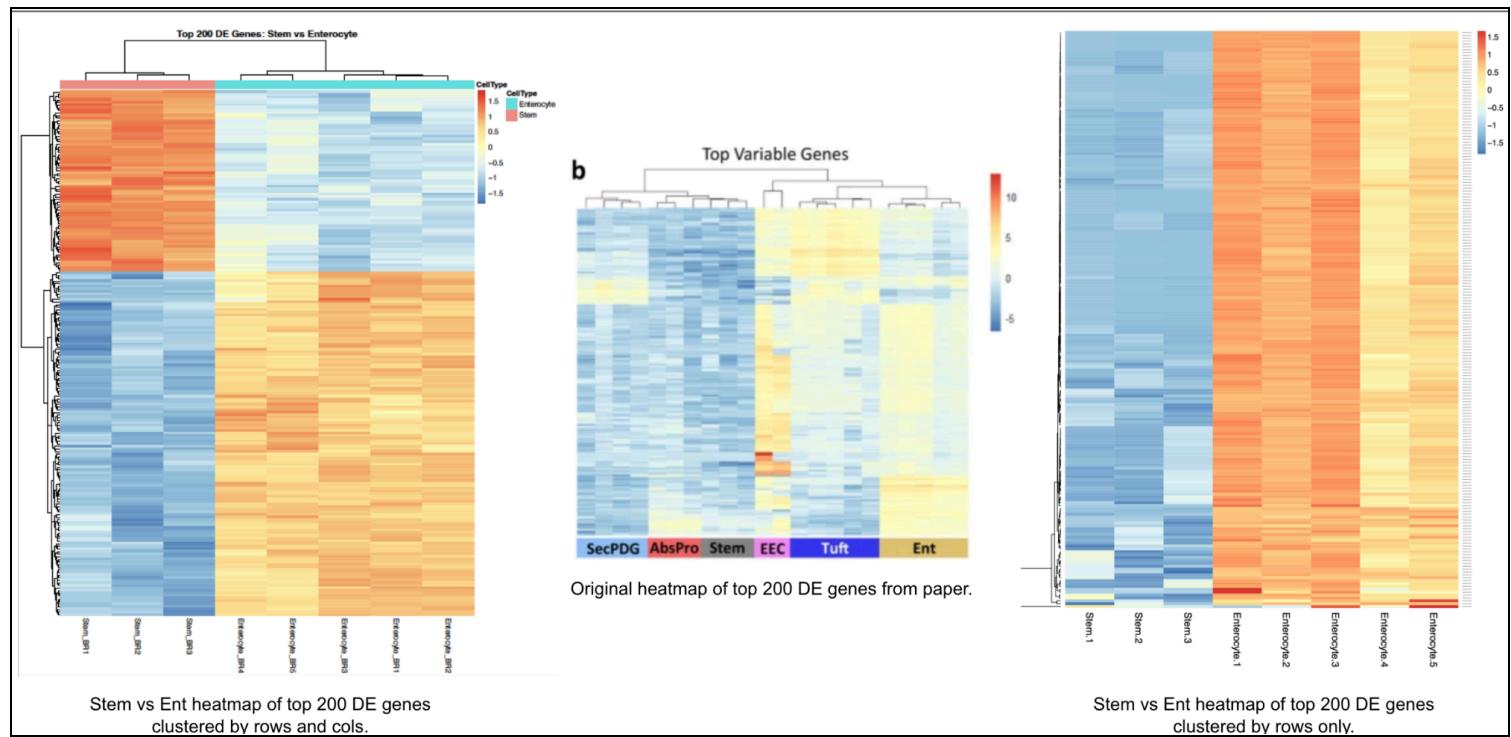
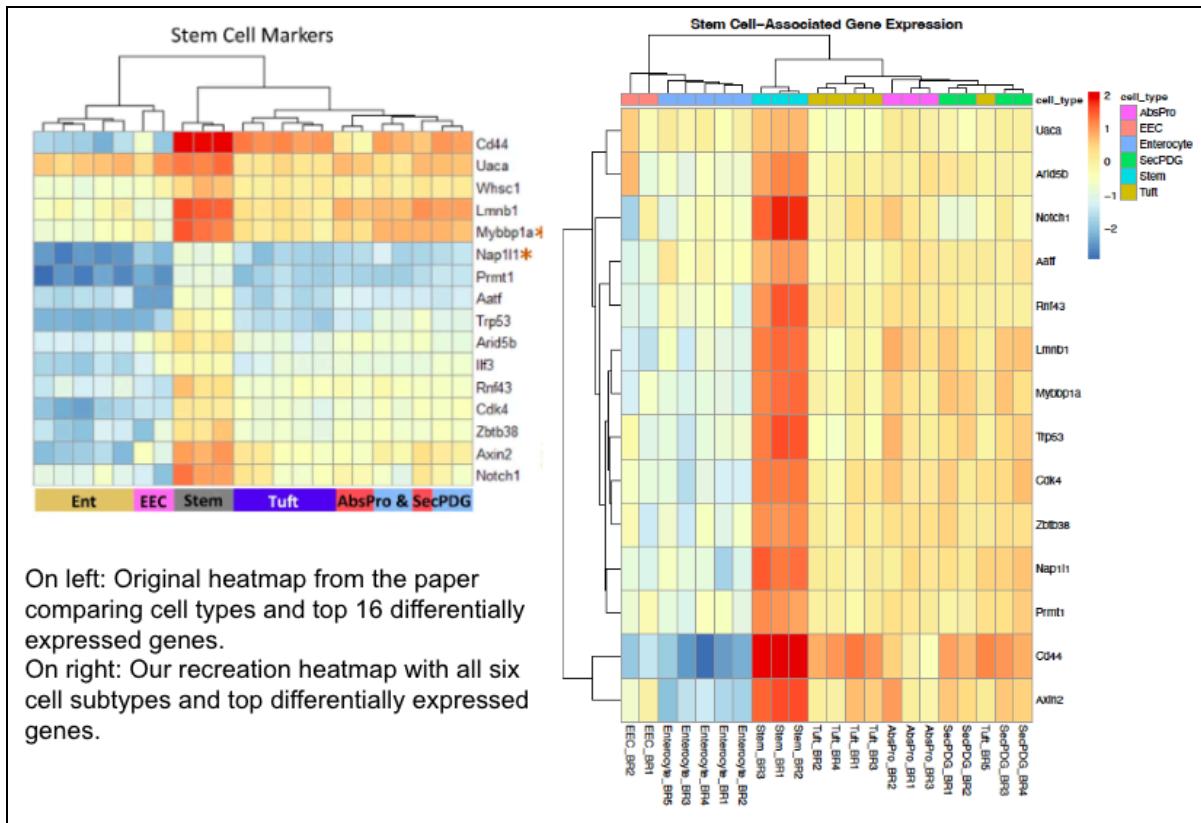
Left: PCA at default parameter of 1. Right: PCA with sliding window adjusted to 2.

Max Read Length Parameter Change. The second parameter we adjusted for the PCA plot was maximum read length, which we set to 85 bp. This changed the plot very slightly, most noticeably by shifting the two EEC samples closer together. Overall, the minimal changes suggested the original full-length reads had good integrity, but the trimming did improve clustering for the enteroendocrine cell subtype.

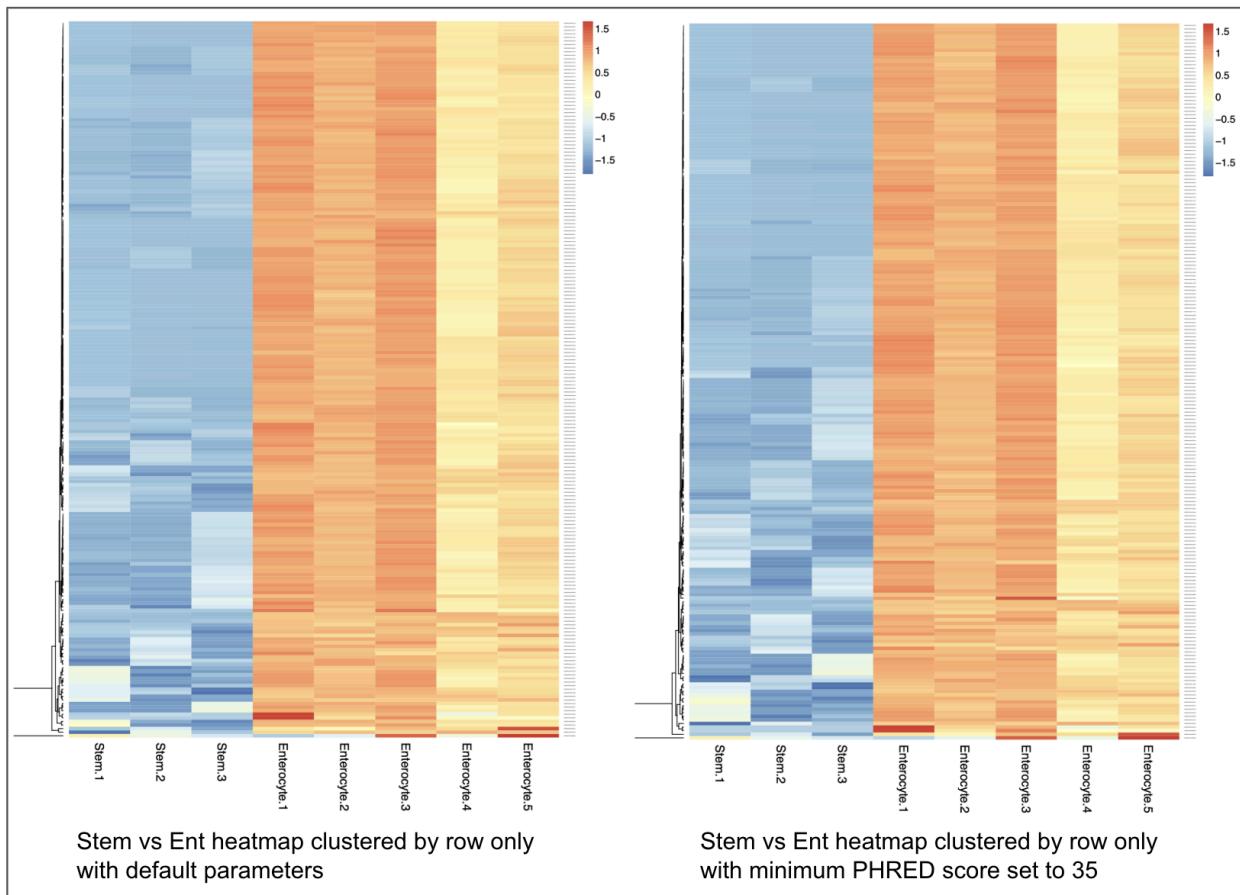


Heatmap

There were three heatmaps in this paper, each segregated by cell type. The first looked at the top 200 differentially expressed genes in all six cell subpopulations (secretory cells, absorptive progenitors, stem cells, enteroendocrine, tuft, and enterocytes); this was autoscaled with unsupervised clustering. Then, the top 16 genes expressed in stem cells were selected and depicted across all six cell subtypes again. We recreated this heatmap to the best of our ability below. Next, the researchers selected the genes that were most changed in the non-stem cell populations of secretory cells and absorptive progenitors. These 17 genes are depicted in stem cells, SecPDG, and AbsPro in the final heatmap. For our purposes, we chose to recreate the heatmap comparing the top 200 differentially expressed genes, but focusing on stem cells and enterocyte populations. We have one Stem vs Ent heatmap that was clustered according to rows and columns, and one clustered only by row; this clustering variation accounts for differences between the default parameter heatmaps. Focusing only on the Stem and Ent columns from the paper's heatmap, our heatmaps can be considered satisfactory recreations.

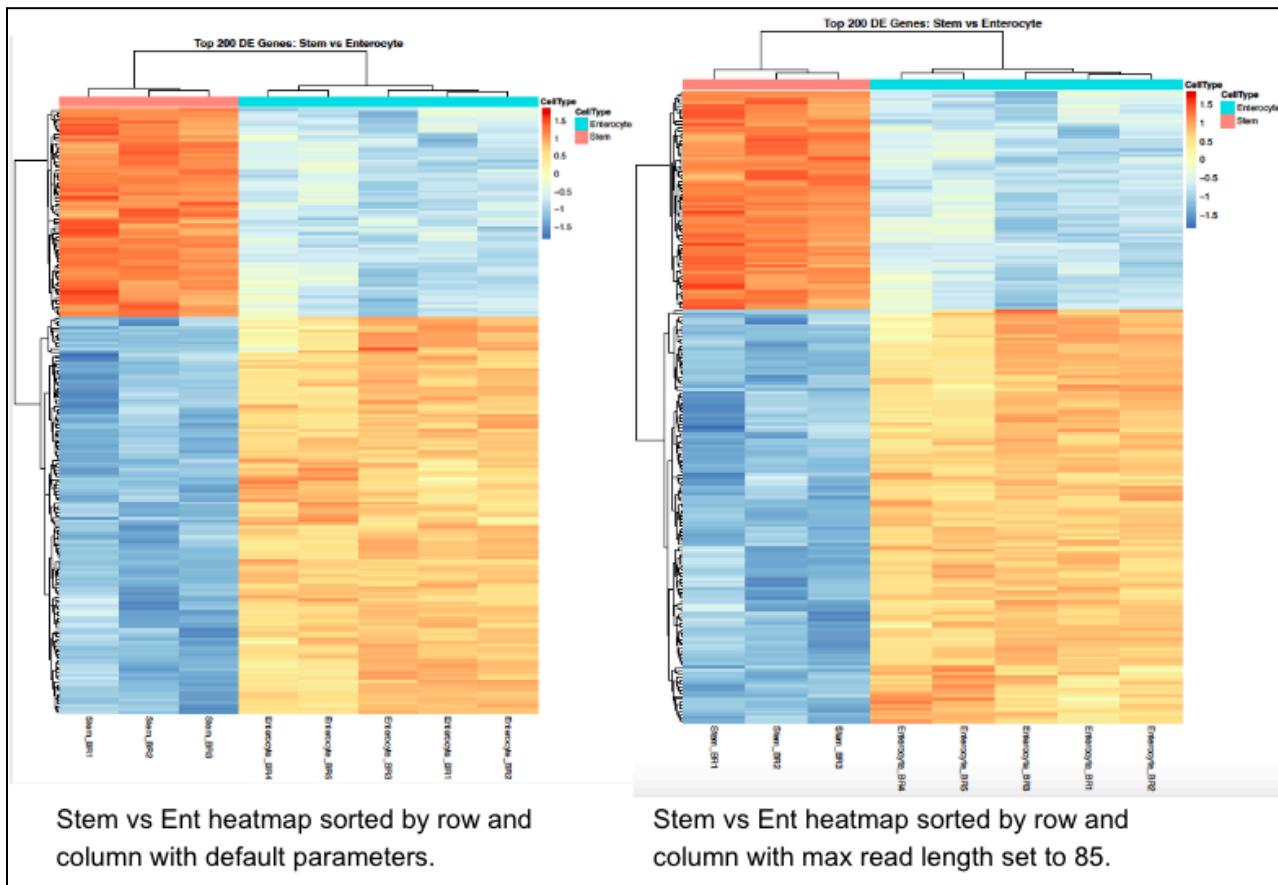


Min PHRED Score Parameter Change. Increasing the minimum PHRED score from 20 to 35 showed very minor changes in the heatmap of the top 200 differentially expressed genes in stem cells vs enterocytes. Since the parameter change eliminates only lower quality reads, the top 200 differentially expressed genes remain biologically consistent across both pipelines. The only difference we see is brighter expression in some areas of the modified PHRED score heatmap due to the elimination of lower quality reads and therefore statistical noise. This is a satisfactory result, as the goal of increasing the minimum PHRED score parameter is to refine results rather than drastically alter them.



Max Read Length Parameter Change. Setting the maximum read length to 85 base pairs did not obviously change the results in our heatmap for top 200 differentially expressed genes in stem cells vs enterocytes. Because the original data was already clean and reliable, the majority of clustering and overall patterns remained the same in the modified max length heatmap. However, the more stringent parameter did likely reduce noise and increase integrity of our

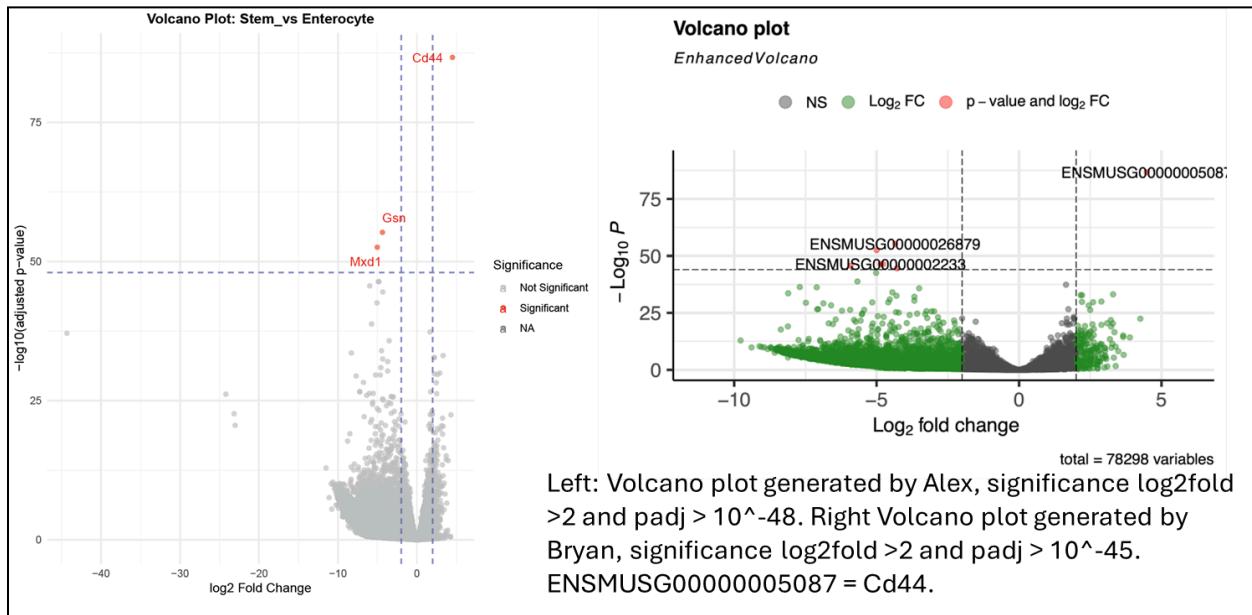
results, and there are some stripes in the heatmap that have shifted around to reflect the more defined clustering.



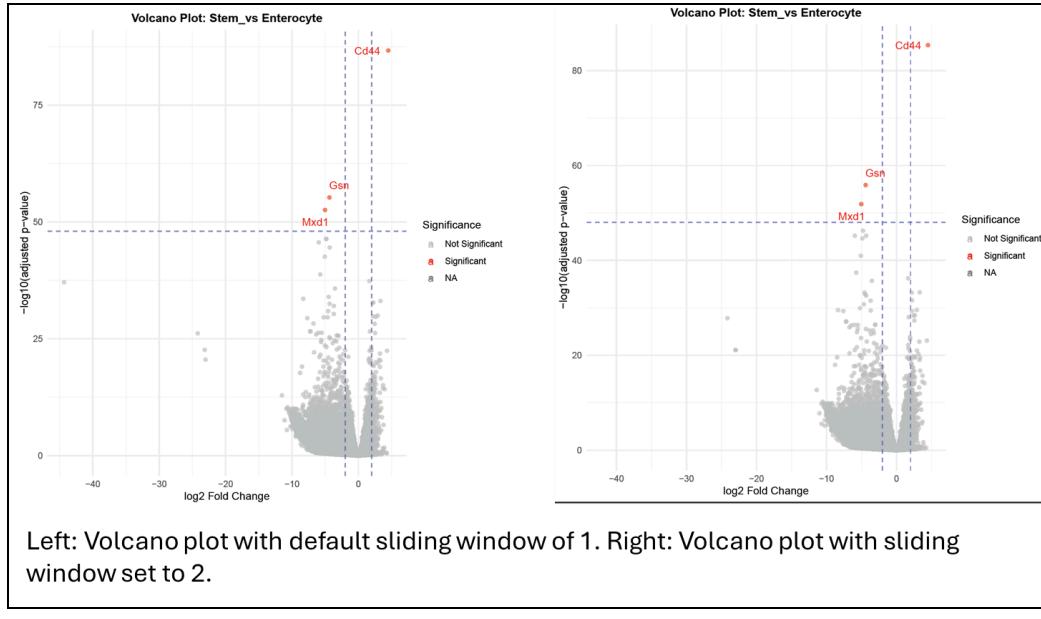
Volcano Plot

The authors did not publish any volcano plots with this paper, however our in-group comparisons of the volcano plots showed good agreement. We identified the same down and upregulated genes using similar criteria for significance when comparing stem to enterocytes. Importantly, when comparing stem cells to enterocytes *Cd44* is the most upregulated gene. Given that the authors used FACS targeting *Cd44* to sort stem cells, we expect that this gene should be highly enriched. This is what our volcano plots show. We set a stringent threshold for our volcano plots ($\log_{2}\text{fold} > 2$ and $\text{padj} < 10^{-48}$ for Alex's plot and 10^{-45} for Bryan's plot). Genes that passed this threshold are biologically relevant for the comparison between stem cells and enterocytes. For example, *Mxd1* was one of the top genes we identified as

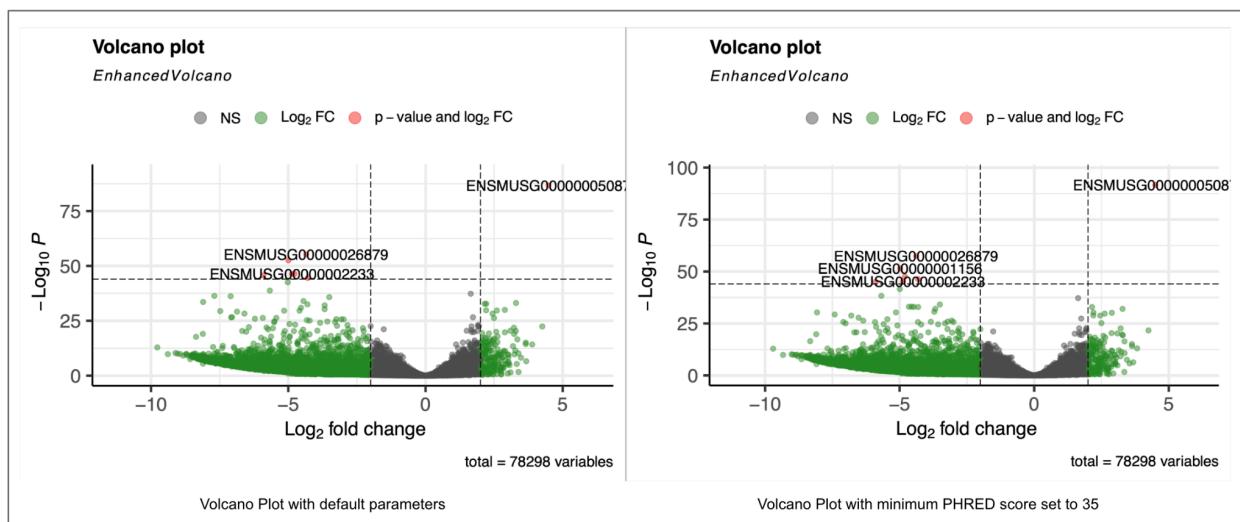
downregulated. *Mxd1* is a competitive antagonist for the transcription factor *Myc*.⁴ *Myc* is important for intestinal stem cells to maintain their fate.³ Since this gene is enriched in enterocytes it could potentially antagonize *Myc* to prevent enterocytes from transcribing pro-stem genes.



Sliding Window Parameter Change. Increasing the window size from 1 to 2 showed minimal effects on the volcano plots. There are slight differences in the overall shape of the data for genes deemed not significant, however genes identified as significant remained so with the parameter change. This is likely due to our stringent threshold for significance.



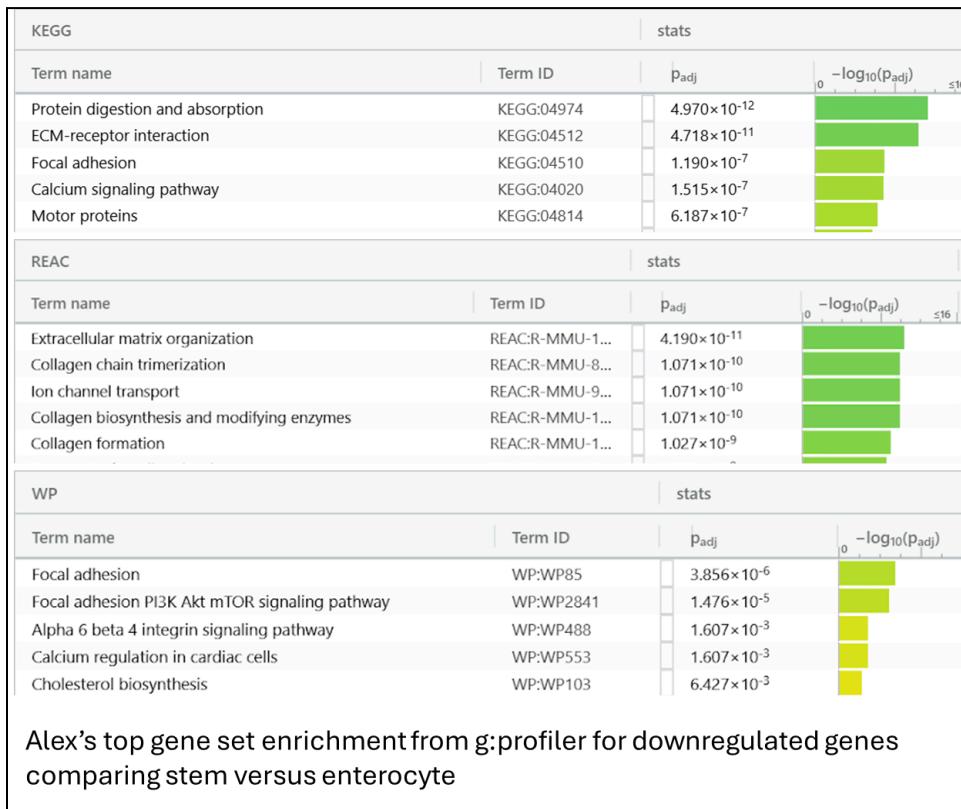
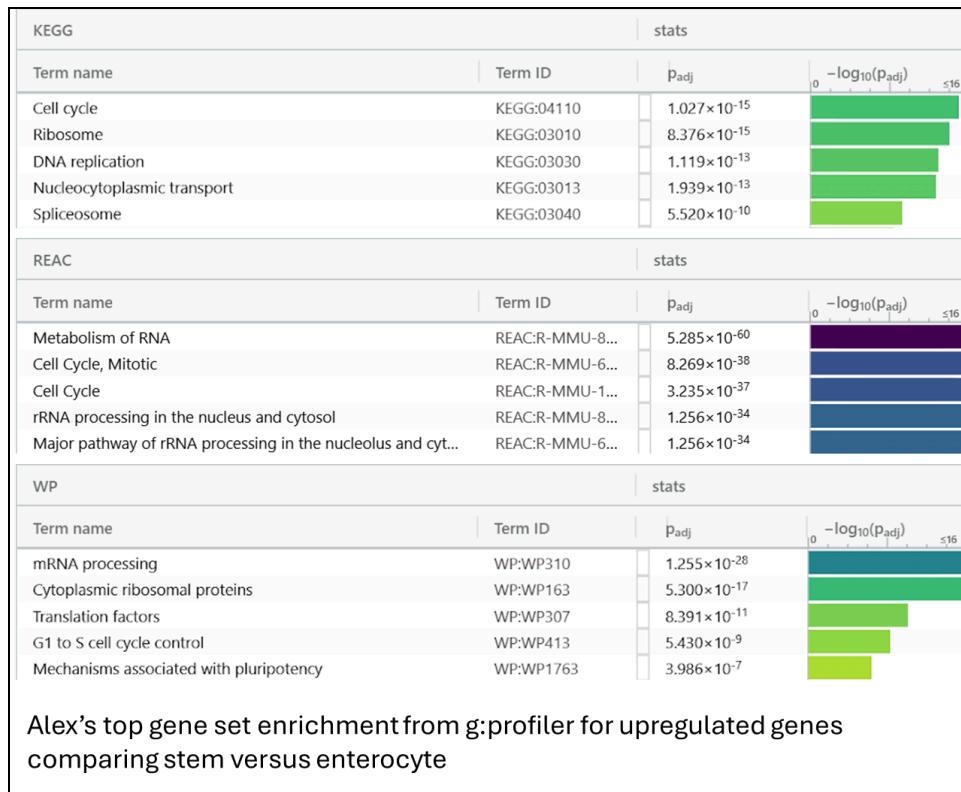
Min PHRED Score Parameter Change. Increasing the minimum PHRED score from 20 to 35 showed moderate effects on significant genes in the volcano plot. The most significant gene, *Cd44*, saw a significant decrease in its adjusted p-value. Log2FC values remained consistent for significant genes across both plots. The decrease in adjusted p-value can be attributed to the high quality filtering strengthening true biological signals by removing noise, while the consistent Log2FC is further proof that the parameter change did not alter any true signals. In short, the volcano plot with minimum PHRED score increased to 35 showed similar results with greater statistical confidence.

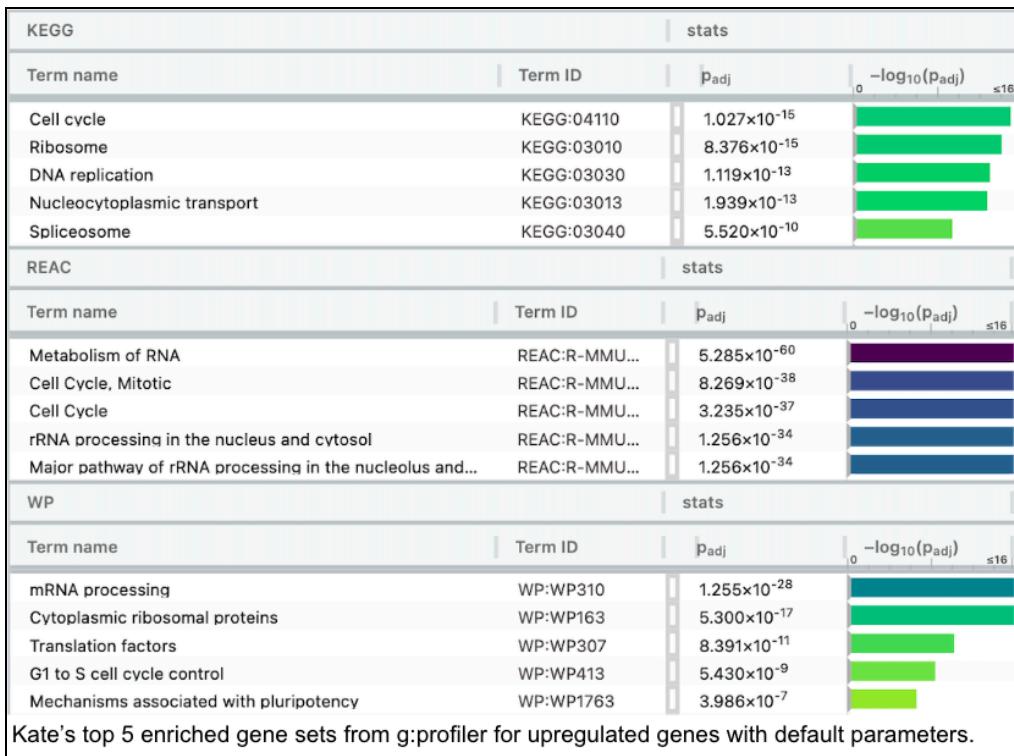


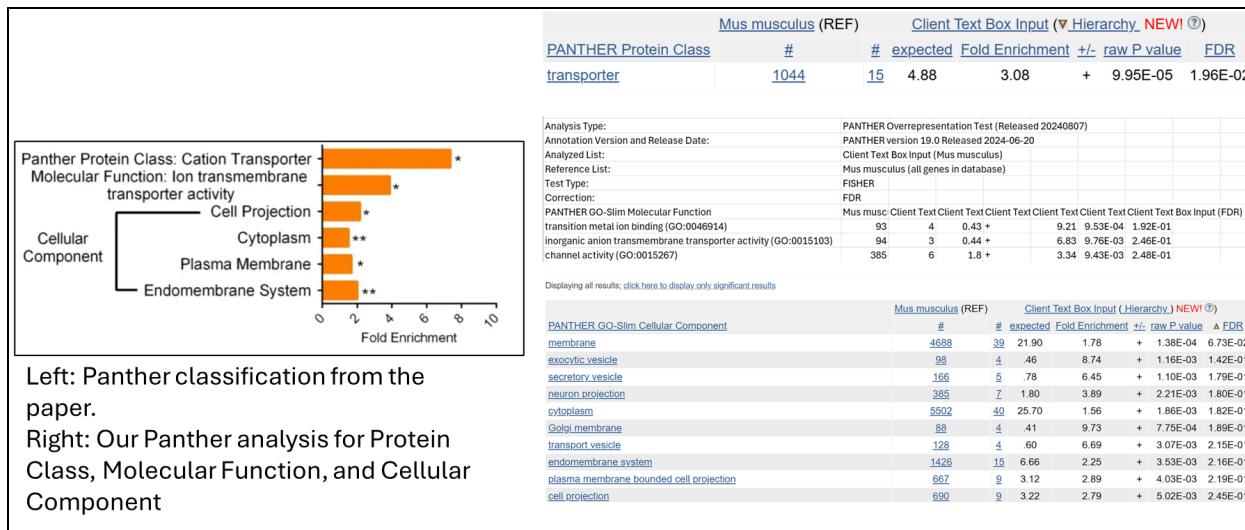
Functional Analysis:

G:profiler

Using default HPC parameters and the significance threshold outlined in the paper of $p_{adj} < 0.01$ and > 50 mean normalized counts, we ran g:profiler to look at KEEG, Reactome, and WikiPathways for both up- and downregulated genes identified comparing stem cells to enterocytes. Significance in g:profiler was set to benjamini-hochberg FDR <0.05 . Both group members had good agreement in g:profiler results. Importantly, the pathways downregulated in this comparison indicated that known pathways for enterocytes were enriched, such as focal adhesions.¹ In addition to our separate g:profiler analyses we also endeavored to recreate GO pathway analysis from our reference paper. In this analysis, the authors used Panther classification system for 107 most upregulated genes in all cells except stem cells. We were able to somewhat recapitulate the author's results using Panther, however the text lacked detail on specifically how they assessed what qualified as top genes in this data. Our results were similar for Panther Protein Class indicating that Transporter was the top term. For molecular function and cellular component we got similar results to the paper, such as Cytoplasm, however for these two categories our results failed the FDR <0.05 test. We attribute these differences to the lack of detail regarding how to determine the top genes to be used as well as how to appropriately use the Panther classification system.

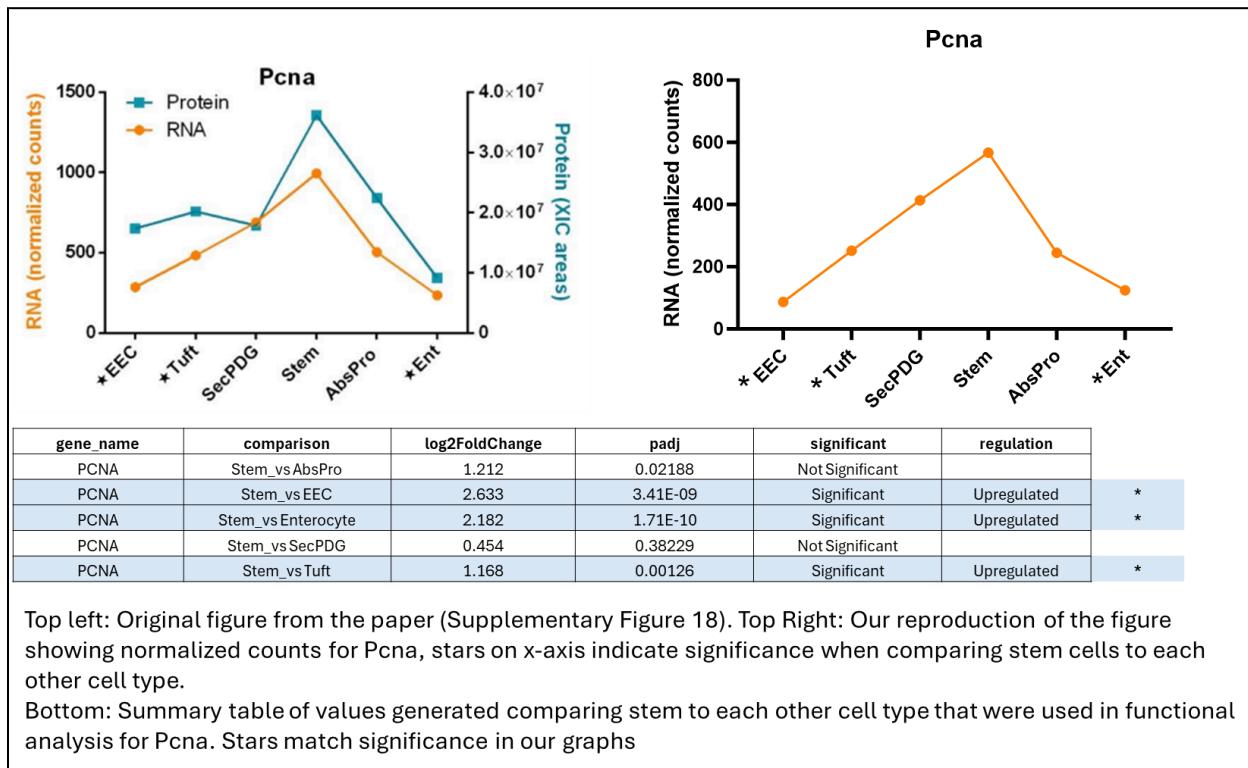




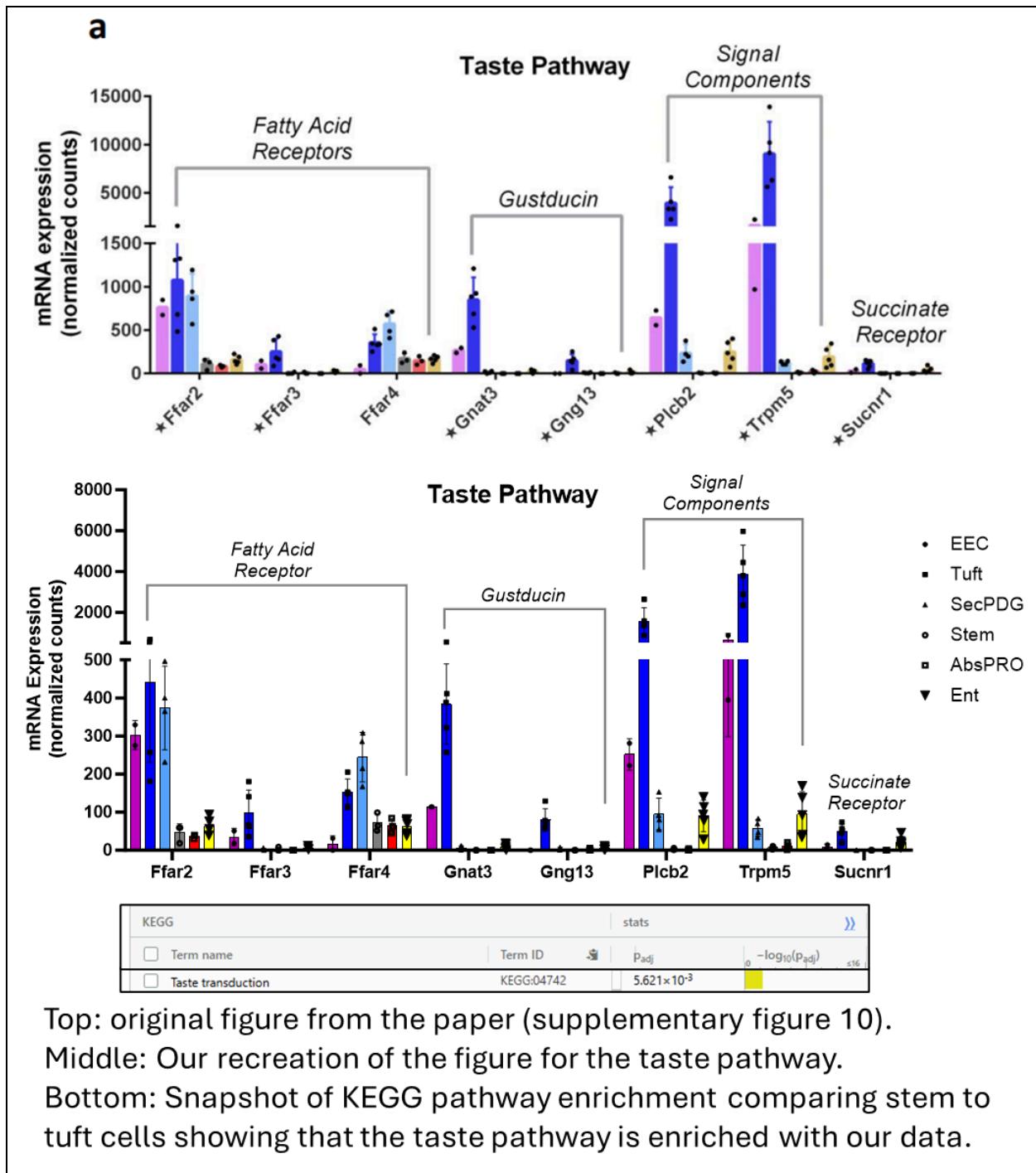


Validation of functional analysis: To further validate that our pathway enrichment results are consistent with what was published in the paper we reproduced two figures from the paper related to specific pathways that appeared in our g:profiler results. The first is based on the upregulation of cell cycle genes in stem cells when compared to enterocytes. The authors of the paper stated, “*RNA markers of proliferation (Mki67, Perna, and Mcms) are highest in stem cells, but interestingly, their protein products are readily detectable in differentiated cells, thus highlighting inconsistencies between mRNA and protein biomarkers of proliferation (Supplementary Fig. 18).*” Based on this text, we pulled expression data from DESeq2 for the gene *Pcna* as well as significance for *Pcna* comparing stem to each of the other cell types. The significance data for stem versus enterocyte is the same data that was used in our pathway enrichment analysis and the comparisons to other cell types was done in an identical manner to the enterocyte comparison. If our data is matching the paper we should be able to reproduce the figure as well as the significance. The figure below shows that we achieved both of these goals indicating that our data does match pathway analysis from the paper with significance defined as $padj < 0.01$ and > 50 mean normalized counts. It is important to note that the scaling on the y-axis for our version of the figure is different from the original. We suspect this is due to our handling of technical replicates in our analysis compared to the paper (first discussed in the PCA section of the report). In our analysis we only kept one technical replicate from the RNAseq runs because we were unsure as to how the authors handled these replicates - nowhere in their text do they indicate how technical replicates were addressed. However, if the authors merged these replicates together before downstream analysis this could lead to the overall larger

normalized counts as we see when comparing their figure to ours. Despite this our figure still recapitulates the original with good accuracy.

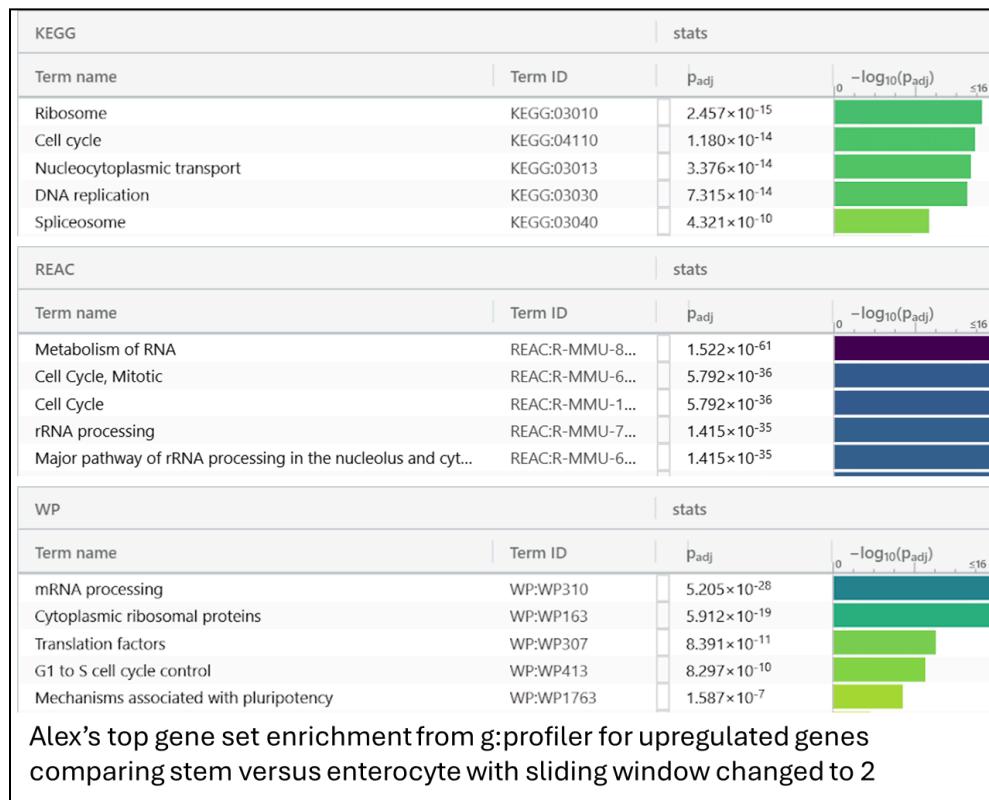


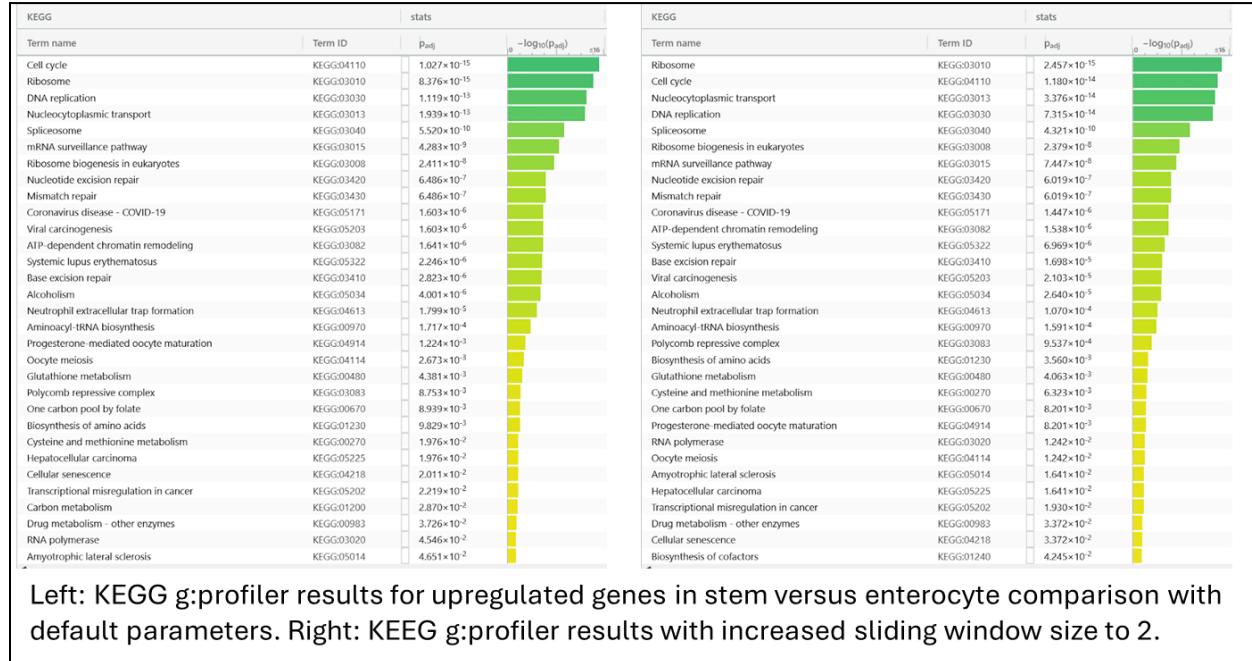
The second analysis we did to validate our pathway enrichment data was done comparing stem cells to tuft cells. This comparison is different from our main comparison of stem versus enterocyte, however we chose this because the authors indicated that the Taste pathway is enriched in tuft cells when compared to stem cells. We did this comparison of stem versus tuft in the same manner as our main comparison of stem versus enterocyte with significance for the genes included in the pathway analysis as $\text{padj} < 0.01$ and >50 mean normalized counts. We were able to reproduce the figure with good accuracy - the only exception is that the scaling on the y-axis is different as discussed above. In conjunction with reproducing the counts for the Taste pathway we also used g:profiler to show that our data comparing stem to tuft cells does show enrichment for the taste pathway for downregulated genes as expected. This comparison was done in an identical manner as to our main comparison of stem versus enterocyte. The results indicate that our pathway enrichment for the main comparison is sound.



Sliding Window Parameter Change. As with our other analyses in this report, increasing the sliding window size from 1 to 2 showed minimal effects on g:profiler results. When comparing stem versus enterocytes at default parameters, most enriched KEGG pathways that were upregulated had some relation to the cell cycle or protein synthesis such as DNA replication, cell cycle, and ribosome. The increased window slightly changed the order of the most enriched

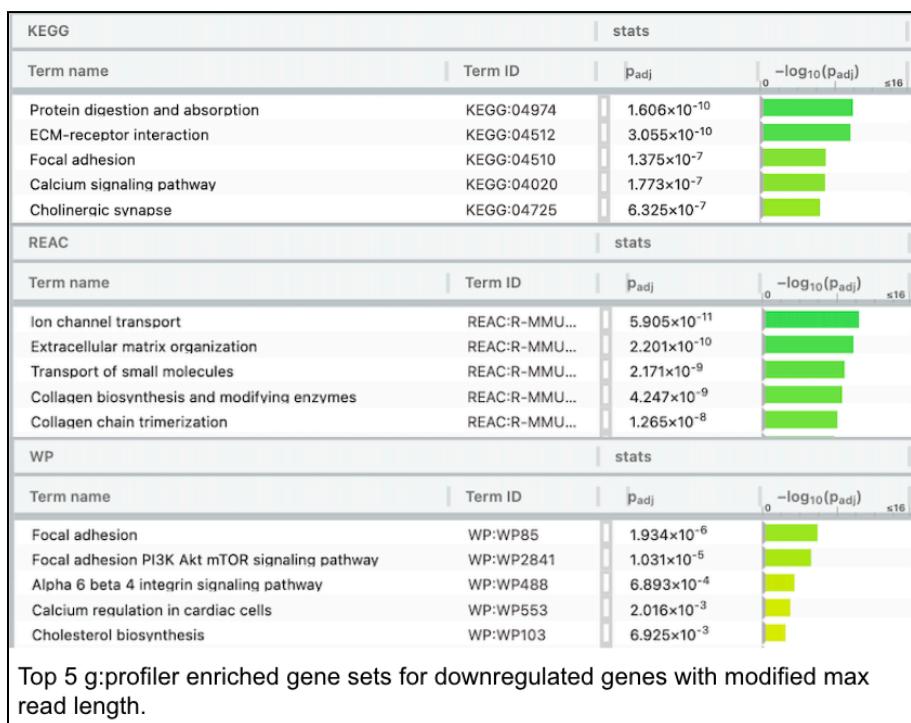
pathways. With default parameters, cell cycle is the top hit followed by ribosome for KEGG. However, with the increased window this is flipped and ribosome is the top hit followed by cell cycle. These are minor changes and the overall interpretation of the data is similar. This indicates that despite the change in trimming, top pathways are still being detected accurately.





Left: KEGG g:profiler results for upregulated genes in stem versus enterocyte comparison with default parameters. Right: KEGG g:profiler results with increased sliding window size to 2.

Max Read Length Parameter Change. Modifying the maximum read length by capping at 85 bp slightly changed our g:profiler data. Within each of the gene groupings (KEGG, Reactome, or WikiPathways), the main pathways did not change, but the Padj score tended to shift slightly when max read length was modified. However, the Padj tenth-power values were the same, so the overall significance was not largely impacted. Modifying max read length increased our confidence in our initial results because the main pathways were consistently up or downregulated.





GSEA

There was no direct comparison we could make to the reference paper for GSEA, however, we were able to separately generate two GSEA runs with default parameters that resemble each other as much as possible. Both runs showed 807 enriched gene sets: 291 upregulated in stem cells, and 516 upregulated in enterocytes. The ranked lists of enriched gene sets in enterocytes were very similar, with the top ranked gene set in both runs being CLASS I PEROXISOMAL MEMBRANE PROTEIN IMPORT. The ranked lists of enriched gene sets in stem cells showed far less consistency between runs. It is important to note that GSEA uses permutation testing to generate p-values and normalized enrichment scores (NES). If a seed is not set for permutation testing, each run of GSEA will likely run different permutations, slightly altering the final results. This fact explains the minor differences between runs. The greater discrepancy between the ranked lists of enriched gene sets in stem cells can be attributed to stem cells' relatively low RNA count. Lower sample size, where sample size in this case is RNA count, causes weaker statistical power and therefore greater variance between different runs of permutation testing within GSEA.

Enrichment in phenotype: 0 (3 samples)	Enrichment in phenotype: 0 (3 samples)
<ul style="list-style-type: none"> • 291 / 807 gene sets are upregulated in phenotype 0 • 280 gene sets are significant at FDR < 25% • 45 gene sets are significantly enriched at nominal pvalue < 1% • 45 gene sets are significantly enriched at nominal pvalue < 5% • Snapshot of enrichment results • Detailed enrichment results in html format • Detailed enrichment results in TSV format (tab delimited text) • Guide to interpret results 	<ul style="list-style-type: none"> • 291 / 807 gene sets are upregulated in phenotype 0 • 279 gene sets are significant at FDR < 25% • 51 gene sets are significantly enriched at nominal pvalue < 1% • 51 gene sets are significantly enriched at nominal pvalue < 5% • Snapshot of enrichment results • Detailed enrichment results in html format • Detailed enrichment results in TSV format (tab delimited text) • Guide to interpret results
Enrichment in phenotype: 4 (5 samples)	Enrichment in phenotype: 4 (5 samples)
<ul style="list-style-type: none"> • 516 / 807 gene sets are upregulated in phenotype 4 • 105 gene sets are significantly enriched at FDR < 25% • 21 gene sets are significantly enriched at nominal pvalue < 1% • 52 gene sets are significantly enriched at nominal pvalue < 5% • Snapshot of enrichment results • Detailed enrichment results in html format • Detailed enrichment results in TSV format (tab delimited text) • Guide to interpret results 	<ul style="list-style-type: none"> • 516 / 807 gene sets are upregulated in phenotype 4 • 100 gene sets are significantly enriched at FDR < 25% • 16 gene sets are significantly enriched at nominal pvalue < 1% • 49 gene sets are significantly enriched at nominal pvalue < 5% • Snapshot of enrichment results • Detailed enrichment results in html format • Detailed enrichment results in TSV format (tab delimited text) • Guide to interpret results

GSEA Index with default parameters generated by
Bryan

GSEA Index with default parameters generated by
Kate

GSE										GSE										
Index Files					Index Files					Index Files					Index Files					
	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size
1	REACTOME:REGULATION_OF_ARNA	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:CLASS_I_PEROXISOMAL_MEMBRANE_PROTEIN_IMPORT	Enterocyte	18	4.81	4.81	0.000	0.000	0.010	1002	High-TPL, HighP
2	REACTOME:PROCESSING_OF_ATPIFENONE_PHE_PHRN	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:SYNTHESIS_OF_BILE_ACIDS_AND_BILE_SALTS_VIA_7ALPHA_HYDROXYCHOLESTEROL	Enterocyte	22	0.70	0.40	0.000	0.000	0.419	6416	High-TPL, HighP
3	REACTOME:ARAC_ASYNTHESIS	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:Negative_Regulation_of_ALD	Enterocyte	11	0.73	0.44	0.004	0.000	0.700	5403	High-TPL, HighP
4	REACTOME:POX_DEPENDENT_LAND_PSIK_ARNA_EXCISION_PATH	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:SYNTHESIS_OF_BILE_ACIDS_AND_BILE_SALTS	Enterocyte	32	0.89	0.41	0.003	0.000	0.895	6416	High-TPL, HighP
5	REACTOME:ARAC_CARNITINE	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:FATTY_ACIDS	Enterocyte	23	0.68	0.41	0.008	0.000	0.895	3403	High-TPL, HighP
6	REACTOME:PROLIFERATION_OF_THE_WCLC	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:Olfactory_SIGNALING_PATHWAY	Enterocyte	72	0.82	0.30	0.000	0.000	0.941	3376	High-TPL, HighP
7	REACTOME:ARAC_DNA_SENSE_GLOMULOSIS	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:WPSI_LIGAND_BINDING_AND_TRANSPORTING	Enterocyte	23	0.84	0.30	0.009	0.000	0.998	14827	High-TPL, HighP
8	REACTOME:ACTIVE_RESPONSE_PICOR_N1_3HPI_ARAC_ASYNTHESIS	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:ARAC_LIGAND_BINDING_RECEP	Enterocyte	43	0.81	0.20	0.000	0.021	1.000	94007	High-TPL, HighP
9	REACTOME:PROCESSIVE_SYNTHESIS_ON_THE_C_STAND_OF_THE_TELOMERE	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_P1_3K_CASCADE_PATH	Enterocyte	21	0.83	0.20	0.001	0.000	0.994	10024	High-TPL, HighP
10	REACTOME:PROCESSIVE_SYNTHESIS_ON_THE_C_STAND_OF_THE_TELOMERE	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_FORMATION_OF_THE_COIFFED_ENVELOPE	Enterocyte	20	0.83	0.20	0.001	0.000	1.000	13044	High-TPL, HighP
11	REACTOME:ARMED_ATTACK	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:GLUTAMATE_METABOLISM_AND_GLYCINE_DECARBOXYLATION	Enterocyte	18	0.84	0.20	0.009	0.000	1.000	12000	High-TPL, HighP
12	REACTOME:FORMATION_OF_ANNEAUATED_HETEROCYCLES_FOLIC_ACID	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:PHOSPHOGLYCIC_C_MEDIATED CASCADE_PATH	Enterocyte	16	0.86	0.20	0.022	0.000	1.000	10252	High-TPL, HighP
13	REACTOME:FORMATION_OF_AN_ARAC_HOMOMERIC_JUNCTION	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_FOLP_LIGAND_BINDING_AND_ACTIVATION	Enterocyte	17	0.88	0.21	0.027	0.000	1.000	10324	High-TPL, HighP
14	REACTOME:ARAC_MERAKS_R3_32_3_ALPHOMERICAS	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_NEURONS_AND_NEUROGEN	Enterocyte	21	0.81	0.21	0.015	0.000	1.000	9596	High-TPL, HighP
15	REACTOME:INITION_OF_ARAC_EQUIVARIANT_EQ	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	17	0.85	0.20	0.008	0.000	1.000	13400	High-TPL, HighP

GSE										GSE										
Index Files					Index Files					Index Files					Index Files					
	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size	File Type	File ID	File Name	File Size
1	REACTOME:HYDROXYLATION_OF_ALCOHOLS_BY_ALCOHOL_OXIDASE	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:CLAS_I_PEROXISOMAL_MEMBRANE_PROTEIN_IMPORT	Enterocyte	19	0.47	0.20	0.000	0.000	0.239	1029	High-TPL, HighP
2	REACTOME:HYDROXYLATION_OF_ALCOHOLS_BY_ALCOHOL_OXIDASE	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:HYDROXYLATION_OF_ALCOHOLS_BY_ALCOHOL_OXIDASE	Enterocyte	20	0.70	0.20	0.021	0.000	0.240	10299	High-TPL, HighP
3	REACTOME:LARSEN_STEDMAN	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:Negative_Regulation_of_ALD	Enterocyte	11	0.73	0.20	0.029	0.000	0.240	10297	High-TPL, HighP
4	REACTOME:INVOLVEMENT_OF_NUCLEAR_JUNQOLYLIC_ACID_REG	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:FATTY_ACIDS	Enterocyte	22	0.89	0.20	0.029	0.000	0.240	12495	High-TPL, HighP
5	REACTOME:TRANSPORTER_REGULATION_BY_EQT	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME:SYNTHESIS_OF_ALC_ACIDS_AND_BILE_SALTS	Enterocyte	20	0.86	0.20	0.029	0.000	0.240	10244	High-TPL, HighP
6	REACTOME:ARAC_EQC	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	20	0.81	0.20	0.027	0.000	0.240	12704	High-TPL, HighP
7	REACTOME:POX_DEPENDENT_LONG_FAITH_ARNA_EXCISION_PATH	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_WPSI_LIGAND_BINDING_AND_TRANSPORTING	Enterocyte	23	0.84	0.20	0.027	0.000	0.240	14827	High-TPL, HighP
8	REACTOME:HIGH_EQC_ARAC_ASYNTHESIS	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	20	0.85	0.20	0.028	0.000	0.240	12495	High-TPL, HighP
9	REACTOME:THE_POSTTRANSLATIONAL_ACTIVITY_OF_ALC_IS_REQUIRED_FOR_THE_ONSET_OF_ANAPHASE_BY_MITOCHONDRIONAL_CHECKPOINT_COMPONENTS	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	19	0.81	0.20	0.028	0.000	0.240	10244	High-TPL, HighP
10	REACTOME:PROLIFERATION_OF_THE_ARC_E	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	20	0.76	0.20	0.028	0.000	0.240	10244	High-TPL, HighP
11	REACTOME:PROCESSING_OF_ATPIFENONE_PHE_PHRN	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	20	0.80	0.20	0.028	0.000	0.240	10244	High-TPL, HighP
12	REACTOME:ATP_ALTERNATIVE_UPLOA	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	17	0.76	0.20	0.028	0.000	0.240	12701	High-TPL, HighP
13	REACTOME:FORMATION_OF_BENDENEASSOCIATED_HETEROBROMOM_FOLIC_ACID	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	20	0.84	0.20	0.028	0.000	0.240	12495	High-TPL, HighP
14	REACTOME:PROCESSIVE_EQC	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	14	0.81	0.20	0.028	0.000	0.240	10244	High-TPL, HighP
15	REACTOME:PROCESSIVE_EQC_ON_THE_C_STAND_OF_THE_TELOMERE	Stem	20	20M	20	20M	20M	20M	20	20	REACTOME_EQC	Enterocyte	20	0.84	0.20	0.028	0.000	0.240	12495	High-TPL, HighP

Ranked enriched gene set lists generated in GSEA with default parameters. Top left: Bryan, stem. Top right: Bryan, enterocyte. Bottom left: Kate, stem. Bottom right: Kate, enterocyte.

Min PHRED Score Parameter Change. Increasing the minimum PHRED score from 20 to 35 exhibited moderate impacts on both GSEA generated index files and ranked enriched gene set lists. The number of gene sets enriched in stem cells decreased from 291 to 287, while the number of significantly enriched gene sets in stem cells increased from 45 to 53. The number of gene sets enriched in enterocytes increased from 516 to 520, while the number of significantly enriched gene sets in enterocytes remained about the same. The changes across runs were relatively small due to the high quality of the data, and only a small fraction of reads being filtered out by the increased PHRED threshold. The decrease in enriched gene sets in stem cells and increase in enterocytes can be attributed to each cell's RNA content. The low RNA content of stem cells leads to more low quality associated reads that get filtered out by the stricter PHRED threshold, while the high RNA content of enterocytes benefits from the reduced noise after stricter filtering. This filtering of lower quality reads associated with stem cells explains the increase in significantly enriched gene sets in stem cells. Further, results showed consistency between ranked enriched gene set lists for enterocytes, and relatively high discrepancy between ranked enriched gene set lists for stem cells. The top two enriched gene sets in enterocytes remained the same. These were: CLASS I PEROXISOMAL MEMBRANE PROTEIN IMPORT and SYNTHESIS OF BILE ACIDS AND BILE SALTS VIA 7ALPHA HYDROXYCHOLESTEROL. These results can also be attributed to the RNA content of stem cells and enterocytes, and the discarding of more stem cell associated reads because of their

relatively low RNA count compared to enterocytes.

Enrichment in phenotype: 0 (3 samples)

- 291 / 807 gene sets are upregulated in phenotype 0
 - 280 gene sets are significant at FDR < 25%
 - 45 gene sets are significantly enriched at nominal pvalue < 1%
 - 45 gene sets are significantly enriched at nominal pvalue < 5%
 - **Snapshot** of enrichment results
 - Detailed [enrichment results in html](#) format
 - Detailed [enrichment results in TSV](#) format (tab delimited text)
 - [Guide](#) to interpret results

Enrichment in phenotype: 4 (5 samples)

- 516 / 807 gene sets are upregulated in phenotype 4
 - 105 gene sets are significantly enriched at FDR < 25%
 - 21 gene sets are significantly enriched at nominal pvalue < 1%
 - 52 gene sets are significantly enriched at nominal pvalue < 5%
 - **Snapshot** of enrichment results
 - Detailed [enrichment results in html](#) format
 - Detailed [enrichment results in TSV](#) format (tab delimited text)
 - **Guide to interpret results**

GSEA Index with default parameters

Enrichment in phenotype: 0 (3 samples)

- 287 / 807 gene sets are upregulated in phenotype **0**
 - 277 gene sets are significant at FDR < 25%
 - 53 gene sets are significantly enriched at nominal pvalue < 1%
 - 53 gene sets are significantly enriched at nominal pvalue < 5%
 - **Snapshot** of enrichment results
 - Detailed [enrichment results in html](#) format
 - Detailed [enrichment results in TSV](#) format (tab delimited text)
 - [Guide to interpret results](#)

Enrichment in phenotype: 4 (5 samples)

- 520 / 807 gene sets are upregulated in phenotype 4
 - 118 gene sets are significantly enriched at FDR < 25%
 - 19 gene sets are significantly enriched at nominal pvalue < 1%
 - 52 gene sets are significantly enriched at nominal pvalue < 5%
 - **Snapshot** of enrichment results
 - Detailed [enrichment results in html](#) format
 - Detailed [enrichment results in TSV](#) format (tab delimited text)
 - [Guide to interpret results](#)

GSEA Index with minimum PHRED score set to 35

	NAME	KEGG ID	EC NUMBER	PUBMED REFERENCES	INTERACTOME REFERENCES	GO REFERENCES	KEGG REFERENCES	LEADERBOARD REFERENCES
Pathway: Bile Acid Synthesis								
1	REACTOME: CLASS I PEROXISOMAL MEMBRANE PROTEIN IMPORT	Details...	19 -0.91 -1.40	0.000	0.000	0.010	14932	lsgp111%;nsgp10%;
2	REACTOME: SYNTHESIS OF BILE ACIDS AND BILE SALTS VIA TAPAHA HYDROXYCHOLESTEROL	Details...	22 -0.70 -1.40	0.000	0.019	0.019	6416	lsgp111%;nsgp11%;
3	REACTOME: NEGATIVE REGULATION OF F1AT	Details...	15 -0.73 -1.40	0.004	0.009	0.709	5473	lsgp111%;nsgp11%;
4	REACTOME: SYNTHESIS OF BILE ACIDS AND BILE SALTS	Details...	32 -0.81 -1.40	0.003	0.085	0.616	5416	lsgp111%;nsgp11%;
5	REACTOME: FATTY ACIDS	Details...	22 -0.81 -1.40	0.000	0.005	0.885	13495	lsgp111%;nsgp11%;
6	REACTOME: OLFACTORY SIGNALING PATHWAY	Details...	72 -0.82 -1.30	0.000	0.000	0.941	13764	lsgp111%;nsgp11%;
7	REACTOME: WNT LEADERSHIP, INGURNESS, AND TRAFICKING	Details...	26 -0.84 -1.30	0.000	0.016	0.998	14937	lsgp111%;nsgp11%;
8	REACTOME: AMINE-LIGAND BINDING RECEPTORS	Details...	43 -0.81 -1.30	0.000	0.022	1.000	18467	lsgp111%;nsgp11%;
9	REACTOME: PI_3K CASCADE FORM	Details...	21 -0.83 -1.30	0.041	0.000	1.000	16234	lsgp111%;nsgp11%;
10	REACTOME: FORMATION OF THE CONFINED ENVELOPE	Details...	50 -0.86 -1.30	0.000	0.000	1.000	17021	lsgp111%;nsgp11%;
11	REACTOME: GLYCOLYSE, METABOLISM AND GLYCINE DEGRADATION	Details...	18 -0.84 -1.30	0.009	0.000	1.000	12556	lsgp111%;nsgp11%;
12	REACTOME: PHOSPHATE-REGULATED LIPIDACE FORM	Details...	14 -0.82 -1.20	0.000	0.000	1.000	16234	lsgp111%;nsgp11%;
13	REACTOME: FGFR1 LIGAND BINDING AND ACTIVATION	Details...	17 -0.85 -1.31	0.007	0.000	1.000	16234	lsgp111%;nsgp11%;
14	REACTOME: NEUROGENES AND NEUROGENES	Details...	81 -0.81 -1.31	0.015	0.009	0.998	6560	lsgp111%;nsgp11%;
15	REACTOME: EXOCYTOSIS	Details...	17 -0.85 -1.30	0.000	0.000	1.000	13495	lsgp111%;nsgp11%;

Top: Ranked enriched gene set lists generated in GSEA with default parameters. Bottom: Ranked enriched gene set lists generated in GSEA with minimum PHRED score increased to 35. Left: Stem cells. Right: Enterocytes.

Max Read Length Parameter Change. Modifying max read length by capping at 85 bp slightly lowered the number of gene sets enriched in GSEA for stem cells (291 in default to 288), but increased in enterocytes (516 in default to 519). This is possibly because there is more mRNA in enterocytes so the natural robustness is higher than in stem cells, which have less mRNA.

Thus, increasing stringency can reduce the number of enriched gene sets in stem cells because the stem cells are more susceptible to changes in parameters. We would also expect some changes in GSEA after capping max read length because trimming some reads impacted alignment and thus expression levels, and these subtle changes can shift genes around on the rank list. GSEA is based on this rank list so even small changes may show up in the gene sets that were enriched. For example, in the stem cell GSEA data below, “postmitotic nuclear pore complex” is most enriched before and after modifying maximum read length, but the second most enriched gene set is different (“butyrate response factor” vs “PCNA dependent base excision repair”). Although the order changes, other enriched gene sets like “lagging strand synthesis” and “formation of senescence-associated heterochromatin foci” are in the top 15 for both GSEA results, which further confirms the raw data had high integrity.

GSEA results summary with default parameters (stem cell on top, enterocyte below)	GSEA results summary with modified max read length (stem cell on top, enterocyte below)
<p>Enrichment in phenotype: 0 (3 samples)</p> <ul style="list-style-type: none"> 291 / 807 gene sets are upregulated in phenotype 0 279 gene sets are significant at FDR < 25% 51 gene sets are significantly enriched at nominal pvalue < 1% 51 gene sets are significantly enriched at nominal pvalue < 5% Snapshot of enrichment results Detailed enrichment results in html format Detailed enrichment results in TSV format (tab delimited text) Guide to interpret results <p>Enrichment in phenotype: 4 (5 samples)</p> <ul style="list-style-type: none"> 516 / 807 gene sets are upregulated in phenotype 4 100 gene sets are significantly enriched at FDR < 25% 16 gene sets are significantly enriched at nominal pvalue < 1% 49 gene sets are significantly enriched at nominal pvalue < 5% Snapshot of enrichment results Detailed enrichment results in html format Detailed enrichment results in TSV format (tab delimited text) Guide to interpret results 	<p>Enrichment in phenotype: 0 (3 samples)</p> <ul style="list-style-type: none"> 288 / 807 gene sets are upregulated in phenotype 0 267 gene sets are significant at FDR < 25% 46 gene sets are significantly enriched at nominal pvalue < 1% 46 gene sets are significantly enriched at nominal pvalue < 5% Snapshot of enrichment results Detailed enrichment results in html format Detailed enrichment results in TSV format (tab delimited text) Guide to interpret results <p>Enrichment in phenotype: 4 (5 samples)</p> <ul style="list-style-type: none"> 519 / 807 gene sets are upregulated in phenotype 4 103 gene sets are significantly enriched at FDR < 25% 20 gene sets are significantly enriched at nominal pvalue < 1% 52 gene sets are significantly enriched at nominal pvalue < 5% Snapshot of enrichment results Detailed enrichment results in html format Detailed enrichment results in TSV format (tab delimited text) Guide to interpret results

GSEA results for Stem Cell

	GS follow link to MSigDB	GS DETAILS	SIZE	ES	NES	NOM p-val	FDR q-val	PWER p-val	RANK AT MAX	LEADING EDGE
1	REACTOME_POSTMITOTIC_NUCLEAR_PORE_COMPLEX_NPC_REFORMATION	Details...	22	0.90	3.37	0.000	0.000	3662		tags=>8%, list=4%, signal=>7%
2	REACTOME_BUTYRATE_RESPONSE_FACTOR_1_BRF1_BINDS_AND_DESTABILIZES_MRNA	Details...	17	0.85	3.32	0.000	0.000	4800		tags=>7%, list=5%, signal=>6%
3	REACTOME_LAGGING_STRAND_SYNTHESIS	Details...	20	0.92	3.28	0.000	0.000	2729		tags=>9%, list=5%, signal=>8%
4	REACTOME_INITIATION_OF_NUCLEAR_ENVELOPE_NE_REFORMATION	Details...	17	0.79	3.22	0.000	0.000	3237		tags=>8%, list=5%, signal=>6%
5	REACTOME_TRANSITIONAL_REGULATION_BY_EER	Details...	23	0.78	3.14	0.000	0.000	4852		tags=>7%, list=5%, signal=>8%
6	REACTOME_MRNA_DECAY_BY_3_TO_5_EXORIBONUCLEASE	Details...	16	0.90	3.09	0.000	0.000	3009		tags=>8%, list=5%, signal=>7%
7	REACTOME_PCPA_DEPENDENT_LONG_PATCH_BASE_EXCISION_REPAIR	Details...	21	0.87	3.05	0.000	0.000	2564		tags=>8%, list=5%, signal=>9%
8	REACTOME_KSRP_KHSPB BINDS AND DESTABILIZES MRNA	Details...	15	0.91	3.01	0.000	0.000	4800		tags=>7%, list=5%, signal=>9%
9	REACTOME_INHIBITION_OF_THE_PROTEOLYTIC_ACTIVITY_OF_APOLY_C_REQUIRED_FOR_THE_ONSET_OF_ANAPHASE_BY_MITOTIC_SPINDLE_CHECKPOINT_COMPONENTS	Details...	21	0.80	2.94	0.000	0.000	4420		tags=>7%, list=5%, signal=>8%
10	REACTOME_PHOSPHORYLATION_OF_THE_APOLY_C	Details...	20	0.78	2.92	0.000	0.000	5023		tags=>8%, list=5%, signal=>8%
11	REACTOME_PROCESSING_OF_INTRONLESS_PRE_MRNAS	Details...	20	0.90	2.90	0.000	0.000	5668		tags=>7%, list=5%, signal=>10%
12	REACTOME_FGFR_ALTERNATIVE_SPLICING	Details...	17	0.78	2.77	0.000	0.000	4217		tags=>7%, list=5%, signal=>7%
13	REACTOME_FORMATION_OF_SENESCENCE_ASSOCIATED_HETEROCHROMATIN_FOCI_SAHF	Details...	16	0.86	2.62	0.000	0.000	3319		tags=>7%, list=5%, signal=>6%
14	REACTOME_PROCESSIVE_SYNTHESIS_ON_THE_LAGGING_STRAND	Details...	15	0.91	2.57	0.000	0.000	2729		tags=>7%, list=5%, signal=>6%
15	REACTOME_PROCESSIVE_SYNTHESIS_ON_THE_C_STRAND_OF_THE_TELOMERE	Details...	18	0.83	2.57	0.000	0.000	2471		tags=>7%, list=4%, signal=>7%

GSEA results for Stem Cell with modified max read length

	GS follow link to MSigDB	GS DETAILS	SIZE	ES	NES	NOM p-val	FDR q-val	PWER p-val	RANK AT MAX	LEADING EDGE
1	REACTOME_POSTMITOTIC_NUCLEAR_PORE_COMPLEX_NPC_REFORMATION	Details...	22	0.89	3.19	0.000	0.000	3777		tags=>82%, list=7%, signal=>8%
2	REACTOME_PCPA_DEPENDENT_LONG_PATCH_BASE_EXCISION_REPAIR	Details...	21	0.87	3.12	0.000	0.000	2589		tags=>76%, list=5%, signal=>8%
3	REACTOME_RHOBTB2_GTPASE_CYCLE	Details...	22	0.85	3.25	0.000	0.000	5305		tags=>82%, list=6%, signal=>9%
4	REACTOME_INHIBITION_OF_THE_PROTEOLYTIC_ACTIVITY_OF_APOLY_C_REQUIRED_FOR_THE_ONSET_OF_ANAPHASE_BY_MITOTIC_SPINDLE_CHECKPOINT_COMPONENTS	Details...	21	0.80	3.09	0.000	0.000	5006		tags=>81%, list=5%, signal=>9%
5	REACTOME_PROCESSIVE_SYNTHESIS_ON_THE_LAGGING_STRAND	Details...	15	0.91	3.09	0.000	0.000	2767		tags=>7%, list=5%, signal=>9%
6	REACTOME_LAGGING_STRAND_SYNTHESIS	Details...	20	0.92	3.05	0.000	0.000	2767		tags=>90%, list=5%, signal=>9%
7	REACTOME_MISMATCH_REPAIR	Details...	15	0.87	2.91	0.000	0.000	2333		tags=>7%, list=4%, signal=>7%
8	REACTOME_FORMATION_OF_SENESCENCE_ASSOCIATED_HETEROCHROMATIN_FOCI_SAHF	Details...	16	0.88	2.89	0.000	0.000	3374		tags=>91%, list=4%, signal=>10%
9	REACTOME_KSRP_KHSPB BINDS AND DESTABILIZES MRNA	Details...	15	0.91	2.83	0.000	0.000	4852		tags=>7%, list=4%, signal=>9%
10	REACTOME_MRNA_DECAY_BY_3_TO_5_EXORIBONUCLEASE	Details...	16	0.90	2.83	0.000	0.000	5319		tags=>88%, list=5%, signal=>7%
11	REACTOME_PROTEIN_HYDROXYLATION	Details...	18	0.74	2.72	0.000	0.000	4749		tags=>67%, list=4%, signal=>7%
12	REACTOME_TRISTETRAPROLIN_TTP_ZFP36_BINDS_AND_DESTABILIZES_MRNA	Details...	17	0.71	2.68	0.000	0.000	4852		tags=>71%, list=5%, signal=>7%
13	REACTOME_INITIATION_OF_NUCLEAR_ENVELOPE_NE_REFORMATION	Details...	17	0.79	2.66	0.000	0.000	3243		tags=>65%, list=4%, signal=>9%
14	REACTOME_CONVERSION_FROM_APOLY_C_CD20_TO_APOLY_C_CDH1_IN_LATE_ANAPHASE	Details...	20	0.73	2.54	0.000	0.000	5006		tags=>75%, list=4%, signal=>8%
15	REACTOME_TELOMERE_EXTENSION_BY_TELOMerase	Details...	20	0.71	2.52	0.000	0.000	5626		tags=>65%, list=4%, signal=>7%

GSEA results for Enterocyte (default parameters)

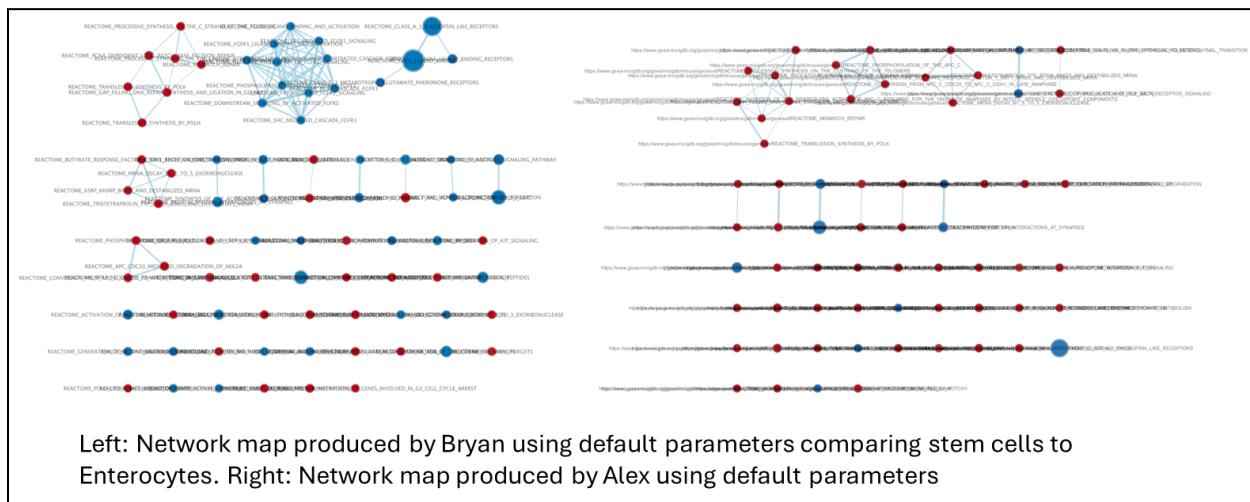
	GS follow link to MSigDB	GS DETAILS	SIZE	ES	NES	NOM p-val	FDR q-val	PWER p-val	RANK AT MAX	LEADING EDGE
1	REACTOME_CLASS_I_PEROXISOMAL_MEMBRANE_PROTEIN_IMPORT	Details...	19	-0.81	-1.66	0.000	0.000	0.009	1693	tags=<1%, list=<1%, signal=<1%
2	REACTOME_SYNTHESIS_OF_BILE_ACIDS_AND_BILE_SALTS_VIA_TALPHA_HYDROXYCHOLESTEROL	Details...	22	-0.70	-1.45	0.000	0.002	0.580	6415	tags=>36%, list=11%, signal=>4%
3	REACTOME_NEGATIVE_REGULATION_OF_FLTs	Details...	15	-0.72	-1.43	0.012	0.004	0.779	5493	tags=>15%, list=10%, signal=>10%
4	REACTOME_FATTY_ACIDS	Details...	22	-0.68	-1.42	0.005	0.004	0.798	1349	tags=>75%, list=2%, signal=>7%
5	REACTOME_SYNTHESIS_OF_BILE_ACIDS_AND_BILE_SALTS	Details...	32	-0.66	-1.41	0.003	0.005	0.854	6416	tags=>36%, list=11%, signal=>10%
6	REACTOME_Olfactory_SIGNALING_PATHWAY	Details...	72	-0.62	-1.40	0.000	0.007	0.917	13764	tags=>40%, list=24%, signal=>2%
7	REACTOME_WNT_LIGAND_BIOGENESIS_AND_TRAFFICKING	Details...	25	-0.64	-1.34	0.017	0.019	1.000	14637	tags=>95%, list=26%, signal=>7%
8	REACTOME_EICOSANOIDS	Details...	17	-0.65	-1.32	0.033	0.024	1.000	13495	tags=>71%, list=24%, signal=>10%
9	REACTOME_AMINE_LIGAND_BINDING_RECEPORS	Details...	43	-0.61	-1.32	0.004	0.028	1.000	18488	tags=>74%, list=39%, signal=>10%
10	REACTOME_PHOSPHOLIPASE_C_MEDIANDED CASCADE_FGFR1	Details...	16	-0.65	-1.32	0.038	0.028	1.000	16334	tags=>75%, list=29%, signal=>10%
11	REACTOME_FORMATION_OF_THE_CORNIFIED_ENVELOPE	Details...	90	-0.58	-1.32	0.000	0.026	1.000	17021	tags=>65%, list=20%, signal=>6%
12	REACTOME_SHC_MEDIANDED CASCADE_FGFR1	Details...	39	-0.64	-1.32	0.035	0.027	1.000	16333	tags=>65%, list=29%, signal=>10%
13	REACTOME_NEUROKININS_AND_NEUROGLIGINS	Details...	31	-0.61	-1.31	0.017	0.029	1.000	9598	tags=>39%, list=17%, signal=>4%
14	REACTOME_FGFR1_LIGAND_BINDING_AND_ACTIVATION	Details...	17	-0.65	-1.31	0.003	0.028	1.000	16334	tags=>76%, list=29%, signal=>10%
15	REACTOME_PI_3K CASCADE_FGFR1	Details...	21	-0.63	-1.30	0.042	0.023	1.000	16334	tags=>57%, list=29%, signal=>8%

GSEA results for Enterocyte with modified max read length

	GS follow link to MSigDB	GS DETAILS	SIZE	ES	NES	NOM p-val	FDR q-val	PWER p-val	RANK AT MAX	LEADING EDGE
1	REACTOME_CLASS_I_PEROXISOMAL_MEMBRANE_PROTEIN_IMPORT	Details...	19	-0.62	-1.67	0.000	0.000	0.005	1685	tags=<1%, list=<1%, signal=<1%
2	REACTOME_SYNTHESIS_OF_BILE_ACIDS_AND_BILE_SALTS_VIA_TALPHA_HYDROXYCHOLESTEROL	Details...	22	-0.70	-1.44	0.001	0.003	0.605	6372	tags=>36%, list=11%, signal=>4%
3	REACTOME_FATTY_ACIDS	Details...	22	-0.68	-1.42	0.001	0.005	0.796	13408	tags=>73%, list=24%, signal=>7%
4	REACTOME_NEGATIVE_REGULATION_OF_FLTs	Details...	15	-0.71	-1.41	0.014	0.009	0.877	6432	tags=>13%, list=10%, signal=>15%
5	REACTOME_SYNTHESIS_OF_BILE_ACIDS_AND_BILE_SALTS	Details...	32	-0.66	-1.41	0.003	0.006	0.864	8377	tags=>25%, list=11%, signal=>2%
6	REACTOME_Olfactory_SIGNALING_PATHWAY	Details...	72	-0.62	-1.38	0.000	0.010	0.973	13728	tags=>69%, list=21%, signal=>9%
7	REACTOME_PI_3K CASCADE_FGFR1	Details...	21	-0.66	-1.36	0.012	0.013	0.994	15122	tags=>57%, list=27%, signal=>10%
8	REACTOME_FGFR1_LIGAND_BINDING_AND_ACTIVATION	Details...	17	-0.67	-1.35	0.029	0.015	0.996	15122	tags=>76%, list=27%, signal=>14%
9	REACTOME_SHC_MEDIANDED CASCADE_FGFR1	Details...	20	-0.66	-1.36	0.012	0.015	0.996	15122	tags=>69%, list=27%, signal=>12%
10	REACTOME_PHOSPHOLIPASE_C_MEDIANDED CASCADE_FGFR1	Details...	16	-0.67	-1.34	0.037	0.017	0.989	15122	tags=>75%, list=27%, signal=>12%
11	REACTOME_WNT_LIGAND_BIOGENESIS_AND_TRAFFICKING	Details...	25	-0.64	-1.34	0.017	0.019	1.000	14540	tags=>69%, list=25%, signal=>10%
12	REACTOME_PI_3K CASCADE_FGFR2	Details...	22	-0.63	-1.32	0.029	0.028	1.000	15122	tags=>69%, list=27%, signal=>7%
13	REACTOME_EICOSANOIDS	Details...	17	-0.65	-1.31	0.042	0.026	1.000	13408	tags=>71%, list=24%, signal=>9%
14	REACTOME_FRS_MEDIANDED_FGFR1_SIGNALING	Details...	22	-0.69	-1.31	0.019	0.028	1.000	15122	tags=>55%, list=27%, signal=>7%
15	REACTOME_AMINE_LIGAND_BINDING_RECEPORS	Details...	43	-0.60	-1.31	0.002	0.026	1.000	18363	tags=>74%, list=33%, signal=>10%

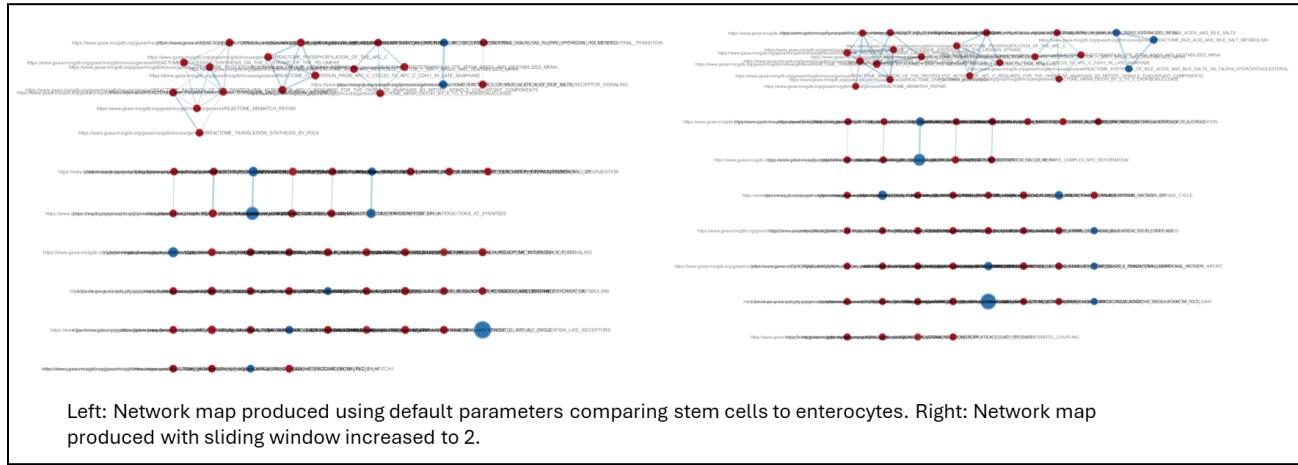
Enrichment Mapping:

There was no direct comparison we could make to the reference paper for enrichment mapping, however, each group member was able to generate enrichment maps that somewhat resemble each other. We attribute the differences between runs to how GSEA runs the initial analysis that we then used to make enrichment maps in Cytoscape. Due to the use of permutation testing and random seeding in GSEA, we would expect differences in which enriched pathways are most significant. This will affect the downstream enrichment map visualization. However, the in-group comparison between network maps are still similar with nodes correlating between separate analyses. For instance, each member using default parameters produced an enrichment map with seven upregulated nodes connected with one of the nodes labeled as REACTOME_TRANSLATION_SYNTHESIS_BY_POLK.

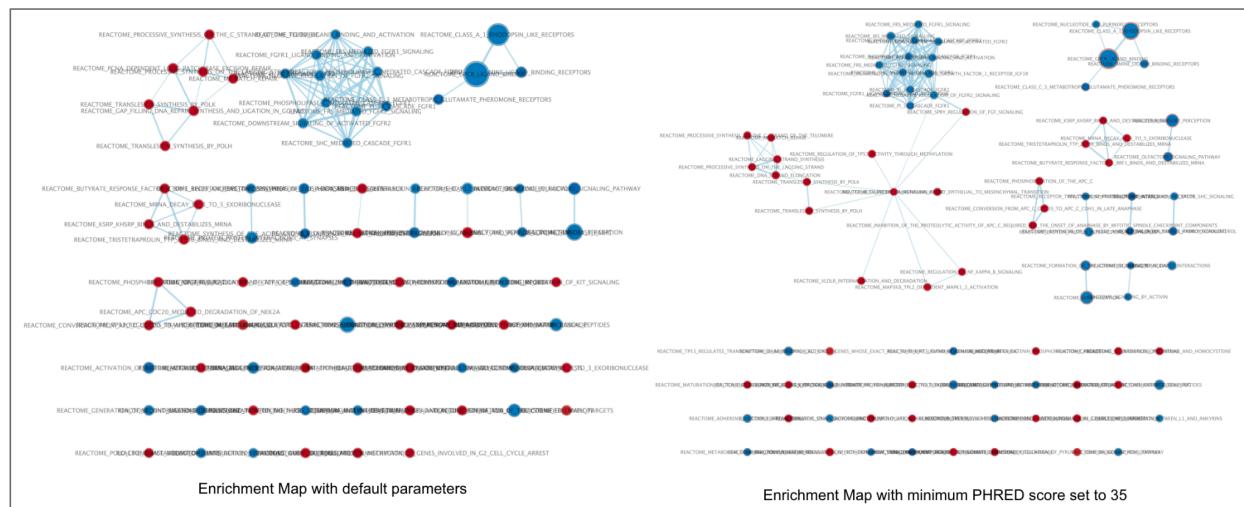


Sliding Window Parameter Change. Adjusting the sliding window during trimming from 1 to 2 slightly altered the final network maps. However, they remained similar overall. In general, the increased window produced more connected nodes with common nodes shared between the two analyses. For instance, with default parameters the left of the network map contains seven nodes connected that are upregulated with one of the nodes being REACTOME_MISMATCH_REPAIR. With the increased window, this portion of the network map increased from seven nodes to nine while still retaining the REACTOME_MISMATCH_REPAIR

node. These results are expected since the increased window should trim less reads, and retain more, leading to an increased number of nodes in our downstream enrichment maps.



Min PHRED Score Parameter Change. Increasing the minimum PHRED score parameter from 20 to 35 showed minor, yet significant differences between enrichment maps. The map with default parameters exhibited two major clusters — one of gene sets upregulated in stem cells, and one of gene sets upregulated in enterocytes. In this map, the clusters were lone standing, and not connected to any other clusters or gene sets. In the enrichment map with the minimum PHRED score parameter increased to 35, many previously unconnected gene sets became intertwined, and the two major clusters became indirectly connected through this new network. This increase in connections is likely due to the elimination of lower quality reads. With the removal of lower quality reads, previously misrepresented or ambiguous connections due to noise now show clearer, more reliable associations.



Methods:

Raw RNAseq data for our report was generated from SRA#SRP242699 (GSE GSE143915). Raw reads were downloaded from the European Nucleotide Archive and stored on The University of Arizona's HPC. Slurm batch jobs were used to run Fastp (v0.23.2), FastQC (v0.11.9), STAR (v2.7.3a), and MultiQC (v1.20) using the HPC. Parameter changes were done at the level of Fastp and consisted of - changing the default minimum PHRED score to 35, increasing the sliding window to 2, and decreasing the maximum read length to 85. STAR annotations were done using mouse genome gencode.vM36.primary_assembly.annotation.gtf. Featurecounts (v2.0.3) was performed on STAR output .bam files for both default and modified parameters. Featurecounts output .txt files were loaded into R (v4.5) and downstream analysis was performed in RStudio (Build 563 2024.12.1). Gene IDs were converted into gene names using biomaRt (v2.64.0) in Rstudio using mmusculus_gene_ensembl dataset. DESeq2 (v1.48.0) analysis was performed in Rstudio using default parameters. Significance was determined using the reference papers published values of padj < 0.01 and > 50 mean counts. rlog (blind = false) was used to transform DESeq2 data for visualization. PCA plots were generated from rlog data using Rstudio's plotPCA function. PCA plots were generated using top 100 genes based on variance. Heatmaps were generated using Rstudio's pheatmap function using rlog transformed data. Volcano plots were generated in Rstudio using enhancedvolcano plot function with rlog transformed data. Significance for volcano plots was set to log2fold change >2 and padj <10^-48 for Alex or <10^-45 for Bryan.

All functional analyses were performed comparing stem cells to enterocytes and significance was determined using padj <0.01 and > 50 normalized count. G:profiler web interface (version e112_eg59_p19_25aa4782, database updated on 03/02/2025) with g:GOSt was used to assess pathway enrichment for KEEG, Reactome, and WikiPathways with significance defined as Benjamini-Hochberg FDR <0.05. Panther (v19.0) web interface was used for Panther pathways classification. Gene set enrichment analysis was performed using GSEA (v4.4.0) on rlog transformed data using m2.cp.reactome.v2024.1. Enrichment maps were generated in Cytoscape (v3.10.3) by directly exporting GSEA data into Cytoscape using GSEA's enrichment map visualization tool.

The final SRR accession numbers for each biological replicate used in our analysis are SRR10913313 (Stem, BR#1), SRR10913315 (Stem, BR#2), SRR10913317 (Stem, BR#3), SRR10913319 (AbsPro, BR#1), SRR10913321 (AbsPro, BR#2), SRR10913323 (AbsPro, BR#3), SRR10913325 (SecPDG, BR#1), SRR10913328 (SecPDG, BR#2), SRR10913331 (SecPDG, BR#3), SRR10913334 (SecPDG, BR#4), SRR10913337 (Tuft, BR#1), SRR10913340 (Tuft, BR#2), SRR10913343 (Tuft, BR#3), SRR10913346 (Tuft, BR#4), SRR10913349 (Tuft, BR#5), SRR10913352 (Enterocyte, BR#1), SRR10913354 (Enterocyte, BR#2), SRR10913356 (Enterocyte, BR#3), SRR10913358 (Enterocyte, BR#4), SRR10913360 (Enterocyte, BR#5), SRR10913362 (EEC, BR#1), SRR10913364 (EEC, BR#2).

AI was used to assist in generating code for HPC data processing and for Rstudio analysis. AI models used were ChatGPT-4, Claude 3.7 Sonnet, and Cursor (v0.46.6).

References:

1. Ashton, G. H., Morton, J. P., Myant, K., Phesse, T. J., Ridgway, R. A., Marsh, V., Wilkins, J. A., Athineos, D., Muncan, V., Kemp, R., Neufeld, K., Clevers, H., Brunton, V., Winton, D. J., Wang, X., Sears, R. C., Clarke, A. R., Frame, M. C., & Sansom, O. J. (2010). Focal Adhesion Kinase Is Required for Intestinal Regeneration and Tumorigenesis Downstream of Wnt/c-Myc Signaling. *Developmental Cell*, 19(2), 259–269. <https://doi.org/10.1016/j.devcel.2010.07.015>
2. Liu, Y., & Chen, Y.-G. (2020). Intestinal epithelial plasticity and regeneration via cell dedifferentiation. *Cell Regeneration*, 9(1). <https://doi.org/10.1186/s13619-020-00053-5>
3. Pond, K. W., Morris, J. M., Alkhimenok, O., Varghese, R. P., Cabel, C. R., Ellis, N. A., Chakrabarti, J., Zavros, Y., Merchant, J. L., Thorne, C. A., & Paek, A. L. (2022). Live-cell imaging in human colonic monolayers reveals ERK waves limit the stem cell compartment to maintain epithelial homeostasis. *eLife*, 11. <https://doi.org/10.7554/elife.78837>
4. Raffeiner, P., Hart, J. R., García-Caballero, D., Bar-Peled, L., Weinberg, M. S., & Vogt, P. K. (2020). An MXD1-derived repressor peptide identifies noncoding mediators of MYC-driven cell proliferation. *Proceedings of the National Academy of Sciences*, 117(12), 6571–6579. <https://doi.org/10.1073/pnas.1921786117>
5. VanDussen, K. L., Carulli, A. J., Keeley, T. M., Patel, S. R., Puthoff, B. J., Magness, S. T., Tran, I. T., Maillard, I., Siebel, C., Kolterud, Å., Grosse, A. S., Gumucio, D. L., Ernst, S. A., Tsai, Y.-H., Dempsey, P. J., & Samuelson, L. C. (2012). Notch signaling modulates proliferation and differentiation of intestinal crypt base columnar stem cells. *Development*, 139(3), 488–497. <https://doi.org/10.1242/dev.070763>
6. Zhao, R., & Michor, F. (2013). Patterns of Proliferative Activity in the Colonic Crypt Determine Crypt Stability and Rates of Somatic Evolution. *PLoS Computational Biology*, 9(6), e1003082. <https://doi.org/10.1371/journal.pcbi.1003082>
7. Gehart, H., & Clevers, H. (2019). Tales from the crypt: New insights into intestinal stem cells. *Nature Reviews Gastroenterology & Hepatology*, 16(1), 19–34. <https://doi.org/10.1038/s41575-018-0081-y>