
Single cell sample integration

Remi Montagne

Institut Curie

Initiation single cell - 10/20/2023

Introduction

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

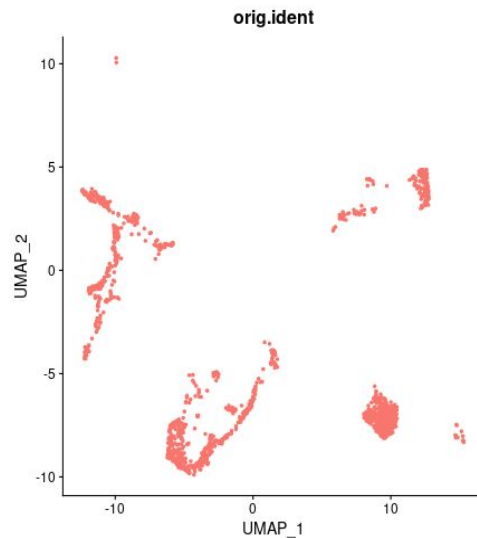
Next step: start getting information

Introduction

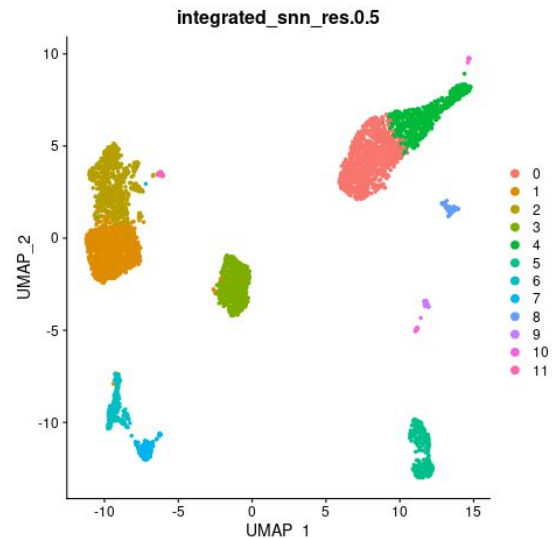
Starting point: normalized, reduced **individual** matrices

Next step: start getting information

→ Visualize the cells



→ Understand what is in the samples (clustering)



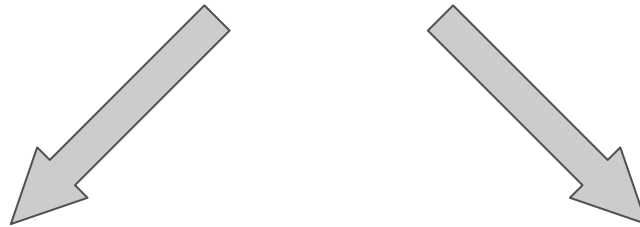
Introduction

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

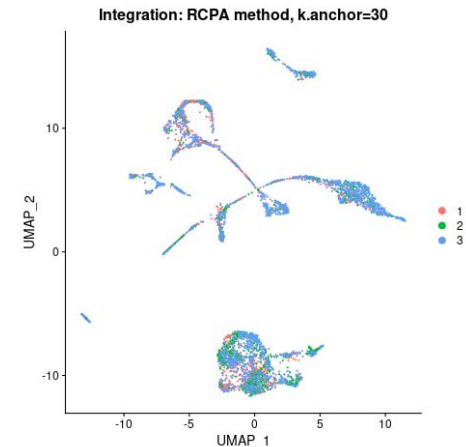
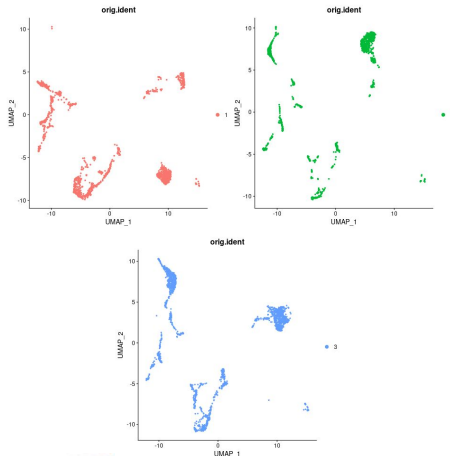
Next step: start getting information

But should we do that

on individual
samples?



On all samples
together ?



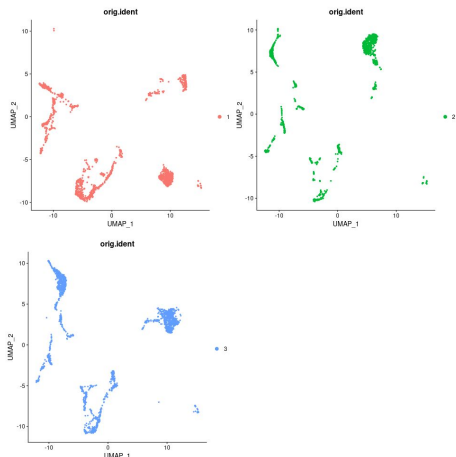
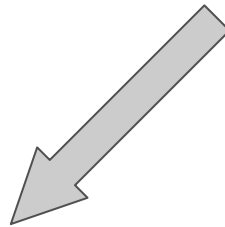
Introduction

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

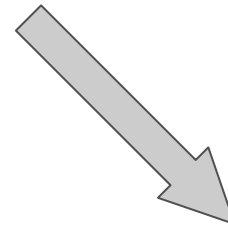
But should we do that



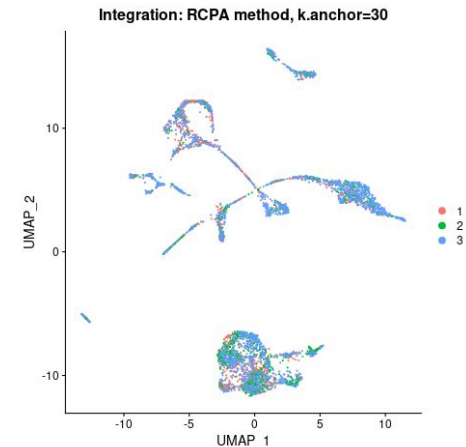
- Allows to work across multiple samples.
- Particularly important for cell populations visualization and identification
- Many cells : helps identifying rare populations



- Overcorrection?



On all samples together ?



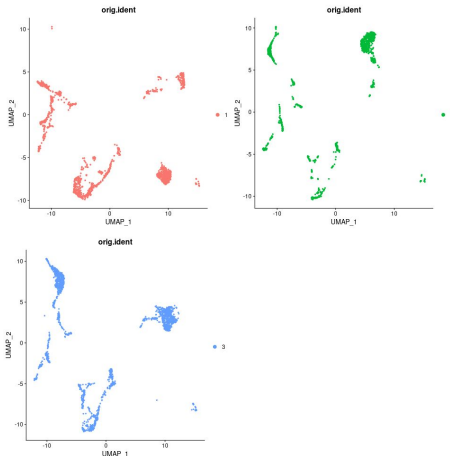
Introduction

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

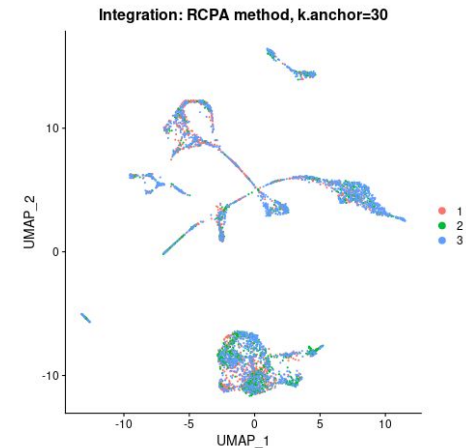
Next step: start getting information

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

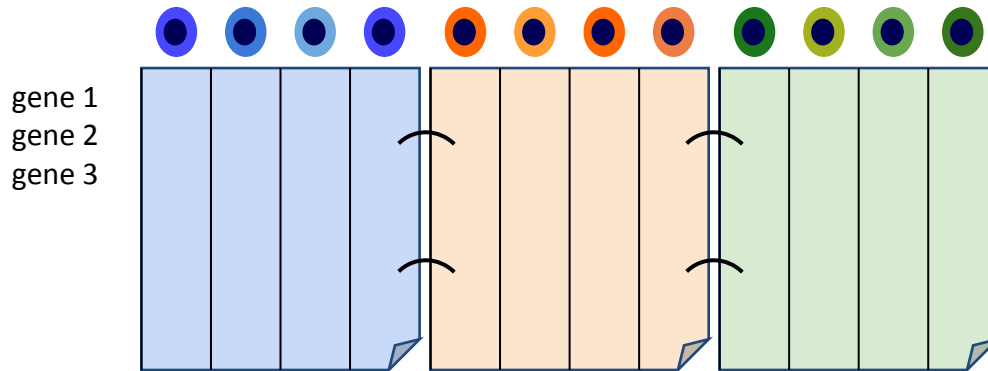
Introduction

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



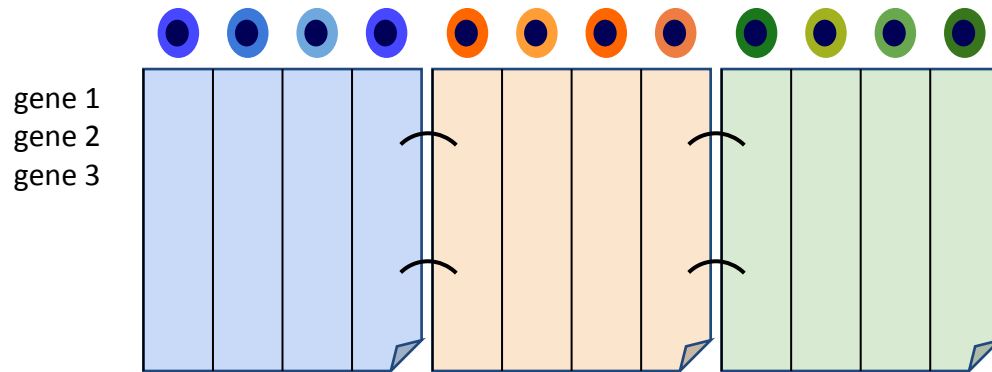
Introduction

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

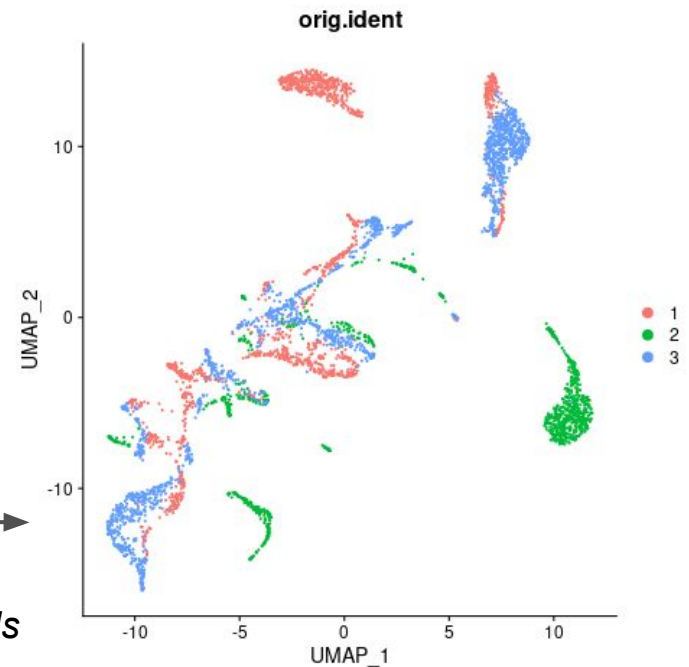
Next step: start getting information

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

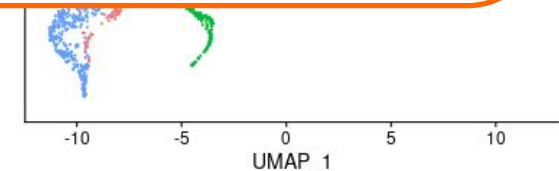
Problem: simple matrix concatenation does not always work

gene 1
gene 2
gene 3

This is a problem of batch effect.

We need a more sophisticated **integration method**

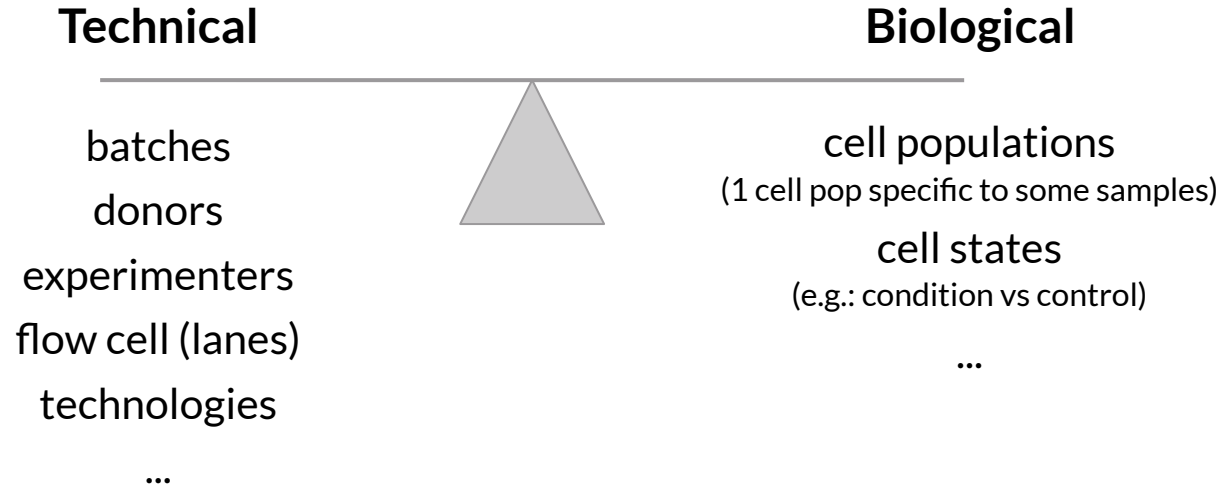
orig.ident



Variability across samples

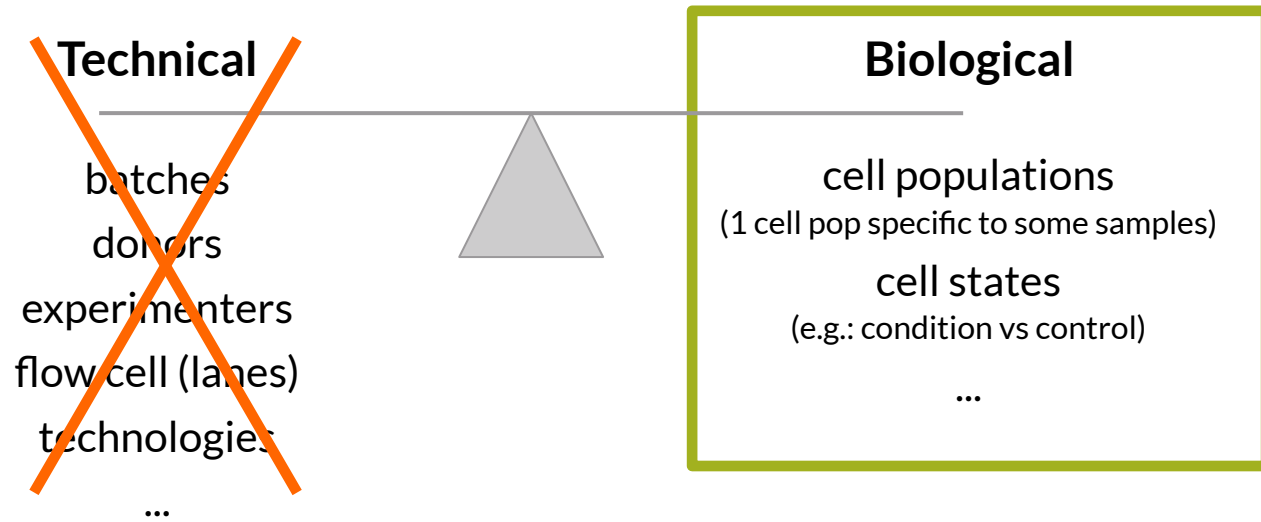
Variability across samples

2 sources of variability across samples



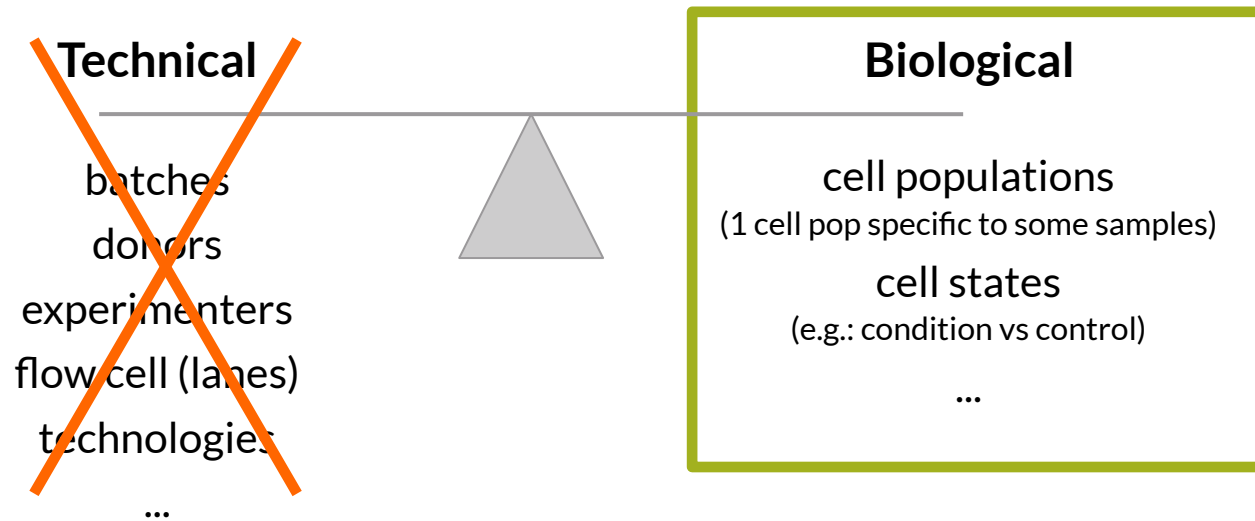
Variability across samples

2 sources of variability across samples



Variability across samples

2 sources of variability across samples



→ Solutions:

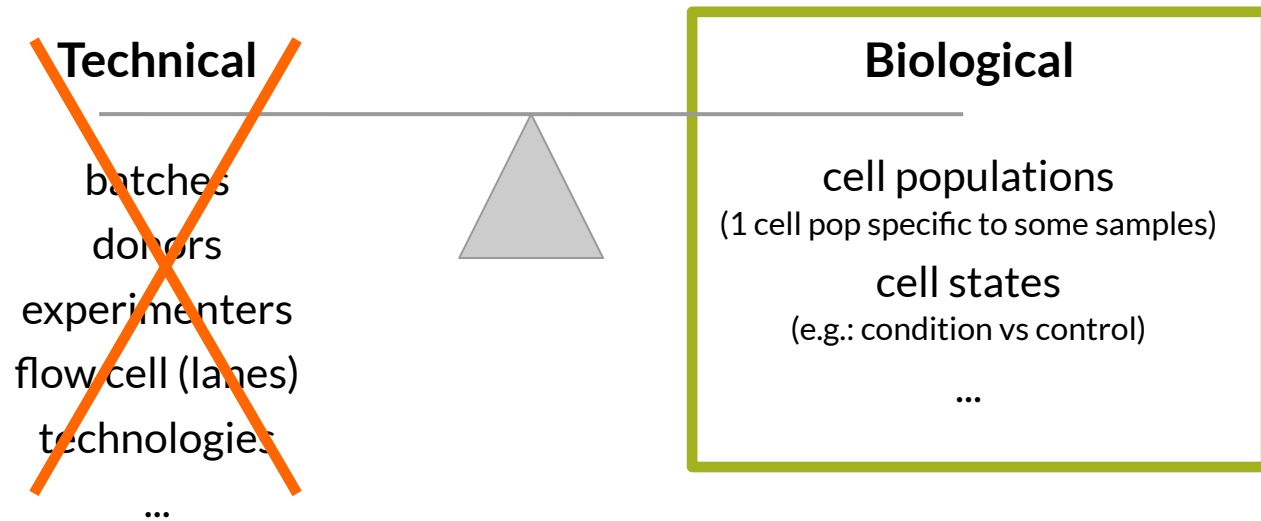
Strategies to avoid factors causing batch effect in the lab

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

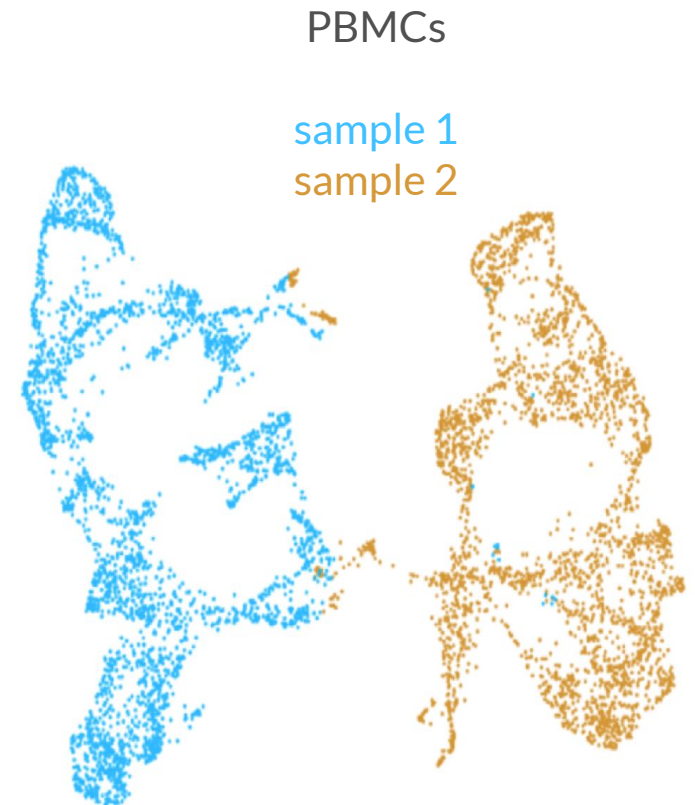
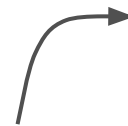
When to integrate

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

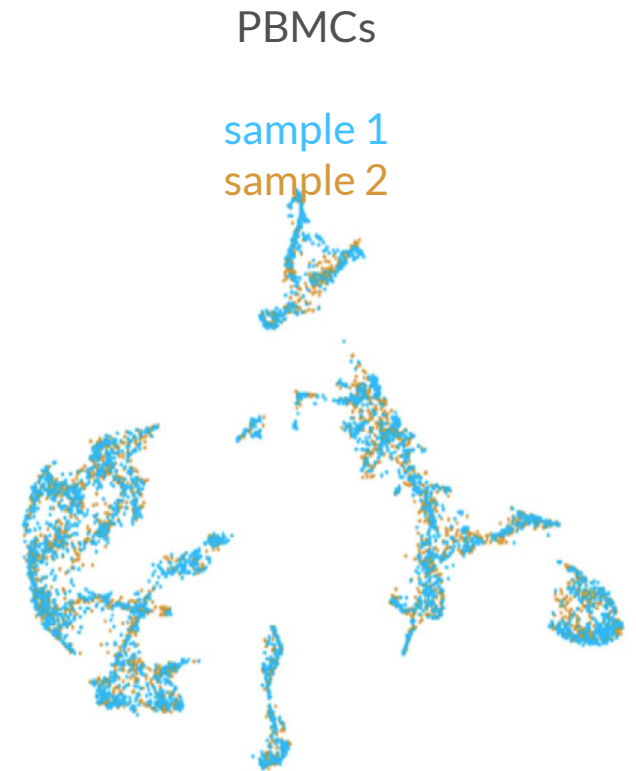
The samples need integration to align cell types/clusters and then identify them correctly



<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



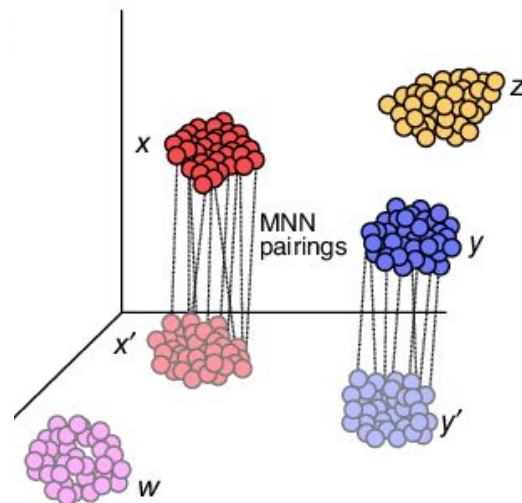
<https://www.10xgenomics.com>

Integration with Seurat

Integration with Seurat

Many methods

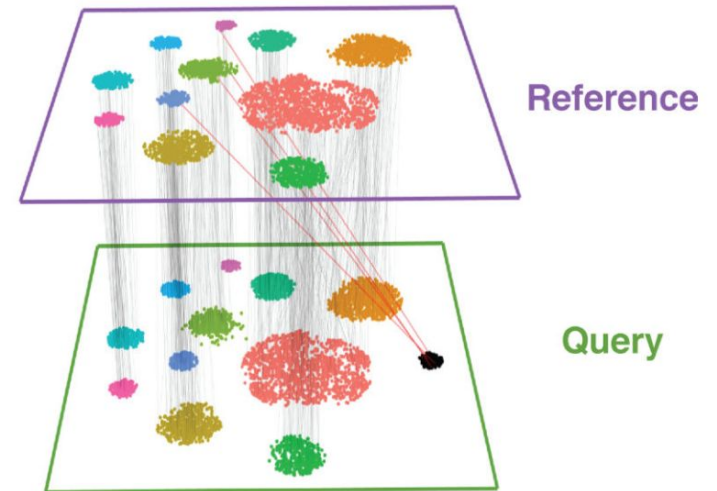
- Over 49 methods (Luecken et al., Nat Methods 2022)
- Seurat integration: group of **similarity-based** methods (most methods)



Integration with Seurat

Principle

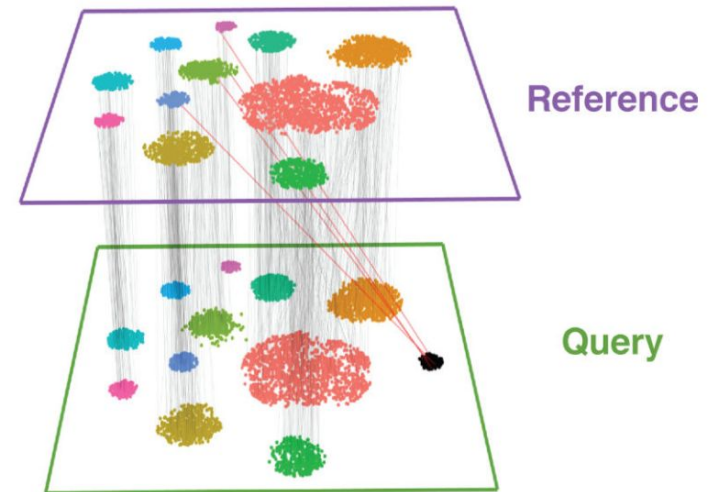
- Integration is always **pairwise**: correct a sample, **the query** to match the expression data of another sample, **the reference**.



Integration with Seurat

Principle

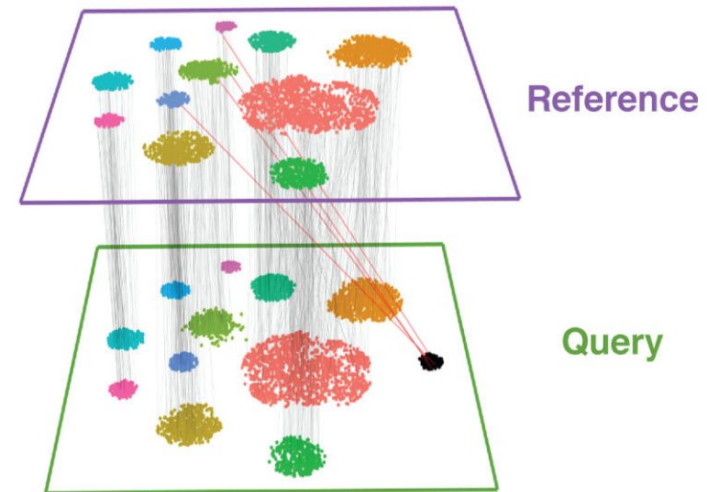
- Integration is always **pairwise**: correct a sample, **the query** to match the expression data of another sample, **the reference**.
- Seurat identifies pairs of close cells across both datasets, the **anchors** (Mutual Nearest Neighbors).



Integration with Seurat

Principle

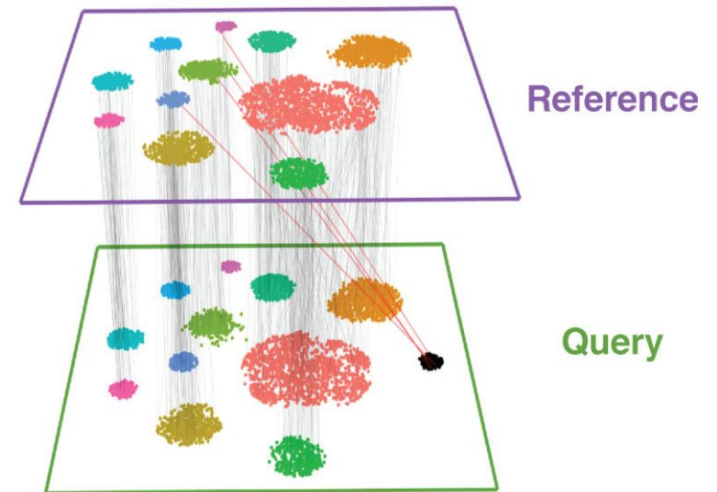
- Integration is always **pairwise**: correct a sample, **the query** to match the expression data of another sample, **the reference**.
- Seurat identifies pairs of close cells across both datasets, the **anchors** (Mutual Nearest Neighbors).
- The difference between them is used to compute a **correction**.



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

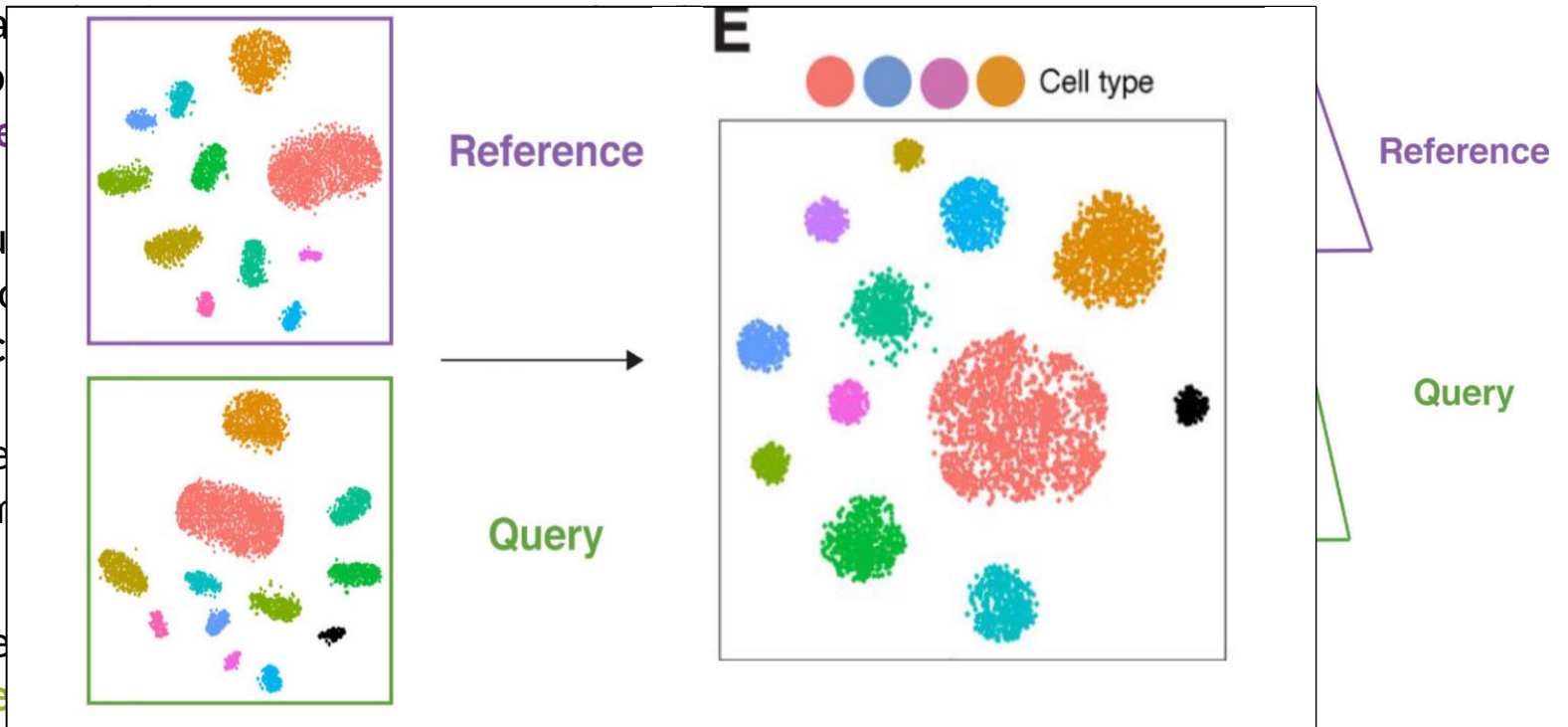
- Integration is always **pairwise**: correct

a sa
exp
refe

- Seu
acro
anc

- The
com

- The
que



Integration with Seurat

Principle

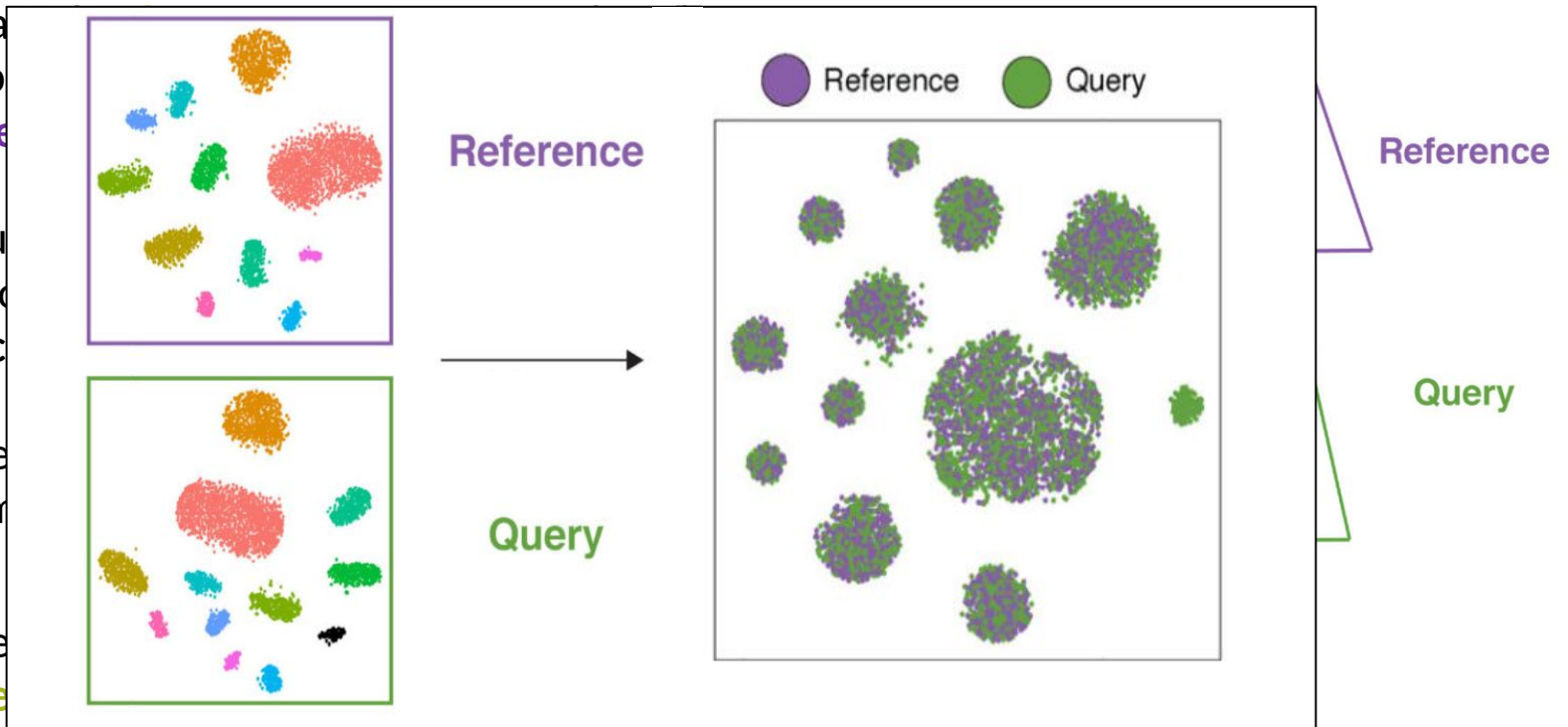
- Integration is always **pairwise**: correct

a sa
exp
refe

- Seu
acro
anc

- The
com

- The
que



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

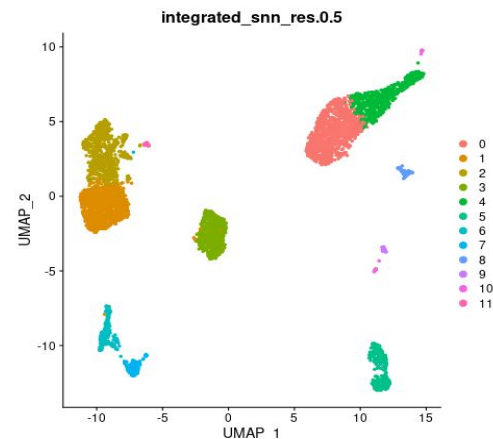
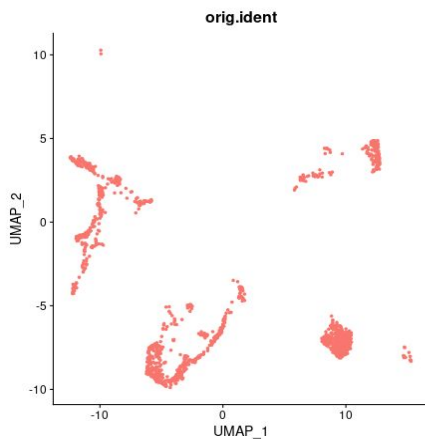


Luecken *et al.*, Nature Methods 2022

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**



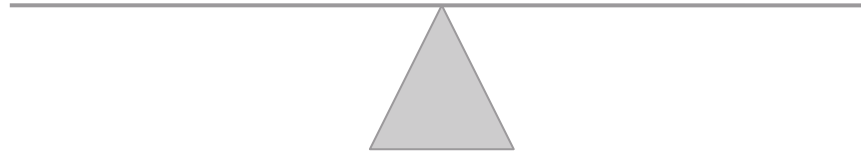
- For differential expression analysis, we go back to raw data

Conclusion

A good integration method

Technical

Biological



- Corrects for technical variability:

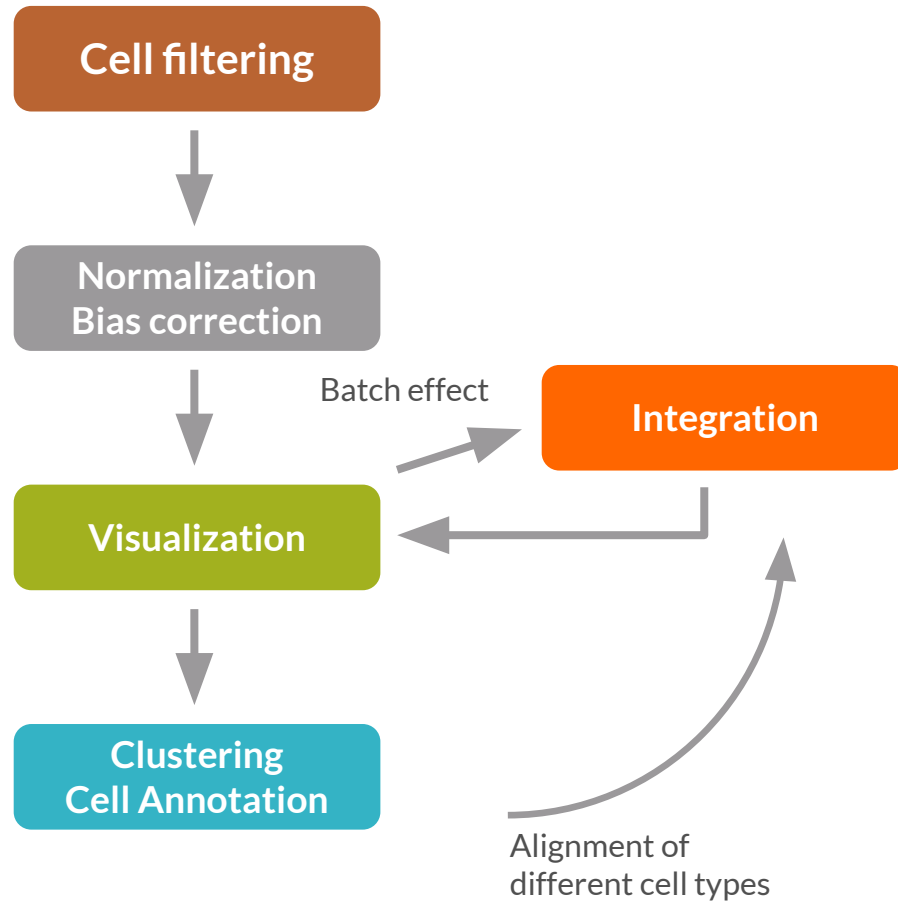
- samples
- donors
- experimenter
- technologies

- 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

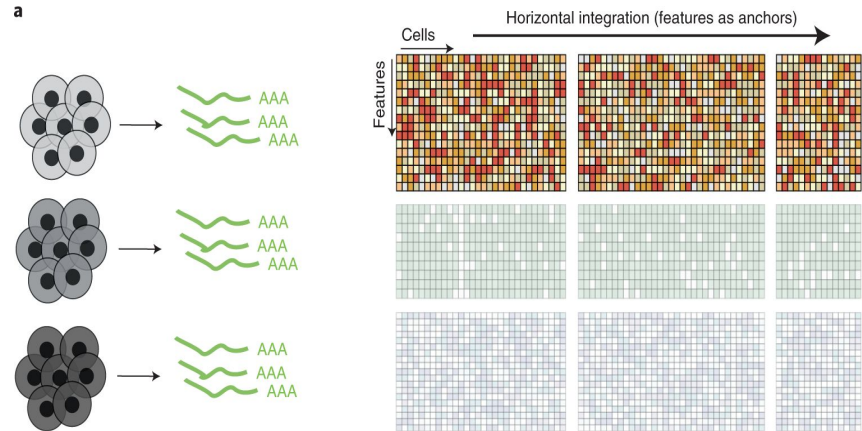


Conclusion

Different types of integrations

- Horizontal: different samples same modality

We saw horizontal integration



Luecken *et al.*, Nat Met 2021

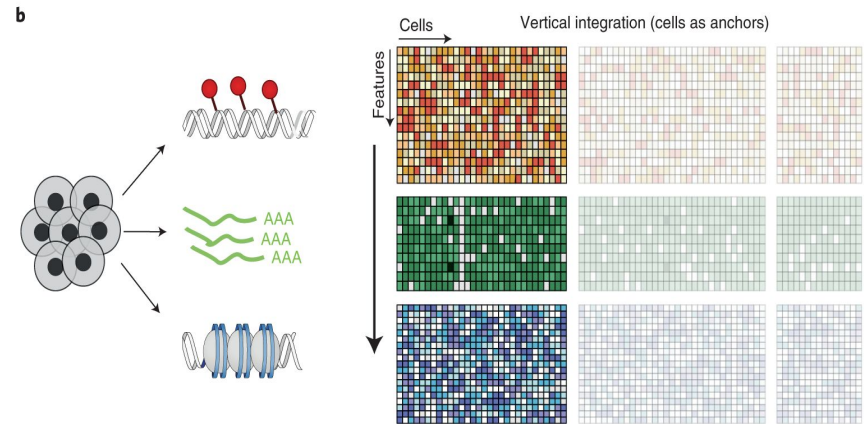
Conclusion

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

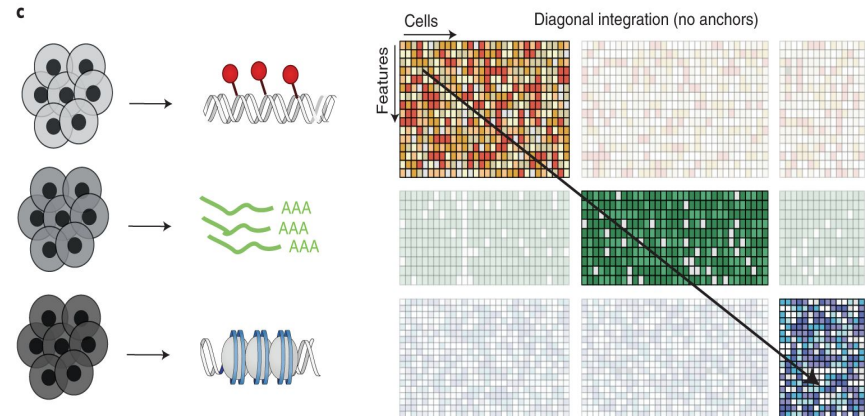
Conclusion

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](#)

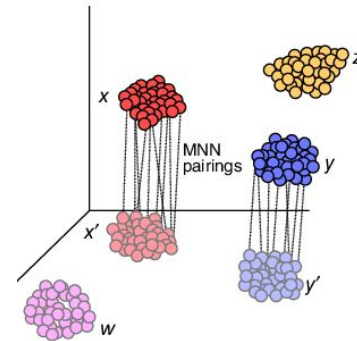
Integration with Seurat

Many methods

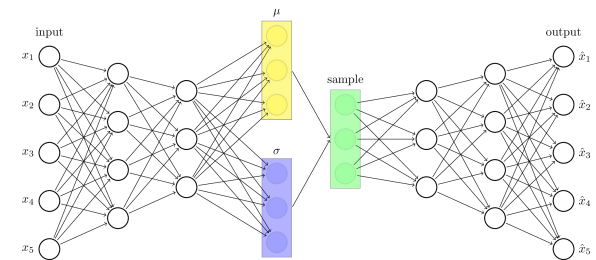
1. Linear decomposition methods

$$A_{m \times n} \approx B_{m \times k} \times C_{k \times n}$$

2. similarity-based (in reduced dimension space)



3. Deep learning



Integration with Seurat

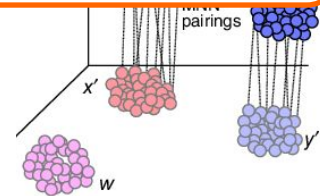
Many methods

1. Linear decomposition methods

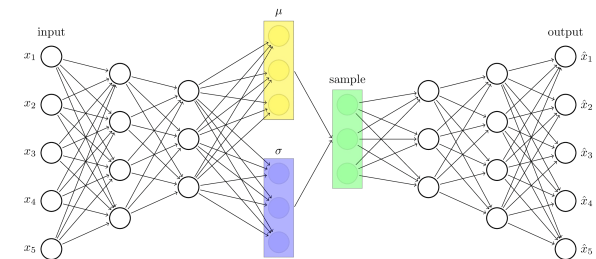
$$\begin{matrix} A \\ m \times n \end{matrix} \approx \begin{matrix} B \\ m \times k \end{matrix} \times \begin{matrix} C \\ k \times n \end{matrix}$$

Over 49 methods (Luecken et al., Nat Methods 2022)

2. similarity-based (in reduced dimension space)



3. Deep learning



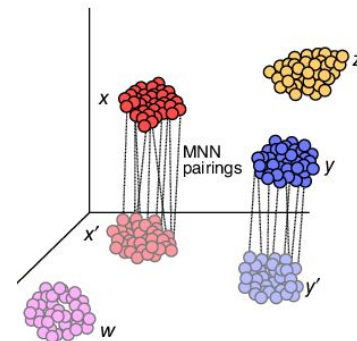
Integration with Seurat

Many methods

1. Linear decomposition methods

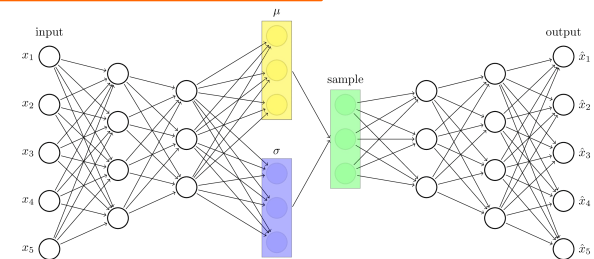
$$A_{m \times n} \approx B_{m \times k} \times C_{k \times n}$$

2. similarity-based (in reduced dimension space)



Seurat

3. Deep learning

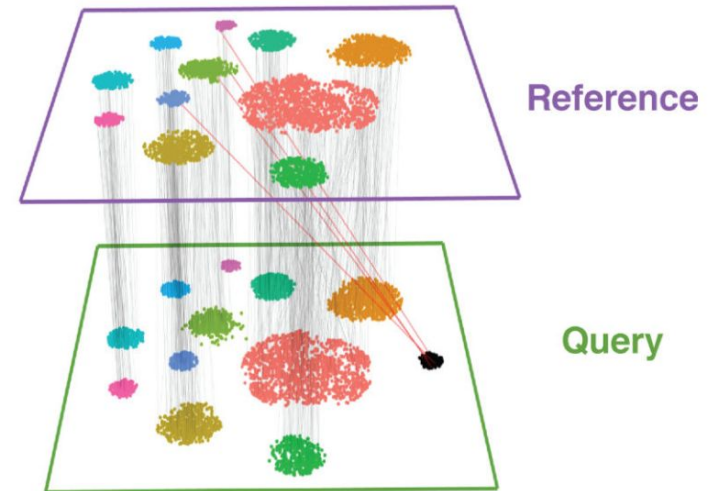


Integration with Seurat

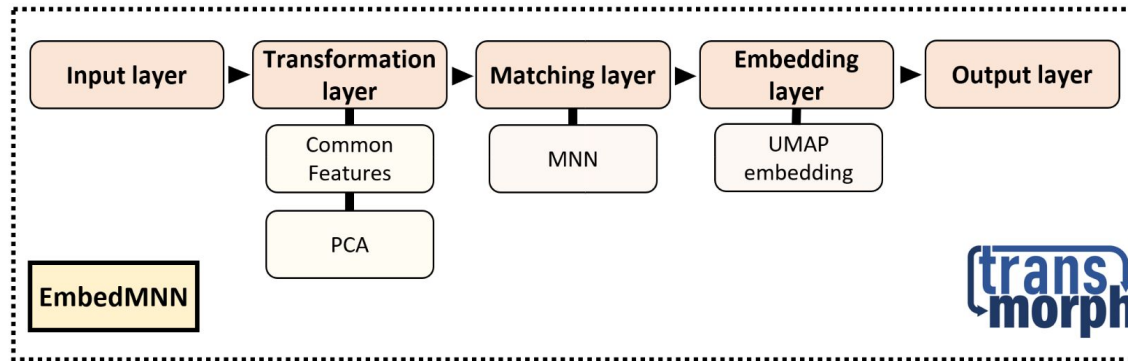
Dimension reduction

- Integration is performed in low dimension space.
- The reduction method is an important parameter

	focus on	when	limits
CCA	Finding highly variable genes between samples	Dataset with many differences	Can overcorrect biological signal
RPCA	Telling signal and noise apart from each other	Less different datasets, huge datasets	Can fail to align populations perfectly
LSI	Identify latent structure of texts (here DNA sequences)	scATAC-Seq	

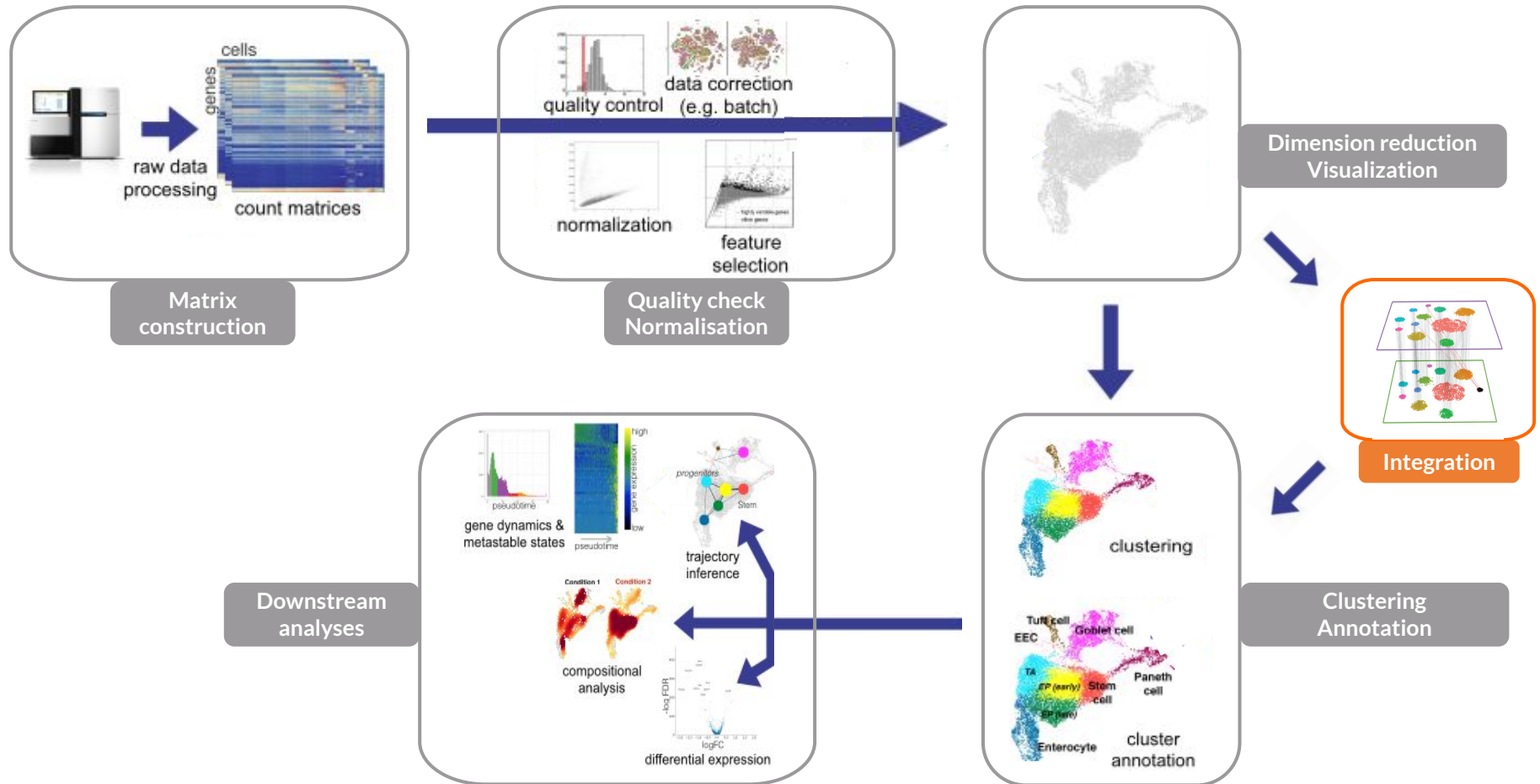


Benchmarking methods



- Developed by A. Fouché, L. Chadoutaud, A. Zinovyev in U900
- Framework breaking down integration algorithms into building blocks
- Allow to combine the building blocks into personalized integration workflow
- Databank for benchmarking

Introduction



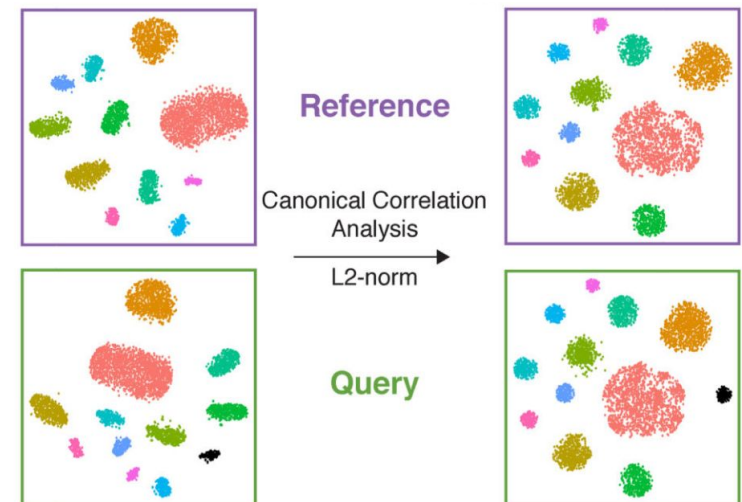
Integration with Seurat



Algorithm

1) Dimension reduction

- Like PCA or UMAP, it projects the cells into a lower space
- The dimension reduction methods used here roughly align similar cells.



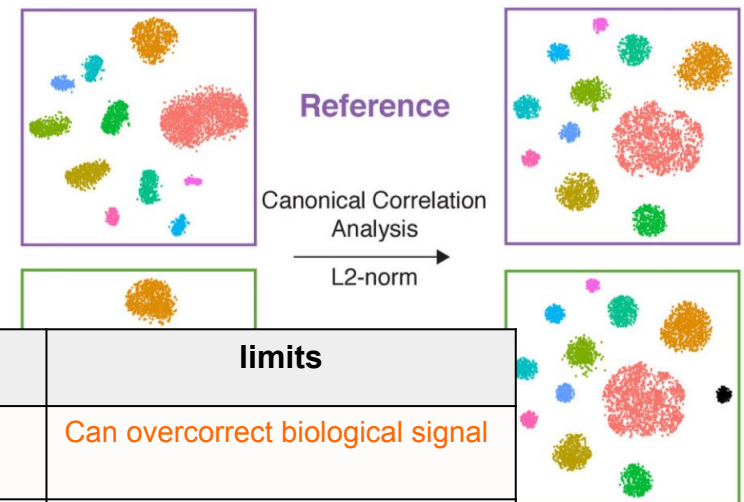
Integration with Seurat



Algorithm

1) Dimension reduction

- Like PCA or UMAP, it projects the cells into a lower space
- The dimension reduction methods used here roughly align similar cells.

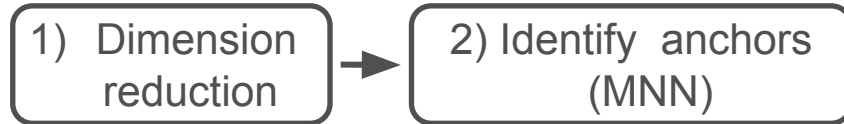


	focus on	when	limits
CCA	Finding highly variable genes between samples	Dataset with many differences	Can overcorrect biological signal
RPCA	Telling signal and noise apart from each other	Less different dataset, huge datasets	Can fail to align populations perfectly
LSI	Identify latent structure of texts (here DNA sequences)	scATAC-Seq	

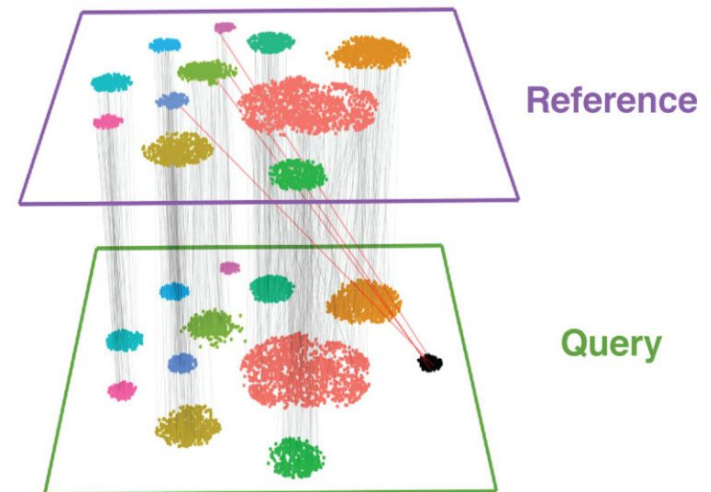
Integration with Seurat



Algorithm



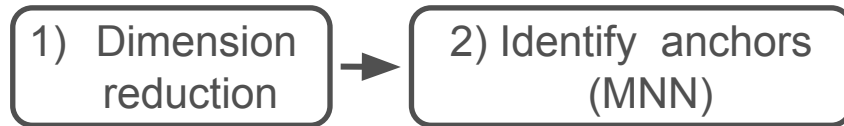
- MNN: Mutual Nearest Neighbors
- In **reference** and **query**, identify 2 cells that are close (neighbors) in terms of euclidean distance: **anchors**
- Identify many anchors



Integration with Seurat



Algorithm



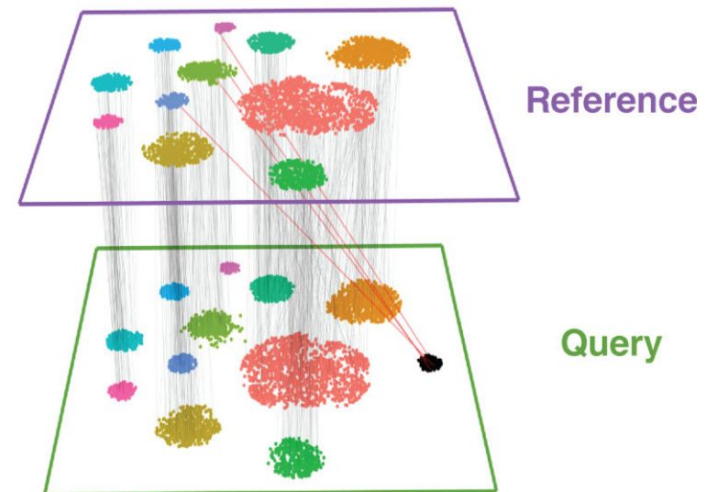
- MNN: Mutual Nearest Neighbors
- In **reference** and **query**, identify 2 cells that are close (neighbors) in terms of euclidean distance: **anchors**
- Identify many anchors
- Note: a cell is represented as a vector

	query	reference
dim 1	3	1
dim 2	5.2	2.1
	.	.
⋮	.	.
	.	.
dim K	4.8	4

Integration with Seurat



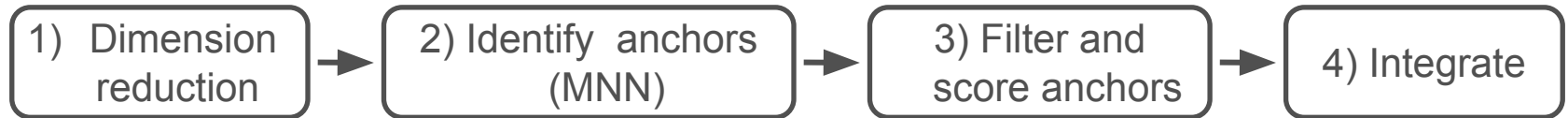
Algorithm



Integration with Seurat

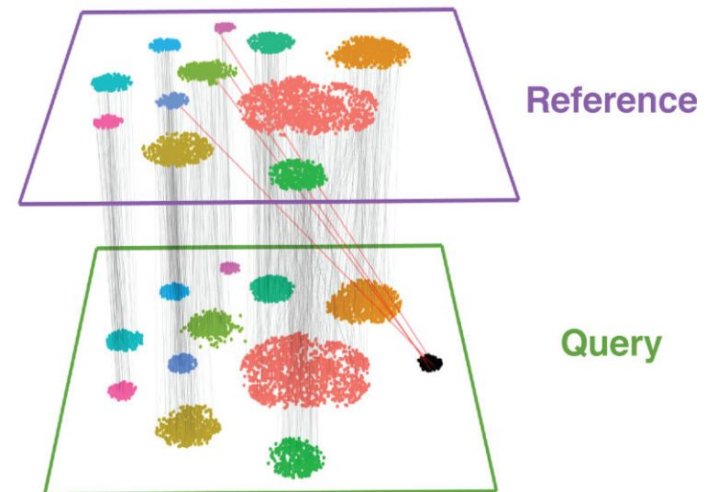


Algorithm



- Deduce correction from anchors

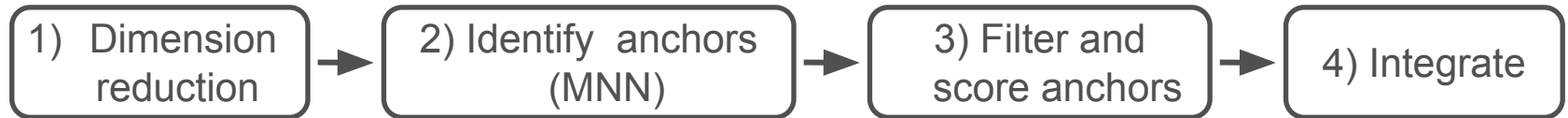
$$\begin{array}{c} \text{correction} \\ \left[\begin{array}{c} 2 \\ 3.1 \\ \vdots \\ 0.8 \end{array} \right] \end{array} = \begin{array}{c} \text{query} \\ \left[\begin{array}{c} 3 \\ 5.2 \\ \vdots \\ 4.8 \end{array} \right] \end{array} = \begin{array}{c} \text{reference} \\ \left[\begin{array}{c} 1 \\ 2.1 \\ \vdots \\ 4 \end{array} \right] \end{array}$$



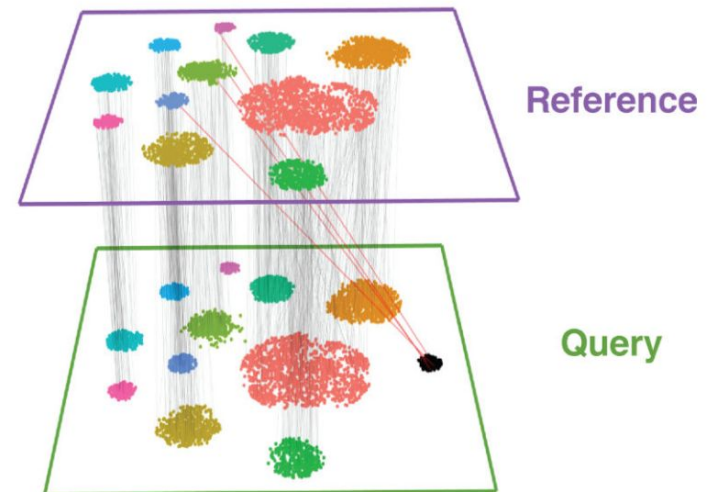
Integration with Seurat



Algorithm



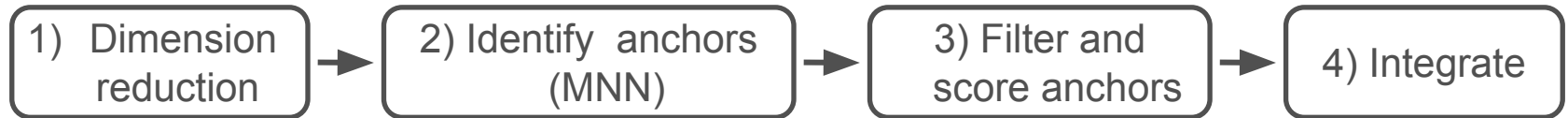
- Deduce correction from anchors
- Apply correction vector to all query cells.



Integration with Seurat



Algorithm



- Deduce correction from anchors
- Apply correction vector to all query cells.

