



Swiss Institute of  
Bioinformatics

SINGLE-CELL TRANSCRIPTOMICS WITH R

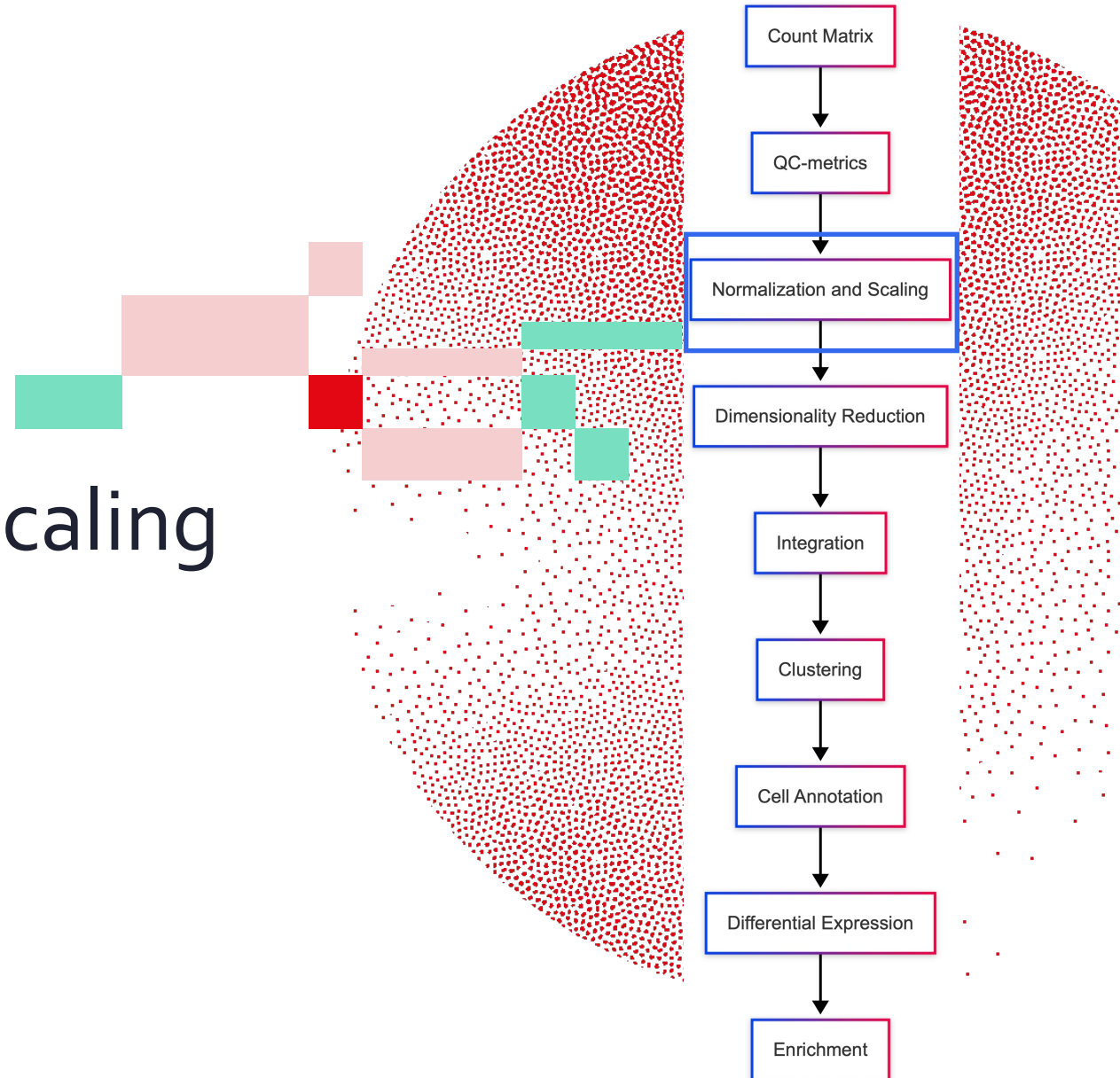
# Normalization and scaling

Deepak Tanwar

March 18-20, 2025

Adapted from previous year courses

Feedback from Geert van Geest



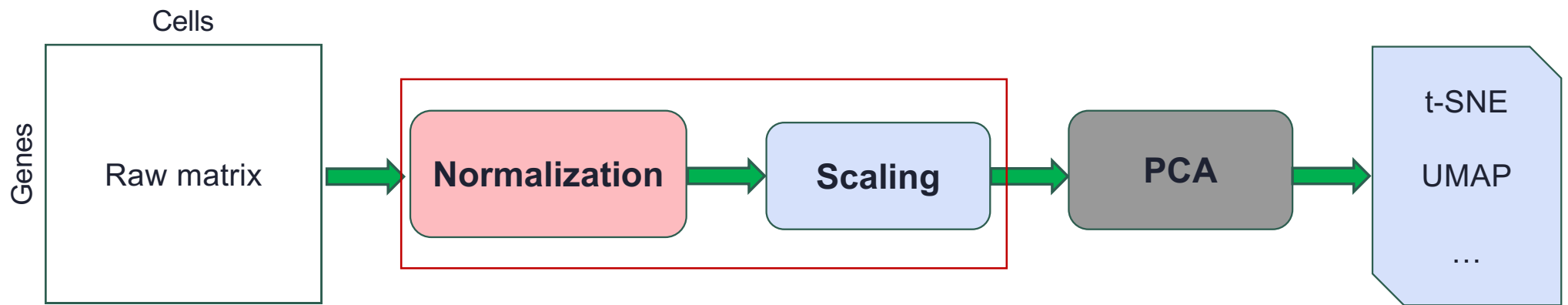
# Learning objectives

Understand the importance of Normalization and Scaling

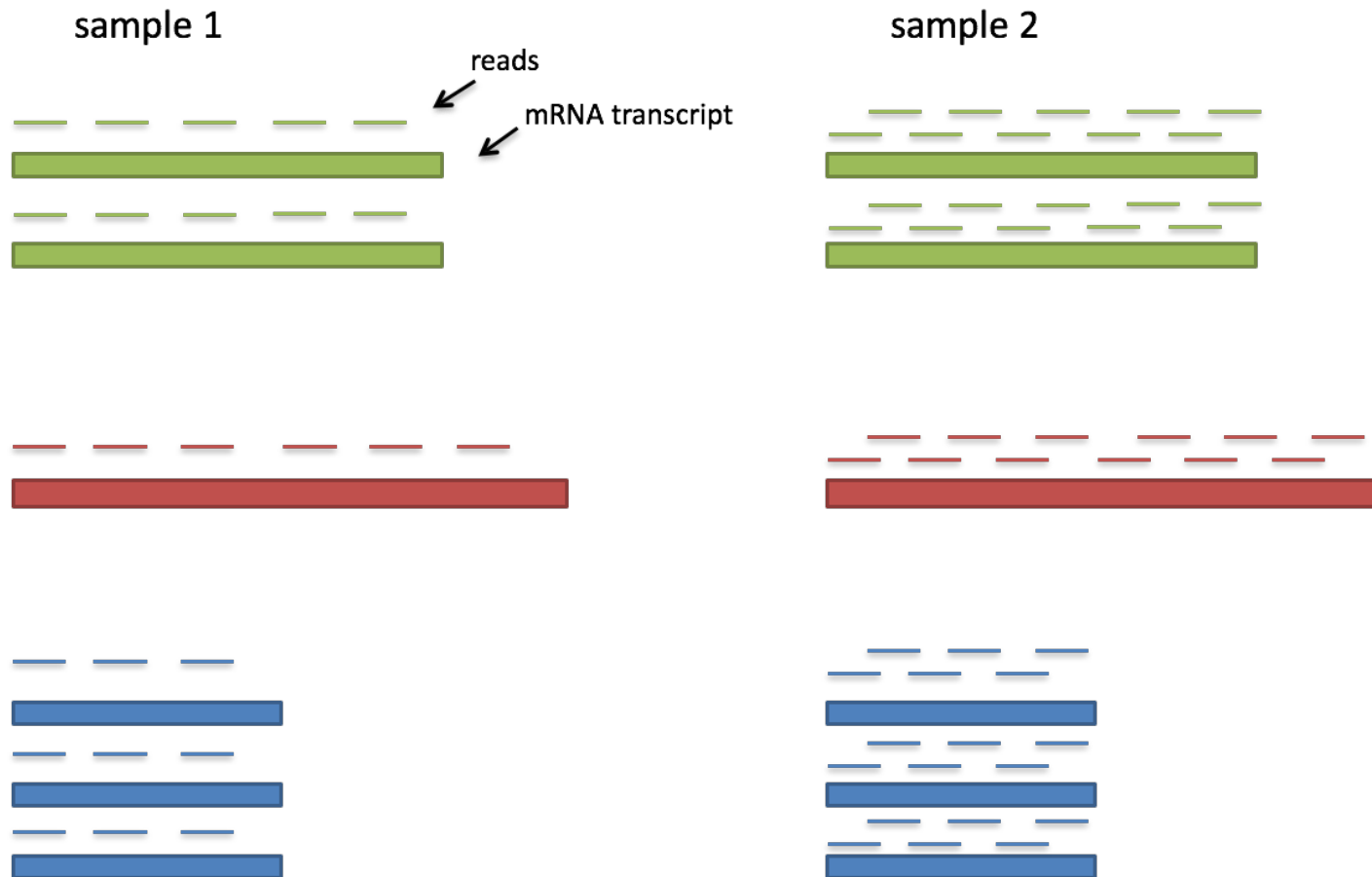
Identify and apply Normalization techniques

Understand Scaling and Transformation

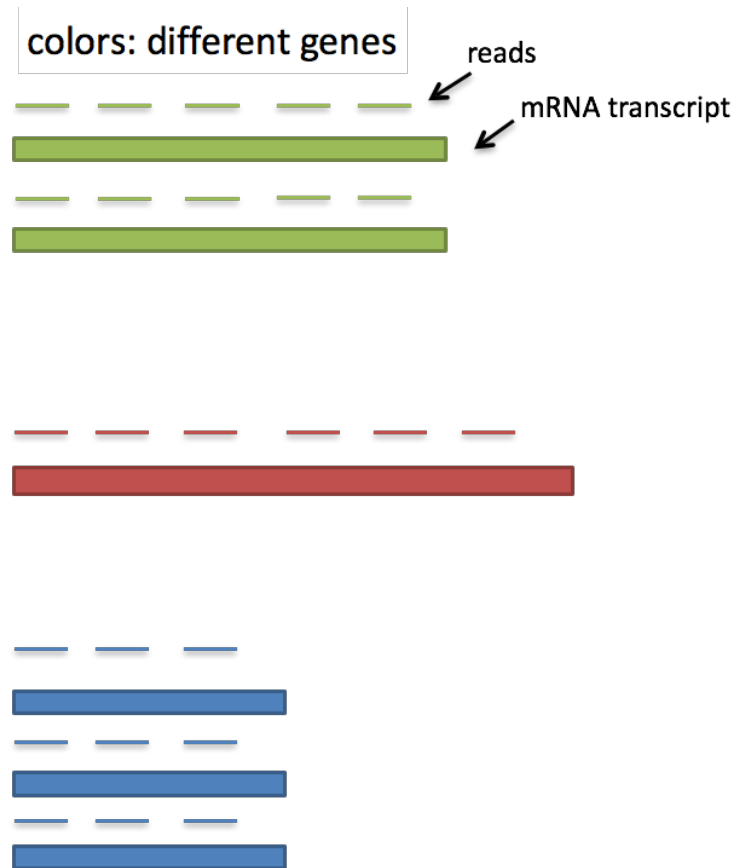
# Normalization and scaling



# Understanding the differences



# Understanding the differences



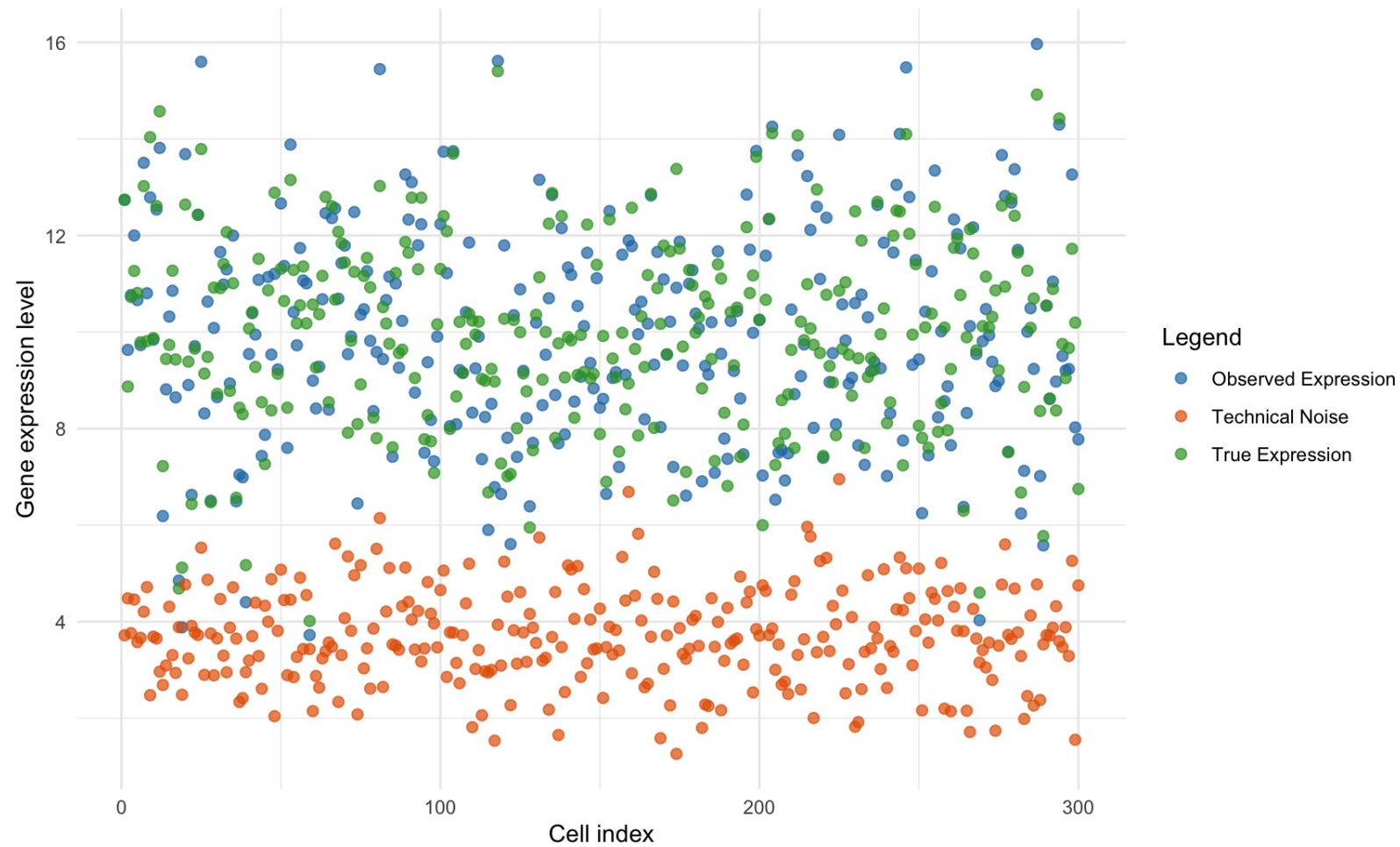
Slide adapted from MI Love: RNA-seq statistical analysis

# Goal of normalization

Remove technical noise while preserving biological signal

- Library size

# Understanding the differences



# Normalization we apply

## UMI (Unique Molecular Identifiers):

- Cells with **extremely high UMI counts** could be **doublets** (two or more cells captured in a single droplet).
- Cells with **very low UMI counts** might be **low-quality or empty droplets**.

## Detected genes:

- A healthy cell will express a moderate number of genes.
- Very **low gene count** could indicate a dead cell or an empty droplet.
- **High gene count** could indicate a doublet.



# Normalization techniques applied in scRNA-seq

## % Mitochondrial UMI:

1. Mitochondrial genes are usually expressed at **low levels**.
2. **High mitochondrial RNA percentage (>10-20%)** indicates **stressed or dying cells**.

## % Ribosomal UMI:

1. High ribosomal content may suggest technical artifacts or certain cell types (e.g., rapidly dividing cells).

## % Globin UMI:

1. In blood samples, **high globin content** comes from red blood cells (RBCs).
2. Filtering out these cells is often necessary when focusing on immune or other cell types.

High UMI count + high gene count → **Doublet suspicion**  
High mitochondrial percentage → **Apoptotic or stressed cell**  
Low gene count + low UMI → **Low-quality or empty droplet**  
High ribosomal percentage → **Potential technical artifact**

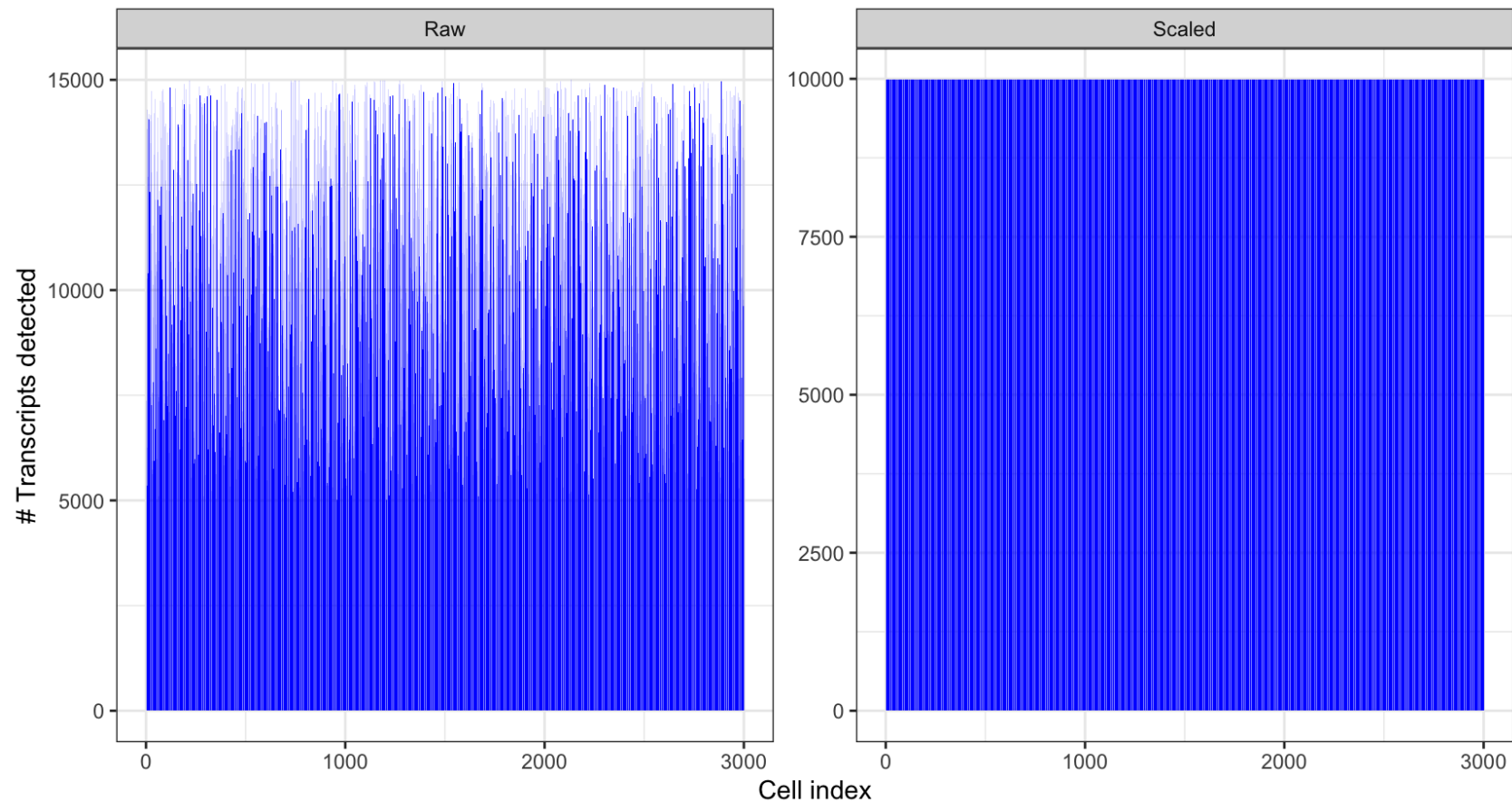
```
seurat_obj <- subset(seurat_obj,  
                     subset = nFeature_RNA > 200 &  
                               nFeature_RNA < 5000 &  
                               nCount_RNA < 20000 &  
                               percent.mt < 10 &  
                               percent.ribo < 40 &  
                               percent.globin < 5)
```

# Scaling

Multiply each UMI count by a cell specific factor to get all cells to have the same UMI counts

*Different cells have different amounts of mRNA; this could be due to differences between cell types or variation within the same cell type depending on how well the chemistry worked in one drop versus another.*

# Scaling: standardize range, mean and variance



# Transformation

- Simple transformations
- Pearson residuals

# Transformation : Simple transformations

Raw data			
	Cell Type A	Cell Type B	$\Delta$
Gene 1	1	2	1
Gene 2	100	200	100
Gene 3	5	25	20
Gene 4	400	800	400
Gene 5	10	60	50

# Transformation : Simple transformations

	Raw data			Log <sub>2</sub> transform		
	Cell Type A	Cell Type B	$\Delta$	Cell Type A	Cell Type B	$\Delta$
Gene 1	1	2	1	0.00	1.00	1.00
Gene 2	100	200	100	6.64	7.64	1.00
Gene 3	5	25	20	2.32	4.64	2.32
Gene 4	400	800	400	8.64	9.64	1.00
Gene 5	10	60	50	3.32	5.91	2.58

# Transformation : Simple transformations

	Raw data			Log <sub>2</sub> transform			Square root transform		
	Cell Type A	Cell Type B	$\Delta$	Cell Type A	Cell Type B	$\Delta$	Cell Type A	Cell Type B	$\Delta$
Gene 1	1	2	1	0.00	1.00	1.00	1.00	1.41	0.41
Gene 2	100	200	100	6.64	7.64	1.00	10.00	14.14	4.14
Gene 3	5	25	20	2.32	4.64	2.32	2.24	5.00	2.76
Gene 4	400	800	400	8.64	9.64	1.00	20.00	28.28	8.28
Gene 5	10	60	50	3.32	5.91	2.58	3.16	7.75	4.58



# Transformation : Simple transformations

- Log transformation
- Square root transformation

$$y_{ij} = f(x_{ij})$$

i: cell

j: gene

# Quiz

Which of the simple transformation methods transform each measurements individually?

1. Log
2. Square root
3. None
4. Both

Transformation : Pearson residuals

$$y_{ij} = w_j * x_{ij}$$

w: weight

How do we choose weights?  
1/sqrt(mean of the gene)

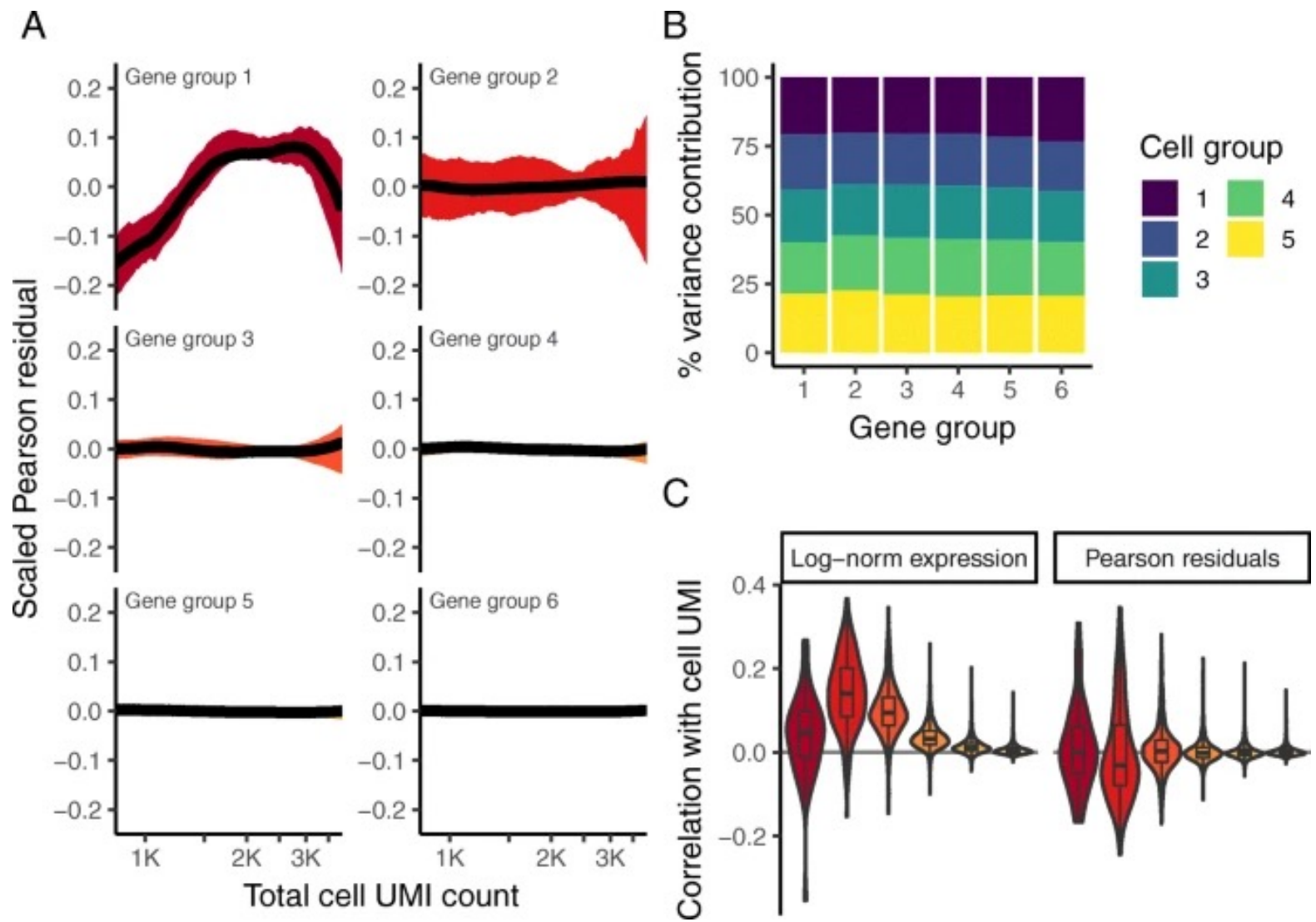
# Transformation : Pearson residuals

Instead of transforming each measurements individually, Pearson residuals apply a weight to all measurements of a gene

# Transformation : Pearson residuals

	Raw data			Log transform			Square root transform			Pearson Residuals		
	Cell A: 75%, Cell B: 25%			Cell A: 75%, Cell B: 25%			Cell A: 75%, Cell B: 25%			Cell A: 75%, Cell B: 25%		
	Cell Type A	Cell Type B	$\Delta$	Cell Type A	Cell Type B	$\Delta$	Cell Type A	Cell Type B	$\Delta$	Cell Type A	Cell Type B	$\Delta$
Gene 1	1.00	2.00	1.00	0.00	1.00	1.00	1.00	1.41	0.41	-0.83	1.44	2.28
Gene 2	100.00	200.00	100.00	6.64	7.64	1.00	10.00	14.14	4.14	-8.33	14.43	22.77
Gene 3	5.00	25.00	20.00	2.32	4.64	2.32	2.24	5.00	2.76	-3.69	6.39	10.08
Gene 4	400.00	800.00	400.00	8.64	9.64	1.00	20.00	28.28	8.28	-16.67	28.87	45.53
Gene 5	10.00	60.00	50.00	3.32	5.91	2.58	3.16	7.75	4.58	-5.87	10.16	16.03

# sctransform



# Summary

**Normalization:** Adjust UMI counts, mitochondrial, ribosomal, globin RNA percentages

**Goal:** Remove technical noise, preserve biological signals

**Scaling:** Standardize range, mean, variance

**Transformations:** Log, square root, Pearson residuals

**Outcome:** Reliable, meaningful scRNA-seq data analysis

