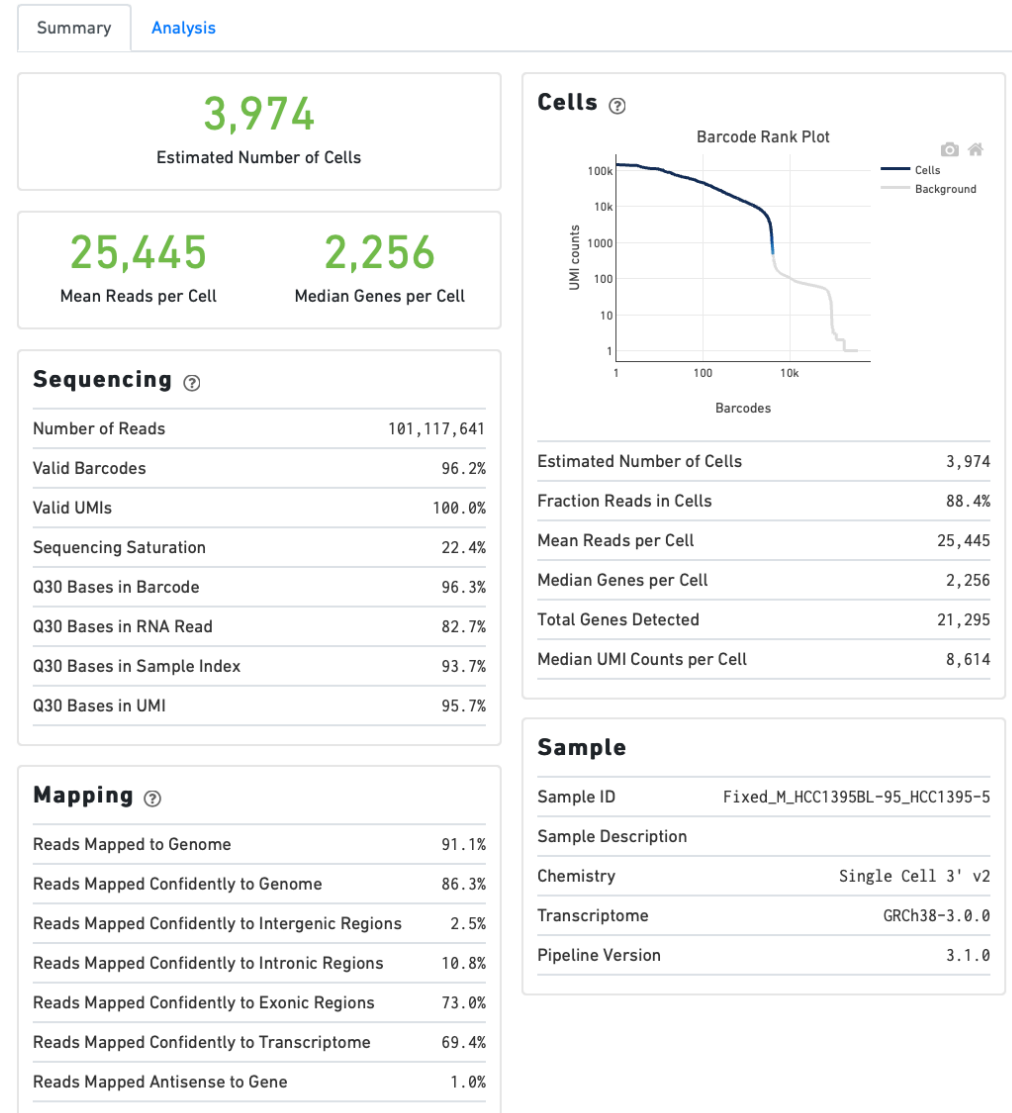


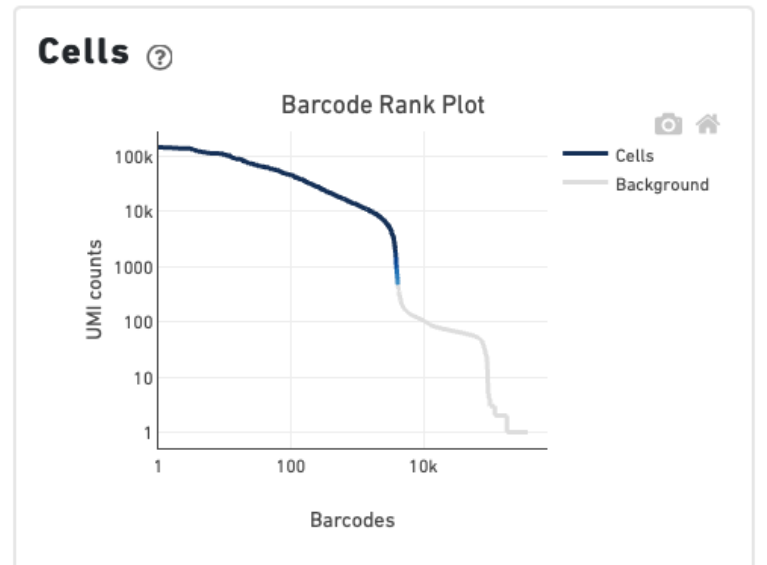
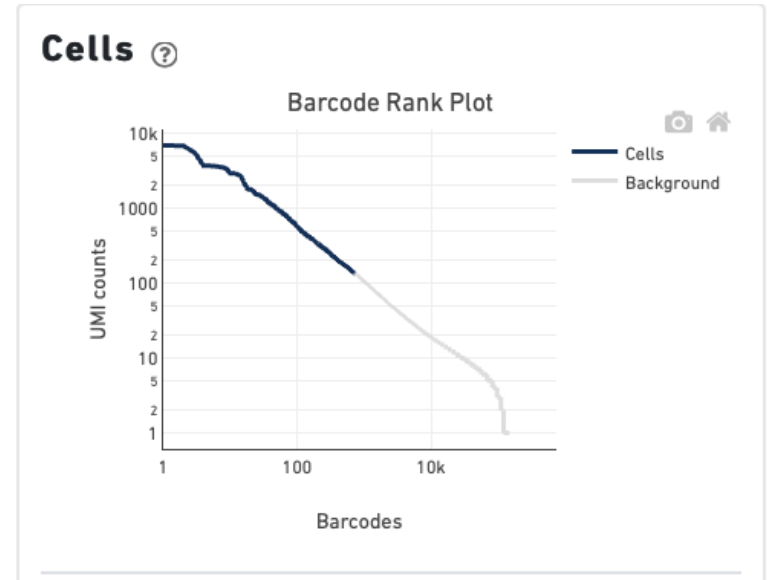
scRNA-seq preprocessing and quality control

Nathan Wong (CCBR) and Vicky Chen (CCR-SF Bioinformatics Group)

- Used to process FASTQ files for 10X samples
- Generates UMI expression matrices, basic sample statistics, and interactive analysis platform

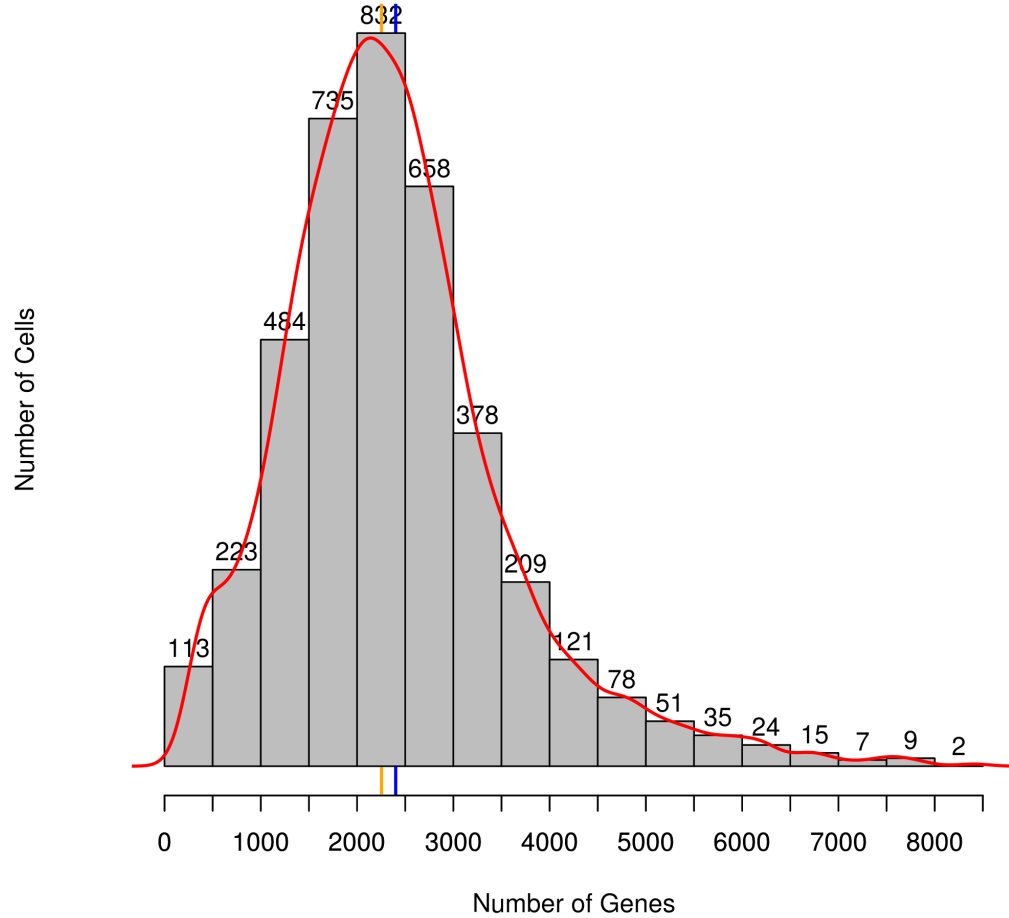


- **Barcode Rank Plot (Knee plot) can be used to determine sample quality**
- **Cell Ranger 3 increased sensitivity for low UMI cell populations**



Fixed_M_HCC1395BL-95_HCC1395-5

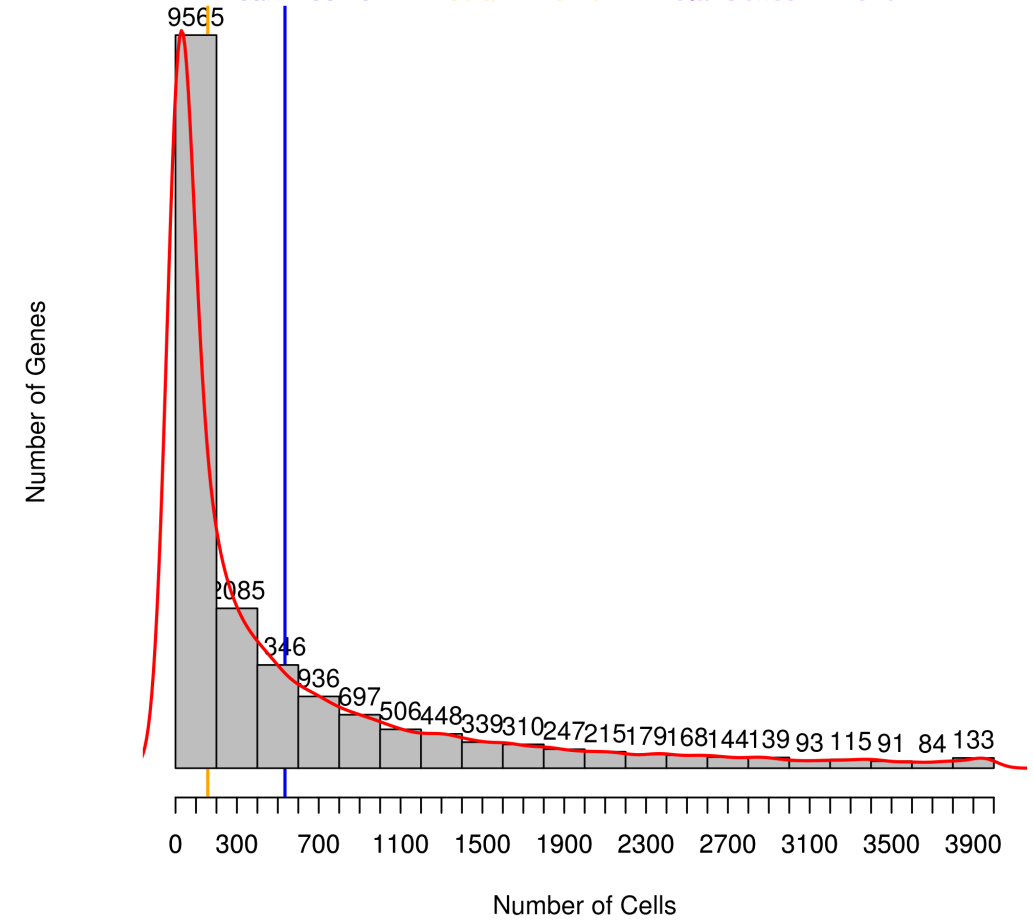
Mean = 2401.28 Median = 2254



Genes per Cell

Fixed_M_HCC1395BL-95_HCC1395-5

Mean = 534.9 Median = 157.5 Total Genes = 17840



Cells per Gene

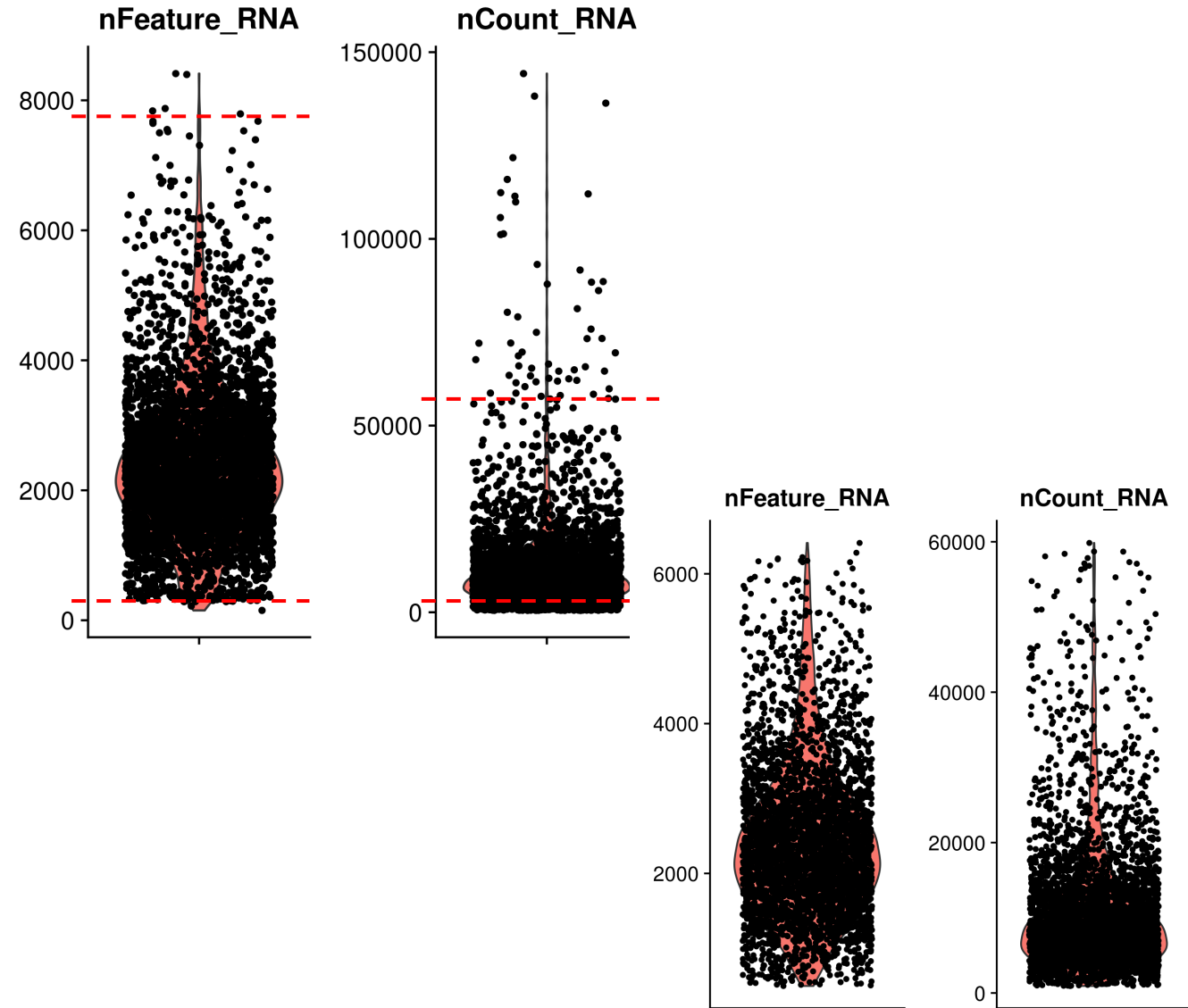
Cell Filtering

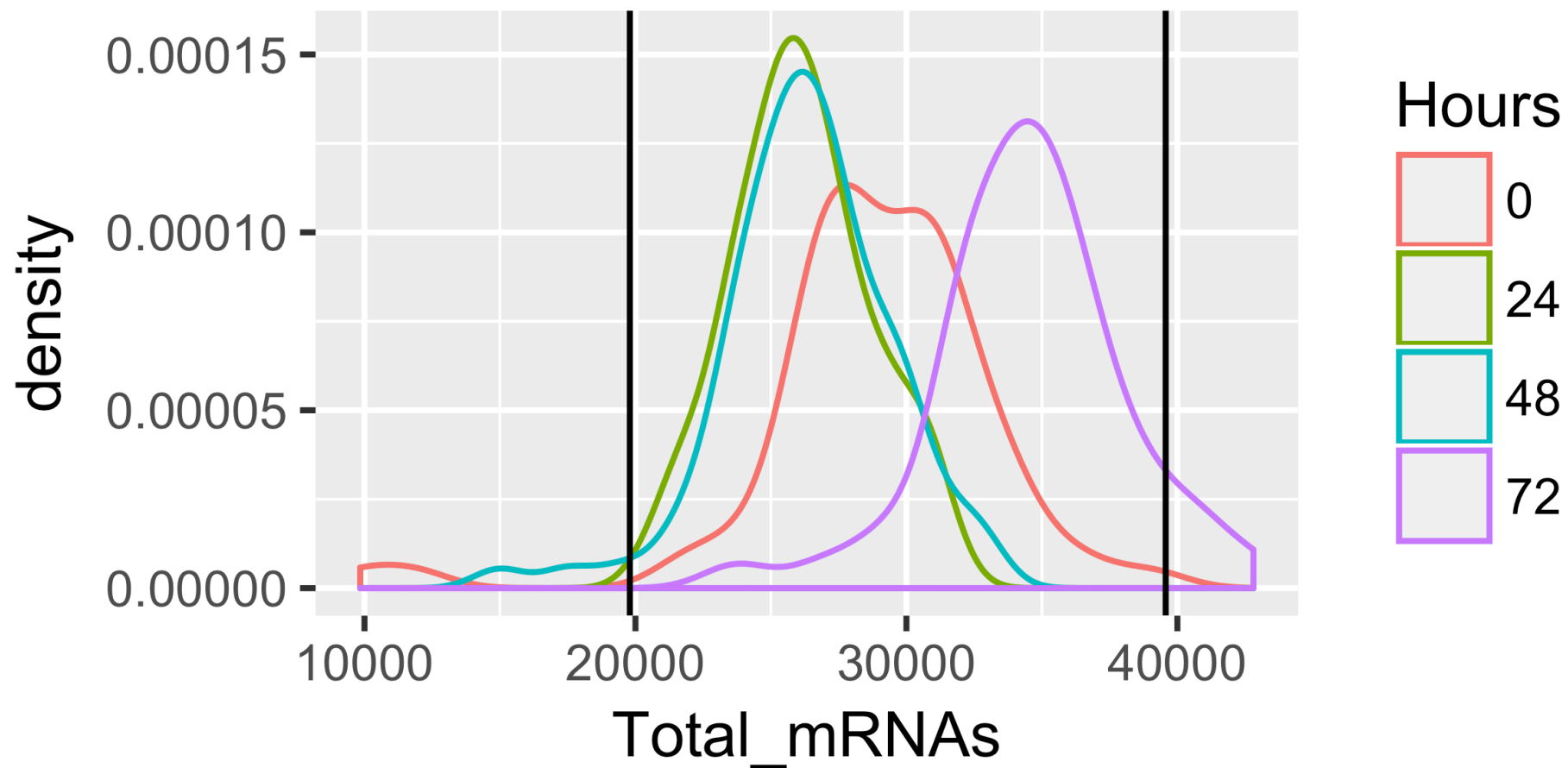
- **Useful because low quality cells or doublets/multiplets might be included in data**
- **Doublet/Multiplets are when more than one cell is captured and labeled with the same cell barcode**
- **Low quality cells include dying cells or cells with broken membranes**
 - Contains lower amounts of genes
 - Has a higher expression of mitochondrial genes

- **Low quality cells or doublets/multiplets might be included in data**
- **Filtering is used to remove the excess noise to have a clean analysis**
- **Stringent filters risk losing useful data**
- **Loose filters risk leaving in noise**

Cell Filtering

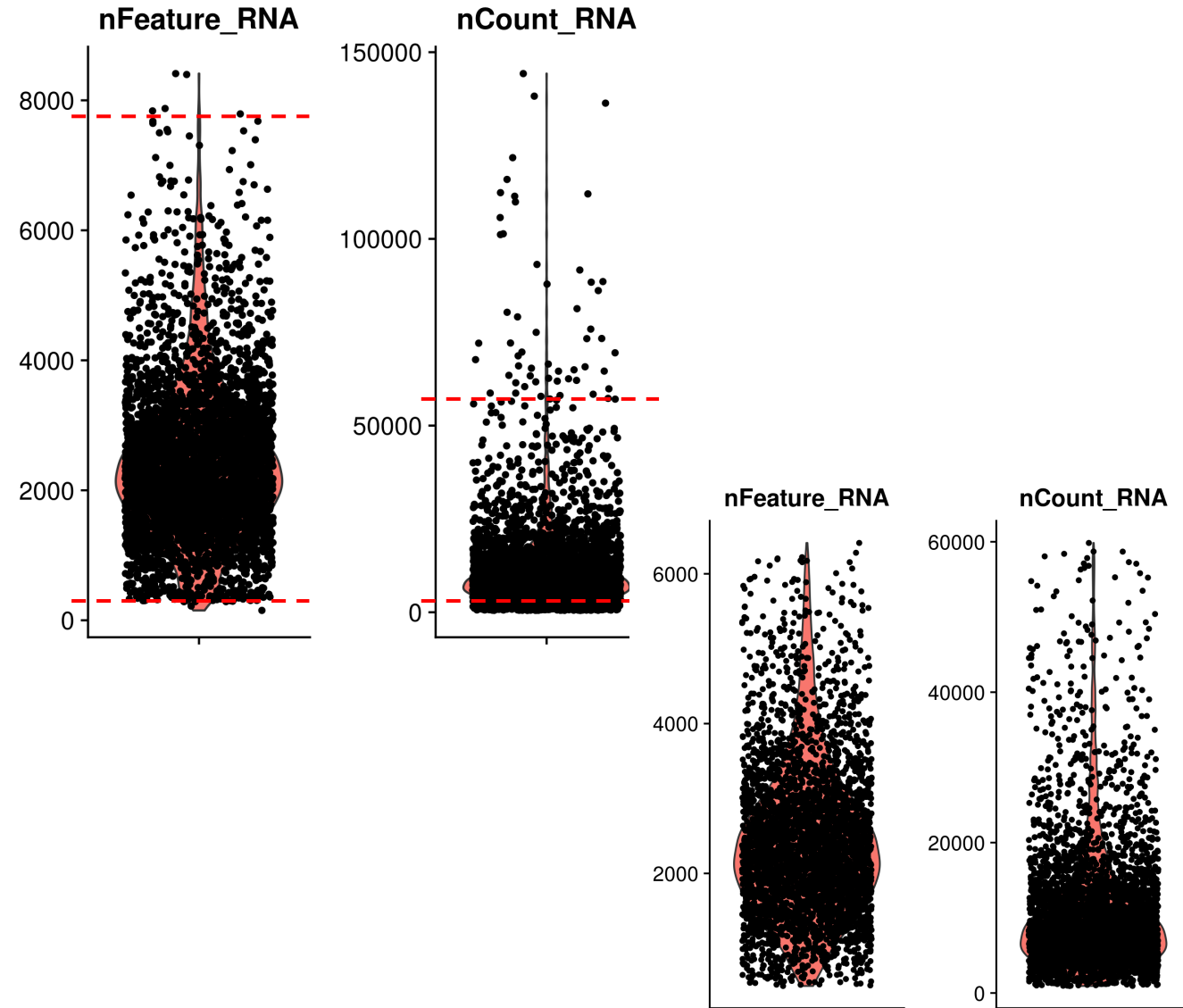
- Different cell types have different expression levels
- Filtering based on UMI count, gene count, and mitochondrial gene expression
- UMI count and gene count filters based on negative binomial distribution
- Other distribution and statistical methods can be used





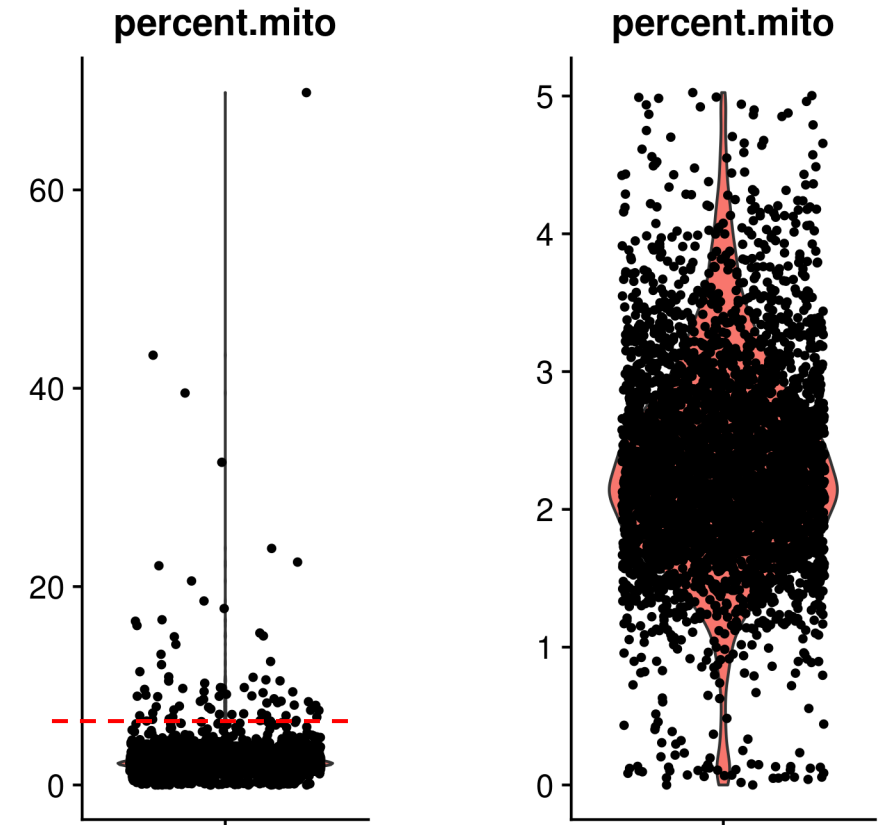
Cell Filtering

- Different cell types have different expression levels
- Filtering based on UMI count, gene count, and mitochondrial gene expression
- UMI count and gene count filters based on negative binomial distribution



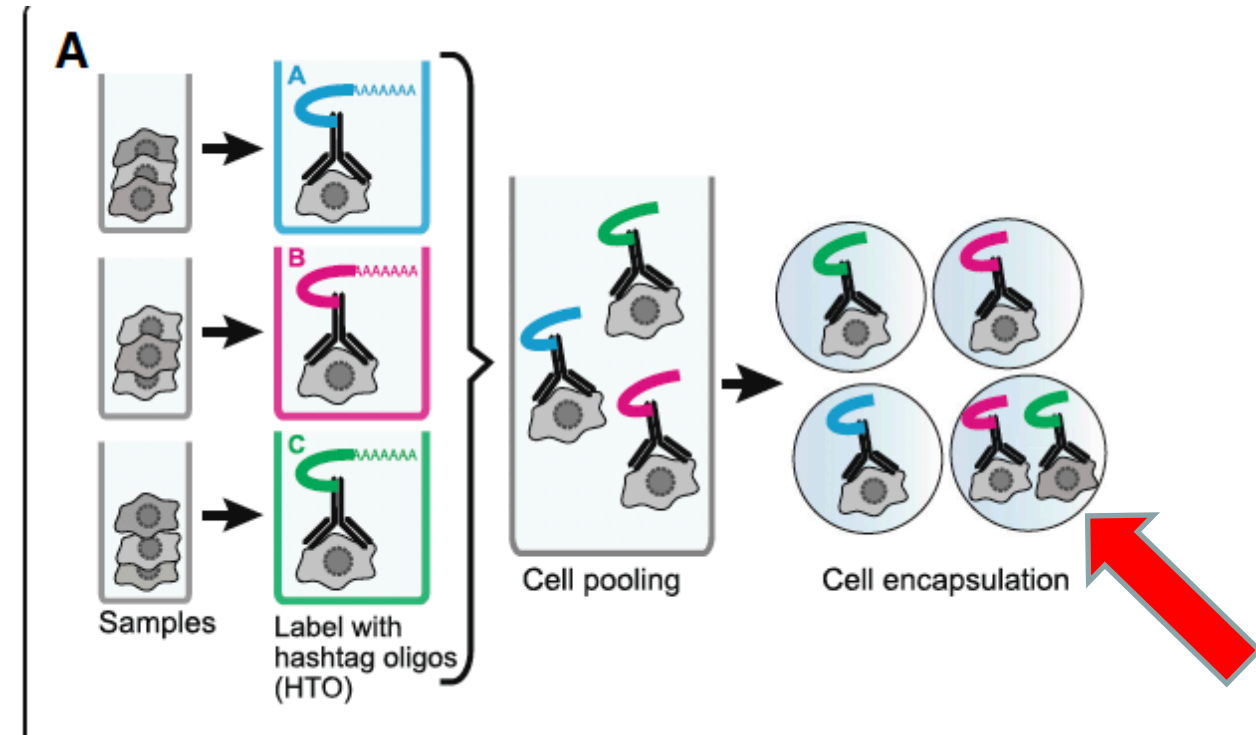
Cell Filtering

- Filtering based on UMI count, gene count, and mitochondrial gene expression
- Mitochondrial gene expression threshold is 4 median absolute deviation above median
- Mitochondrial fraction is linked to cell death, which may influence normalization
- Different cell types have different expression levels



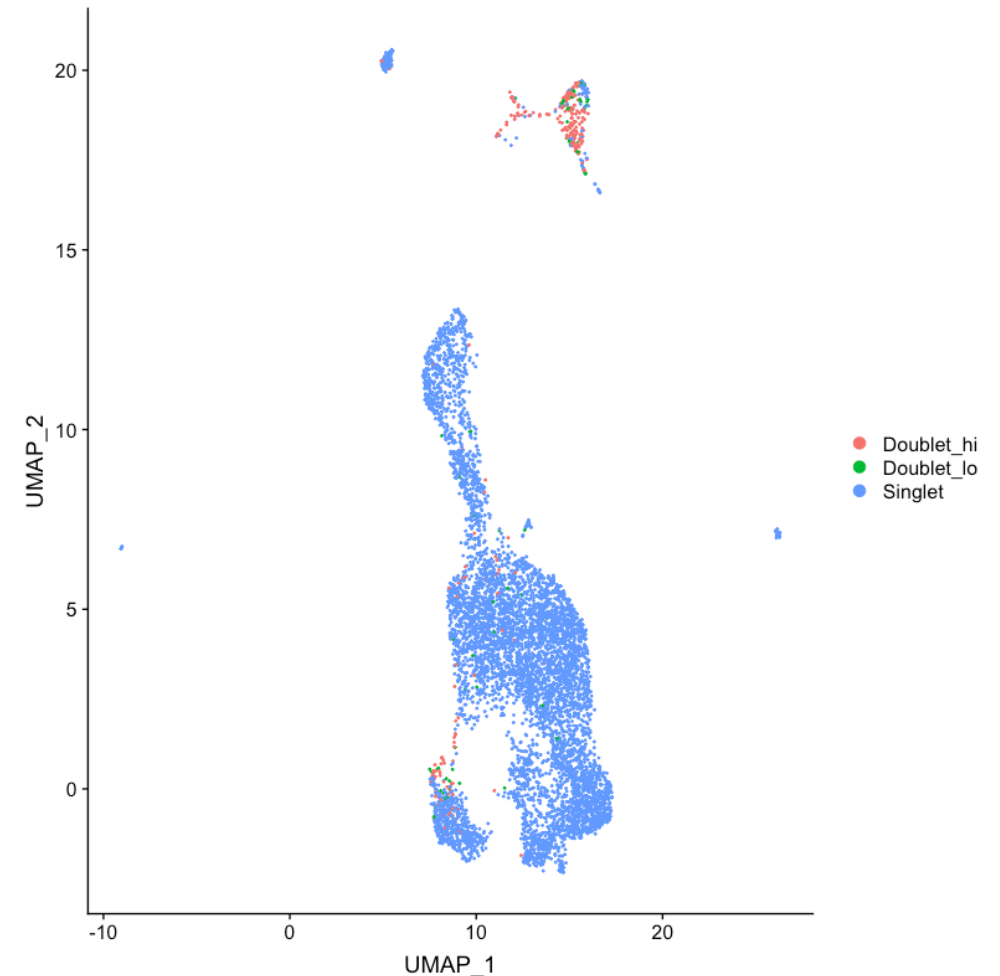
Finding Doublets

- Doublets (or multiplets) are a technical byproduct of single-cell droplet sequencing
- Doublets can interfere with downstream analysis by including high read counts per “cell” and changing cluster identities
- There is no current method to identify transcripts associated with the individual cells in doublets
- Doublets can be homotypic (same cell type) or heterotypic (different cell types)

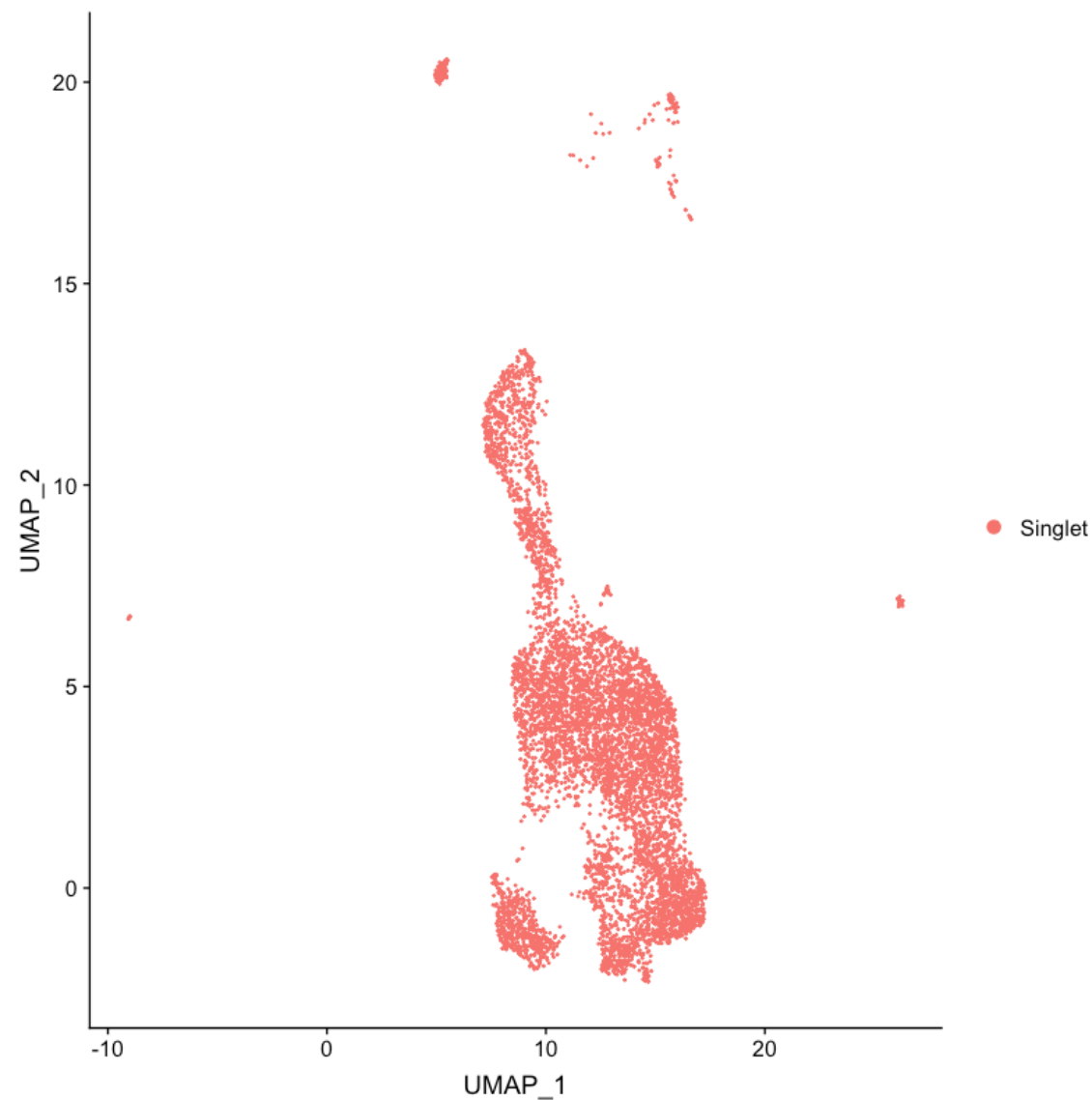
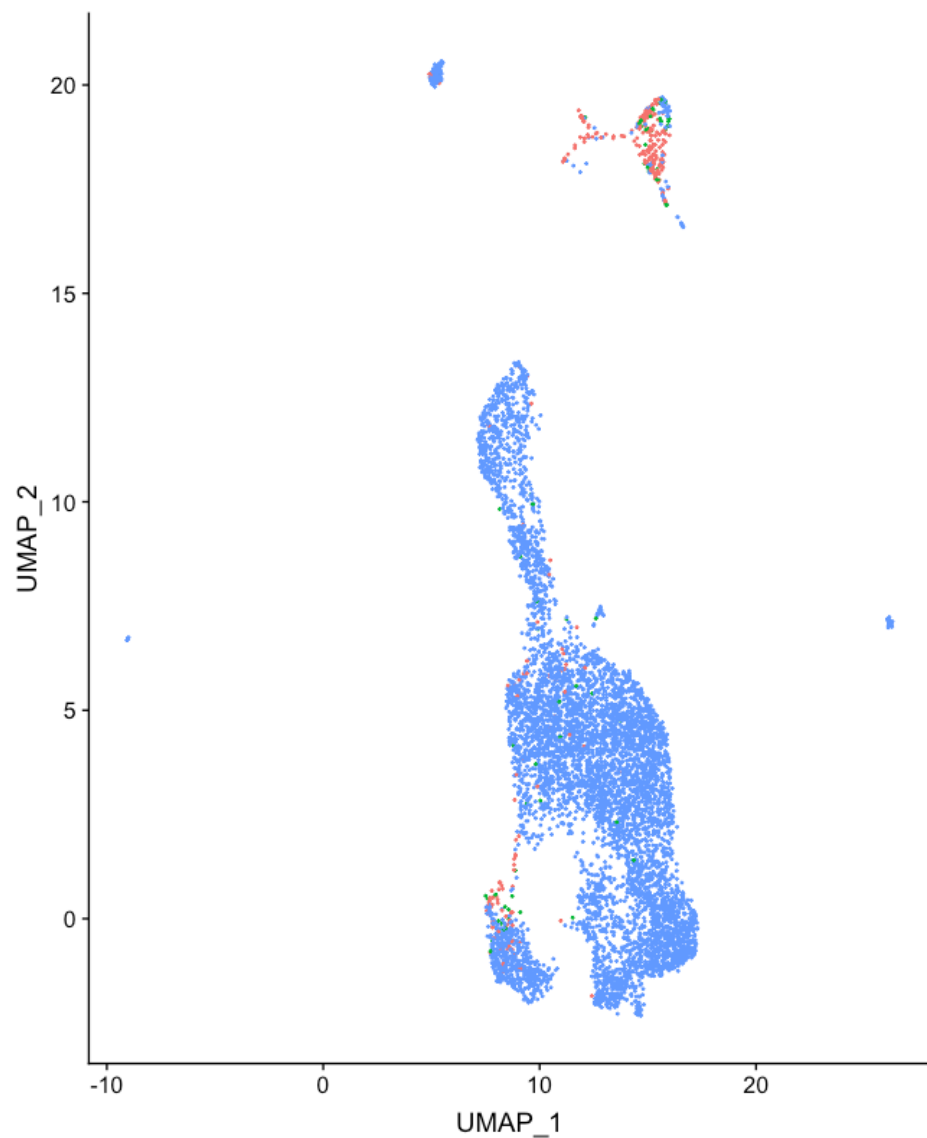


Finding Doublets

- **Statistical removal of doublets:**
 - UMI count and gene count based filters
- **Algorithmic removal of doublets:**
 - DoubletFinder (McGinnis, Murrow and Gartner 2019)
 - Scrublet (Wolock, Lopez, and Klein 2018)
- **The estimated doublet rate as provided by 10x Genomics is:**
 - *Doublet Fraction* = $0.008 \times \frac{n_{Cells}}{1000}$



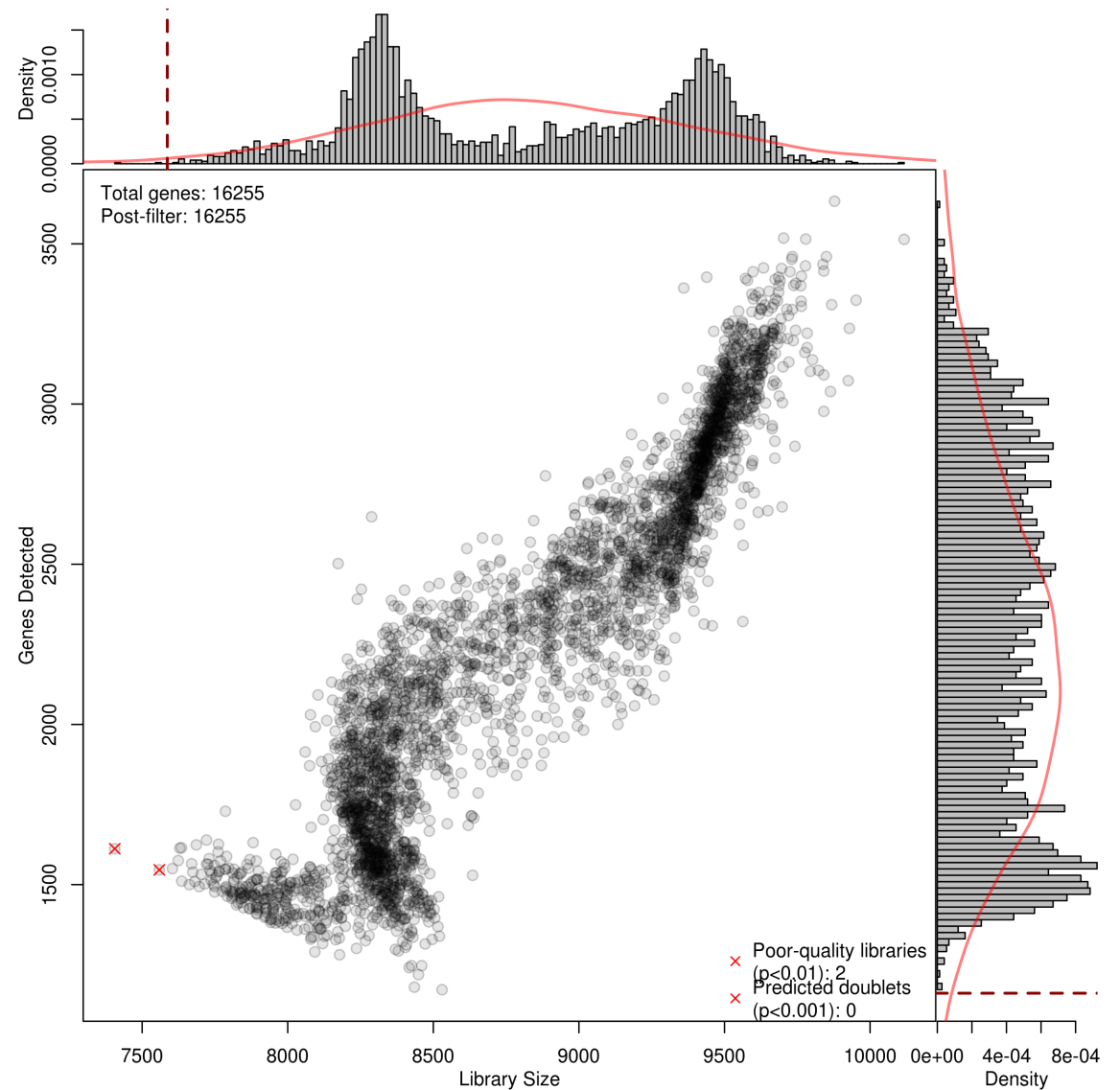
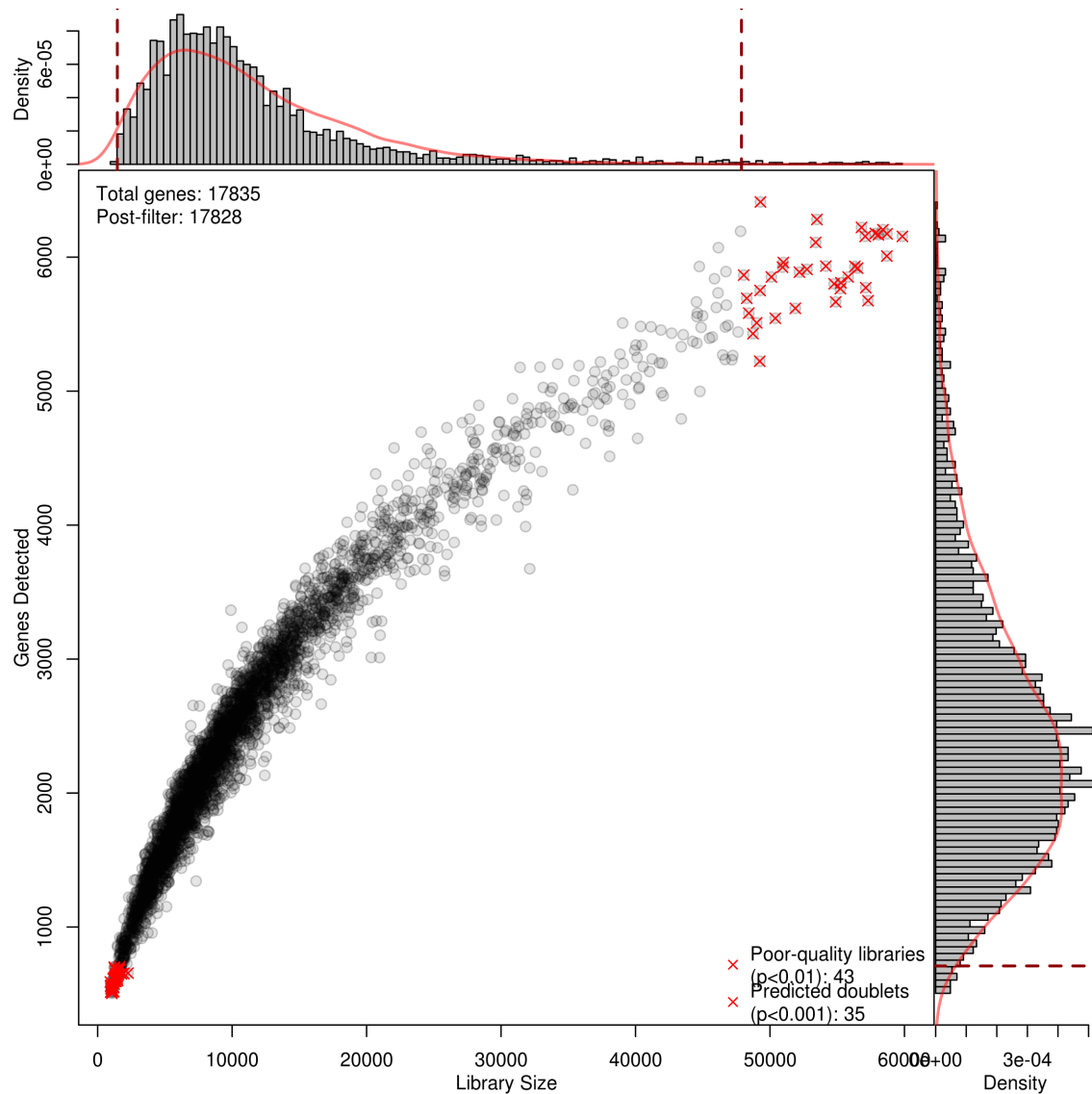
Removal of doublets allows for downstream re-clustering



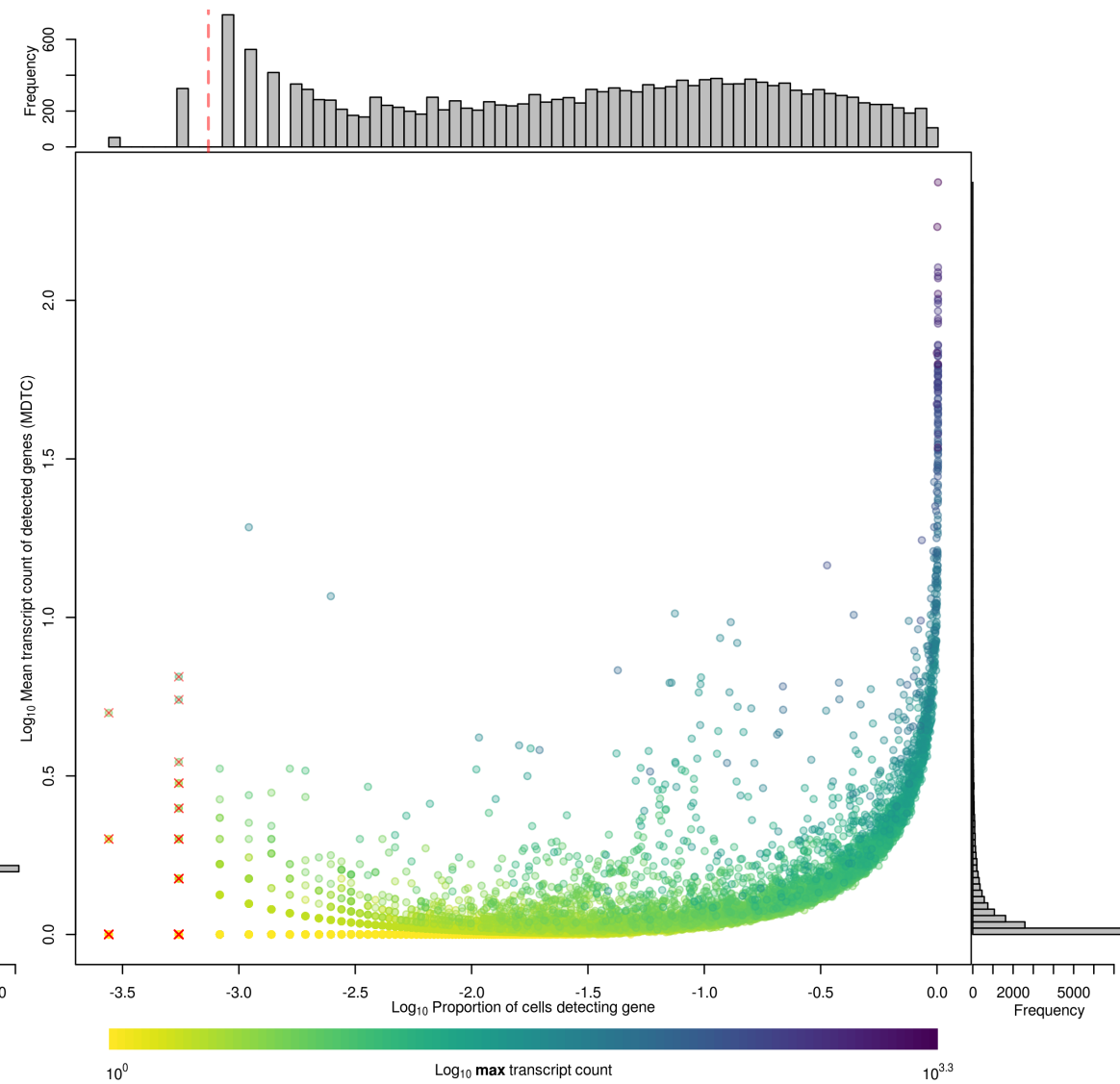
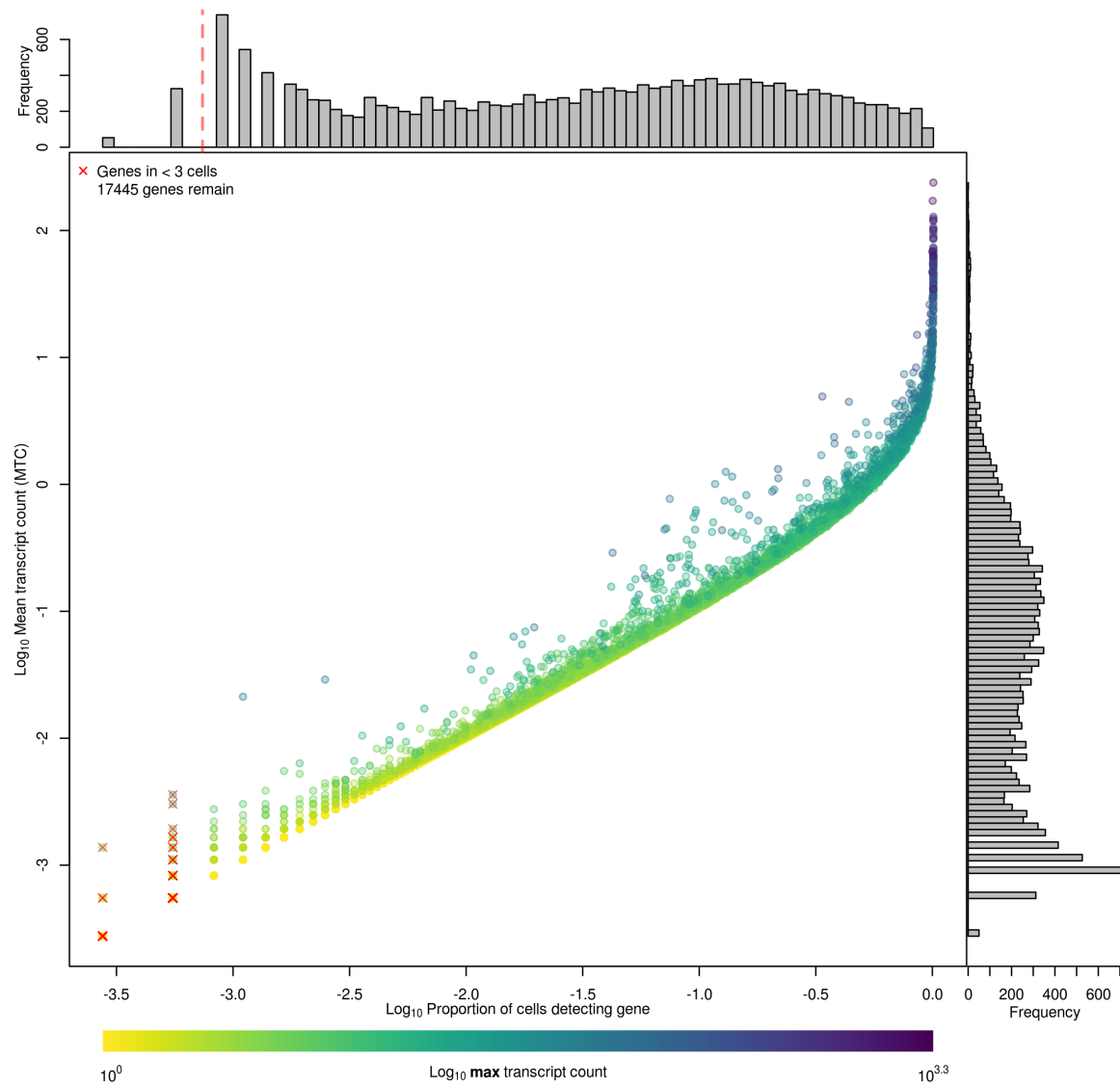
Normalization

- **Aim is to remove technical effects while retaining biological variation**
 - Differences in detected gene expression can be due to sequencing depth of cell
- **Many different normalization techniques available**
- **Seurat has different normalization algorithms available**
 - NormalizeData, ScaleData
 - NormalizeData - Default normalization is log normalize. Each cell divided by total counts, multiplied by scale factor, and natural log transformed
 - ScaleData - Scales and centers features in the data. Can optionally regress out effects of variables (i.e. mitochondrial expression, cell cycle)
 - scTransform - combined NormalizeData, FindVariableFeatures, ScaleData

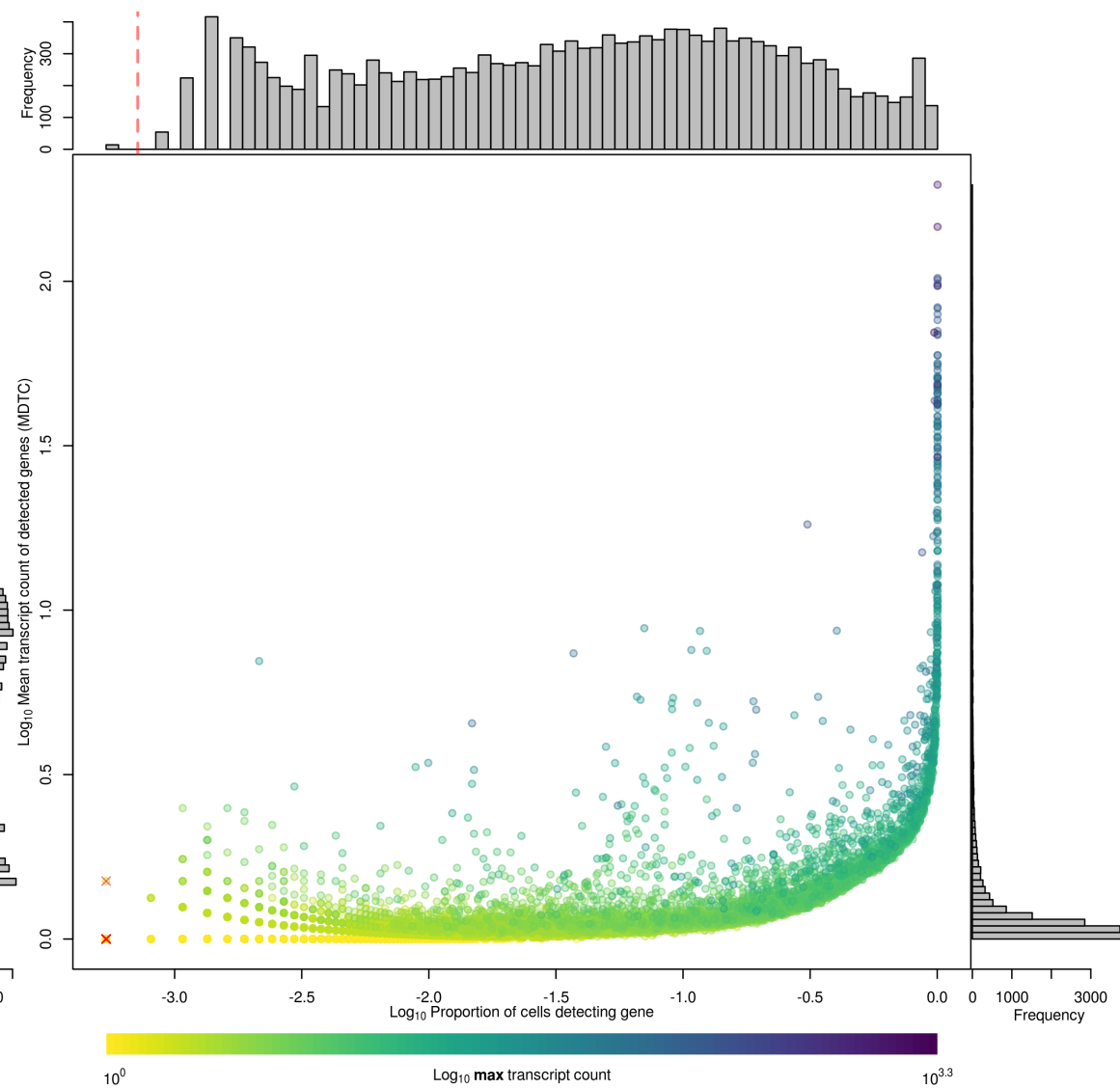
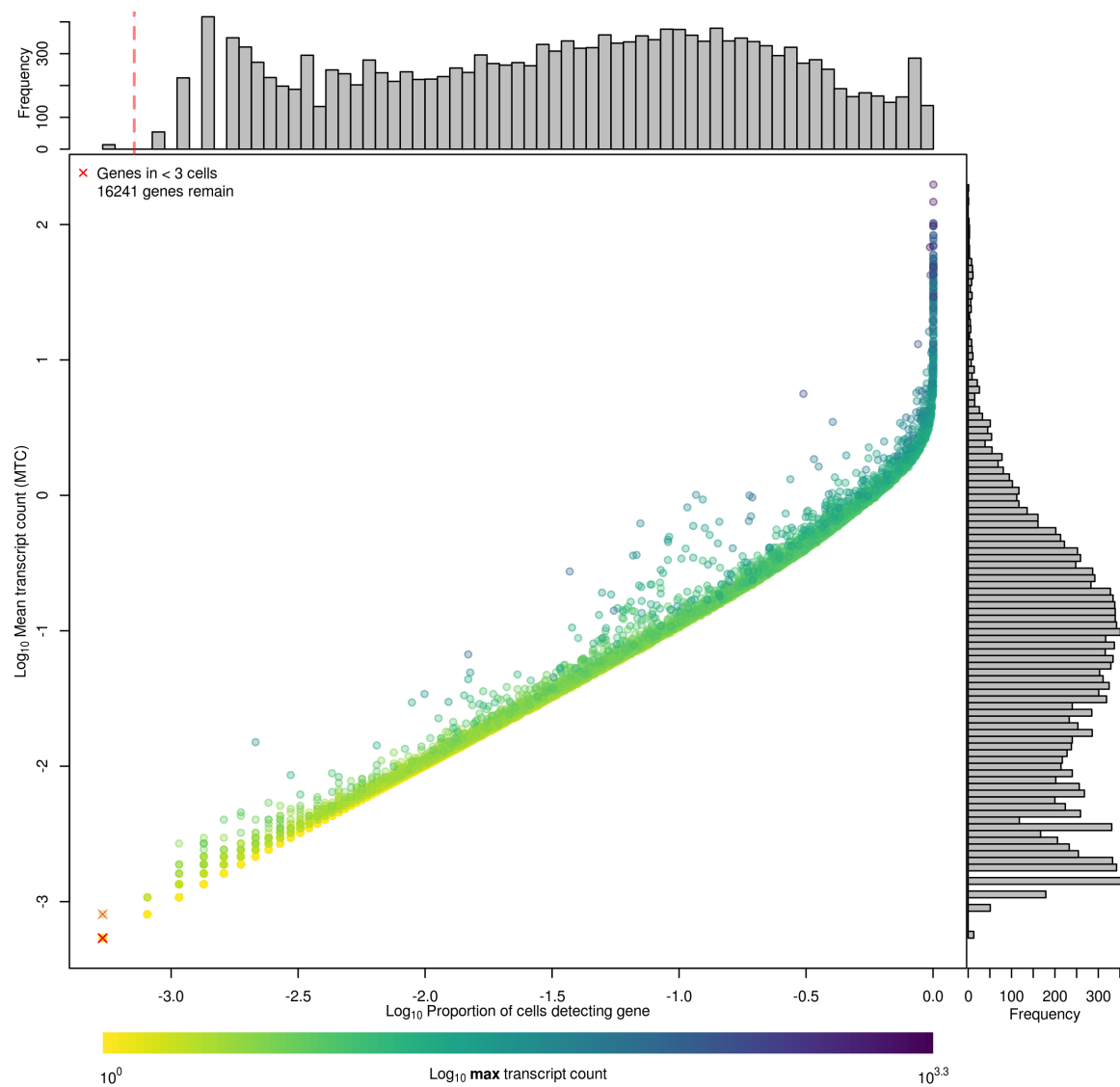
Seurat log Normalize vs scTransform



Expression Plot v2

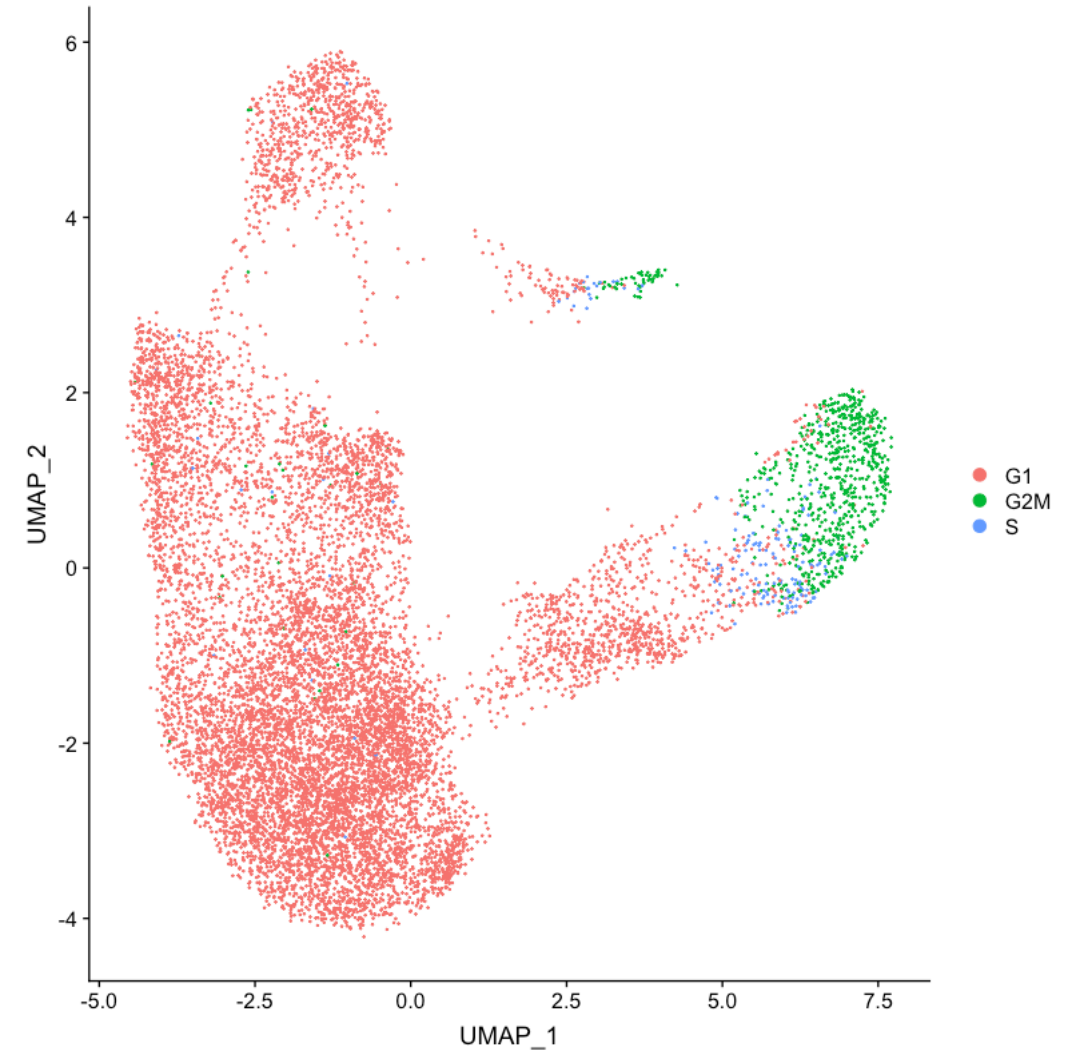


Expression Plot – v3 scTransform



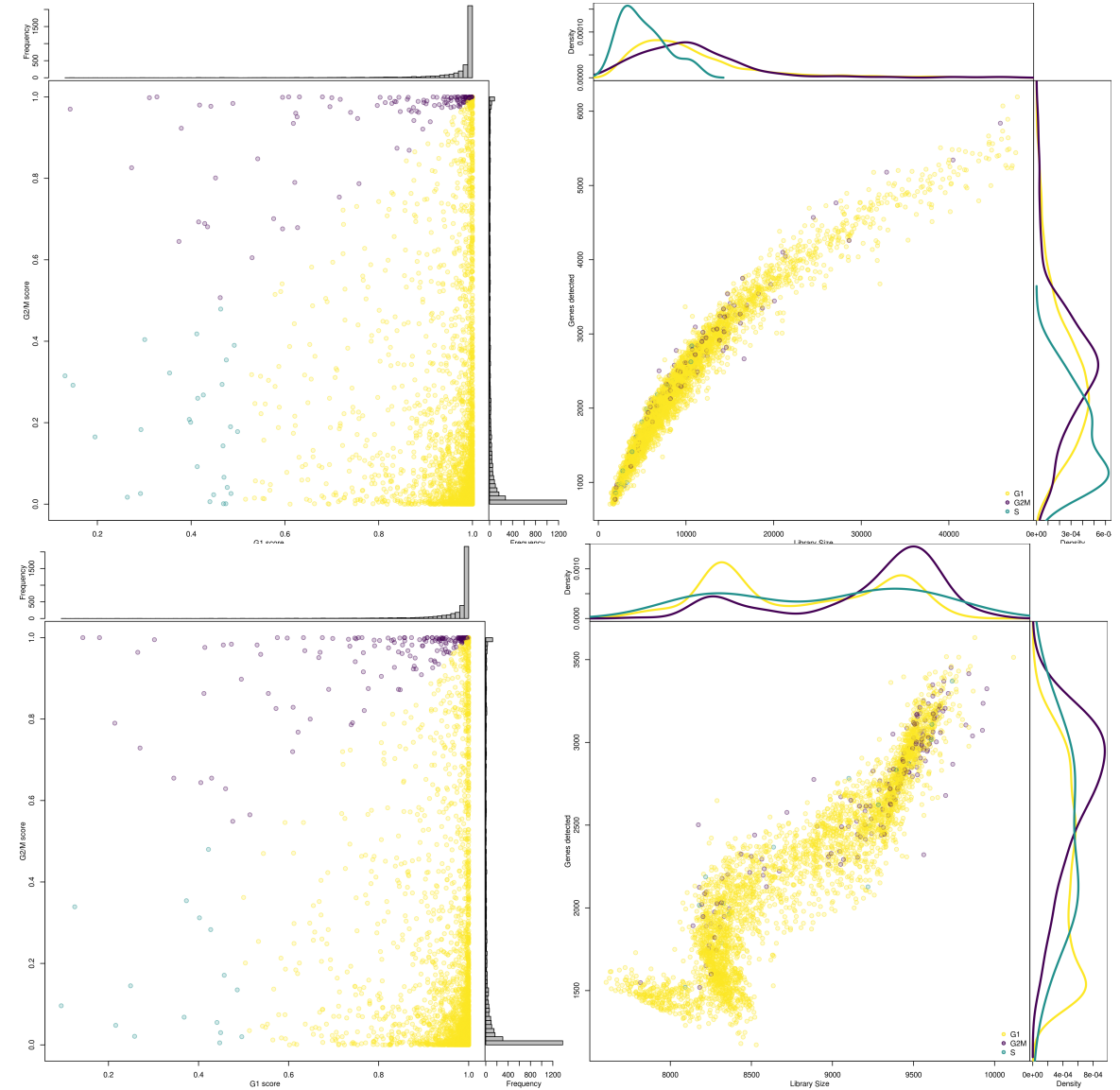
Cell Cycle

- **Cell cycle can introduce bias or obscure differences in expression by cell types**
- **Cell cycle can be identified using available tools, including:**
 - Seurat: CellCycleScoring
 - Scran: Cyclone
- **A variety of tools and techniques are available that can be used to remove effect**
 - ccRemover (Li and Barron 2017)
 - Seurat – ScaleData can be used to regress out effects after labelling cell cycle



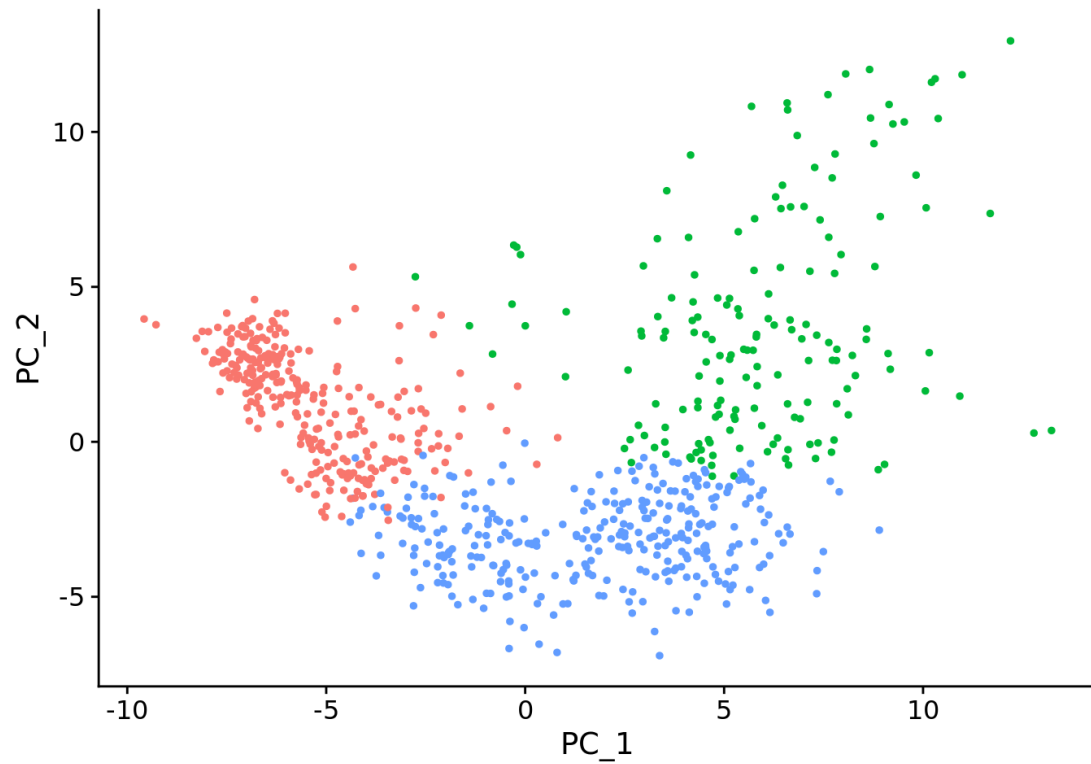
Cell Cycle

- **Cell cycle can introduce bias or obscure differences in expression by cell types**
- **Cell cycle can be identified using available tools, including:**
 - Seurat: CellCycleScoring
 - Scrn: Cyclone
- **A variety of tools and techniques are available that can be used to remove effect**
 - ccRemover (Li and Barron 2017)
 - Seurat – ScaleData can be used to regress out effects after labelling cell cycle

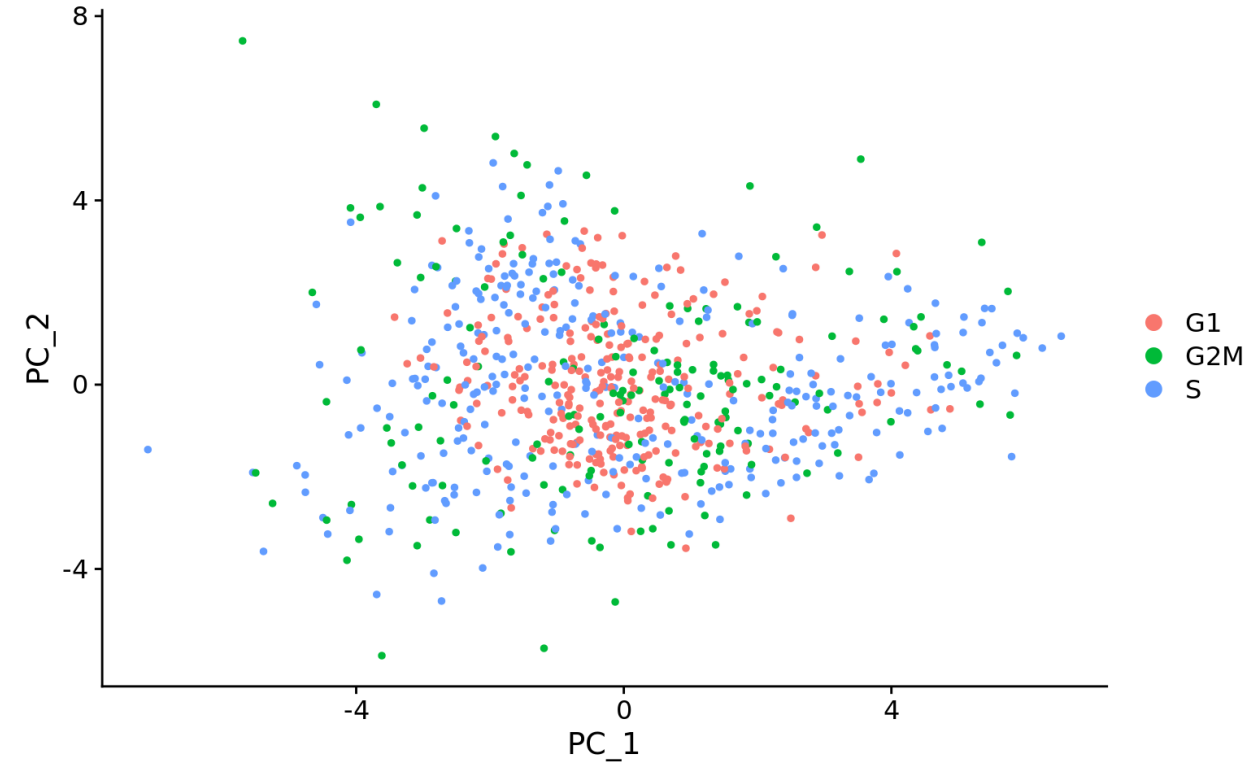


Regressing out cell cycle effects

Prior to Regression



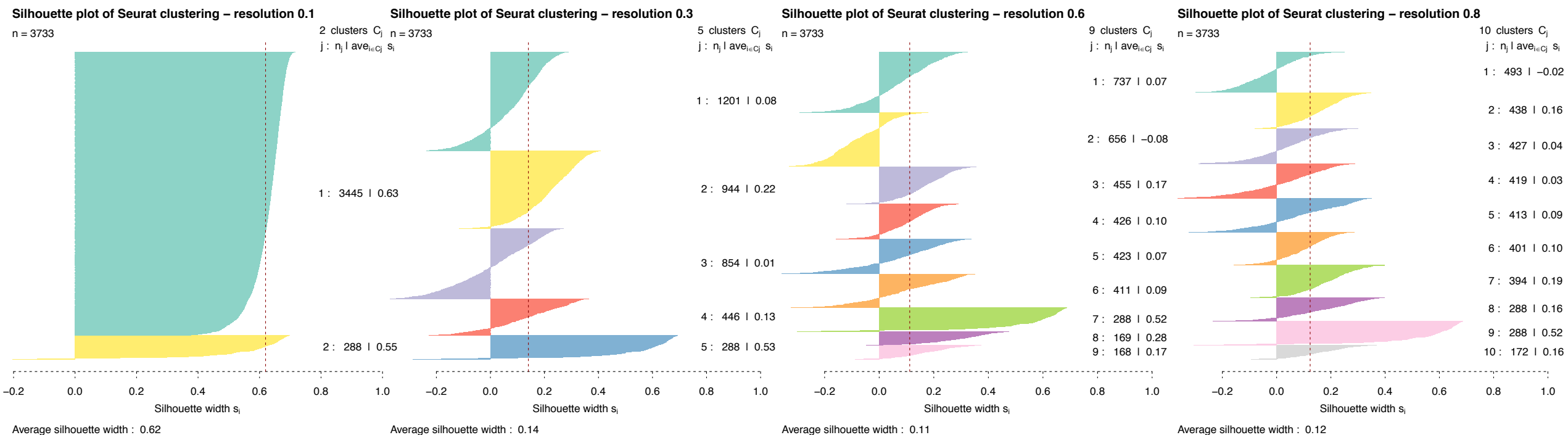
After Regression



Measuring Cluster Quality

- **Different numbers of clusters can be used to group cells within a sample**
- **Can be difficult to determine appropriate number of clusters without prior knowledge**
- **Metrics can be used to measure the quality of the clusters**
 - Silhouette score, Rand index, Davies-Bouldin index
- **Cluster size that results in best score indicates an appropriate number of clusters**


Silhouette Plots – After Seurat Clustering



Imputation

- **Noise and signal dropout are (currently) unavoidable errors in single cell RNA-Seq**
- **Characterized by zero count genes in individual cells**
 - 10x Genomics v3 captures 30-32% of mRNA transcripts per cell
- **Imputation attempts to fill in those zeros based on:**
 - Count distribution
 - Overdispersion
 - Sparsity of the data
 - Noise modeling
 - Gene-gene dependencies

0	2	5.0	3.0	6.0	NaN
1	9	NaN	9.0	0.0	7.0
2	19	17.0	NaN	9.0	NaN
3	7	10.0	3.0	6.0	4.0
4	2	8.0	10.0	NaN	3.0

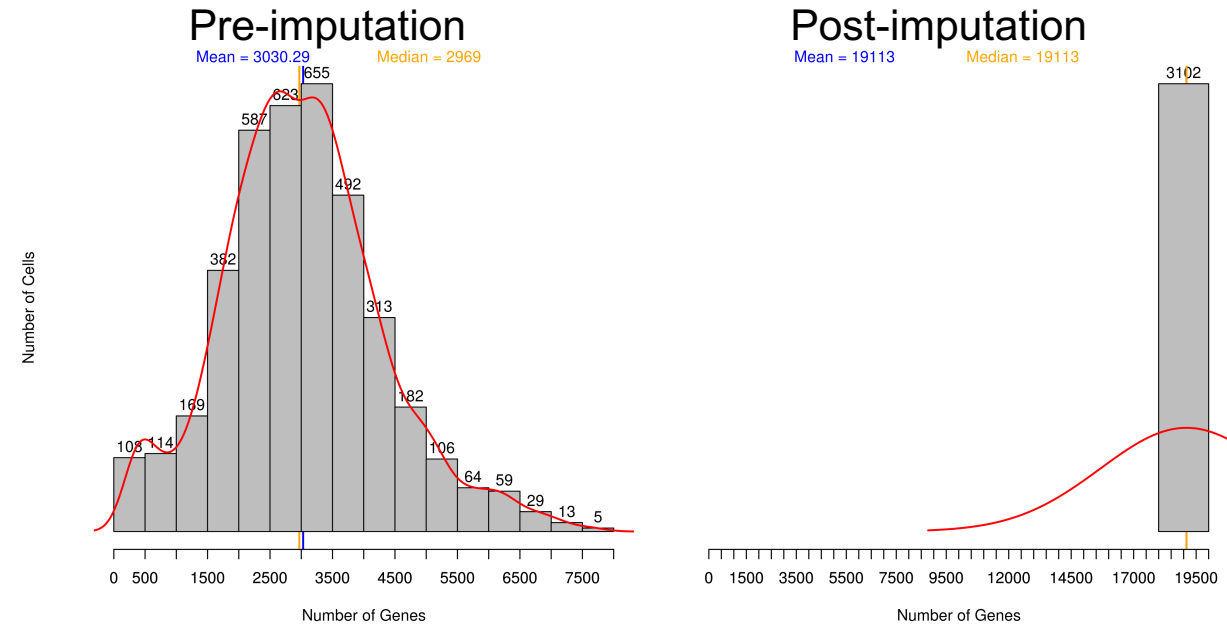


0	2.0	5.0	3.00	6.00	4.666667
1	9.0	10.0	9.00	0.00	7.000000
2	19.0	17.0	6.25	9.00	4.666667
3	7.0	10.0	3.00	6.00	4.000000
4	2.0	8.0	10.00	5.25	3.000000

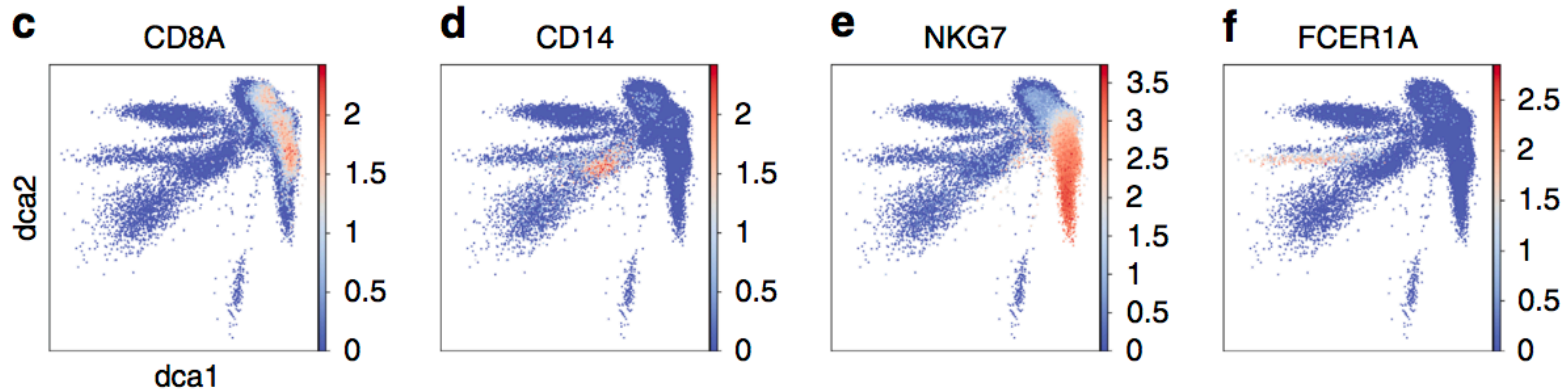
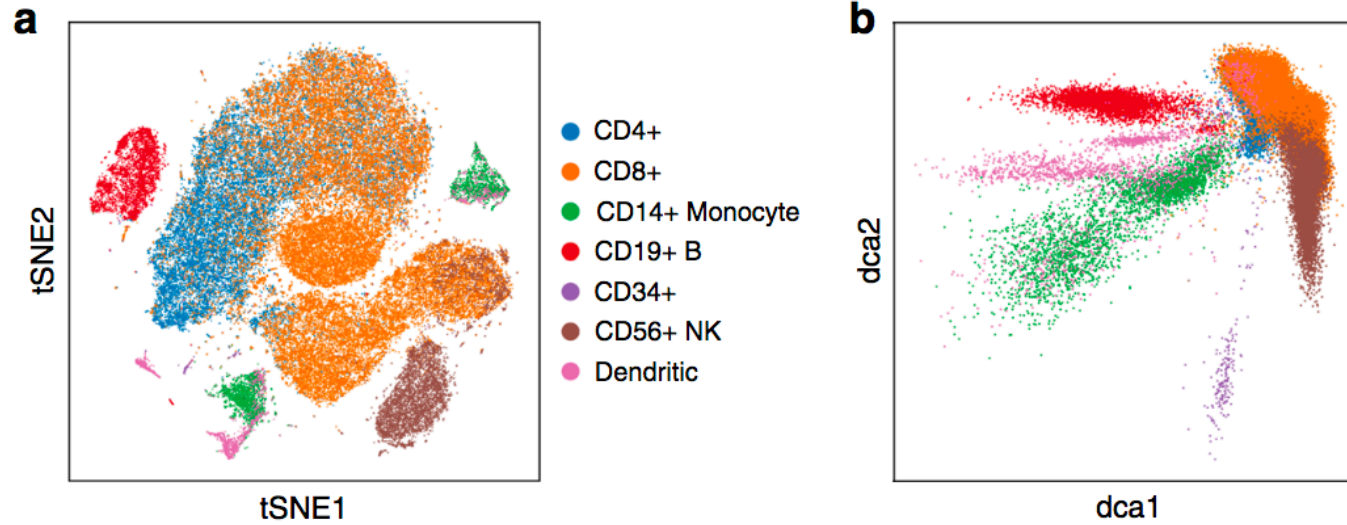
Available imputation tools include:

- **dca (Deep count autoencoder) (Erslan, et al. 2019)**
- **SCRABBLE (Peng, et al. 2019)**
- **SAVER (Huang, et al. 2018)**
- **DrImpute (Gong, et al. 2018)**
- **scImpute (Li and Li 2018)**
- **bayNorm (Tang, et al. 2018)**
- **knn-smooth (Wagner, Yan and Yanai 2018)**
- **MAGIC (van Dijk, et al. 2017)**
- **CIDR (Lin, Troup, and Ho 2017)**

Genes per Barcode (dca)

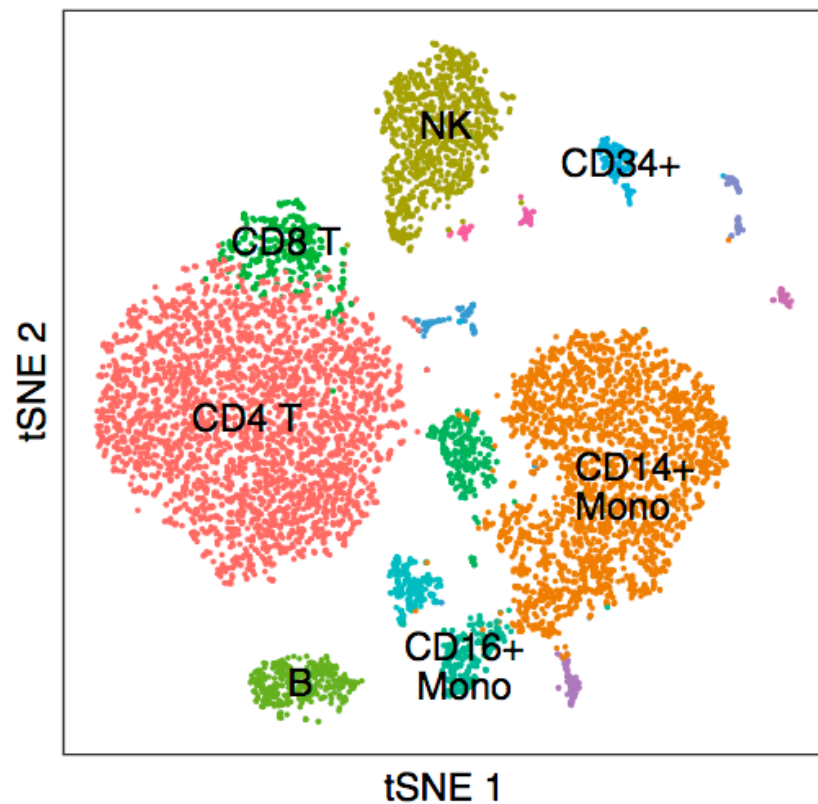


Imputation effects on clusters



Imputation effects on gene expression

a



b

