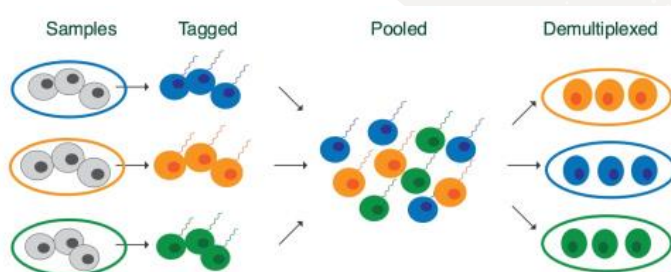# Multiplexing and Normalization of single cell data

Associate Professor Joseph Powell
Director, Garvan-Weizmann Centre for Cellular Genomics
Deputy Director, UNSW Cellular Genomics Futures Institute
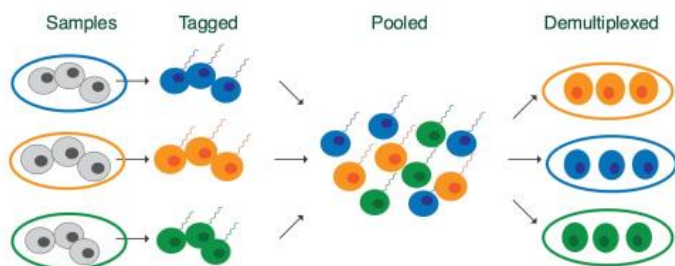
SISG - 2019

---

# What is multiplexing?

- Pooling cells from multiple samples before loading into a single cell capture

# What is multiplexing?

- Pooling cells from multiple samples before loading into a single cell capture



- Lower costs for library preparation
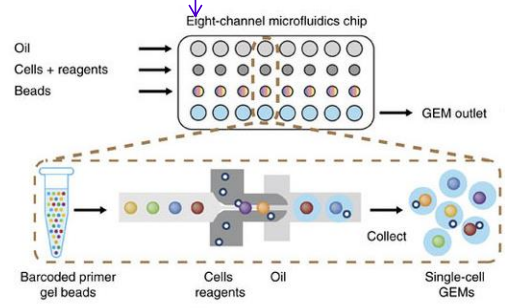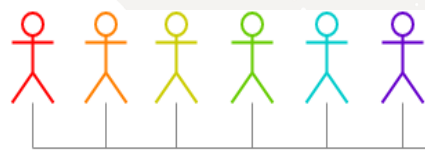- Better doublet identification

# Cell multiplexing methods

- SNP calling and matched genotyping (natural genetic barcoding)

- SNP calling without matched genotypes (natural genetic barcoding)

- Antibody hash tagging

- Lipid- and cholesterol-modified oligonucleotides

4

# Pooling cells for library prep

Garvan Institute
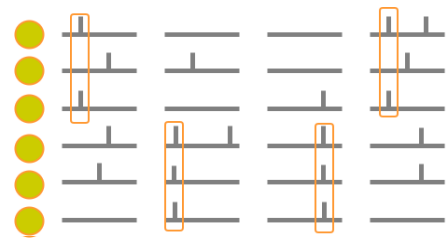of Medical Research

- SNPs
- Anti Bodies
- Lipids

Pooled sample loaded
onto 10X Chromium

Eight-channel microfluidics chip

Oil
Cells + reagents
Beads

GEM outlet

Collect

Barcoded primer
gel beads

Cells
reagents

Oil

Single-cell
GEMs

5

# Multiplexing samples for single cell library prep

Garvan Institute
of Medical Research

Call SNPs from the 3' reads

# Multiplexing samples for single cell library prep

Garvan Institute
of Medical Research

# Multiplexing samples for single cell library prep

Garvan Institute
of Medical Research

A
A
A
a

aa
aa

Aa
Aa
AA

aa
Aa
AA

## Multiplexing samples for single cell library prep



$$P^{(g)}_{SNP} = Pr(g \mid Array_{SNP})$$

Powell et al. Nature Reviews Genetics 2010

## Demuxlet – Kang et al. Nat Biotech 2018



https://www.protocols.io/view/instructional-tutorial-for-using-demuxlet-233gggn

https://github.com/statgen/demuxlet.

10

**Cell Hashing with barcoded antibodies enables multiplexing and doublet detection for single cell genomics**

11 Stoeckius et al. Genome Biology



# scSplit - Genotyping Free Demultiplexing

12

# scSplit - Genotyping Free Demultiplexing



13

# Detecting doublets



Each color comes from a different individual/sample

Demuxlet seems to over estimate doublets

14

# What genotyping to use?

**Garvan Institute** of Medical Research

## Demuxlet optimisation and comparison to popscle

| Assignment | Demuxlet GP | Demuxlet GT | Demuxlet GT (Exon-only) | Popscle GP | Popscle GT | Popscle GT (Exon-only) |
|---|---|---|---|---|---|---|
| AMB | 0 | 1 | 0 | 0 | 1 | 0 |
| DBL | 3726 | 3670 | 4204 | 3691 | 3433 | 3770 |
| SNG | 9517 | 9572 | 9039 | 9552 | 9809 | 9473 |

- Using all imputed SNPs gives the best results
- Genotype-only SNPs still fare better than exon-only filtered imputed SNPs
- Genotype-only SNP runs aren't computationally intensive - need to weigh up benefits of using imputed data

15

# *Scrublet*

**Garvan Institute** of Medical Research

## Identifying doublets with *scrublet*

- *Scrublet*: Identifies neotypic multiplets from scRNA-seq transcriptome data
  - **Publication**: DOI:https://doi.org/10.1016/j.cels.2018.11.005
  - **Website: https://github.com/AllonKleinLab/scrublet**
- Comparing to demuxlet assignments (GT)
  - 9741 agreements, 3502 disagreements
- *Scrublet* cannot assign an individual to a droplet, but we can possibly use the most likely candidate from the *demuxlet* results.

| Software | Number of singlets | Number of doublets |
|---|---|---|
| **Demuxlet** | 9572 | 3671 |
| **Scrublet** | 12818 | 425 |

16

## Combined approach

### Doublet filtering using *demuxlet* and *scrublet*

POAG_scRNA Sample 1
Scrublet Assignments for Demuxlet Doublets



SCRUBLET_PREDICTION   • DBL   • SNG

- Demuxlet is quite certain a cell is a singlet.
- Conversely, it says a cell is just as likely to be a singlet as a doublet.
- Scrublet seems to think the majority of cells are singlets.
- Set threshold high - class as singlets if:
  - Singlet PP >= 0.95
  - Scrublet calls as singlet
  - Assigned to individual with highest singlet PP
- This recovers 3,099 cells
  - **Includes 137 cells from WAB-00069**

Doublet filtering using *demuxlet* and *scrublet*

| Droplet Type | Before | After |
|---|---|---|
| Singlet | 9,517 | 12,616 |
| Doublet | 3,726 | 627 |

- Can we trust this?

17

## No change in cluster identification

### Clusters characterised by Seurat

**Sample 1 Clusters**



- Cells filtered with *ascend*
- Normalisation with Seurat v3's SCtransform
- Clustering with Seurat v3's clustering
- Previous slide shows multiplets identified by demuxlet are clustering in a location that corresponds to clusters 5, 9 and 14.
- Removed multiplets identified by *demuxlet*

18

## Effects on cost

Garvan Institute
of Medical Research

- Costs for generating 1 library 10x (1x chip, 1x reagent)
  - List price:                    $2,830-3,215
  - Doublet rate

| Multiplet Rate (%) | # of Cells Loaded | # of Cells Recovered |
|---|---|---|
| ~0.4% | ~870 | ~500 |
| ~0.8% | ~1700 | ~1000 |
| ~1.6% | ~3500 | ~2000 |
| ~2.3% | ~5300 | ~3000 |
| ~3.1% | ~7000 | ~4000 |
| ~3.9% | ~8700 | ~5000 |
| ~4.6% | ~10500 | ~6000 |
| ~5.4% | ~12200 | ~7000 |
| ~6.1% | ~14000 | ~8000 |
| ~6.9% | ~15700 | ~9000 |
| ~7.6% | ~17400 | ~10000 |

  - Sequencing to 50,000 reads per cell
    - NovaSeq S4          = $0.25
    - NextSeq             = $0.50

- Cost per sample (3,000 cells) = $2,830+750=$3,580

19

## Effects on cost

Garvan Institute
of Medical Research

- Pool cells from 20 samples
- Aim for 20,000 cells = ~**1,000** cells per sample
- Cost per sample
  - Library prep          =$280
  - SNP Chip             =$47
  - Sequencing          =$150
  - **Total**               =$477

- Run all 8 lanes on a 10x chip     =$427
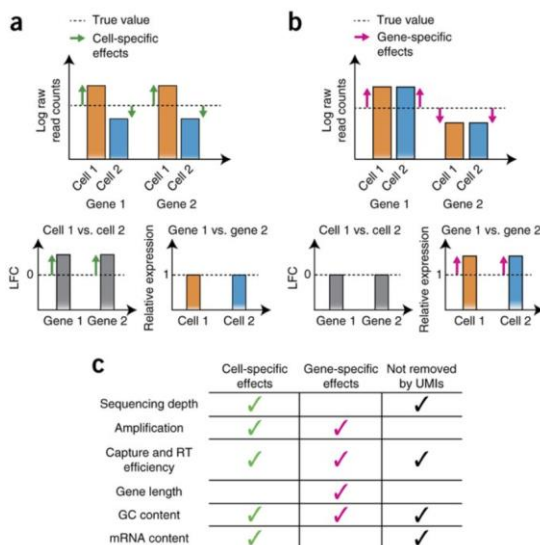- Use BGI sequencer (~$75)        =$352

  - 

20

# Normalization
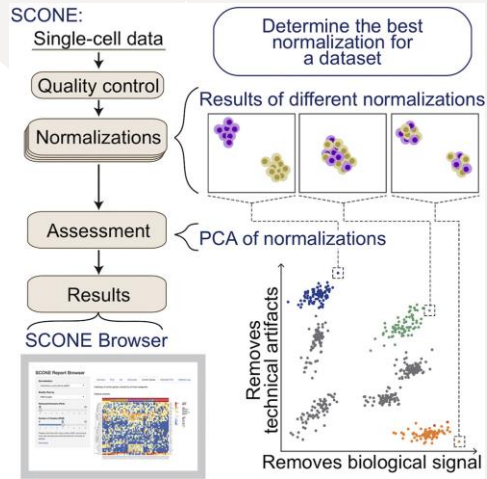
21

---

## Normalization – cell and gene level

This is the next stage of
normalization after the 'batch'
effects have been accounted for



22 Vellejos et al. Nature Methods, 2017

## Performance Assessment and Selection of Normalization Procedures for Single-Cell RNA-Seq

Garvan Institute of Medical Research

Flexible generalized regression method

Implemented in the *scone* package



**23**   Cole et al. Cell Systems 2018

---

## *scone* normalization

Garvan Institute of Medical Research

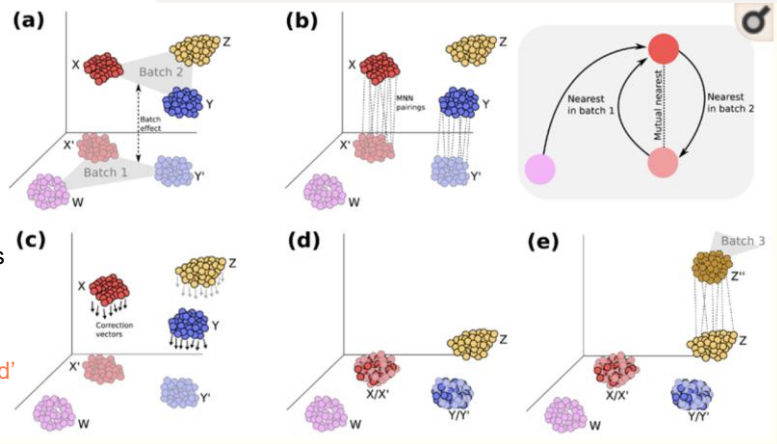$$g(\mathrm{E}[Y\,|\,X, U, W]) = X\beta + U\gamma + W\alpha,$$



**24**

# Normalization by nearest neighbor

Batch correction based on the detection of mutual nearest neighbours (MNNs) in the high-dimensional expression space

Makes the assumption that cells are similar between batches

If there are unique cell types in a sample then they will be 'adjusted' Closer to other cell types.



**25**    Haghverdi et al. Nature Biotechnology 2018

# Resources

- *ascend*:  Senabouth et al. GigaScience 2019 (https://github.com/powellgenomicslab/ascend)

- seurat: Butler *et al*. Nature Biotechnology 2018 (https://github.com/satijalab/seurat)

- scater: McCarthy *et al*. Bioinformatics 2017 (https://bioconductor.org/packages/release/bioc/html/scater.html)

**26**

# QUESTIONS - I

1.  EdgeR or DESeq2 or … ?

2.  Seurat or Monocle or … ?

3.  tSNE or UMAP or … ?

4.  What is the biggest impediment to robust bulk RNA-seq analysis?

5.  What is the biggest impediment to robust scRNA-seq analysis?

# QUESTIONS - II

Garvan Institute
of Medical Research

6. What is the best way to remove Batch effects from scRNA-seq data?

7. Should I aim for more cells or greater read depth?

8. Do I need CITE-seq or other surface protein expression markers?

9. Why are p-values from scRNA-seq comparisons so small?

10.  How do I know if I have screwed up my analysis?

29