

---

# Single cell sample integration

Remi Montagne

Institut Curie

**EBAI 2023 - 11/08/2023**

# Introduction

---

Starting point: normalized, reduced **individual** matrices

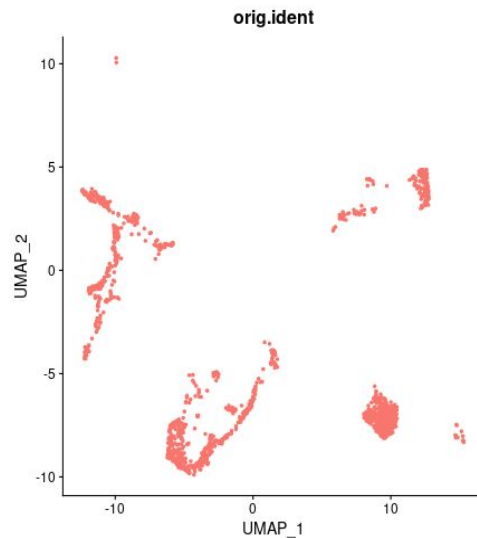
Next step: start getting information

# Introduction

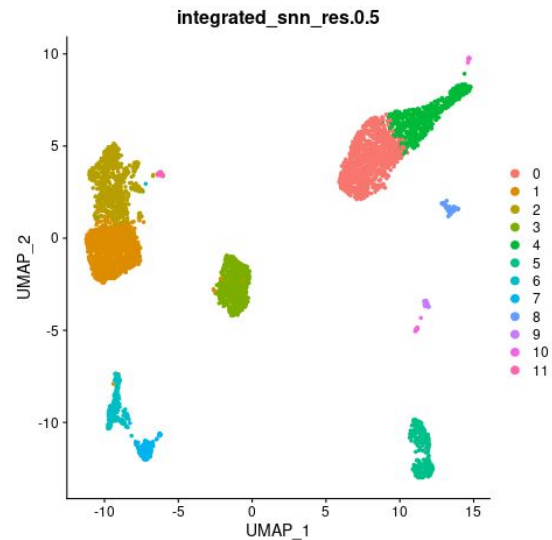
Starting point: normalized, reduced **individual** matrices

Next step: start getting information

→ Visualize the cells



→ Understand what is in the samples (clustering)



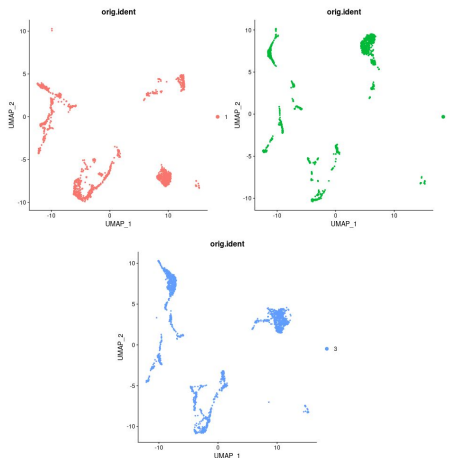
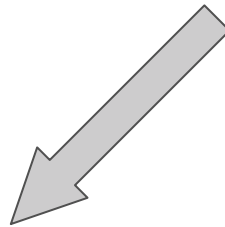
# Introduction

Starting point: normalized, reduced **individual** matrices

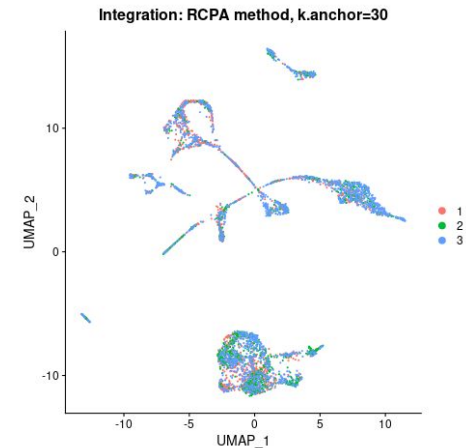
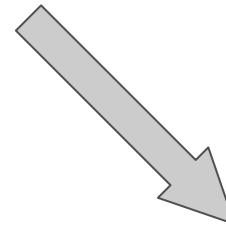
Next step: start getting information

But should we do that

on individual  
samples?



On all samples  
together ?



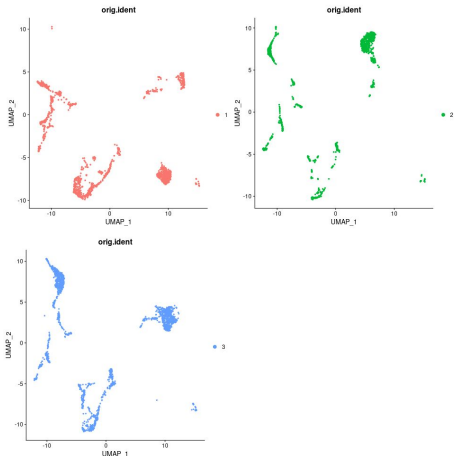
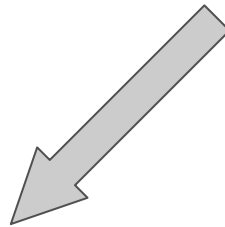
# Introduction

Starting point: normalized, reduced **individual** matrices

Next step: start getting information

But should we do that

on individual  
samples?



- Quick way to have a first look at data
- Repetitive
- Makes more sense to bring replicates together.
- Makes more sense to bring together similar samples (same experiment, organ...)

# Introduction

Starting point: normalized, reduced **individual** matrices

Next step: start getting information

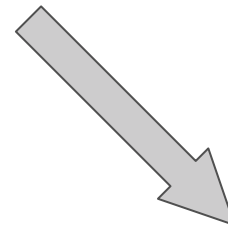
But should we do that



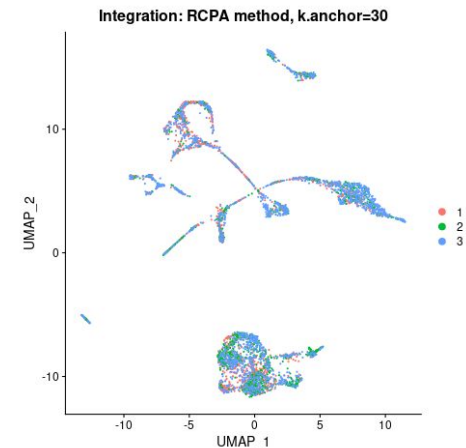
- Allows to work across multiple samples.
- Particularly important for cell populations visualization and identification
- Many cells : helps identifying rare populations



- Overcorrection?



On all samples together ?



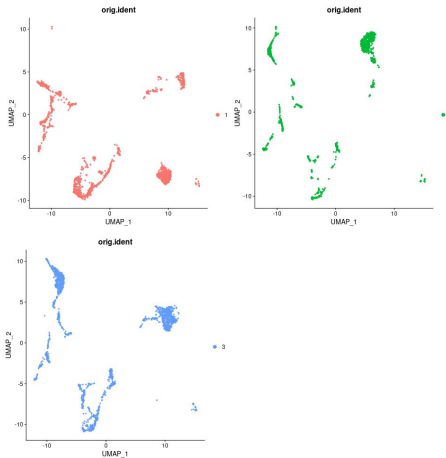
# Introduction

Starting point: normalized, reduced **individual** matrices

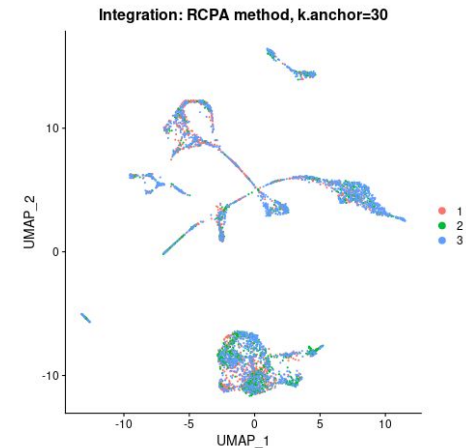
Next step: start getting information

But should we do that

on individual  
samples?



On all samples  
together ?



# Introduction

---

Starting point: normalized, reduced **individual** matrices

Next step: start getting information

We will do that on all samples together



# Introduction

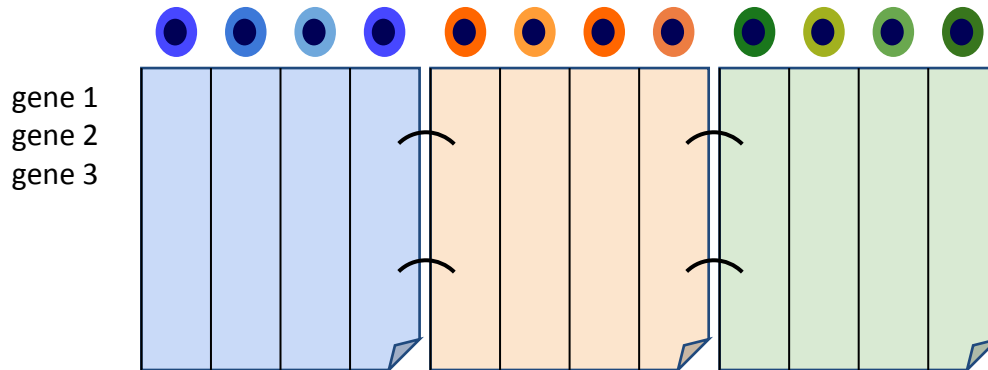
---

Starting point: normalized, reduced **individual** matrices

Next step: start getting information

We will do that on all samples together

**Problem:** simple matrix concatenation does not always work



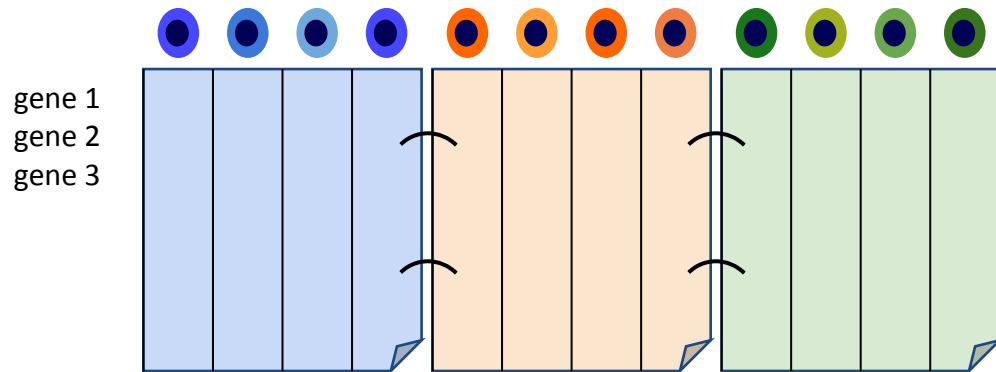
# Introduction

Starting point: normalized, reduced **individual** matrices

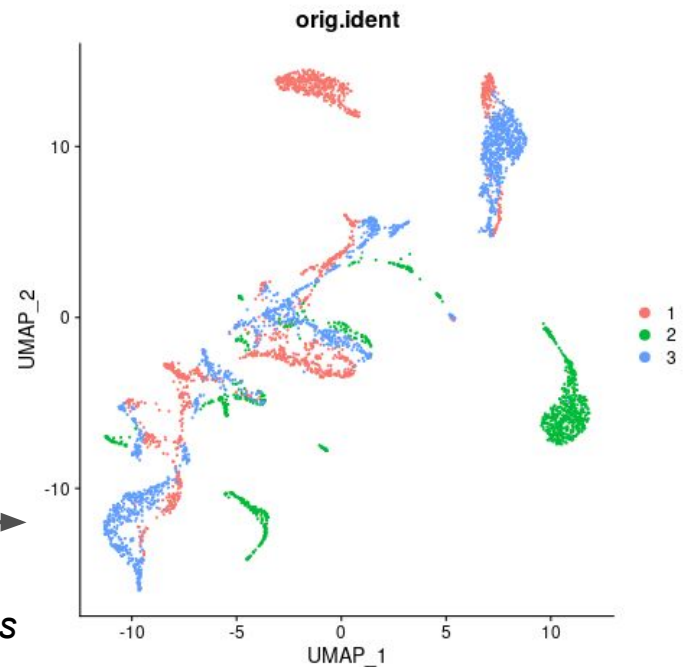
Next step: start getting information

We will do that on all samples together

**Problem:** simple matrix concatenation does not always work



*same model (PBMC), unaligned cells*



# Introduction

Starting point: normalized, reduced **individual** matrices

Next step: start getting information

We will do that on all samples together

**Problem:** simple matrix concatenation does not always work

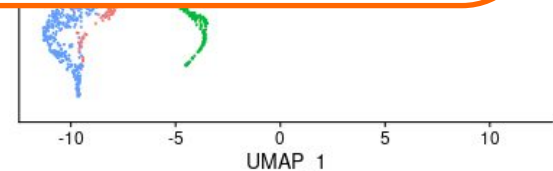
gene 1  
gene 2  
gene 3

This is a problem of batch effect.

We need a more sophisticated **integration method**

orig.ident

1  
2  
3

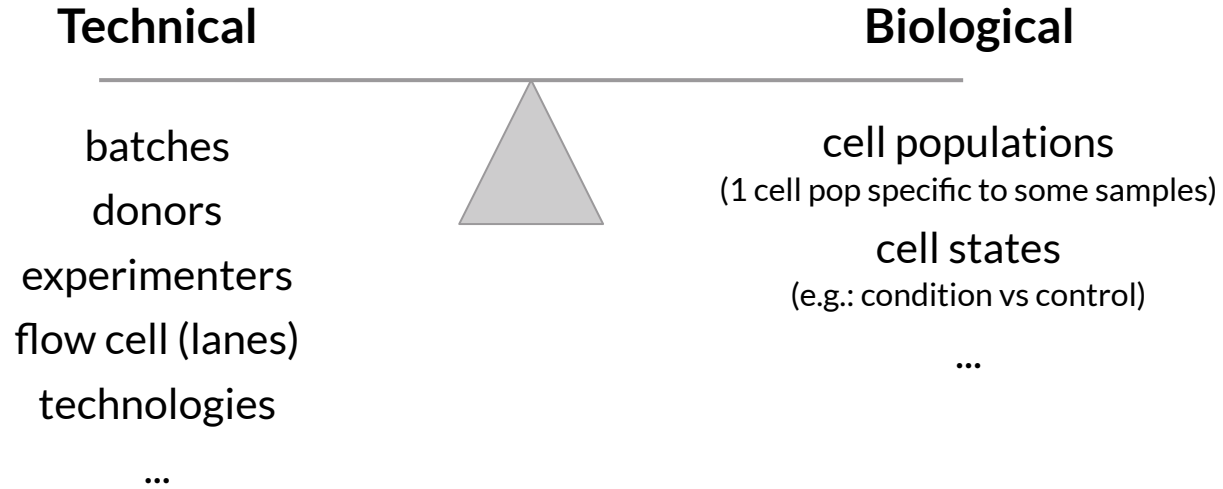


# Variability across samples

# Variability across samples

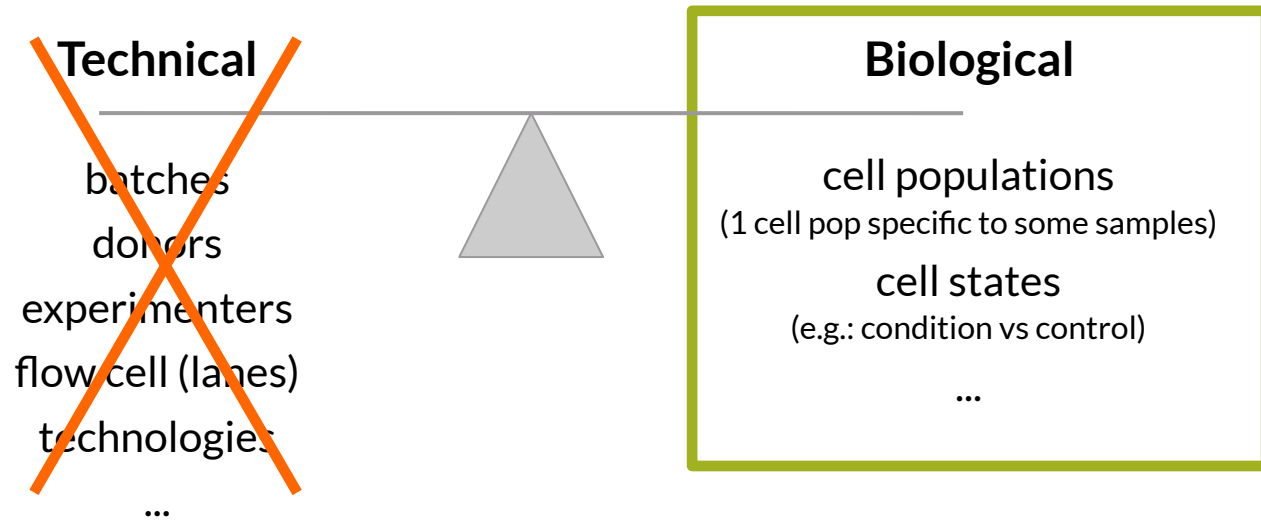
---

## 2 sources of variability across samples



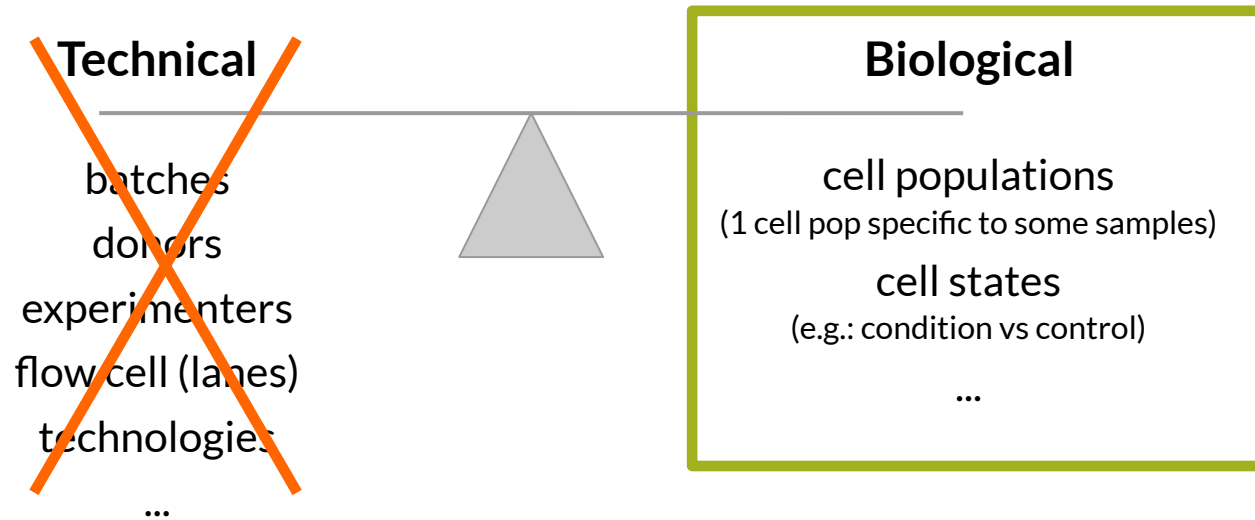
# Variability across samples

## 2 sources of variability across samples



# Variability across samples

## 2 sources of variability across samples



### → Solutions:

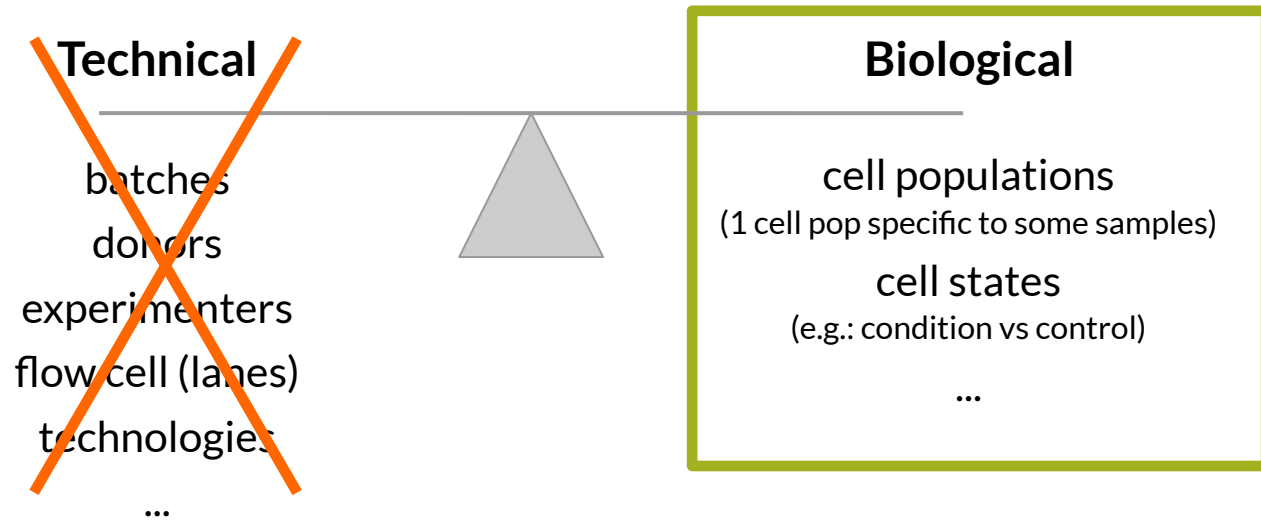
Strategies to avoid factors causing batch effect in the lab

**Solution:** Technical factors that potentially lead to batch effects may be avoided with mitigation strategies in the lab and during sequencing. Examples of lab strategies include: sampling cells on the same day, using the same handling personnel, reagent lots, protocols, reducing PCR amplification bias, and generally using the same equipment. Sequencing strategies can include multiplexing libraries across flow cells. For example, if samples came from two patients, pooling libraries together and spreading them across flow cells can potentially spread out the flow cell-specific variation across samples.

<https://www.10xgenomics.com/resources/analysis-guides/introduction-batch-effect-correction>

# Variability across samples

## 2 sources of variability across samples



→ **Solutions:**

Strategies to avoid factors causing batch effect in the lab

Computational data integration



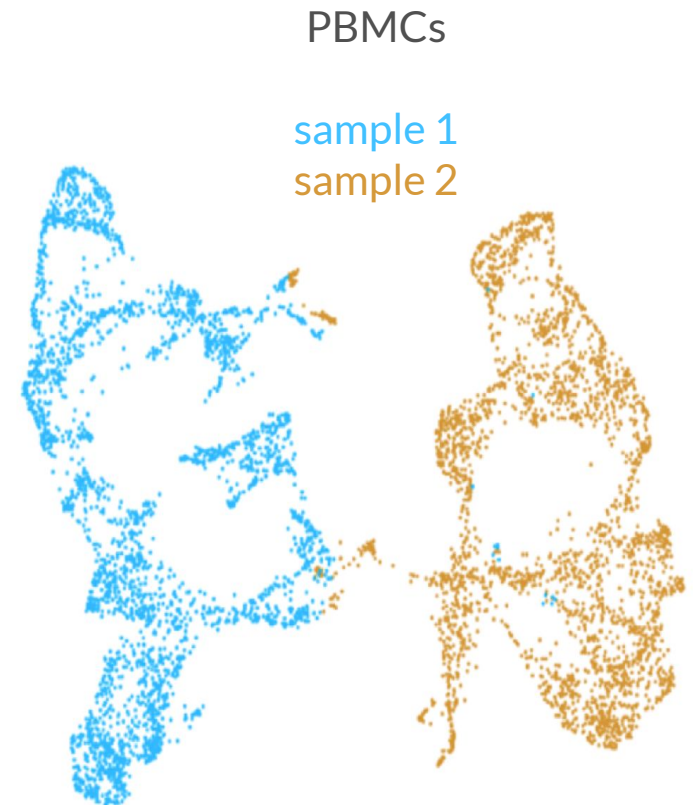
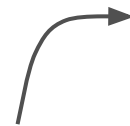
# When to integrate

# When to integrate

- **Integrate** when obvious batch effect between samples, typically seen on low dimension visualization

*In this example, the sample of origin would be a huge bias for clustering*

*The samples need integration to align cell types/clusters and then identify them correctly*

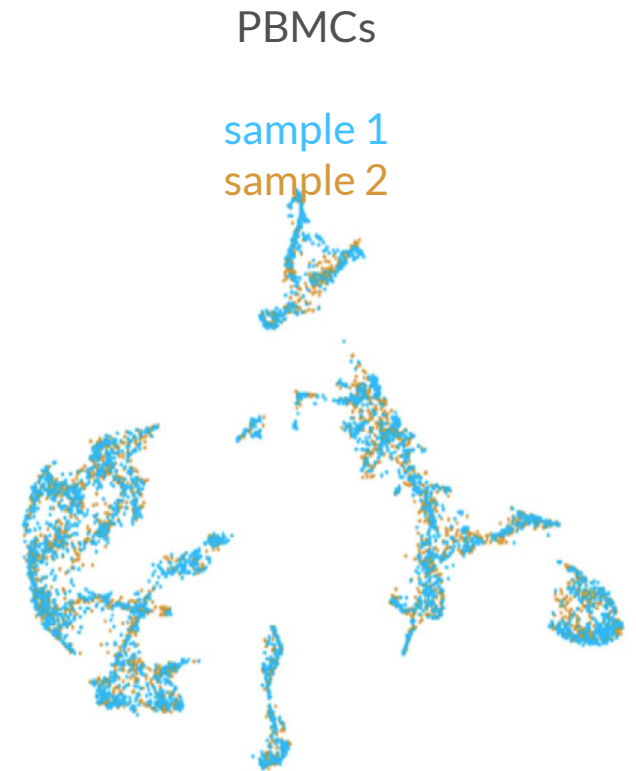


<https://www.10xgenomics.com>

# When to integrate

---

- **Integrate** when obvious batch effect between samples, typically seen on low dimension visualization
- Do not integrate otherwise:  
e.g.: replicates generated in the same time and exactly in the same manner may not need integration



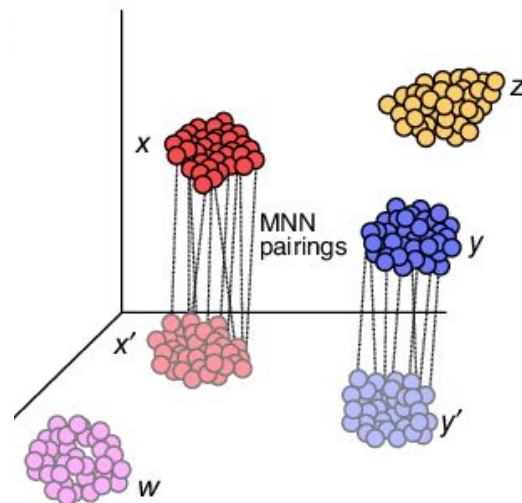
<https://www.10xgenomics.com>

# Integration with Seurat

# Integration with Seurat

## Many methods

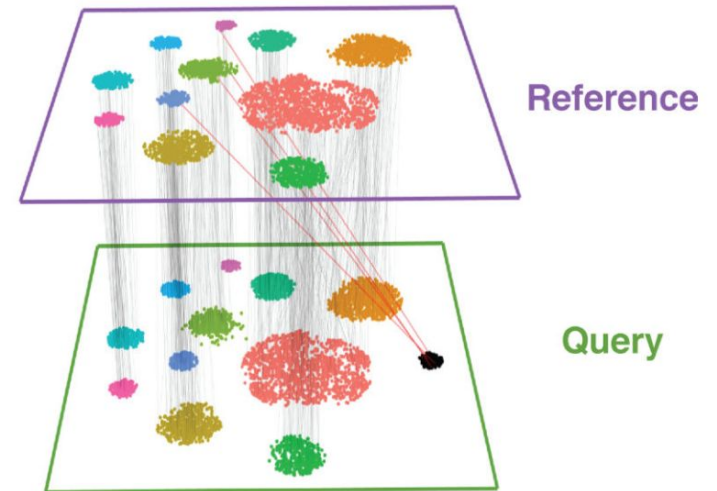
- Over 49 methods (Luecken et al., Nat Methods 2022)
- Seurat integration: group of **similarity-based** methods



# Integration with Seurat

## Principle

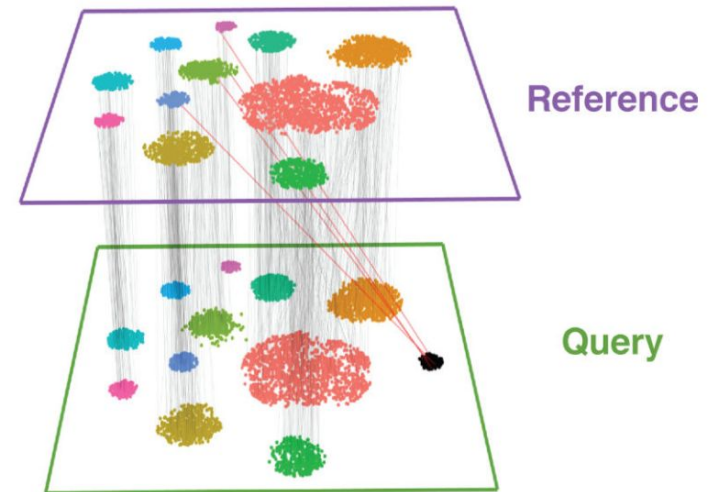
- Integration is always **pairwise**: correct a sample, **the query** to match the expression data of another sample, **the reference**.



# Integration with Seurat

## Principle

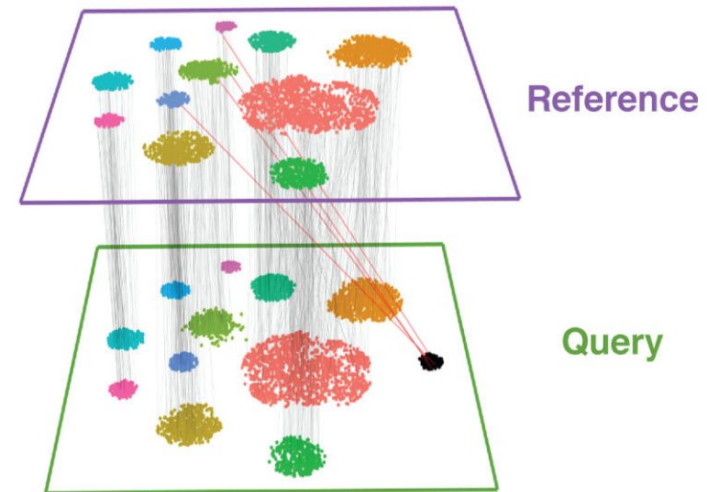
- Integration is always **pairwise**: correct a sample, **the query** to match the expression data of another sample, **the reference**.
- Seurat identifies pairs of close cells across both datasets, the **anchors** (Mutual Nearest Neighbors).



# Integration with Seurat

## Principle

- Integration is always **pairwise**: correct a sample, **the query** to match the expression data of another sample, **the reference**.
- Seurat identifies pairs of close cells across both datasets, the **anchors** (Mutual Nearest Neighbors).
- The difference between them is used to compute a **correction**.

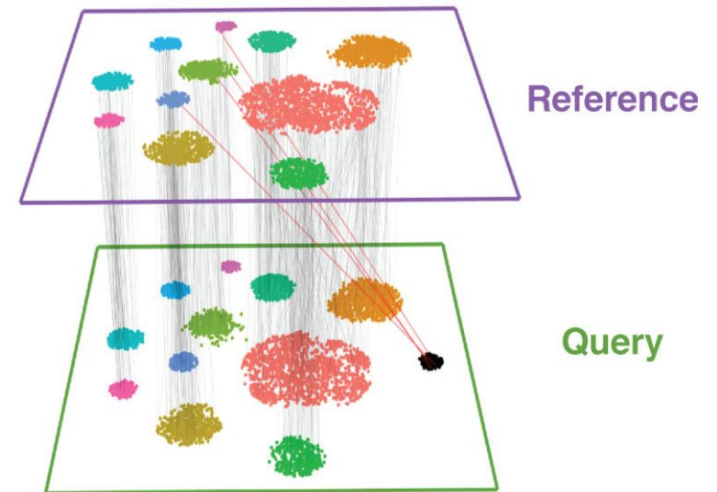




# Integration with Seurat

## Principle

- Integration is always **pairwise**: correct a sample, **the query** to match the expression data of another sample, **the reference**.
- Seurat identifies pairs of close cells across both datasets, the **anchors** (Mutual Nearest Neighbors).
- The difference between them is used to compute a **correction**.
- The correction is used to **align** all the **query** cells on the **reference** cells.



# Integration with Seurat

## Principle

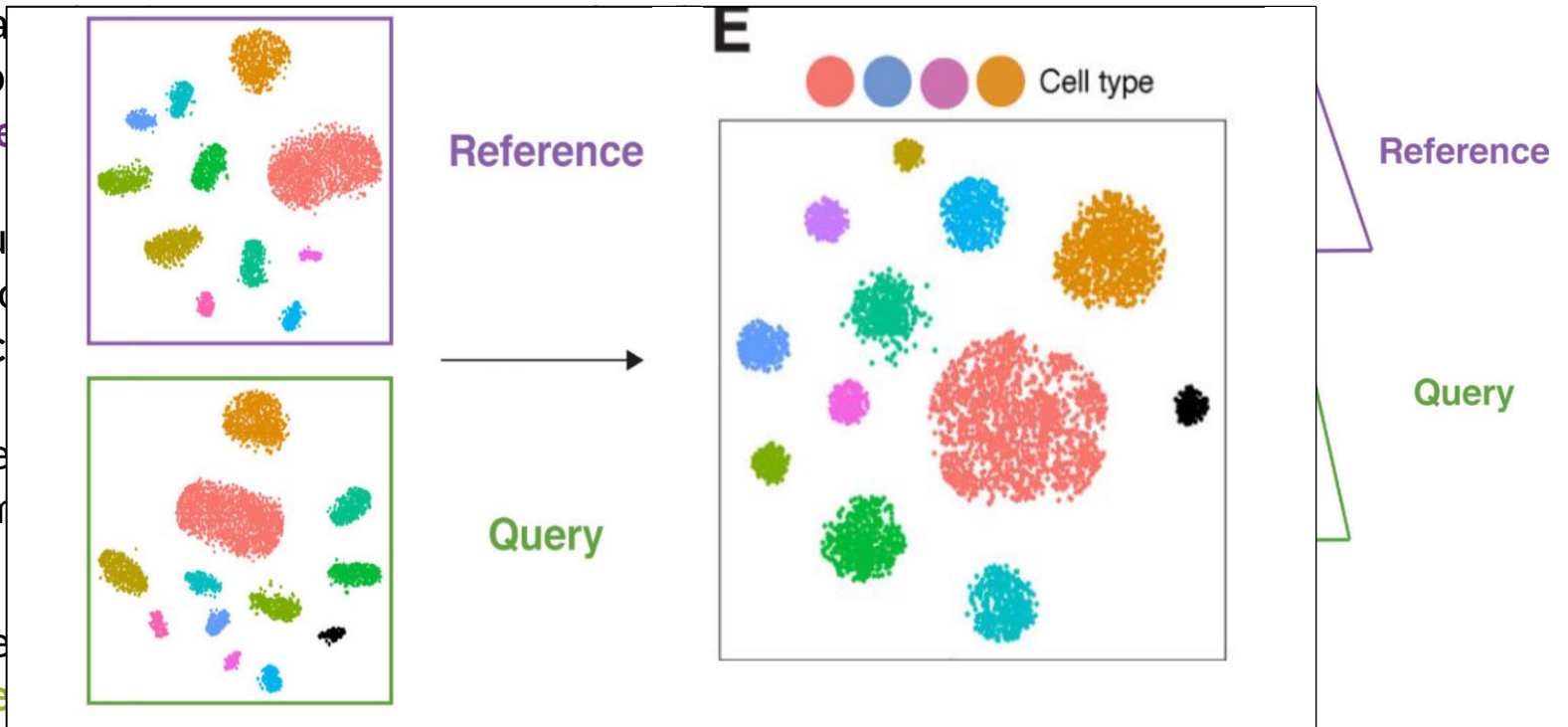
- Integration is always **pairwise**: correct

a sa  
exp  
refe

- Seu  
acro  
anc

- The  
com

- The  
que



# Integration with Seurat

## Principle

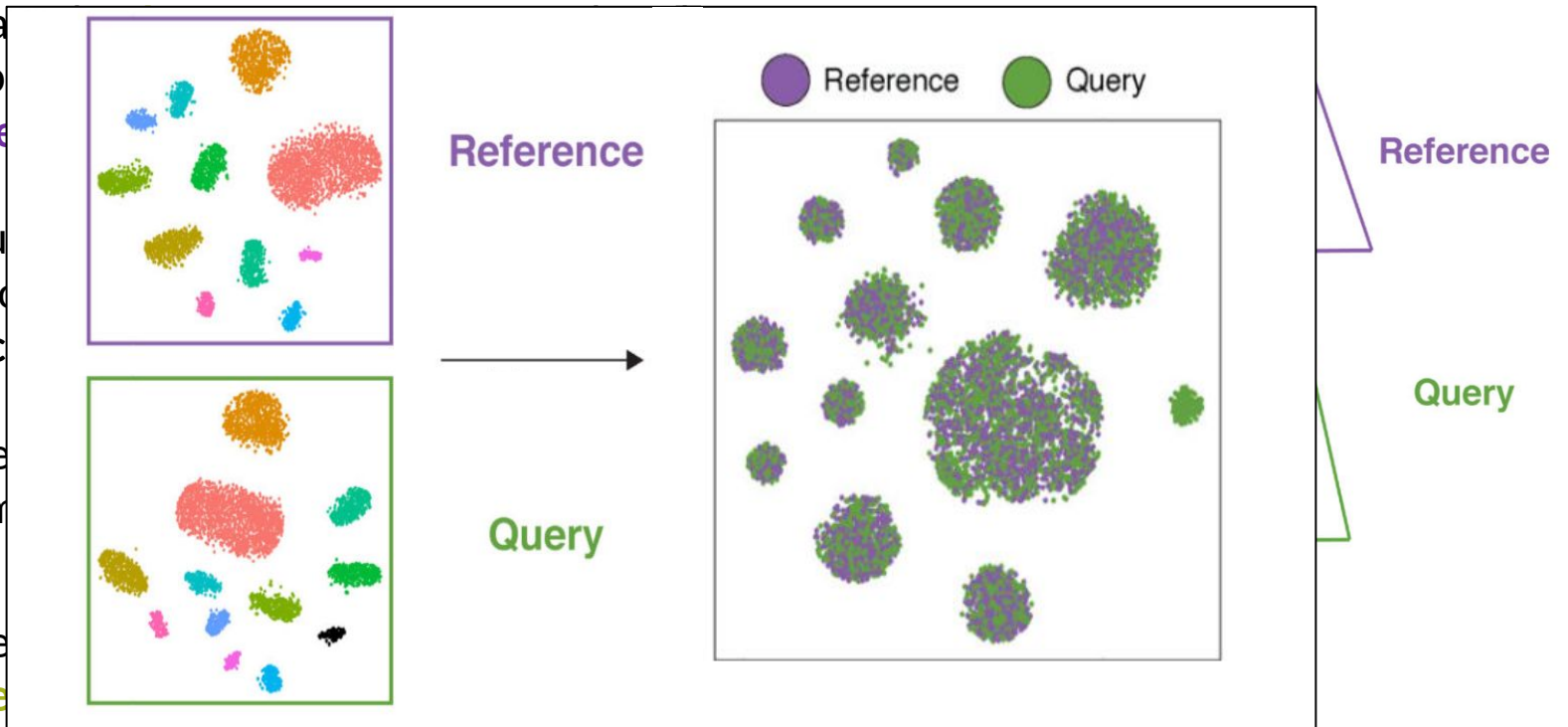
- Integration is always **pairwise**: correct

a sa  
exp  
refe

- Seu  
acro  
anc

- The  
com

- The  
que

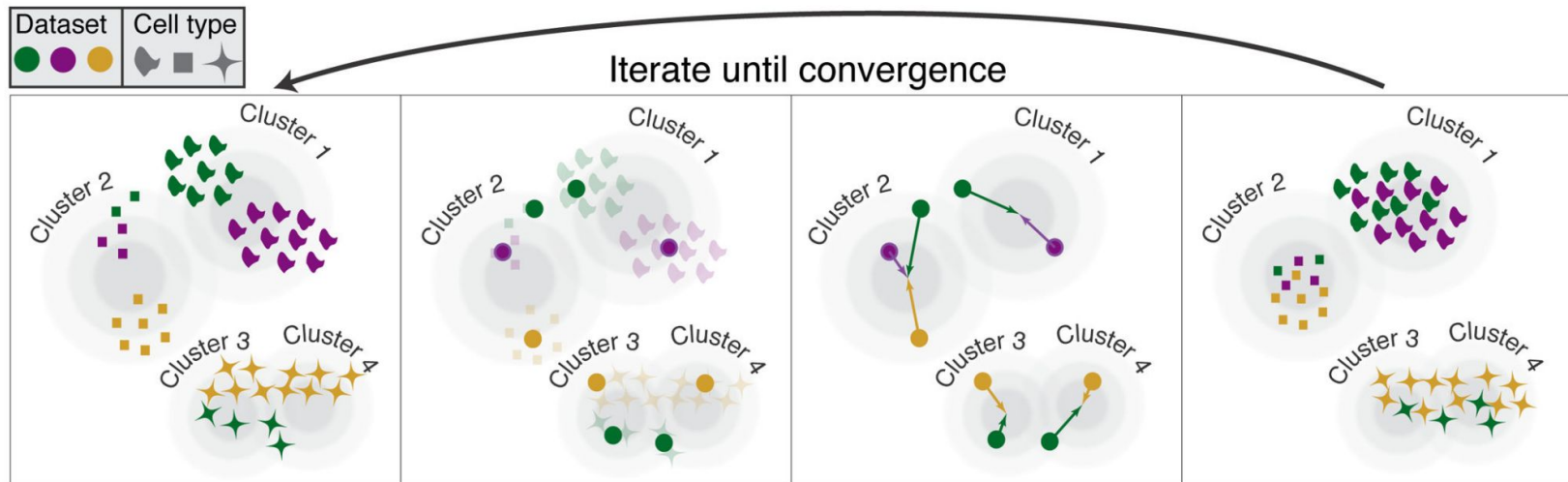


# Integration with Harmony

# Integration with Harmony

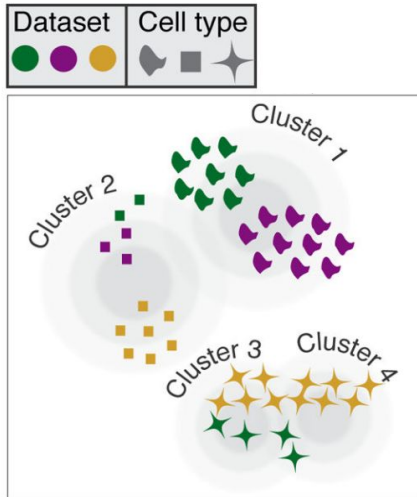
## Many methods

- Harmony integration: Iterative clustering in dimensionally reduced space



# Integration with Seurat

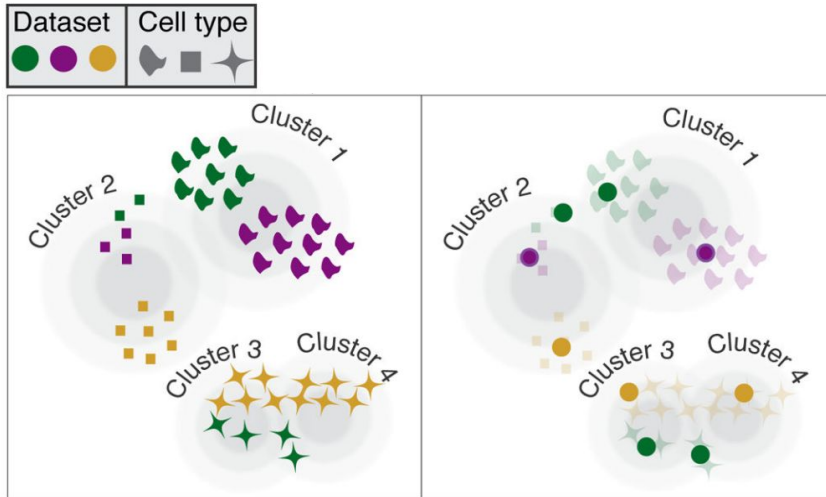
## Principle



- Integration is **not** pairwise: correct all samples in the same time
- Find many small clusters
- Constraint: clusters must contain cell from several samples

# Integration with Seurat

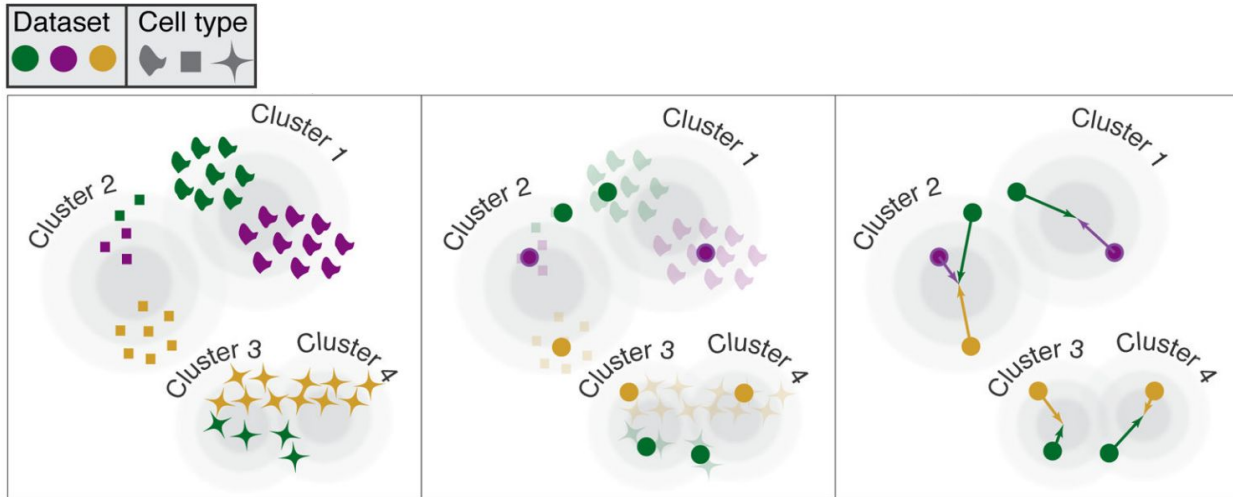
## Principle



- Integration is **not** pairwise: correct all samples in the same time
- Find many small clusters
- Constraint: clusters must contain cell from several samples
- Get cluster centroids (= average position) of each sample.

# Integration with Seurat

## Principle

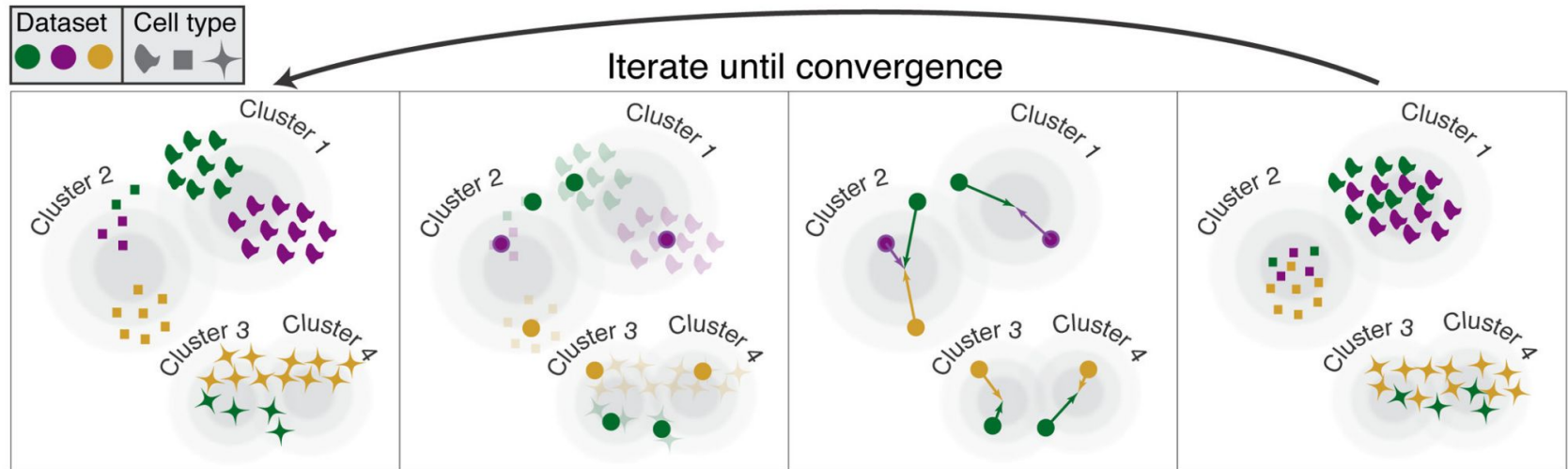


- Integration is **not** pairwise: correct all samples in the same time
- Find many small clusters
- Constraint: clusters must contain cell from several samples
- Get cluster centroids (= average position) of each sample.
- Compute sample corrections for each cluster
- The aim is to get all centroids of the same cluster together



# Integration with Seurat

## Principle



- Integration is **not** pairwise: correct all samples in the same time
- Find many small clusters
- Constraint: clusters must contain cell from several samples
- Get cluster centroids (= average position) of each sample.
- Compute sample corrections for each cluster
- The aim is to get all centroids of the same cluster together
- Apply corrections to cells

# Integration with Seurat

## Principle



Iterate until convergence

This method relies on clustering but the clusters are only used for integration purpose.

Later we will perform clustering for cell population discovery

**not** pairwise: correct all samples in the same time

- Find many small clusters
- Constraint: clusters must contain cell from several samples

centroids (= average position) of each sample.

corrections for each cluster

to cells

- The aim is to get all centroids of the same cluster together

# Benchmarking methods

# Benchmarking methods

**a**

	Considerations	scANVI	Scanorama embed	scVI	FastMNN embed	scGen	Harmony	FastMNN gen	Seurat v3 RPCA	tSNE	Scanorama gene	ComBat	MNN	Seurat v3 CCA	tVAE	Conos	DESC	LIGER	SAUCIE embed	SAUCIE gene
Input	Programming language																			
	Method runs without additional information	✗				✗														
Scib results	Consistent top performer	✓	✓	✓		✓														
	Top method on small/simple tasks		✓		✓	✓	✓													
	Top method on large/complex tasks	✓	✓	✓		✓														
	Top method on ATAC data	—		—			✓											✓		
Task details	Integrates strong batch effects	✓	—	—		✓			—	—				—						
	Top method for recovery cell states or modules	✓	✓								✓	✓	✓							
	Confounding of bio and batch variance	✓	—			✓														
	Top method for trajectories	—	✓	—	✓	✓														
	Method deals with varying compositions											✗								
	Fast method for quick results										✓	✓								
Speed	Scales well to large datasets on CPU	✓	—	✓						✓	✓	—							✓	✓
	Method has GPU support	✓		✓		✓									✓		✓		✓	✓
	Scales well to feature spaces beyond genes														✓	✓				
	Method shows corrected expression					✓		✓	✓		✓	✓	✓	✓						✓
Output	Method gives relative cell embeddings								✗							✗				

Fulfills the criterion    Python  
 Partial fulfillment of criterion    R  
 Does not fulfill criterion

**Seurat v3**    Luecken et al., Nature Methods 2022

A few benchmarks, that do not agree with each other

Büttner et al., Nat. Methods. 2019  
 Chen et al., Nat. Biotechnol 2020  
 Tran et al., Genome Biol. 2020

# Benchmarking methods

---

Do not hesitate to test several methods

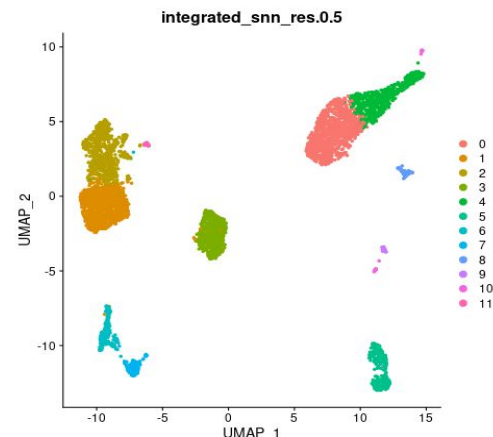
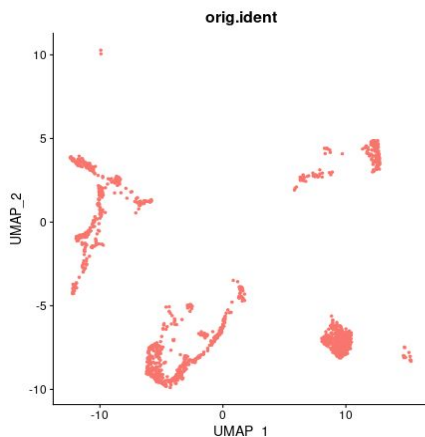


Luecken *et al.*, Nature Methods 2022

What is integration for

# What is integration for

- For computational efficiency, integration is only performed on the most variable genes, not all the genes.
- It is intended for **visualization** and **clustering**



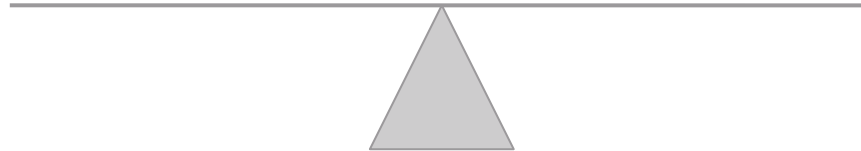
- For differential expression analysis, we go back to raw data

# Conclusion

## A good integration method

**Technical**

**Biological**



- Corrects for technical variability:

- samples
- donors
- experimenter
- technologies

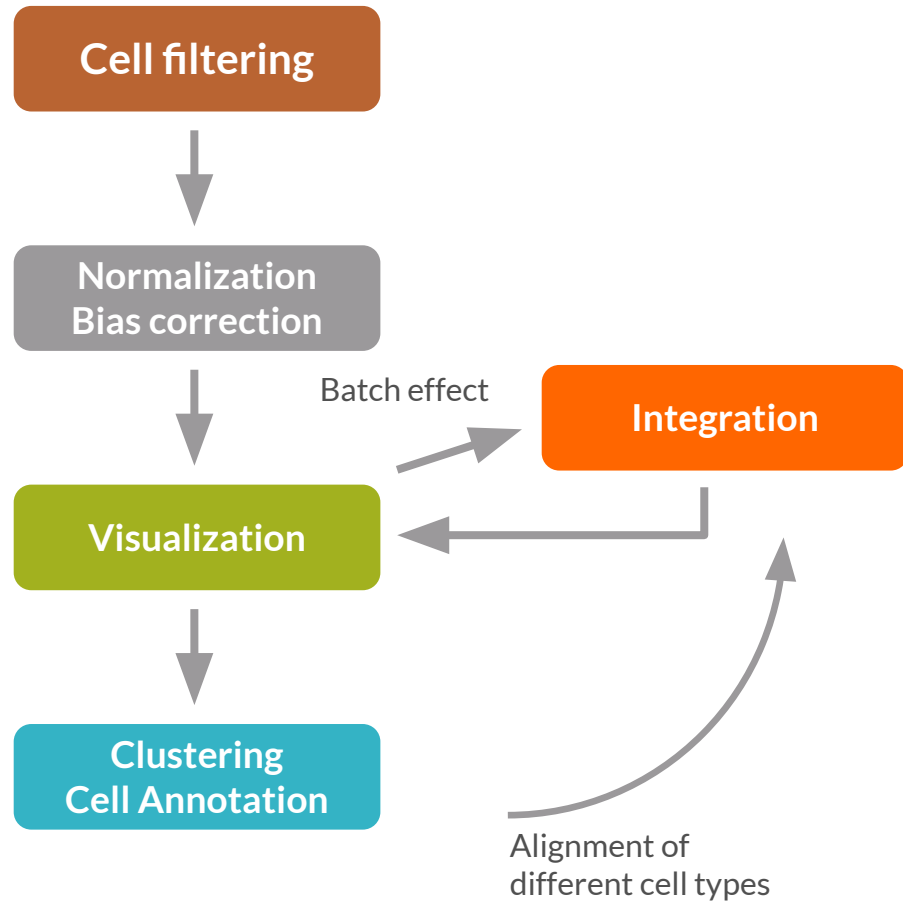
- Preserves biological signal

- cell types across different samples, tissues
- cell trajectories
- differences (cell subtypes, cell states) between condition and control
- population (cell subtypes, cell states) unique to a condition...



# Conclusion

Preparation of the data is not always a linear process

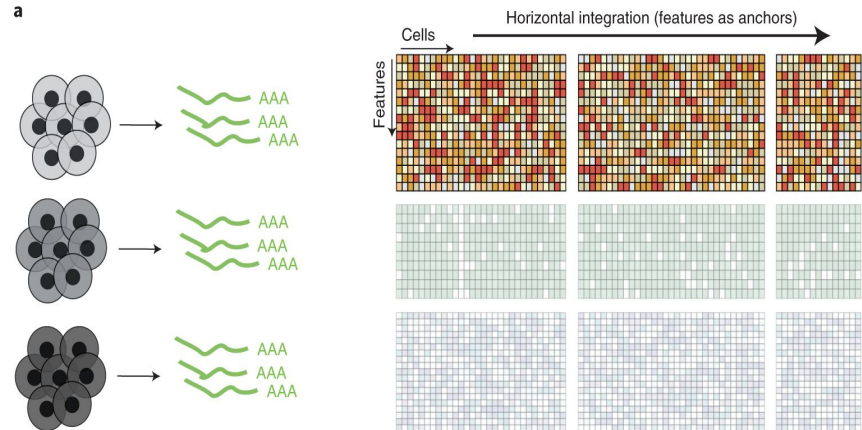


# Conclusion

## Different types of integrations

- Horizontal: different samples same modality

We saw horizontal integration



Luecken *et al.*, Nat Met 2021

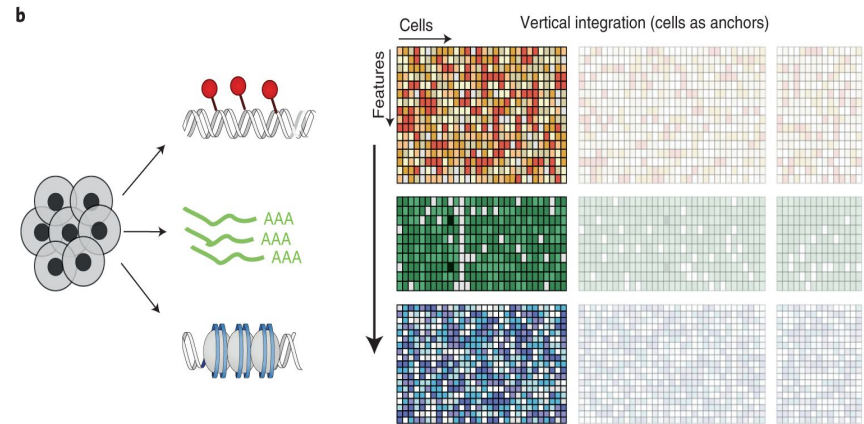
# Conclusion

## Different types of integrations

- Horizontal: different samples same modality

We saw horizontal integration

- Vertical: same sample different modalities (multiomics)



Luecken *et al.*, Nat Met 2021

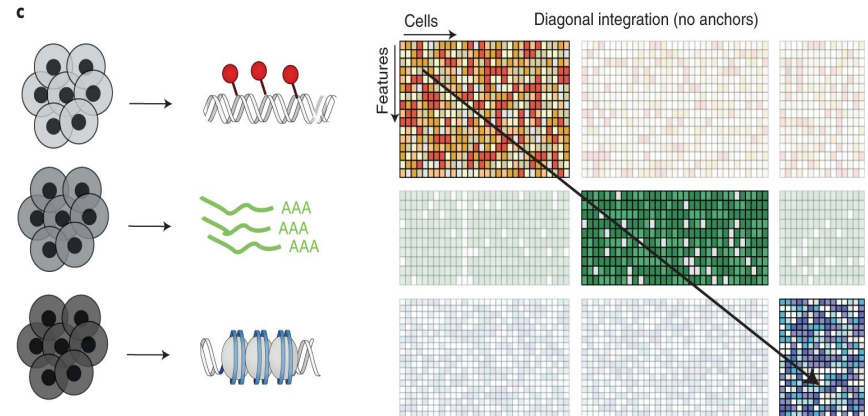
# Conclusion

## Different types of integrations

- Horizontal: different samples same modality

We saw horizontal integration

- Vertical: same sample different modalities (multiomics)
- Diagonal: different samples different modalities



Luecken *et al.*, Nat Met 2021

# Acknowledgements

---

**Parts of this course are inspired by**

The Swiss *Institute of Bioinformatics* course [Single Cell Transcriptomics](#)