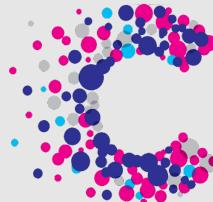




UNIVERSITY OF
CAMBRIDGE



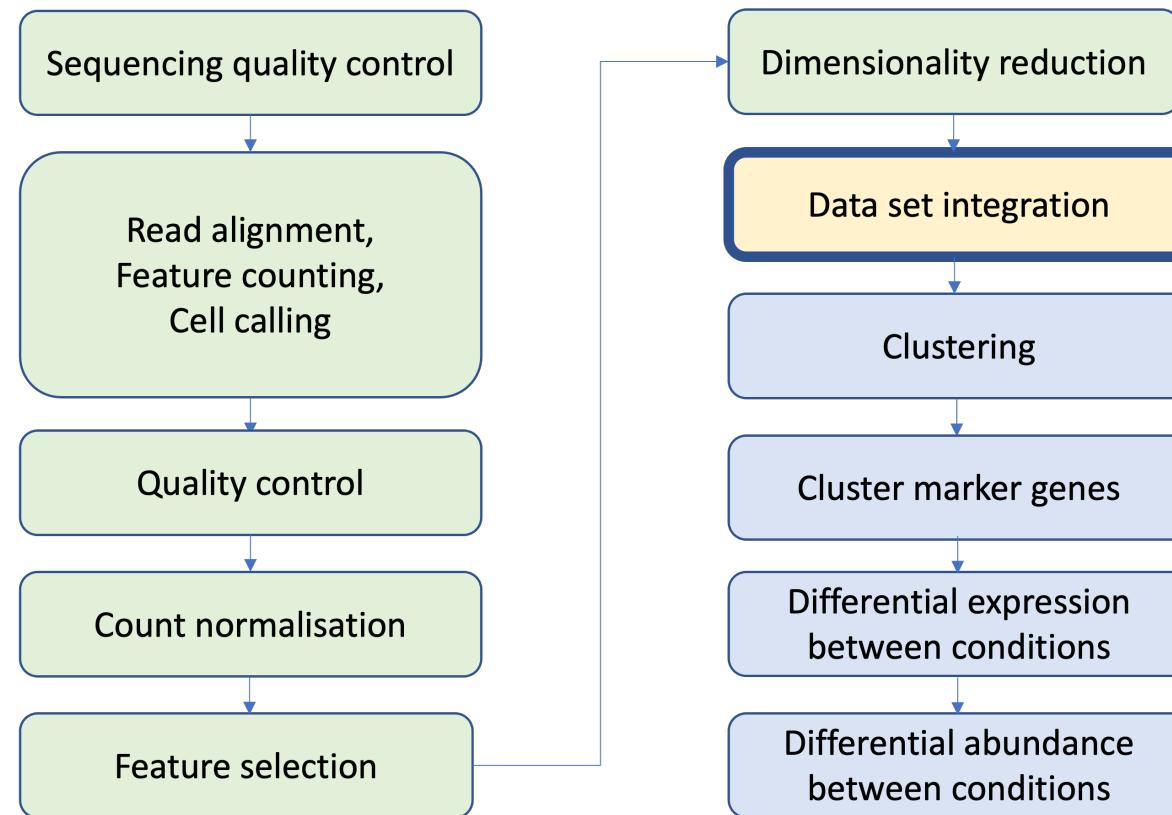
CANCER
RESEARCH
UK

CAMBRIDGE
INSTITUTE

Data Integration and Batch Correction

September 2022

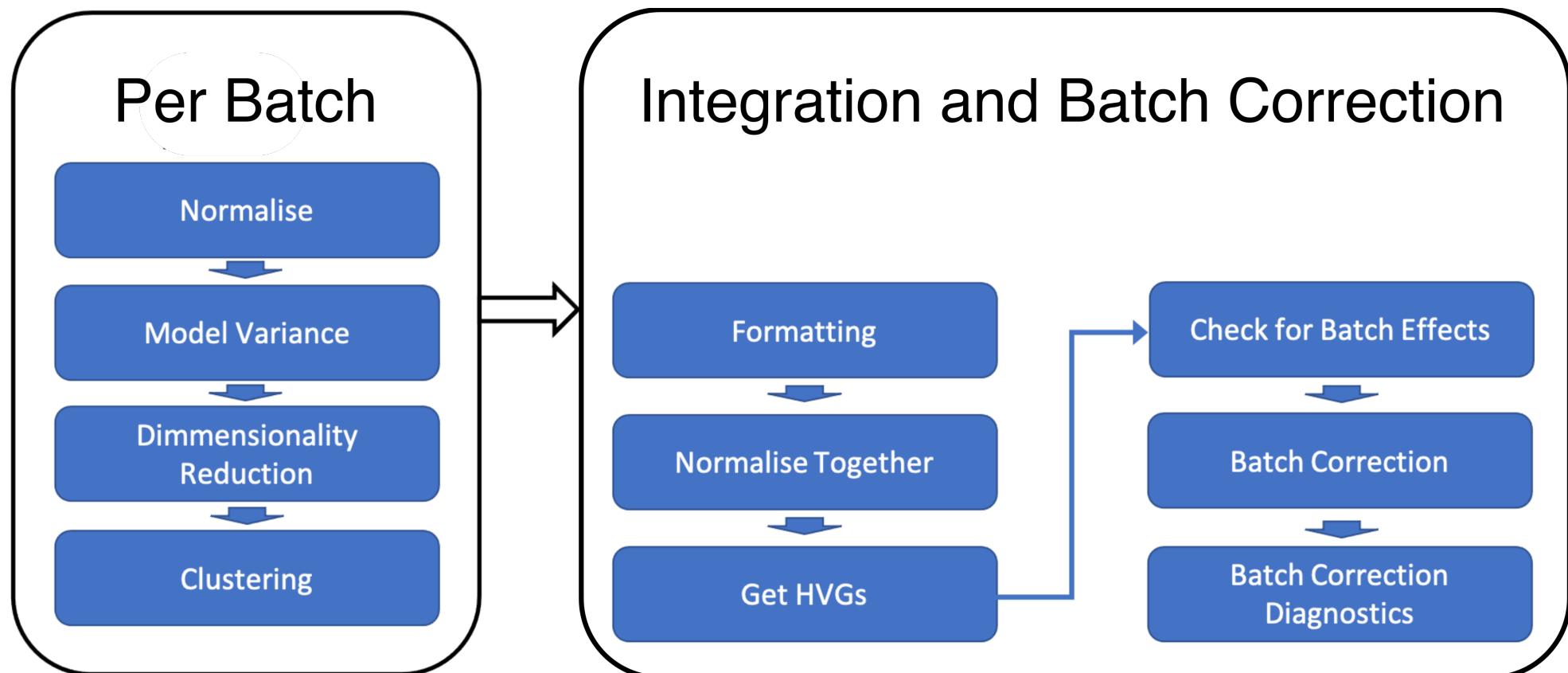
Single Cell RNAseq Analysis Workflow



Why do we need to think about data integration?

- Practicalities of our Experimental Design
- Different 10X runs at different times OR just the same sample run twice
- Obscure real biological changes

Data Integration Workflow



Formatting our data

A few ways our data can be arranged (software-dependent too)

- one large SCE object containing many samples
- many single-sample SCE objects, QC'd in isolation
- multiple large SCE objects with multiple samples

Important we make sure things match up

- Different bioconductor versions
- Different analysts may have formatted things differently

Cellranger aggr

A useful quick look

The screenshot shows a support page from 10X Genomics. At the top, there's a navigation bar with links for Products, Research Areas, Resources, Support, and Company. Below the navigation, a breadcrumb trail shows Support > Single Cell Gene Expression > Software. A search bar and a contact support link are also present. The main content is titled "Setting Up an Aggregation CSV". It instructs users to create a CSV file with a header line containing columns for sample_id and molecule_h5. It provides examples for three samples (LV123, LB456, LP789) with their respective molecule_info.h5 paths. Below this, it says you can either make the CSV in a text editor or Excel. An example Excel spreadsheet is shown with columns A and B. Finally, it shows the resulting CSV text:

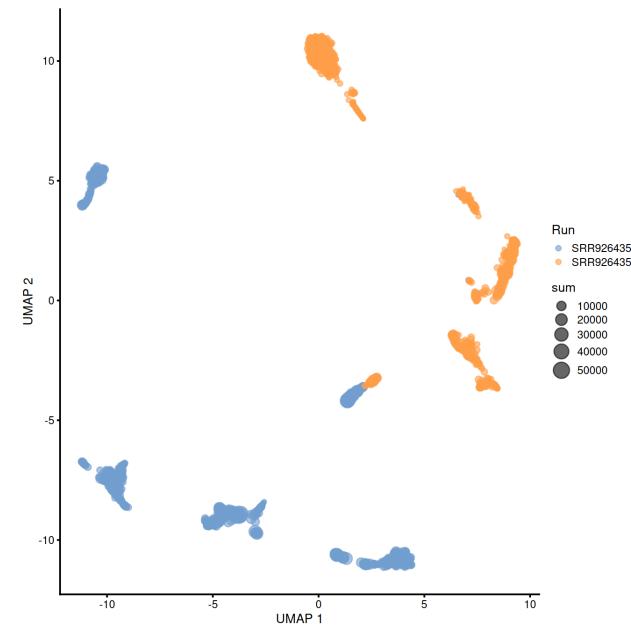
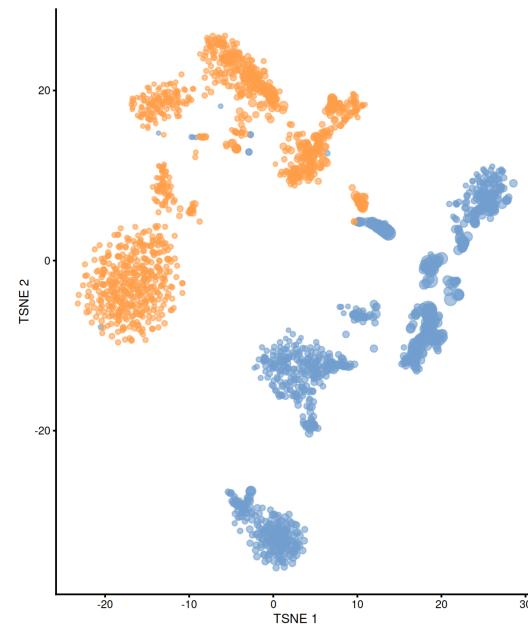
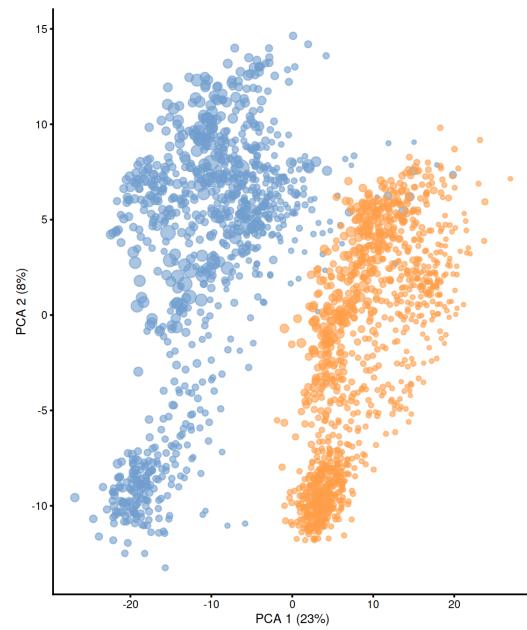
	A	B
1	sample_id	molecule_h5
2	LV123	/opt/runs/LV123/outs/molecule_info.h5
3	LB456	/opt/runs/LB456/outs/molecule_info.h5
4	LP789	/opt/runs/LP789/outs/molecule_info.h5

When you save it as a CSV, the result would look like this:

```
sample_id,molecule_h5
LV123,/opt/runs/LV123/outs/molecule_info.h5
LB456,/opt/runs/LB456/outs/molecule_info.h5
LP789,/opt/runs/LP789/outs/molecule_info.h5
```

Cell Ranger v6.0 (latest)

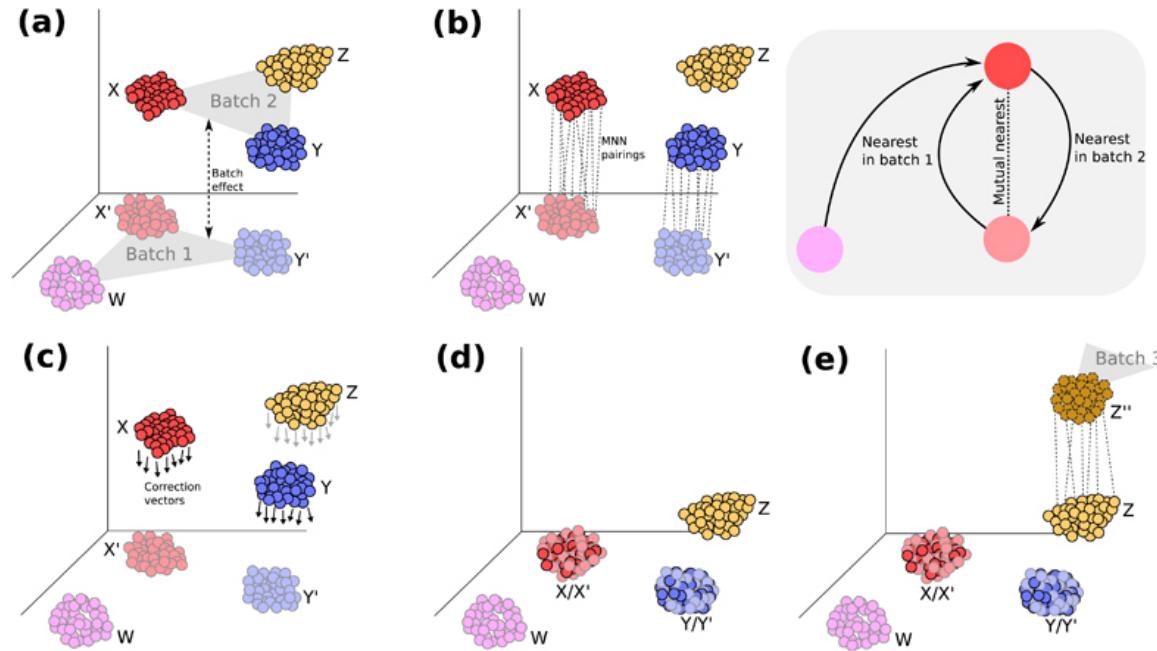
Checking for batch effects



Batch Corrections

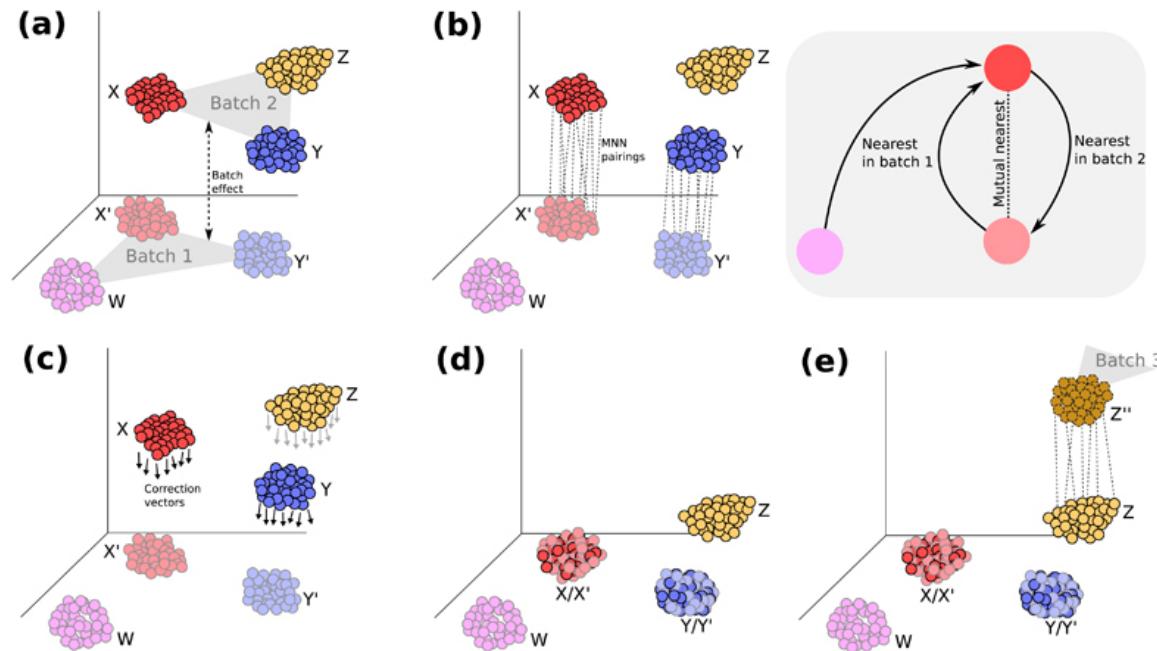
- Gaussian/Linear Regression - removeBatchEffect (limma), comBat (sva), rescaleBatches or regressBatches (batchelor)
- Mutual Nearest Neighbours (MNN) correction - [Haghverdi et al 2018](#)
 - mnnCorrect (batchelor)
 - FastMNN (batchelor)
- And [many more!](#)
 - Different methods may have strengths and weaknesses
 - [Benchmark studies](#) can be used as a reference to choose suitable method

FastMNN ([Haghverdi et al 2018](#))



1. Perform a multi-sample PCA on the (cosine-)normalized expression values to reduce dimensionality.
2. Identify MNN pairs in the low-dimensional space between a reference batch and a target batch.
3. Remove variation along the average batch vector in both reference and target batches.
4. Correct the cells in the target batch towards the reference, using locally weighted correction vectors.
5. Merge the corrected target batch with the reference, and repeat with the next target batch.

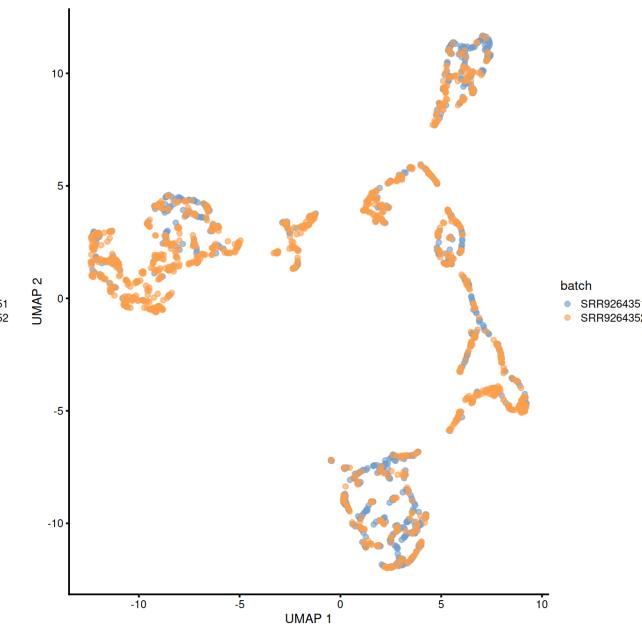
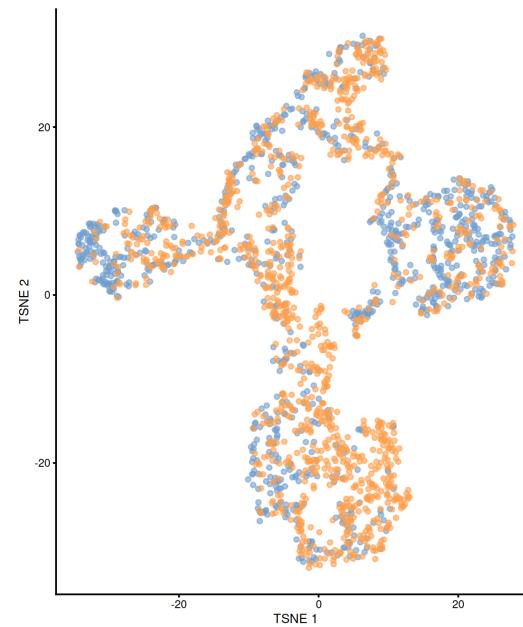
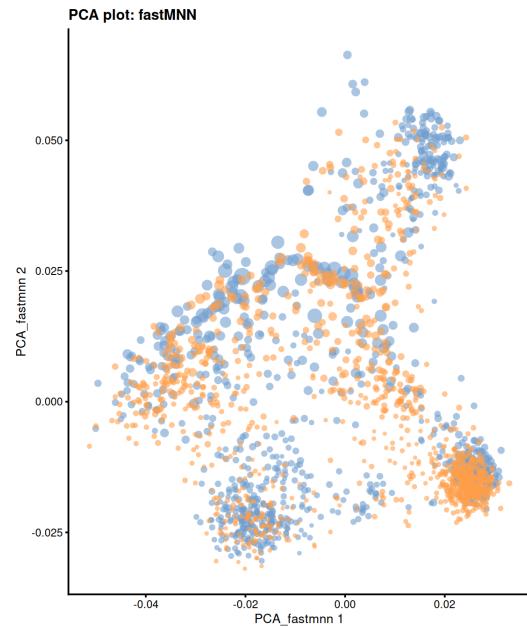
FastMNN ([Haghverdi et al 2018](#))



Assumptions (quoted from the paper):

1. There is at least one cell population that is present in both batches,
2. the batch effect is almost orthogonal [i.e. uncorrelated] to the biological subspace, and
3. the batch-effect variation is much smaller than the biological-effect variation between different cell types

Checking our correction has worked



Checking our correction hasn't over worked

- If you use fastMNN in the absence of a batch effect, it may not work correctly
- It is possible to remove genuine biological heterogeneity
- fastMNN can be instructed to skip the batch correction if the batch effect is below a threshold. You can use the effect sizes it calculates to do this.
- In reality the absence of any batch effect would warrant further investigation.

Using the corrected values

The value in batch correction is that it enables you to see population heterogeneity within clusters/celltypes across batches.

- Also increases the number of cells you have

However the corrected values should not be used for gene based analysis eg. DE/marker detection.

- fastMNN doesn't preserve the magnitude or direction of per-gene expression and may have introduced artificial agreement between batches on the gene level.