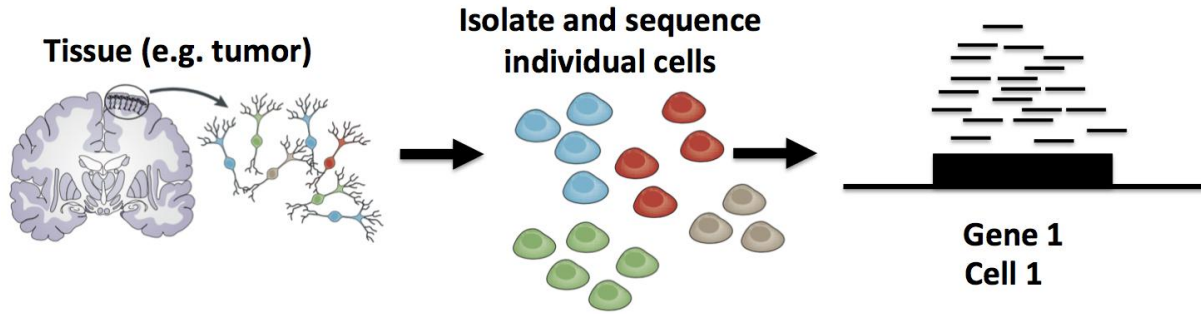# CIDER: an interpretable meta-clustering framework for single-cell RNA-seq data integration and evaluation
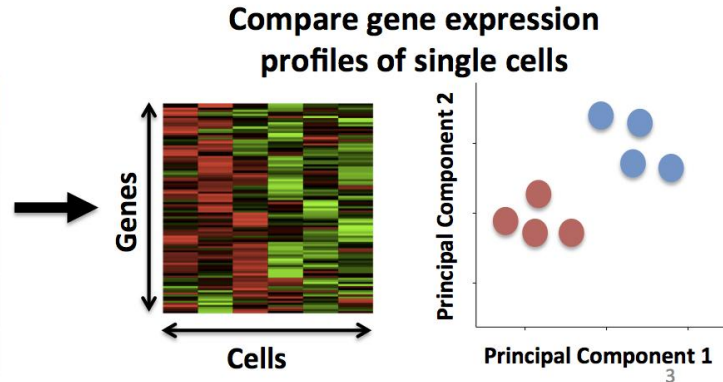
Hu, Ahmed & Yau, 2021
CSE590C  -  2/7/21 (Ayse & Nicasia)
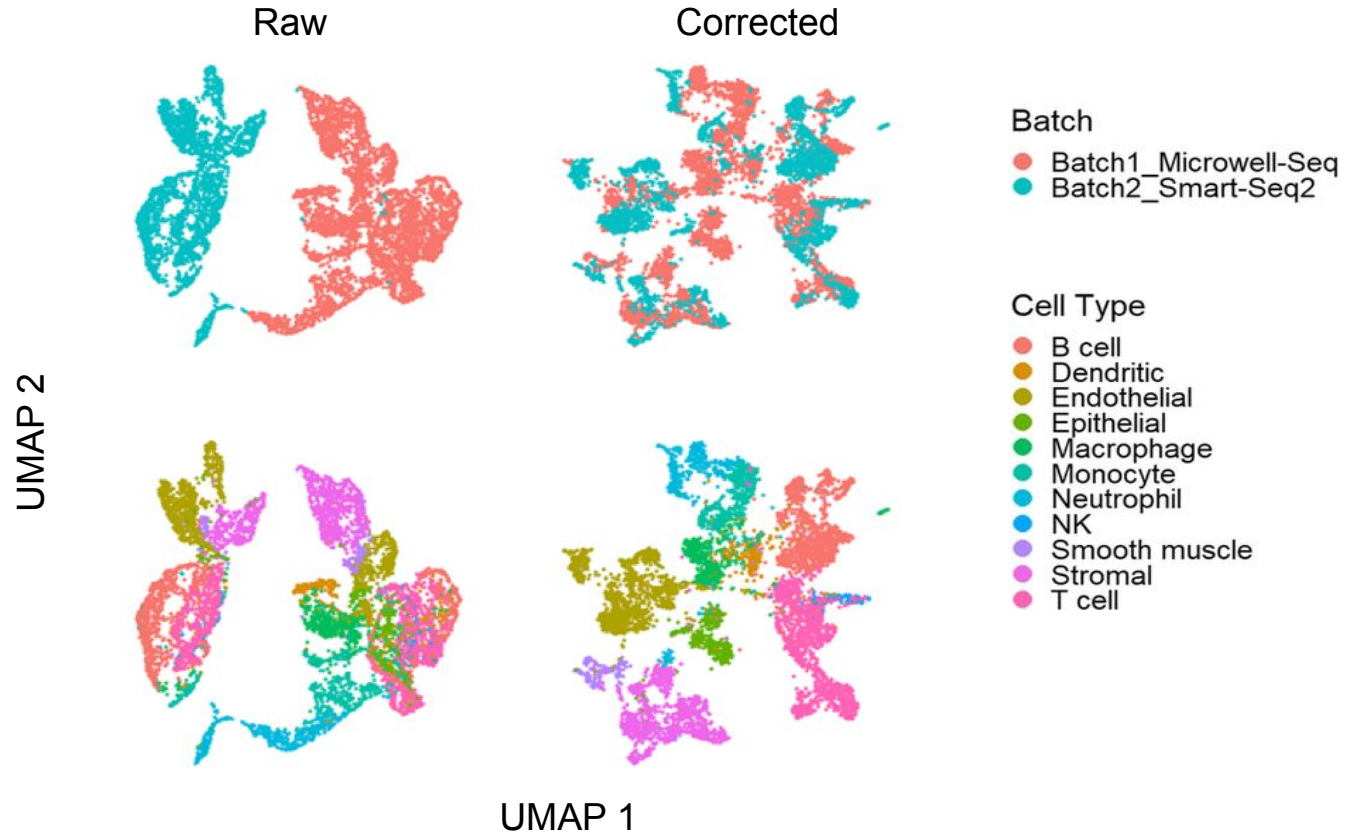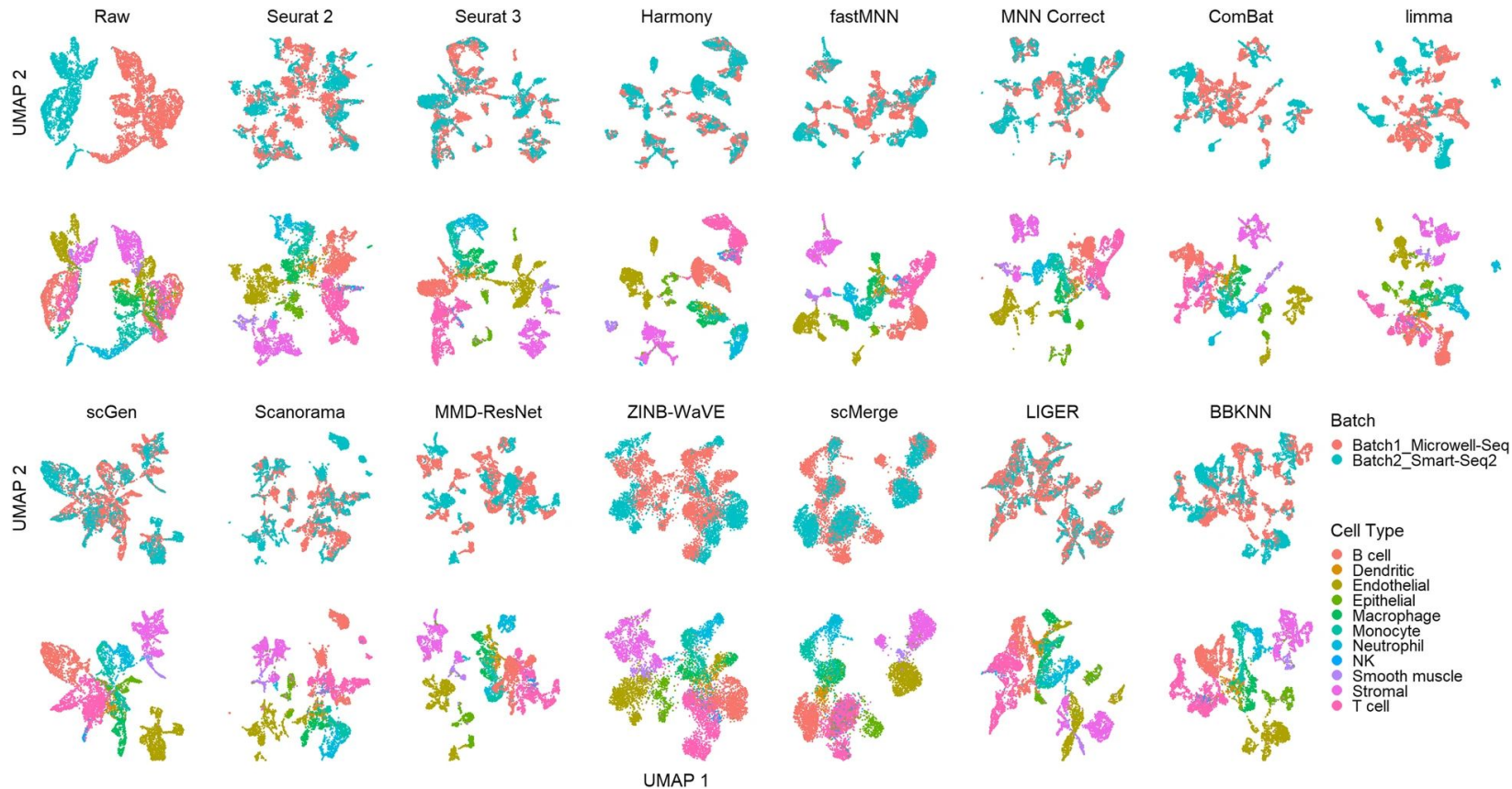
# Single cell RNA sequencing (scRNASeq)



https://www.rna-seqblog.com/top-benefits-of-using-the-technique-of-single-cell-rna-seq/

# scRNASeq - challenges with data integration



Tran, H.T.N., Ang, K.S., Chevrier, M. *et al.* A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol* 21, 12 (2020).

# scRNASeq – current approaches



Tran, H.T.N., Ang, K.S., Chevrier, M. *et al.* A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol* 21, 12 (2020).

# scRNASeq – current approaches

## Clustering Workflows

Can identify cell populations in batch-effect- free datasets

Partition cells by inter-cell distance matrix using PCA or high variance genes (HGVs)

Examples: SC3, RaceID, Seurat v3

**Performance degrades in datasets confounded by batch effects**

## Batch correction + clustering Workflows

Combines batch correction or integration methods and downstream clustering algorithms

Mutual nearest neighbors:
Examples: Monocle3 pipeline, Scanorama, Seurat

Other approaches: Harmony, LIGER, ComBat, Conos

**Performance can vary substantially across data types and scenarios**

# scRNASeq – current approaches

## Limitations

**Bias in initial selection:**

- Integration algorithms work on the low-dimensional representation
- Can be affected by the bias in the initial selection of HVGs and PCs

**Lack of interpretability:**

- Difficult to determine why existing methods drive cells from different batches into the same cluster

**To address these limitations, they introduced CIDER**

# CIDER contributions

1. New similarity metric: Inter-group Differential ExpRession (IDER) → clustering (CIDER)

2. Similar/superior performance compared with other clustering methods for scRNA-Seq data

3. CIDER as a ground-truth-free evaluation metric for other integration methods

# Inter-group Differential ExpRession (IDER) metric

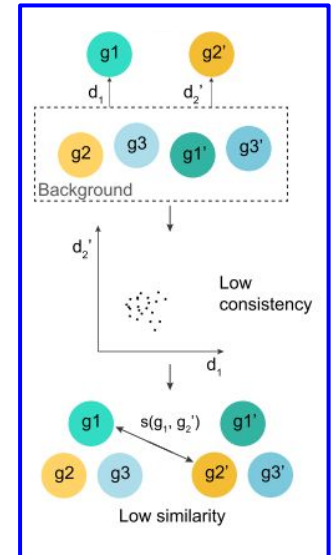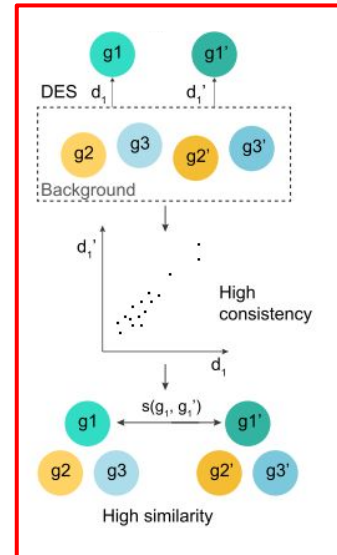Measures similarity between two groups of cells across datasets

IDER for g1 and g1':

1.  Separately, identify differentially expressed genes (DEGs) for g1 and g1' each vs all other groups (limma-trend; can regress out confounders) → d1 and d1' vectors (log2 fold change coeffs for each gene vs background)

2.  IDER(g1,g1') = Pearson r(d1,d1') similarity of the DEG vectors for g1 and g1'



IDER matrix:

|    | g1'  | g2'  | g3'  |
|----|------|------|------|
| g1 | high | low  | low  |
| g2 | low  | high | low  |
| g3 | low  | low  | high |

# Clustering with IDER (CIDER)



**Assumption:** expression level is a linear combination of effects of:

- cluster (of interest)
- batch, donor, platform, etc. (confounders)

**CIDER algorithm:**

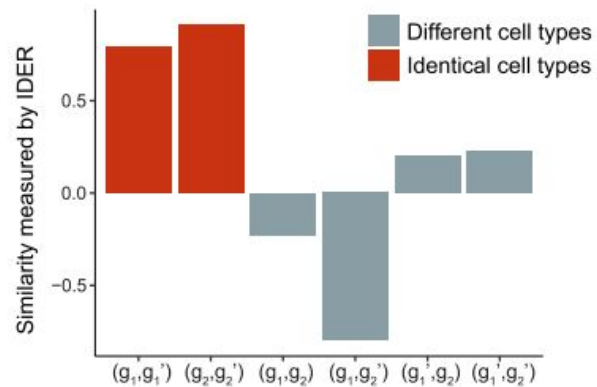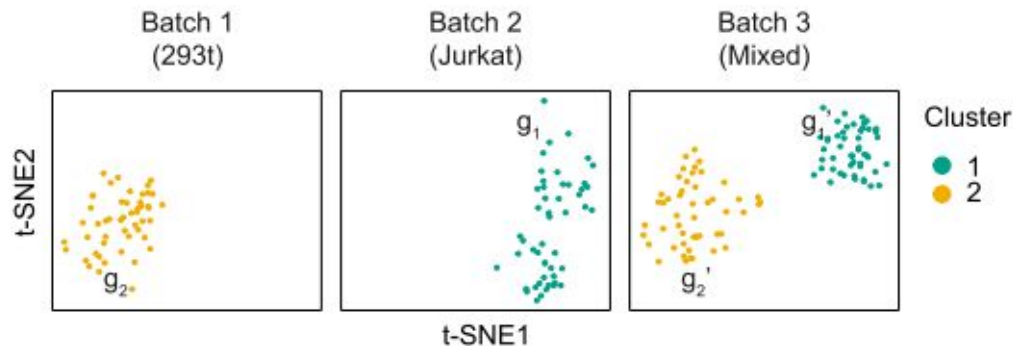1. Within dataset clustering → cluster effect only (confounding effects are constant)
   a. Unsupervised clustering algorithm (e.g., Louvain clustering) → (de novo) dnCIDER
   b. Curated annotations → (assisted) asCIDER

2. Compute IDER similarity matrix across all within-batch clusters to get cross-batch similarity → cluster similar groups across batches
   a. Similarity matrix S → distance matrix (1-S)
   b. Agglomerative clustering with complete linkage
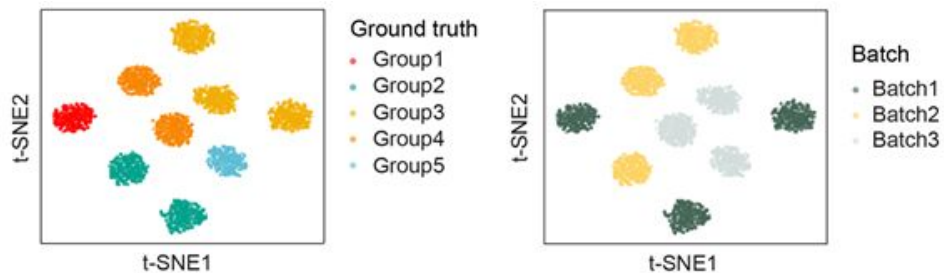3. (optional:) Use limma to regress out confounding effects

# Simple example

**Dataset1:** Batch correction benchmarking dataset (Zheng et al 2017)

1. Only 293T cells
2. Only Jurkat cells
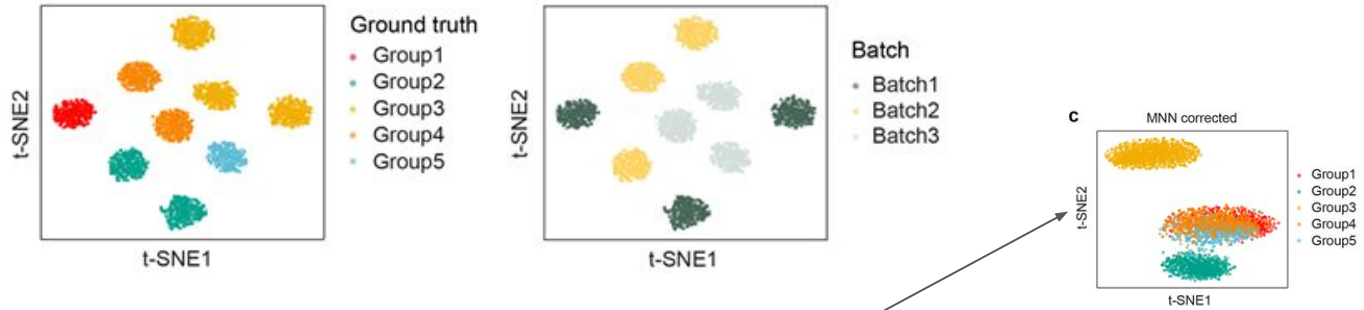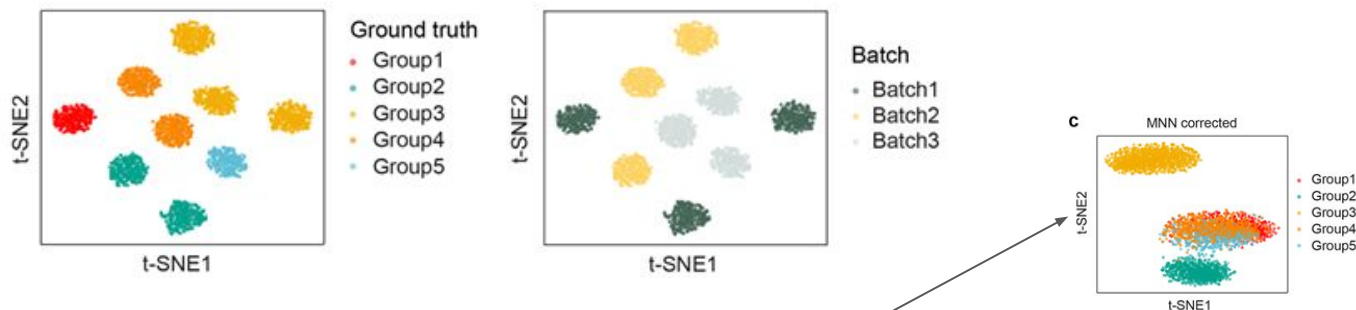3. 1:1 mixture of 293T & Jurkat cells

# Benchmarking with simulated data

- 5 groups across 3 batches with non-identical populations

# Benchmarking with simulated data

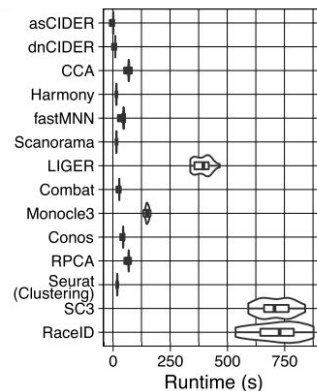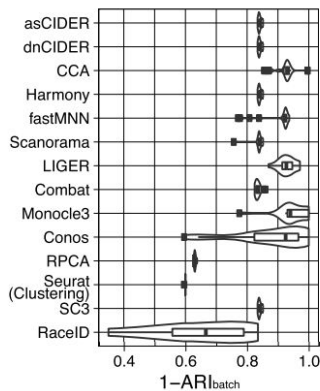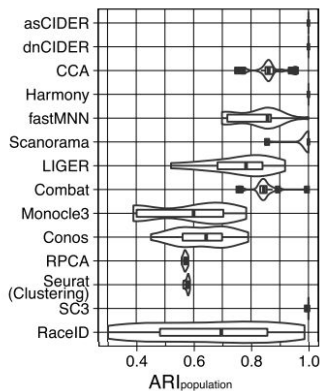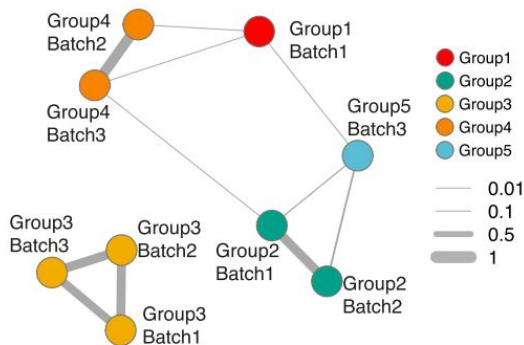- 5 groups across 3 batches with non-identical populations



- Many alternative methods "overcorrect" for batch effects

# Benchmarking with simulated data

- 5 groups across 3 batches with non-identical populations



- Many alternative methods "overcorrect" for batch effects
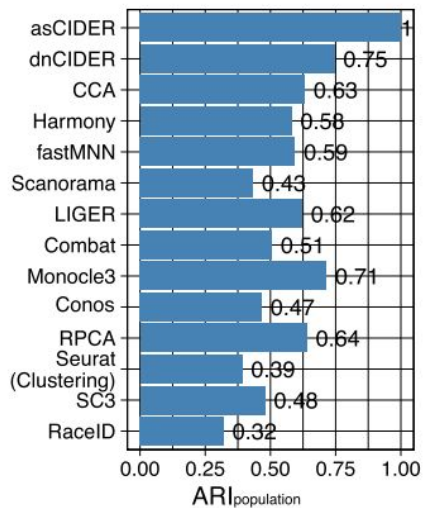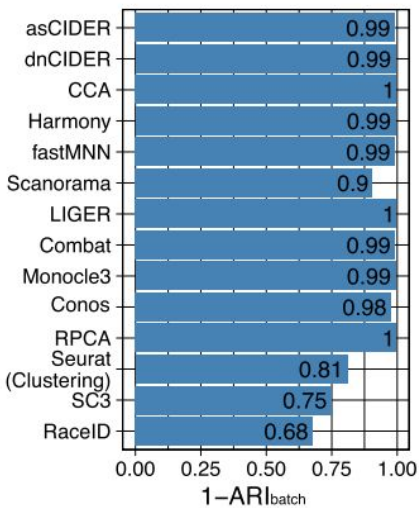
# Benchmarking with real data: PBMCs

**Dataset 3:** human peripheral blood mononuclear cells (PBMCs)

- **9 cell types/subtypes**
- **2 techniques** (10x 3' and 5' single-cell GE) as batches

# Benchmarking with real data: PBMCs

**Dataset 3:** human peripheral blood mononuclear cells (PBMCs)

- 9 cell types/subtypes
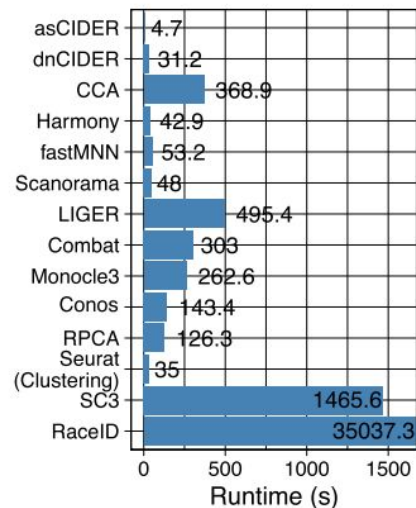- 2 techniques (10x 3' and 5' single-cell GE) as batches
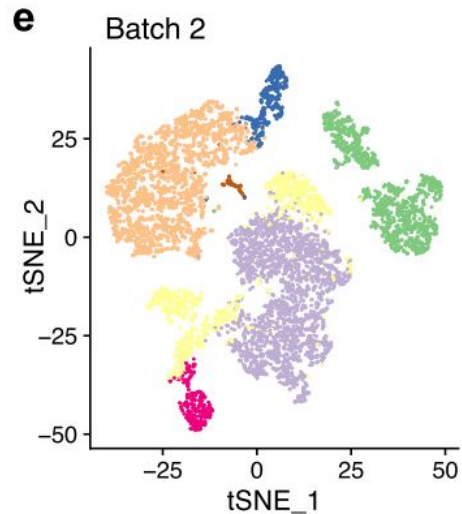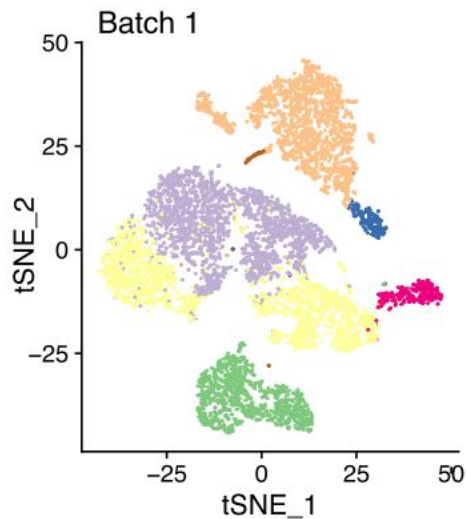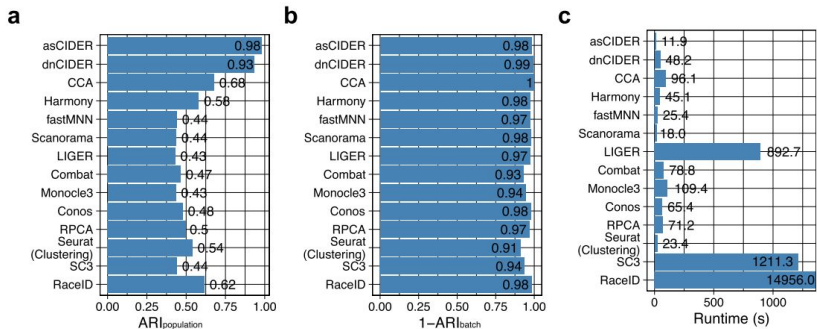
# Benchmarking with real data: PBMCs

**asCIDER could reveal the underlying relationships among initial clusters**
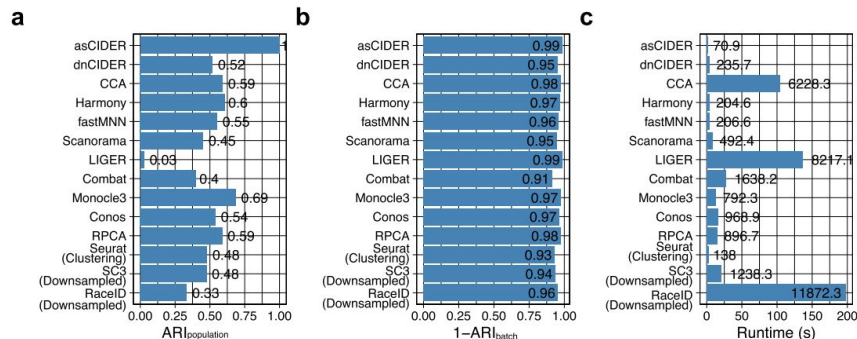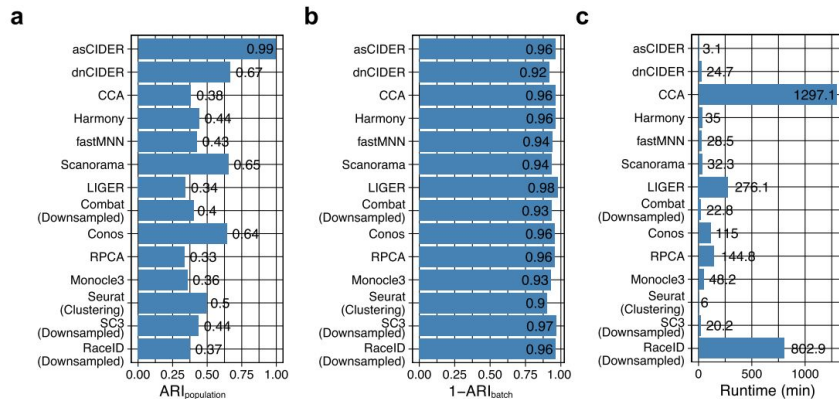
# Benchmarking with real data

## Dataset 4: human and mouse pancreatic data

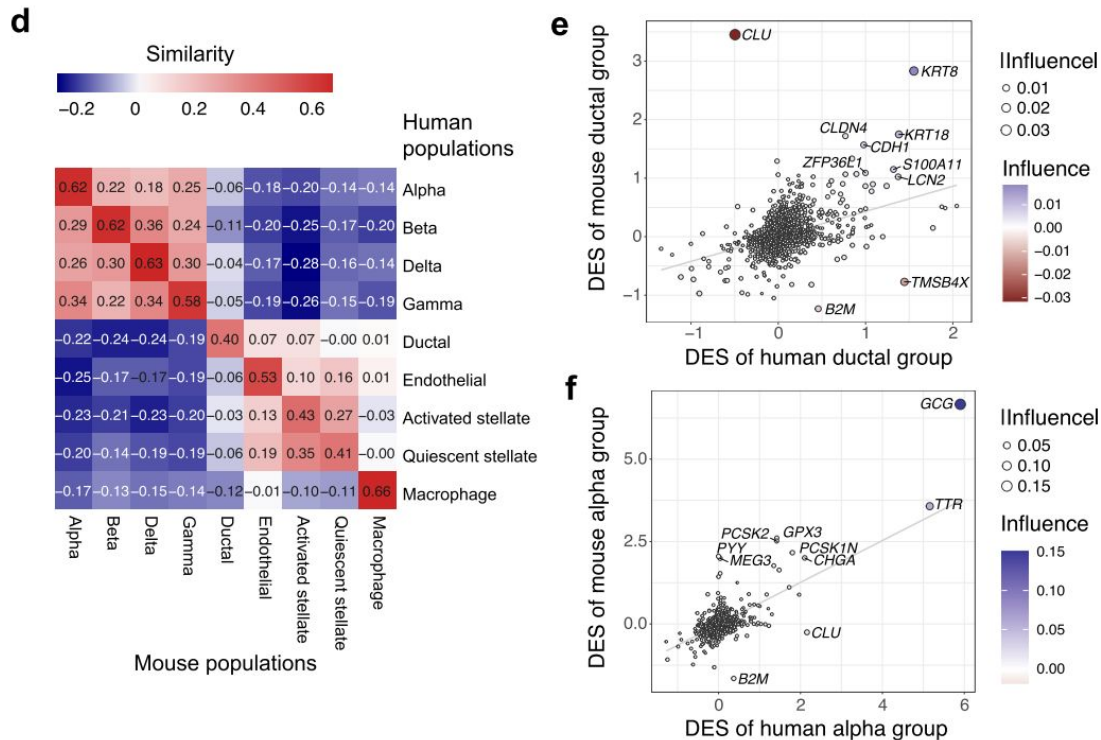## Dataset 5: COVID-19 study
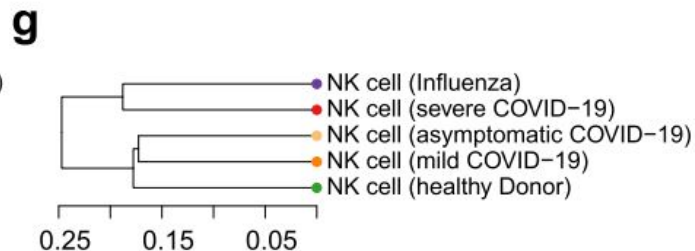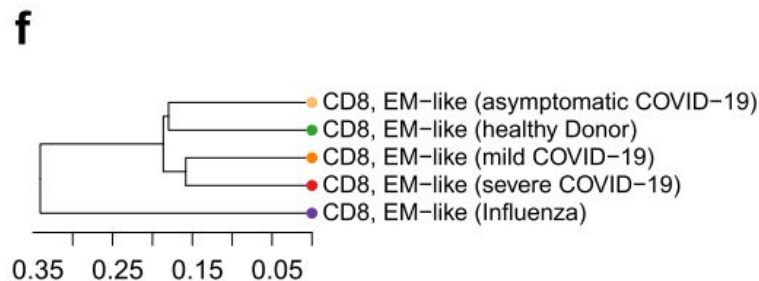
## Dataset 6: breast cancer data

# Benchmarking with real data: human vs mouse pancreatic cells

2 mouse samples, 4 human samples → both species and donor effect

# Benchmarking with real data: COVID-19

PBMCs collected from healthy donors, patients with severe influenza, and patients with various severity of COVID-19 (asymptomatic, mild, and severe)

# Benchmarking with real data: breast cancer
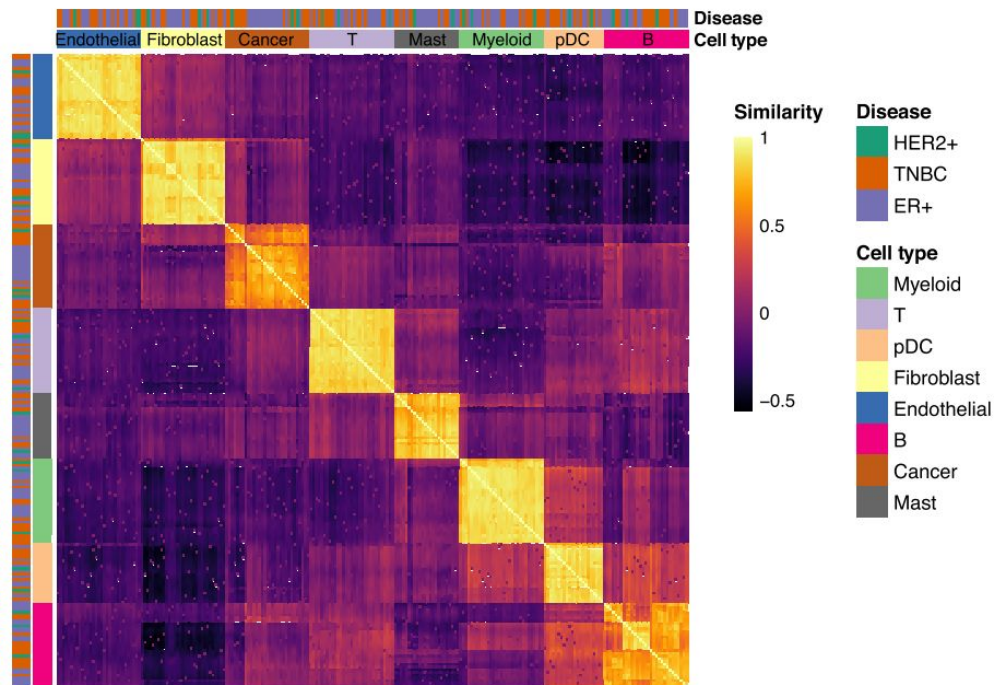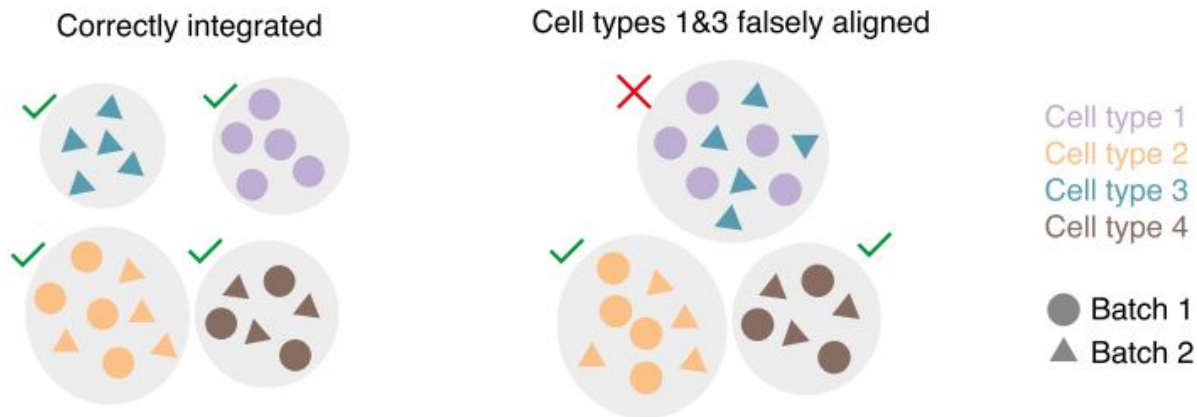
170K cells from 31 breast cancer patients

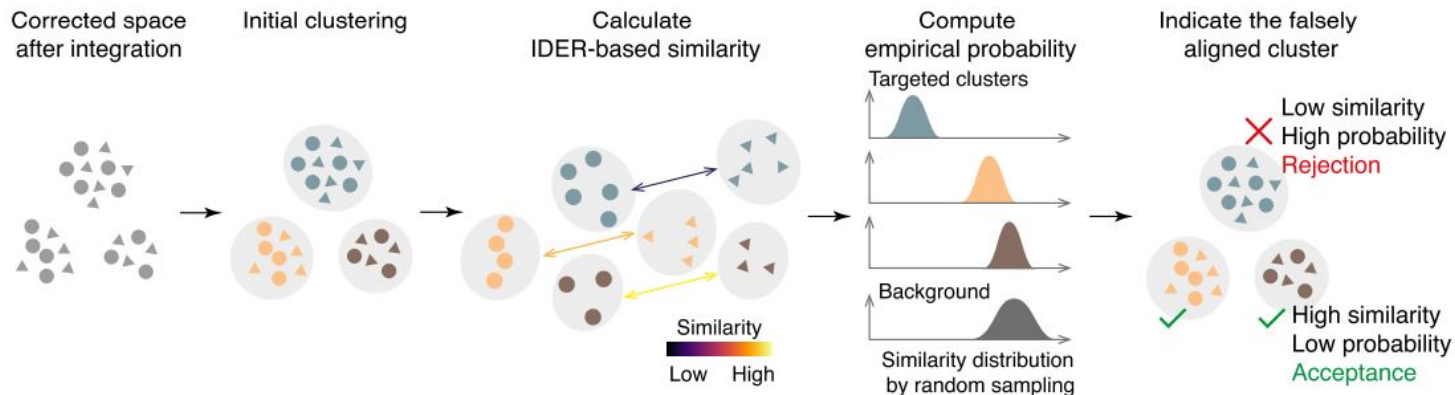Two samples per patient: before and after treatment

# CIDER as a ground-truth-free test metric of integration

- Common issue for integration methods: incorrect alignment - sometimes groups are merged that shouldn't have been
- Other existing metrics require predefined cell populations (e.g., cLISI: Cell-type local inverse Simpson Index)

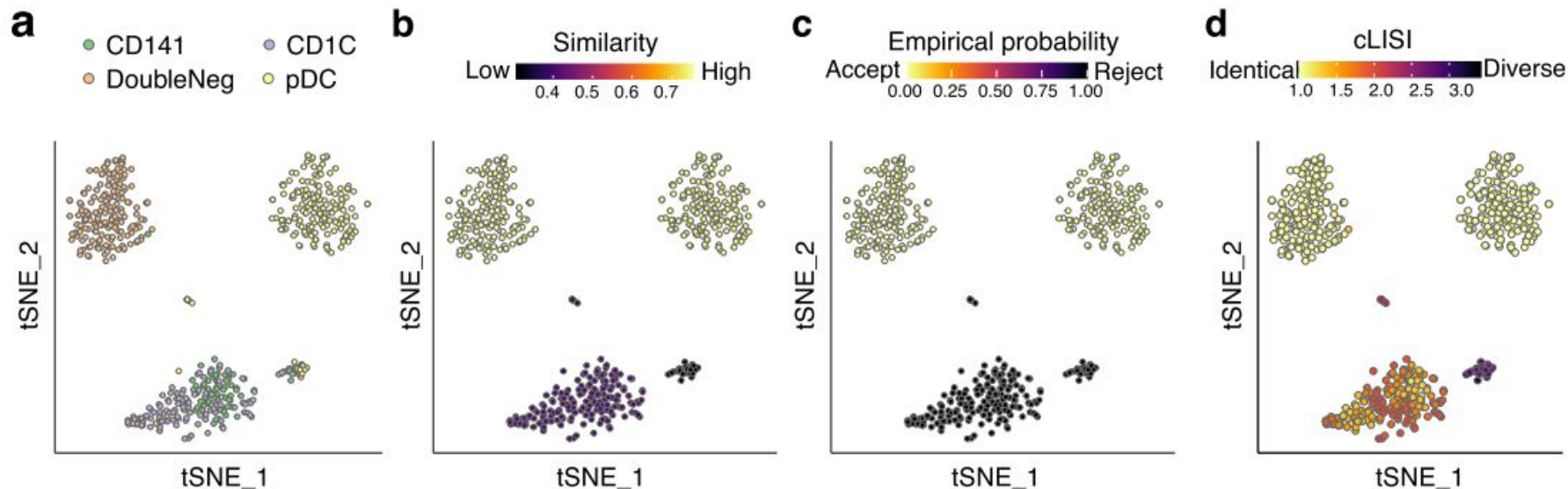# Embedding CIDER into a workflow to evaluate integration:

- Other method: Perform batch correction and learn cross-batch clusters
- Apply IDER metric to cross-batch clusters:
  - For each learned cluster, split by batch
  - Compute IDER similarity for each pair → higher similarity=better integration
  - Compare pairs' similarity to distribution of similarities for random partitions within the cluster
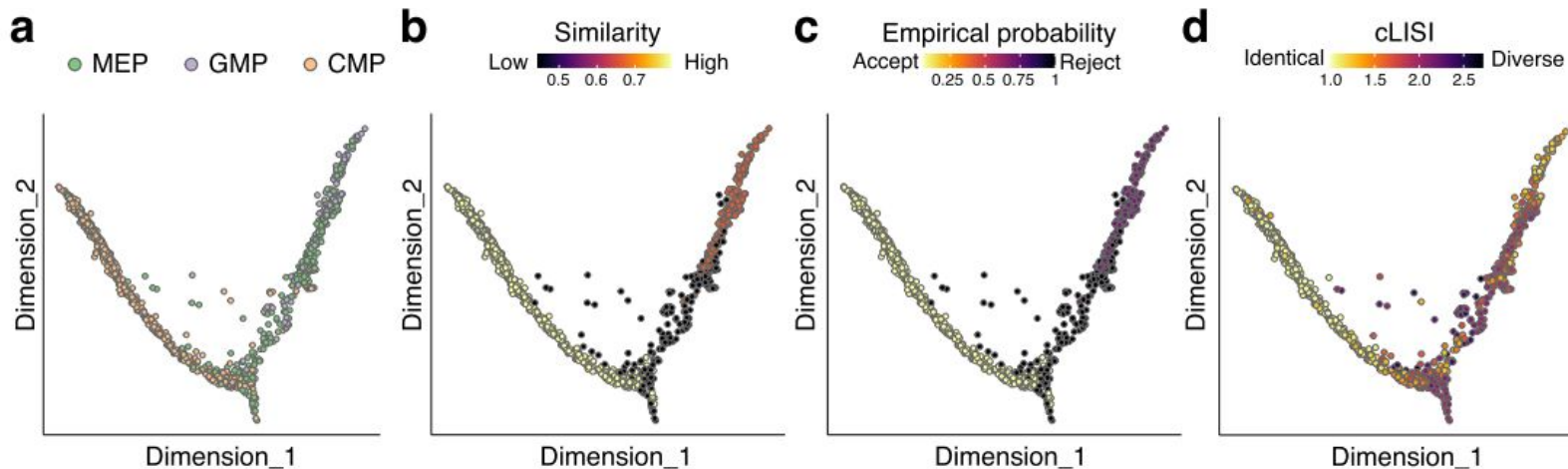
# Using CIDER to evaluate CCA integration on a dendritic cell dataset

- CD141 & CD1C are prone to being merged by batch correction methods
- CIDER has similar results with cLISI (but doesn't require labels to calculate)

# Using CIDER to evaluate mouse hematopoietic progenitor data (continuous data structure)

- Goal: Use CIDER to evaluate local biological heterogeneity without predefined annotations
- Mouse hematopoietic progenitor data (common myeloid, megakaryocyte/erythrocyte, and granulocyte/macrophage progenitor cells) from 2 platforms (MARS-seq and Smart-Seq2)

# Discussion

Summary:

- Introduced IDER, a differentially expressed gene-based similarity metric, which can be used to identify cross-batch clusters
- Both dnCIDER and asCIDER were evaluated on a wide array of benchmarks (dnCIDER was often much better)
- IDER metric can be used to evaluate other batch-correction methods in the absence of ground truth labels

Limitations:

- Developed for scRNA-Seq - currently not designed for multi-modal data
- Linear approach
- Group-level analysis assumes coarse-grained clusters (not continuous data)

# Discussion topics

- Worse performance for dnCIDER vs asCIDER – how do we feel about that, given that one of their presented advantages is not needing labels?

- This space is quite saturated (e.g., all the methods they benchmarked against)
  - What does a new method need to achieve to really be worth using? Did this paper meet that standard?
  - Where should the field go next?

- Circular benchmarks: Most "ground truth" labels are actually the output of clustering methods/previously found gene signatures which are used to identify cell types, so new methods benchmark against these