

# Investigating Inconsistencies in Single-cell RNA Sequencing Data

Pranay V. Singh (Bowdoin College — psingh@bowdoin.edu), Advisors: Juliet Emamaullee, Brittany Rocque, Cameron Goldbeck (University of Southern California)



## INTRODUCTION

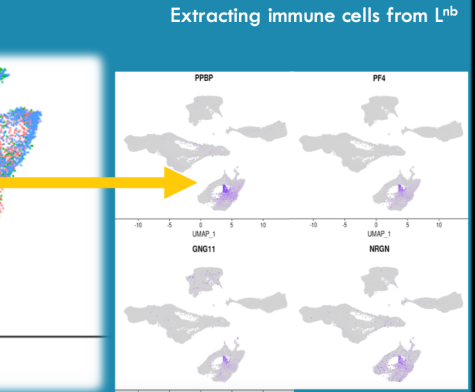
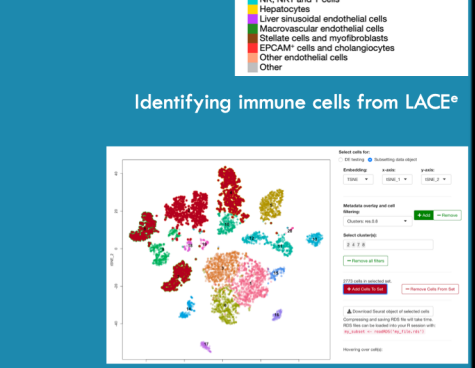
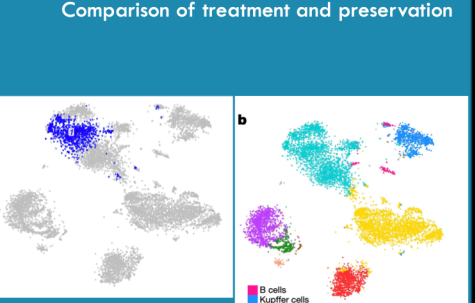
Single-cell RNA sequencing (scRNA-seq) has quickly gained traction, emerging as a robust method for computational genomics. Researchers can use scRNA-seq to better understand how disease manifests in order to find cures, leading to significant and extensive impacts worldwide. In this study, we examined three normal 'human liver' scRNASeq datasets, LCD45<sup>e</sup>, LACE<sup>e</sup>, and L<sup>nb</sup>, to see if different tissue processing techniques could alter scRNASeq outputs. Our analysis suggests a need for a standardized procedure to preserve and treat tissues before sequencing.

## METHODS

We used datasets from three sources: first, Liver atlas, cholangiocyte and endothelial enriched (LACE<sup>e</sup>), second, Liver CD45<sup>+</sup> enriched (LCD45<sup>e</sup>), and third, Non-biased Liver (L<sup>nb</sup>). Data from two of the three sources (LACE<sup>e</sup> and L<sup>nb</sup>) contained all liver cell types and data from the third source (LCD45<sup>e</sup>) contained liver, spleen, and blood immune cells. Since we sought to compare a common set of cells across all three sources, we extracted immune liver cells from each dataset.

Table 3: Comparison of methods across three peer-reviewed single cell RNA sequencing studies of the liver

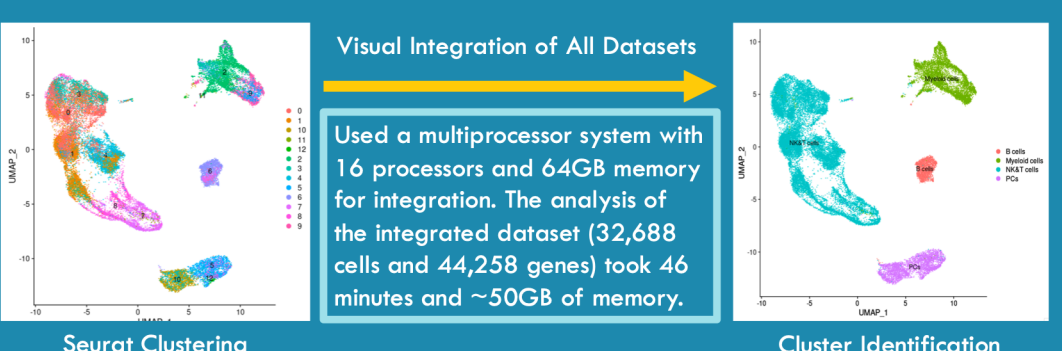
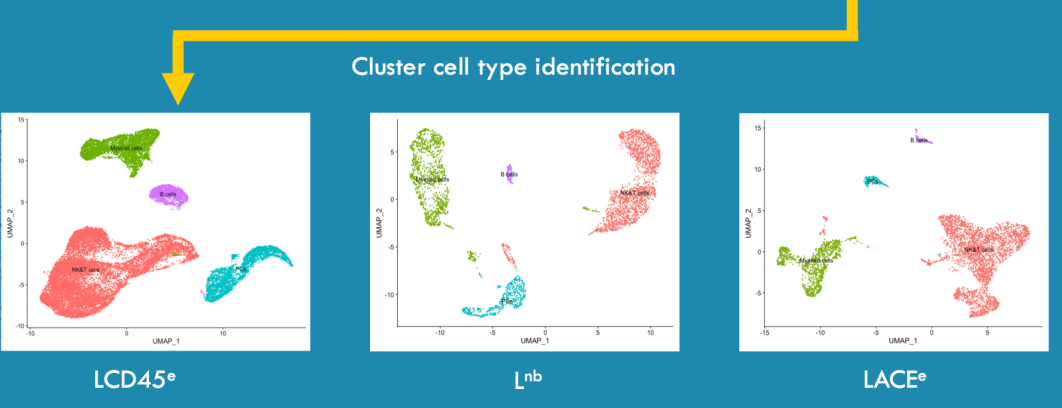
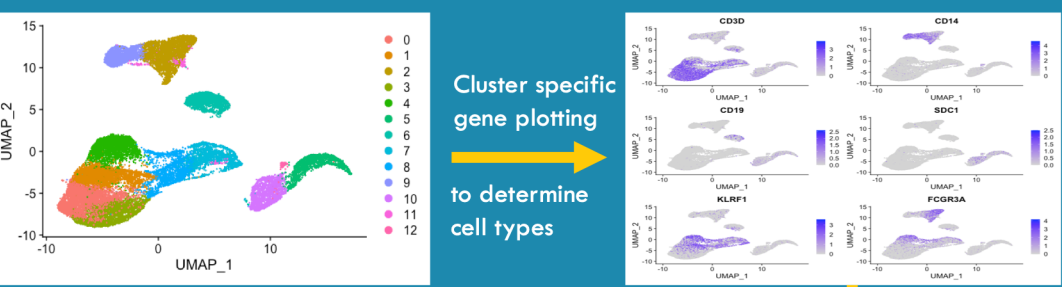
	LACE <sup>e</sup>	L <sup>nb</sup>	LCD45 <sup>e</sup>
Sample	5 liver resection pts (DC, mets or cholangitis)	5 scalable liver L <sup>nb</sup> donors	3 adult liver donors - blood, spleen and liver perfusion
Procedure	HEPES	HTX solution during recovery then HBES + EGCG after oxidative liver isolation	Not addressed
Cell fractionation	HEPES and NPCCs isolated then mixed prior to FACs	None, collagenase protocol to create cell suspension from tissue	Cell filtration and centrifugation and layered onto fluid
Removal of non-viable cells	Gradient centrifugation	Trypan blue exclusion	afflux 450 exclusion
Preservation method	Mix of cryopreserved and fresh	Fresh	Isolated cells fresh
Enrichment of cell population	Yes - FACs	No cell sorting	Yes - FACs
Enrichment for lymphocytes	No	No	Yes - marker CD45
Enrichment for CD45 <sup>+</sup>	No	No	No
Enrichment for B cells	Yes - marker CD19, PECAM	No	No
Enrichment for Cholangiocytes	Yes - marker SOX9	No	No
Removal of low-quality cells	Yes - excluded 42% RING2071 transcripts	Removed >50% mitochondrial content	No
Number of cells with any scRNA-seq	33,372	8,444	70,706
scRNA-seq technique	10x Genomics	10x	10x
Identification of cell types	Scanpy (tSNE) + 1000 nearest + 2, 100 + 25)	Not addressed	Seurat v3
Sample access site	Co-clustered (Seurat)	Co-clustered (Fig 2)	Co-clustered (Fig 3)
Different programs within site	Co-clustered	Not applicable	Not applicable
Cell fractionation	Not addressed	Not removed (with rationale)	Removed



Extraction process for LCD45<sup>e</sup> – split cells into 3 groups, plotted blood/spleen genes to identify liver immune cells (blue)

## METHODS cont.

- We normalized, scaled, clustered, non-linearly dimensionally reduced (UMAP), and graphed each of the datasets and their genes to determine cluster cell types and identified four types: Myeloid, NK&T, B, and Plasma.
- Computed differential gene expression, comparing one cell type to all other cell types within and across datasets, and we found the absolute count gene expression for each cell type, both of which would be used for our analysis.
- Employed three types of evidence: cell proportions, gene expression, and visual integration.



## RESULTS

- Cell proportions differed significantly across sets. B cells were scarce but made up a higher proportion in the LCD45<sup>e</sup> dataset ( $\chi^2$  test,  $p < 0.01$ ). NK and T cells were the largest subset across all datasets but made up a smaller proportion in the L<sup>nb</sup> cell population.
- Pairwise-correlation plots using gene expression counts separated by study demonstrated greater similarity between LCD45<sup>e</sup> and L<sup>nb</sup> ( $R^2 = 0.91$ ), and comparison of LACE<sup>e</sup> to L<sup>nb</sup> and LCD45<sup>e</sup> revealed decreased concordance ( $R^2 = 0.66$ ,  $R^2 = 0.61$ ). When separated by cell types, consistent  $R^2$  values and results were found.
- Volcano plots of differential gene expression across studies revealed the number of differentially expressed genes with significant fold change values ( $p < 0.01$ ,  $0.8 < FC < 1.25$ ). Across cell types, we found an average of 59.75, 145, and 117.25 significantly differentially expressed genes for L<sup>nb</sup> vs LCD45<sup>e</sup>, LCD45<sup>e</sup> vs LACE<sup>e</sup>, and LACE<sup>e</sup> vs L<sup>nb</sup>, respectively, resulting in relatively low differential gene expression for L<sup>nb</sup> and LCD45<sup>e</sup> when compared to the other two combinations.
- We found the top 10 differentially expressed genes in each cell type for each data set and collected the unique genes. Heatmaps of the 12-13 unique genes perceptibly depict relative differences of gene expression across datasets. From left to right, the order of columns is LACE<sup>e</sup>, LCD45<sup>e</sup>, L<sup>nb</sup>.

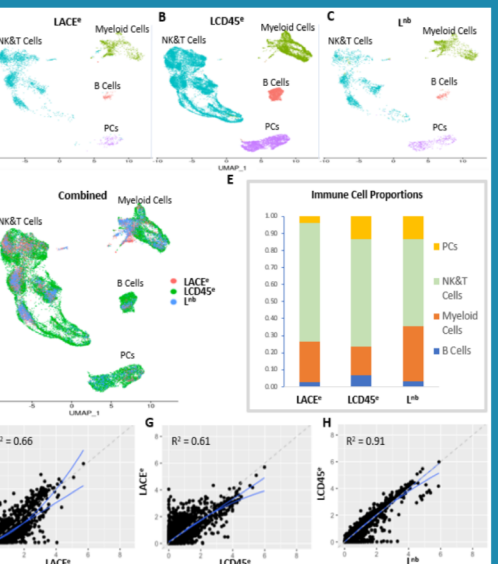
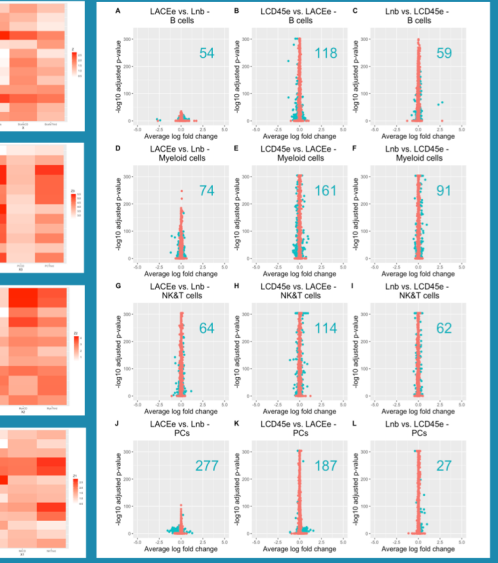


Figure 1: A-C. UMAP plots of individual studies. D. UMAP plot showing co-clustering across single cell datasets. E. Immune cell proportions across the 3 studies. F-H. Gene expression correlation plots with solid lines showing linear and quadratic regressions. Dashed line representing idealized relationship.



- Visual integration of the three datasets reveal that LACE<sup>e</sup> and L<sup>nb</sup> cluster with themselves within LCD45<sup>e</sup>, indicating that each datasets' cells cluster together.

## RESULTS cont.

- Using the gene expression results, we found the top 100 most expressed genes of L<sup>nb</sup> and compared the other two datasets to L<sup>nb</sup>, measuring gene overlap. On average, 76.5% LCD45<sup>e</sup> genes and 52.5% LACE<sup>e</sup> genes matched with L<sup>nb</sup>, demonstrating greater similarity between LCD45<sup>e</sup> and L<sup>nb</sup>.
- Comparing the top 100 differentially expressed genes between the studies by cell type, L<sup>nb</sup> and LCD45<sup>e</sup> consistently matched the most—average was 70.75% of genes matching across all four cell types, while L<sup>nb</sup> and LACE<sup>e</sup> matched at an average of 54.5%, and LACE<sup>e</sup> and LCD45<sup>e</sup> matched the worst at 49.5%, further demonstrating similarity between LCD45<sup>e</sup> and L<sup>nb</sup>.

## CONCLUSION

- Differences in cell proportions and gene expression could be due to variability in treating tissue. For scRNA-seq technique, LACE<sup>e</sup> used mCEL-seq2 while L<sup>nb</sup> and LCD45<sup>e</sup> used 10x Genomics. Variability in preservation techniques, such as cryopreservation or using various solutions to put samples in, could dehydrate cells, contributing to an increase of cells adhering together. Although the livers across these studies are not biologically identical, biological differences between the studies are expected to be comparable across all three datasets.
- There is strong evidence for L<sup>nb</sup> and LCD45<sup>e</sup> resulting in more similarity than L<sup>nb</sup> and LACE<sup>e</sup> or LCD45<sup>e</sup> and LACE<sup>e</sup>; therefore, the notable differences we see between LACE<sup>e</sup> and the other two sources are likely due to preanalytical variables. Due to the scarce and expensive nature of scRNA-seq, we thus must be careful in the process of preserving liver tissue and collecting their scRNA-seq data. Our research and analysis highlight the need for a standardized procedure for the treatment and preservation of liver cells before sequencing.

Acknowledgement: The authors are grateful to Bowdoin College for sponsoring Pranay's summer internship at USC's Transplant Research Lab.

