# Assessing the impact of cell type deconvolution on differential gene expression analysis of COVID-19

Li Qing Wang, Jana Osea, Erick Navarro, Daniela Yanez
STAT540 Group 2
Apr 6, 2022

# Outline

# INTRODUCTION ●●●

- Mechanism of COVID yet to be fully elucidated
- Cell type information is hidden in bulk tissue sequencing

- Sample size (n = 409):
  - 356 positive, 53 negative

**We hypothesize that including cell type as a covariate will influence differential gene expression analysis results.**

# AIMS ● ● ●

### AIM1

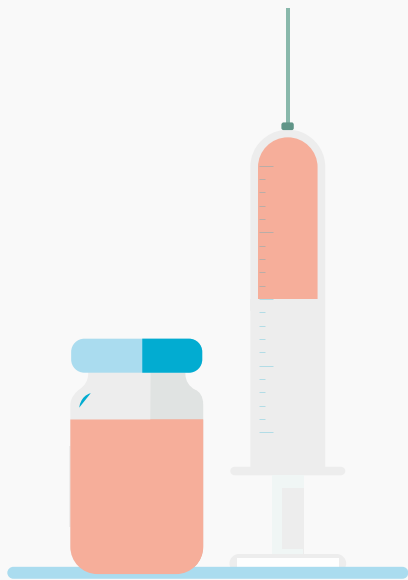Perform cell type deconvolution to the dataset in Lieberman *et al.* paper.

### AIM2

Perform a DE analysis based on infection status and using sex, age and cell type as covariates.

### AIM3

Evaluate the statistical importance of including cell type. Compare our GSEA with the Lieberman *et al.* paper.
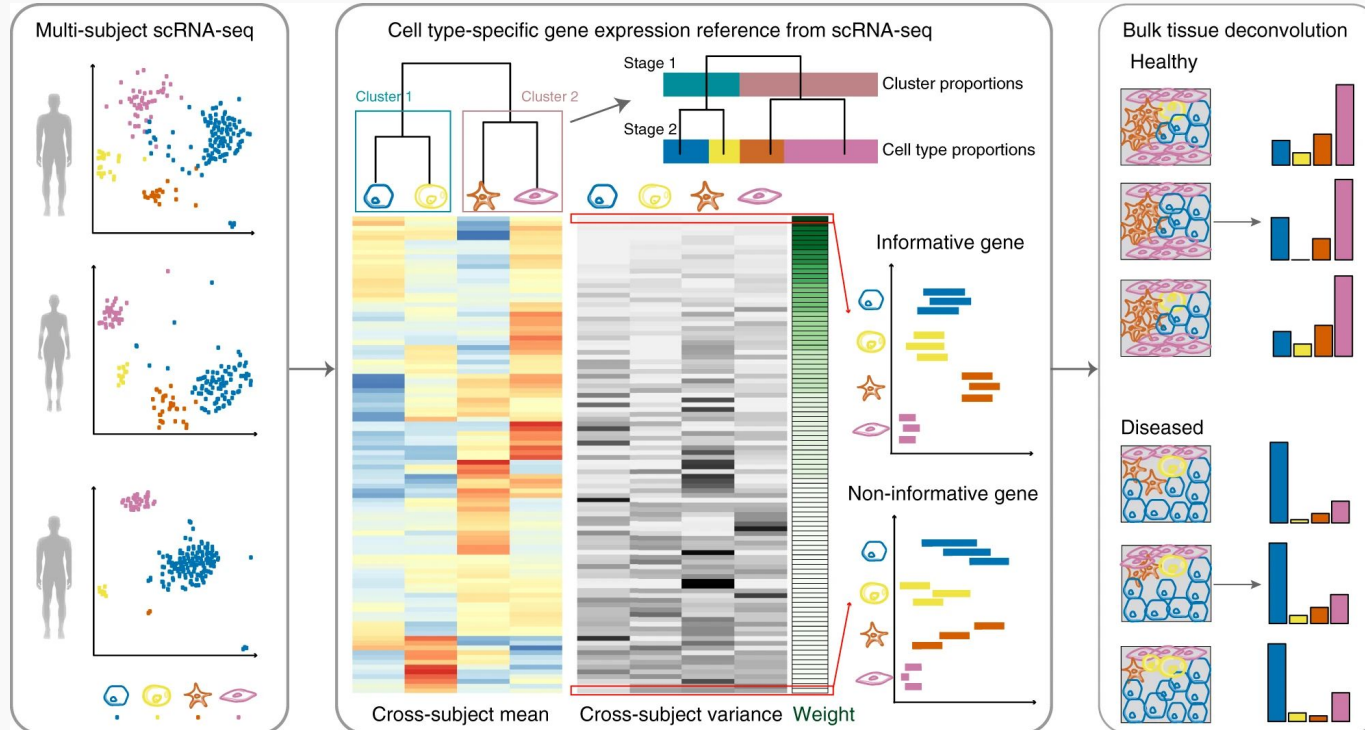
Aim 01

Cell type deconvolution

# Multi-subject Single-cell Deconvolution (MuSiC)
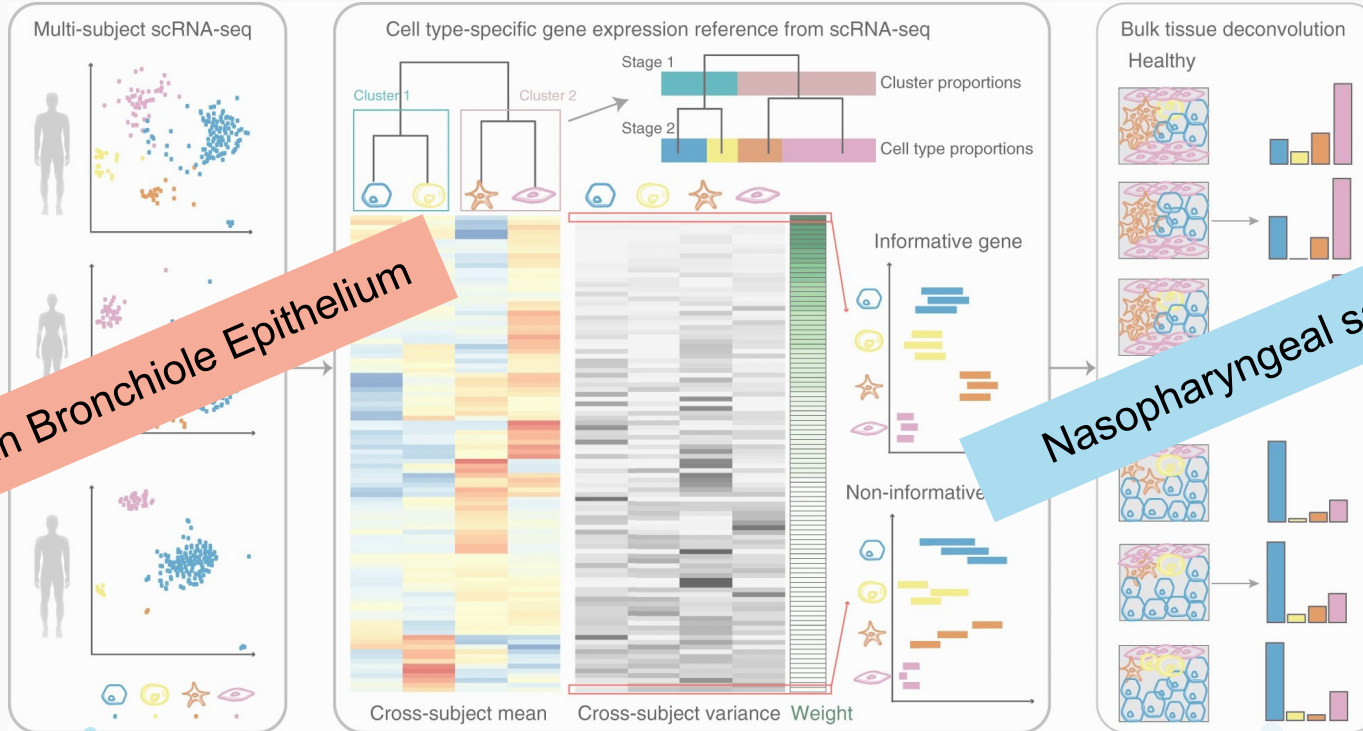
scRNAseq reference

Bulk tissue data

# Multi-subject Single-cell Deconvolution (MuSiC)



scRNAseq reference

Bulk tissue data

Human Bronchiole Epithelium

Nasopharyngeal samples

*Nature Communications (Nat Commun) ISSN 2041-1723 (online)*

# Overview of variables

| Variables | Median or Count (range or %) |
|---|---|
| Age (years) | 53 (2 - 89) |
| Gender (N females) | 217 (53.06%) |
| Covid status (N positive) | 356 (87.04%) |
| Cell type proportions (%) | Ciliated1: 0% (0-0.41%)<br>Goblet: 55.98% (0-100%)<br>FOXN4: 41.58% (0-100%)<br>Basal1: 0% (0-25.72%)<br>Fibroblast: 0% (0-29.38%)<br>Basal3: 0% (0-29.16%) |

# COVID status is associated with cell type



Cell type proportions of COVID positive samples

Cell type proportions of COVID negative samples

- Cell types with low proportions in the overall dataset were dropped

# Aim 02

**Differential gene expression**

# Workflow

Remove lowly expressed genes (<10 counts in total) $\Rightarrow$ Remove samples with missing metadata information $\Rightarrow$ Conduct an exploratory analysis

$\Downarrow$

Downstream analysis $\Leftarrow$ Visualize the results $\Leftarrow$ Perform the DE analysis with DESeq2
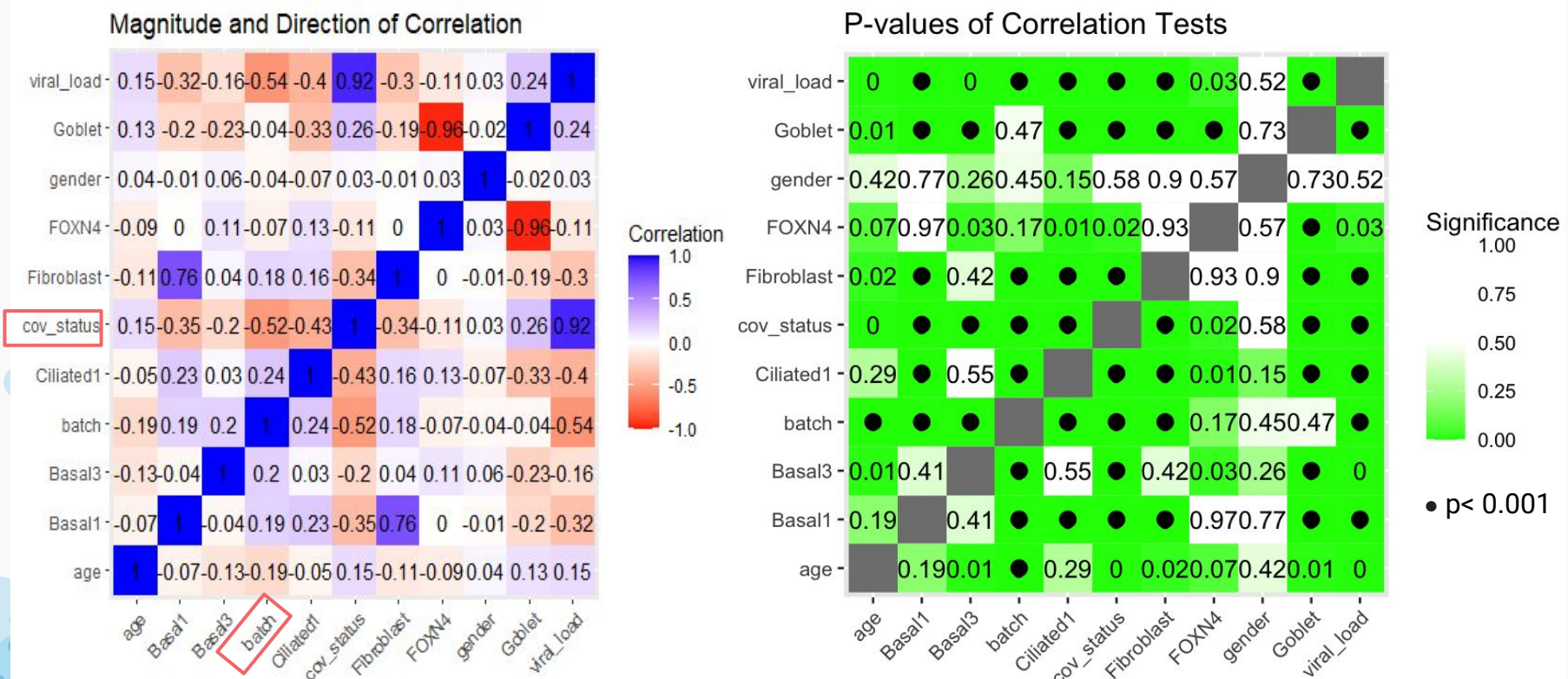
# Samples clustered based on their covid status in the PCA

# Model design

| age | gender | viral_load | cov_status | batch | Ciliated1 | Goblet | FOXN4 | Basal1 | Fibroblast | Basal3 |
|-----|--------|-----------|-----------|-------|-----------|-----------|-----------|--------|-----------|--------|
| 64 | M | 18.88 | pos | I | 0 | 0.7369542 | 0.2630458 | 0 | 0.0000000 | 0 |
| 30 | F | 21.18 | pos | I | 0 | 0.5761618 | 0.4238382 | 0 | 0.0000000 | 0 |
| 47 | M | 24.24 | pos | I | 0 | 0.6791802 | 0.3208198 | 0 | 0.0000000 | 0 |
| 67 | F | 18.91 | pos | G | 0 | 0.5092839 | 0.4907161 | 0 | 0.0000000 | 0 |
| 62 | M | 25.62 | pos | H | 0 | 0.8761618 | 0.1238382 | 0 | 0.0000000 | 0 |
| 52 | F | 25.61 | pos | H | 0 | 0.5874727 | 0.3380129 | 0 | 0.0198715 | 0 |

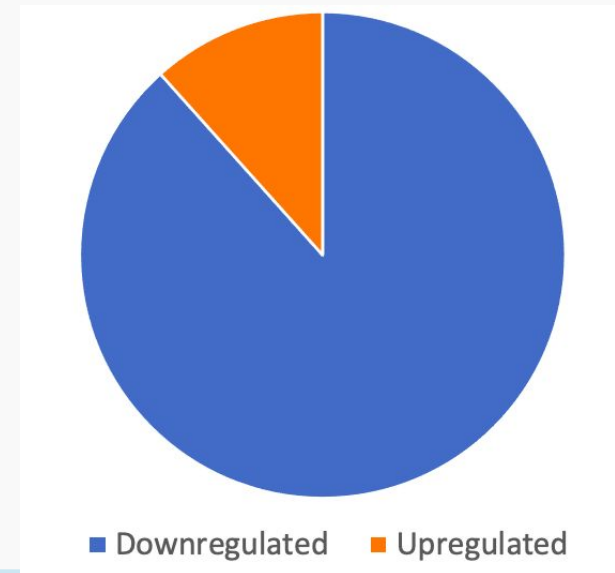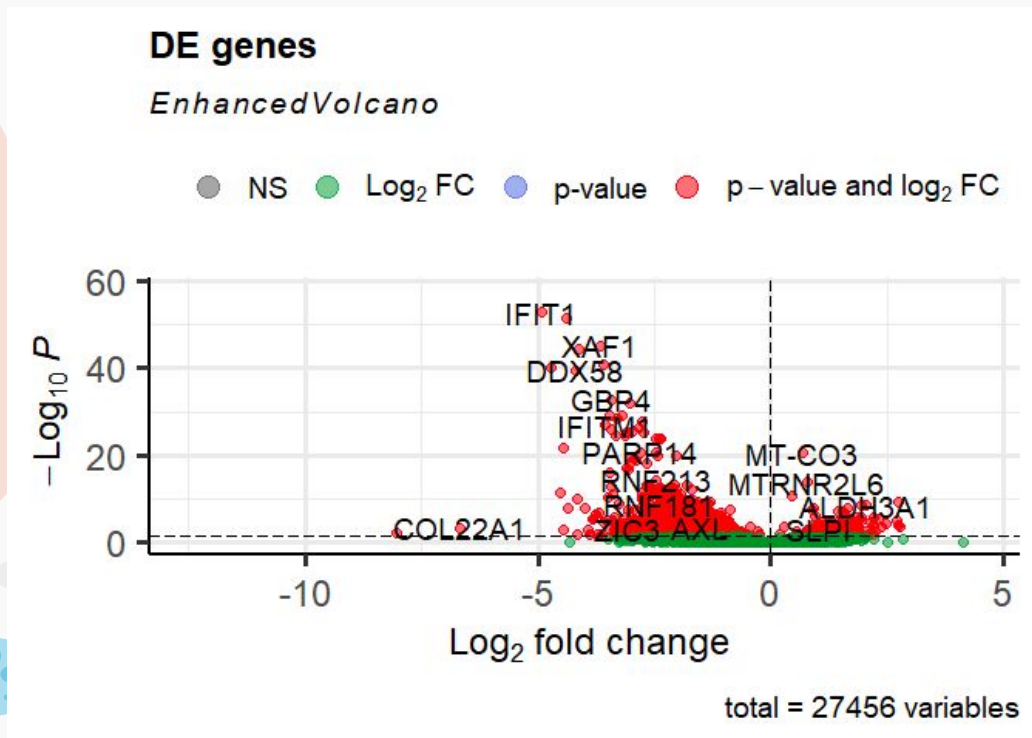# High collinearity between covid status and batch



- **All COVID+ samples were run on the same batches, all COVID- on other batches**

# Model design

design = ~cov_status + age + gender + Ciliated1 + Goblet + FOXN4 + Basal1 + Fibroblast + Basal3

| age | gender | viral_load | cov_status | batch | Ciliated1 | Goblet | FOXN4 | Basal1 | Fibroblast | Basal3 |
|-----|--------|-----------|-----------|-------|-----------|-----------|-----------|--------|-----------|--------|
| 64 | M | 18.88 | pos | I | 0 | 0.7369542 | 0.2630458 | 0 | 0.0000000 | 0 |
| 30 | F | 21.18 | pos | I | 0 | 0.5761618 | 0.4238382 | 0 | 0.0000000 | 0 |
| 47 | M | 24.24 | pos | I | 0 | 0.6791802 | 0.3208198 | 0 | 0.0000000 | 0 |
| 67 | F | 18.91 | pos | G | 0 | 0.5092839 | 0.4907161 | 0 | 0.0000000 | 0 |
| 62 | M | 25.62 | pos | H | 0 | 0.8761618 | 0.1238382 | 0 | 0.0000000 | 0 |
| 52 | F | 25.61 | pos | H | 0 | 0.5874727 | 0.3380129 | 0 | 0.0198715 | 0 |

# Most genes were downregulated

# Results

| ID | baseMean | log2FoldChange | lfcSE | stat | pvalue | padj |
|---|---|---|---|---|---|---|
| CSAG3 | 0.4478138 | -29.930404 | 1.7745408 | -16.86656 | 0 | NA |
| IFIT1 | 449.0871812 | -4.924876 | 0.3075396 | -16.01380 | 0 | 0 |
| OAS3 | 470.0654560 | -4.394149 | 0.2784236 | -15.78224 | 0 | 0 |
| XAF1 | 324.4940444 | -3.648678 | 0.2462750 | -14.81546 | 0 | 0 |
| IFI44L | 351.9585267 | -4.124200 | 0.2811023 | -14.67153 | 0 | 0 |
| OAS2 | 373.5151971 | -3.601094 | 0.2560988 | -14.06135 | 0 | 0 |
| OR52W1 | 0.2509104 | -29.966828 | 2.1357228 | -14.03124 | 0 | NA |
| IFIT2 | 567.2061236 | -4.717428 | 0.3383309 | -13.94324 | 0 | 0 |
| DDX58 | 157.9347320 | -4.202735 | 0.3035097 | -13.84712 | 0 | 0 |
| GBP4 | 162.3396890 | -3.426630 | 0.2705122 | -12.66719 | 0 | 0 |

# Gene Set Enrichment Analysis

1. **Hypergeometric enrichment**
   - MsigDB
2. **Gene set enrichment analysis (GSEA)**
   - MsigDB

   **Broad Molecular Signatures Database (MSigDB) gene sets**
   - a. **Hallmark:** summarize and represent specific well-defined biological processes.
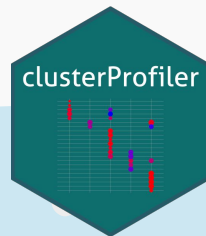
     **Total: 50**

   - b. **Gene ontology [GO]**

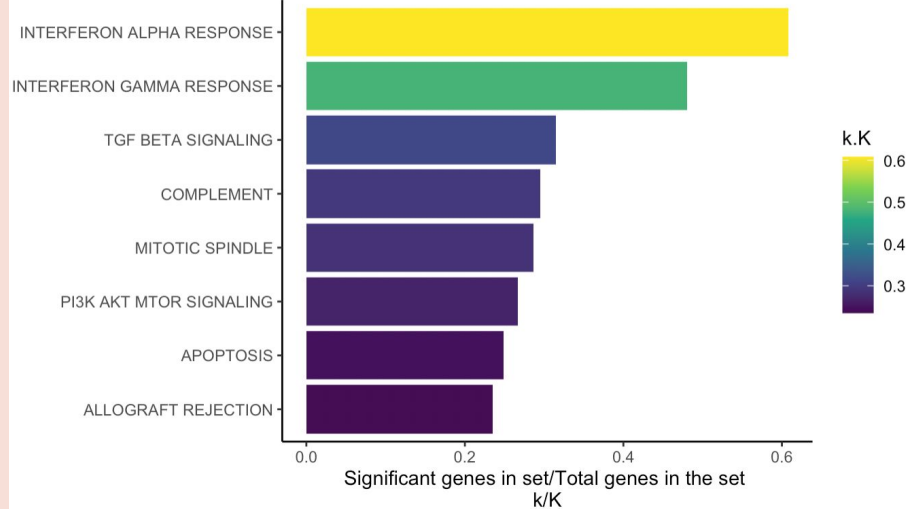   - Biological process: molecular-level activities performed by gene products. **Total: 7481**

3. **Disease enrichment analysis (GSEA)**
   - Cluster profiler
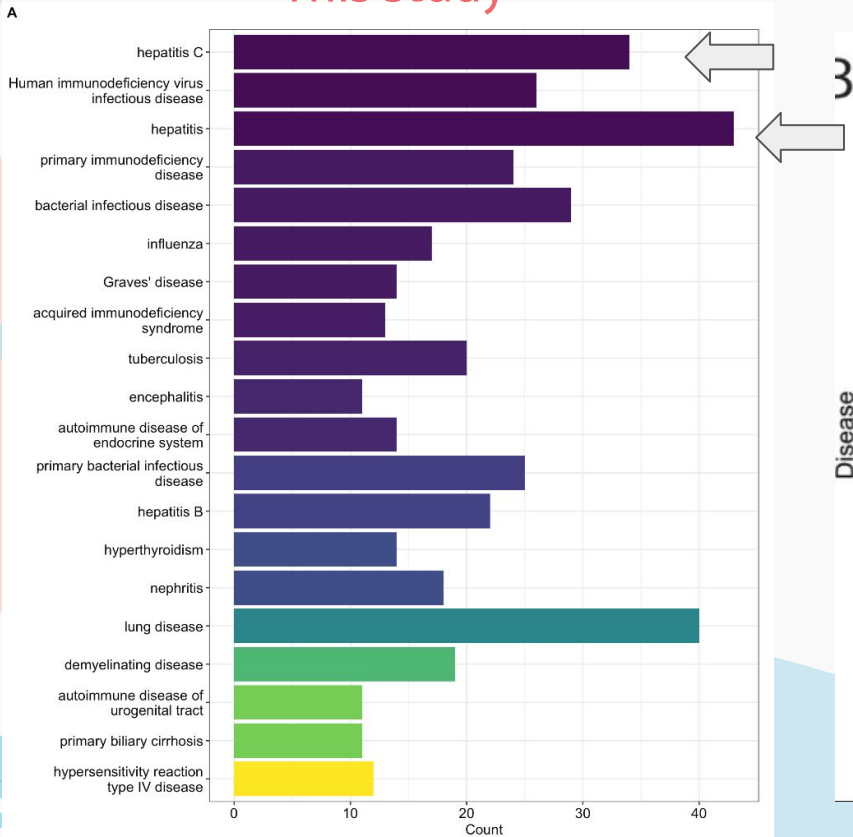   - DOSE: used for DO(Disease Ontology)

# Hallmark Gene Set Results

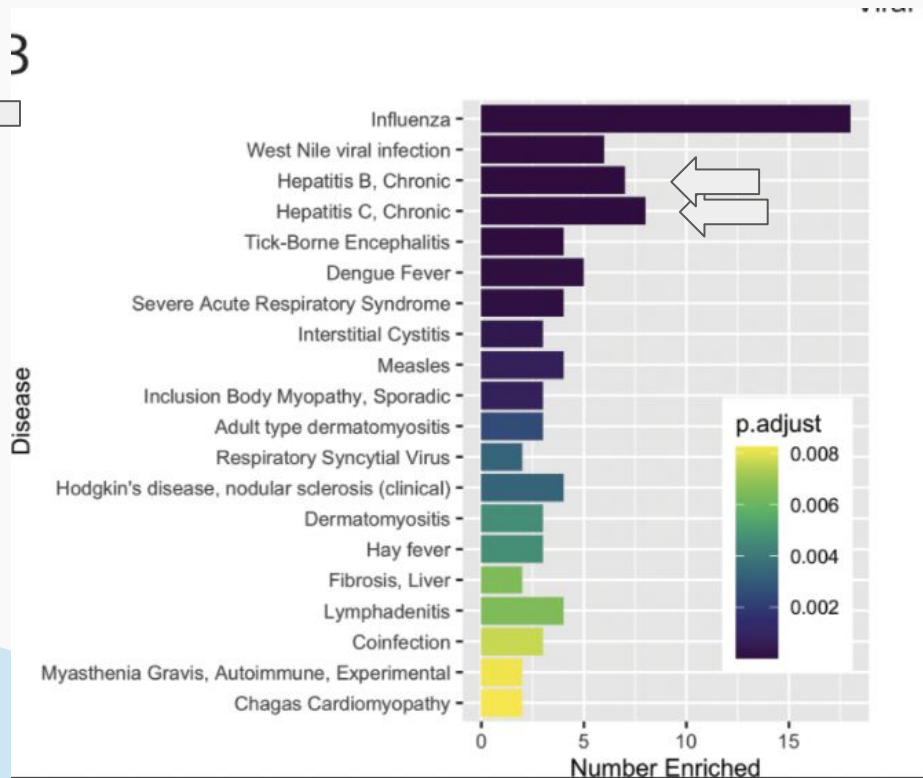Enrichment of signature genes of Interferon Alpha, Interferon Gamma, and Inflammatory Responses
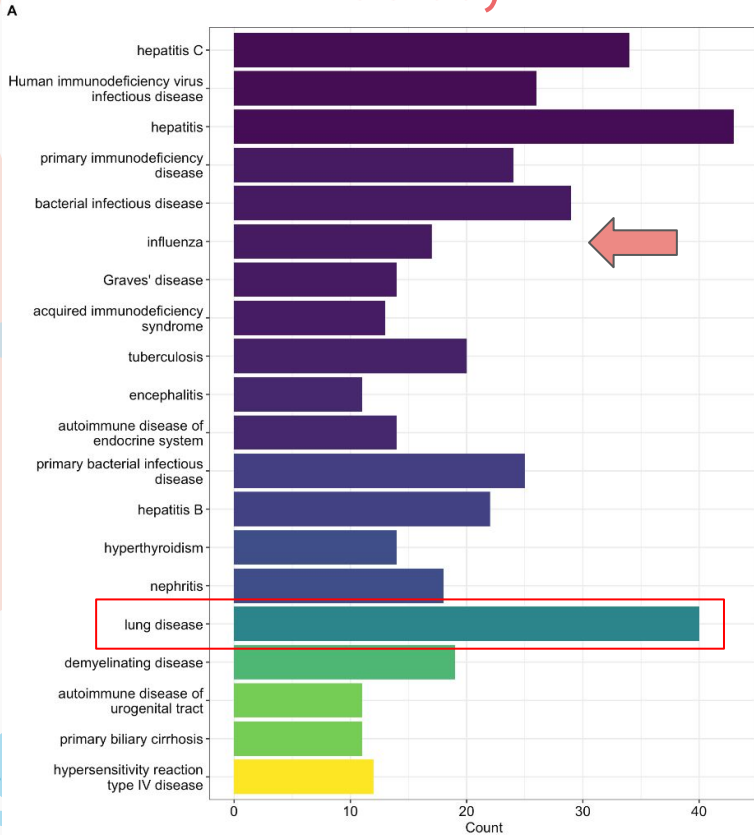
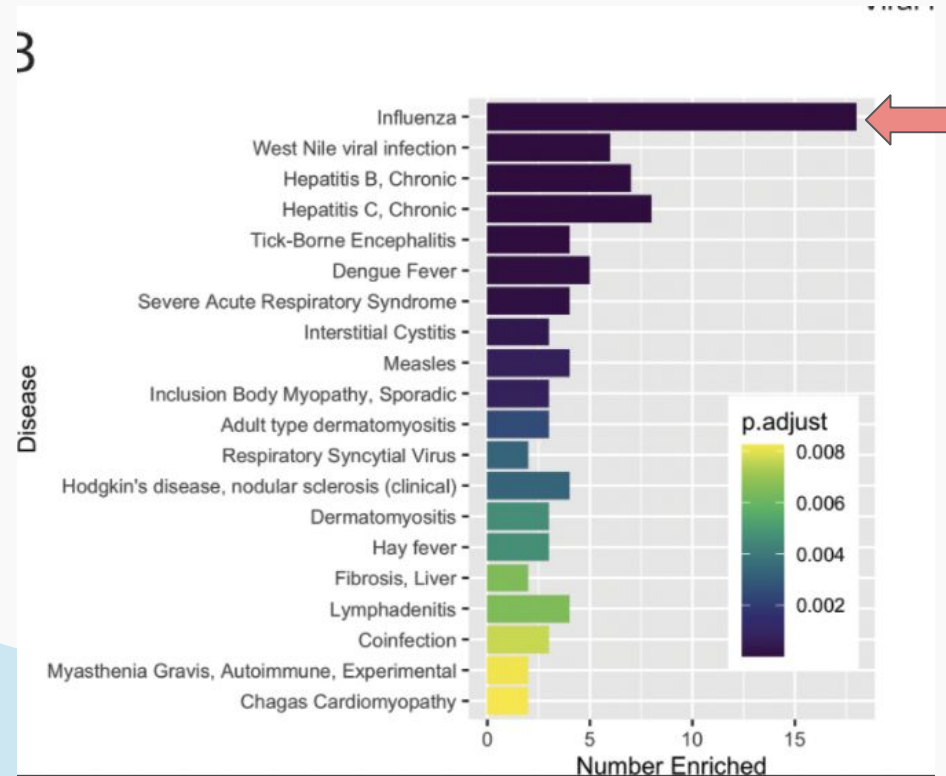# Disease Enrichment Analysis



This study

Liberman's findings

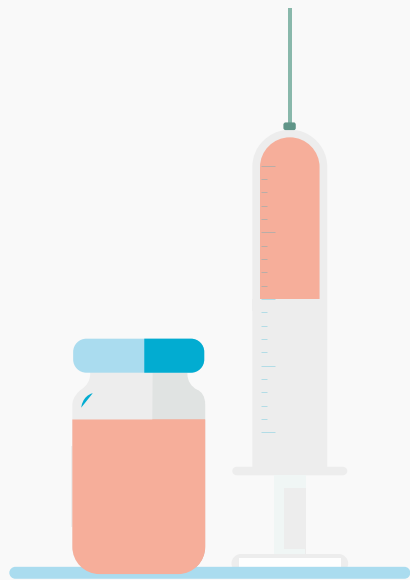# A high degree of overlap with lung disease signature
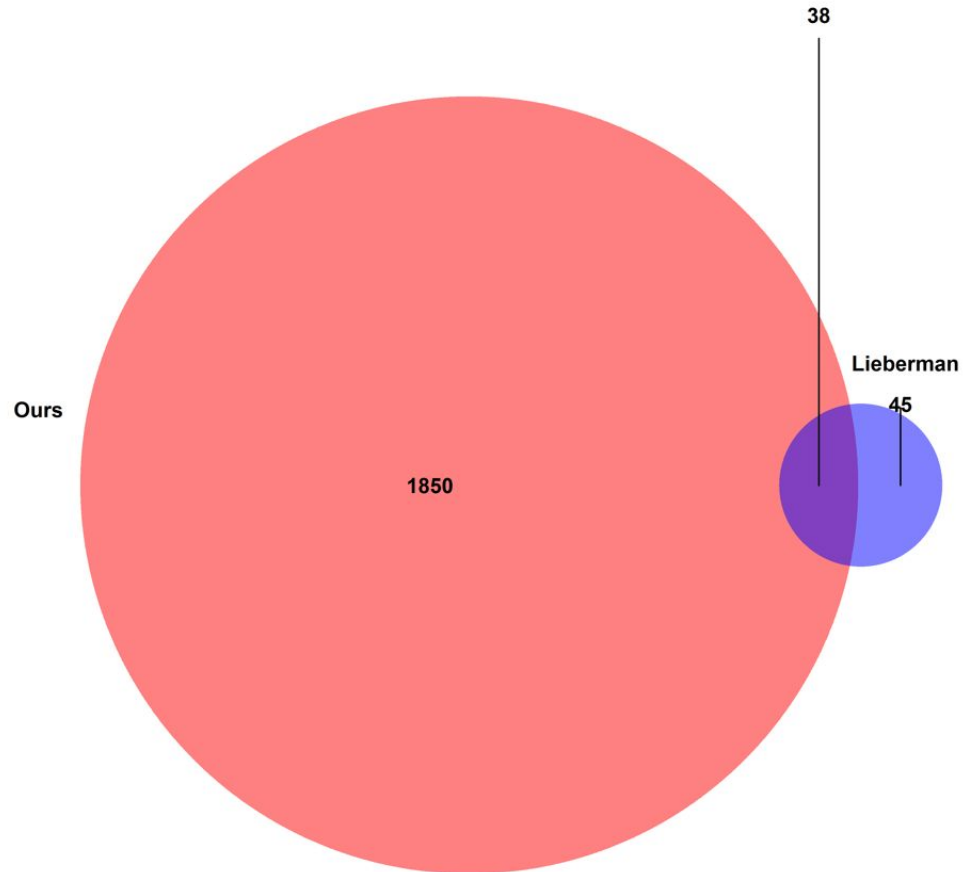
This study

Liberman's findings

# Aim 03

**Quantitative comparison**

1. **How many DE genes overlap between ours and theirs?**
2. **How many of our DE genes were better explained by the cell type model?**

# How many overlapping DE genes?

# Likelihood Ratio Test (LRT)

- Tests the following hypothesis

  H0: covid status + age + gender (reduced model)

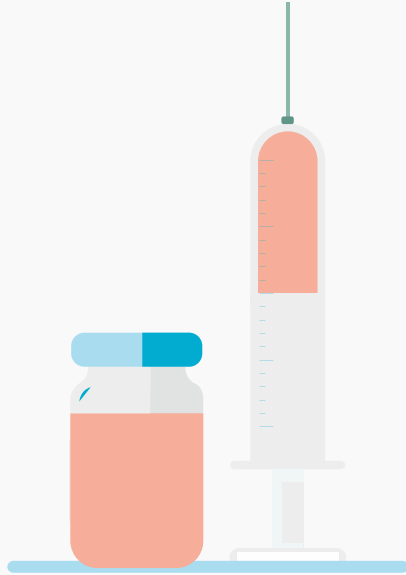  H1: covid status + age + gender + cell types (full model)

- Results
  - 7630 genes with significantly increased likelihood with full model
  - This is 27.7899% of all gene

# How many of our DE genes were better explained by the cell type model?

- Results
  - Out of the 2960 our DE genes, we find 1888 with significantly increased likelihood with full model
  - This is 85.0636% of our findings
- Observation
  - Full model better characterizes significant DE genes

# Conclusion

**Limitations and Final remarks**

# Conclusion

- Limitations
  - Large correlation between batch and Covid status, could not reproduce author's exact model
  - Reference data set is lower airways versus batch data set is upper airways
  - The experiment was not a balanced design (not possible to circumvent statistically)

- Final Remarks
  - Overall, cell type is a possible confounder associated with covid status and gene expression
  - Adding cell type corrects for possible confounding

# References

- Deprez, Marie, Laure-Emmanuelle Zaragosi, Marin Truchi, Christophe Becavin, Sandra Ruiz García, Marie-Jeanne Arguel, Magali Plaisant, et al. 2020. "A Single-Cell Atlas of the Human Healthy Airways." *American Journal of Respiratory and Critical Care Medicine* 202 (12): 1636–45.

- Lieberman, Nicole AP, Vikas Peddu, Hong Xie, Lasata Shrestha, Meei-Li Huang, Megan C Mears, Maria N Cajimat, et al. 2020. "In Vivo Antiviral Host Transcriptional Response to SARS-CoV-2 by Viral Load, Sex, and Age." *PLoS Biology* 18 (9): e3000849.

- Lukassen, Soeren, Robert Lorenz Chua, Timo Trefzer, Nicolas C Kahn, Marc A Schneider, Thomas Muley, Hauke Winter, et al. 2020. "SARS-CoV-2 Receptor ACE 2 and TMPRSS 2 Are Primarily Expressed in Bronchial Transient Secretory Cells." *The EMBO Journal* 39 (10): e105114.

- Wang, Park, X. 2019. "Bulk Tissue Cell Type Deconvolution with Multi-Subject Single-Cell Expression Reference." *Nat Commun* 10 (380).
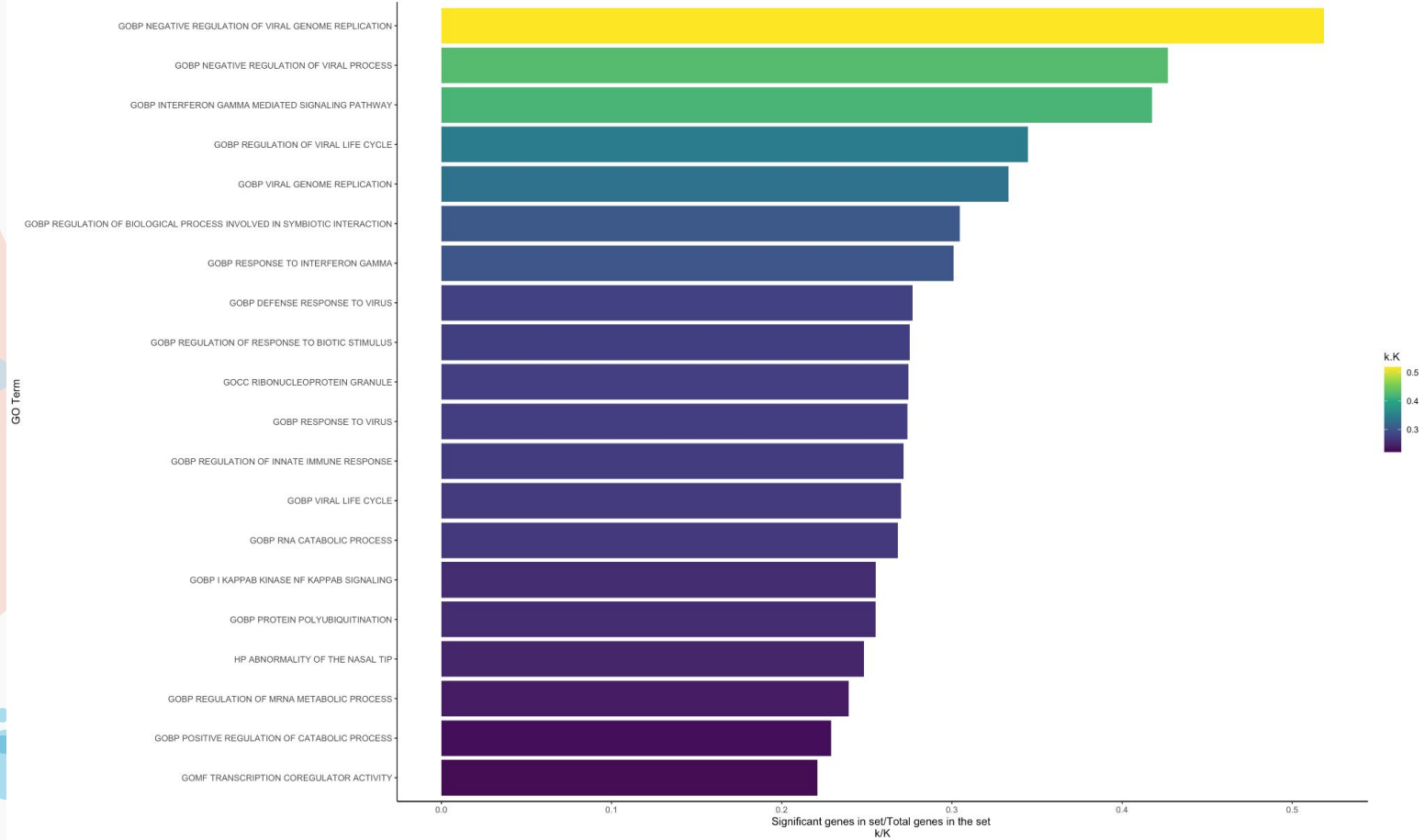
# Thank you!
# Questions?

# Extra slides

Figures that we made but are not part of the main presentation

# Cell type by viral load

# GO Terms

Differentially expressed genes