

MRC

Human
Genetics
Unit



THE UNIVERSITY
of EDINBURGH

The
Alan Turing
Institute



BASiCS: Bayesian Analysis of Single Cell Sequencing data (+scran)

Catalina Vallejos

Chancellor's Fellow, MRC Human Genetics Unit
Fellow, The Alan Turing Institute

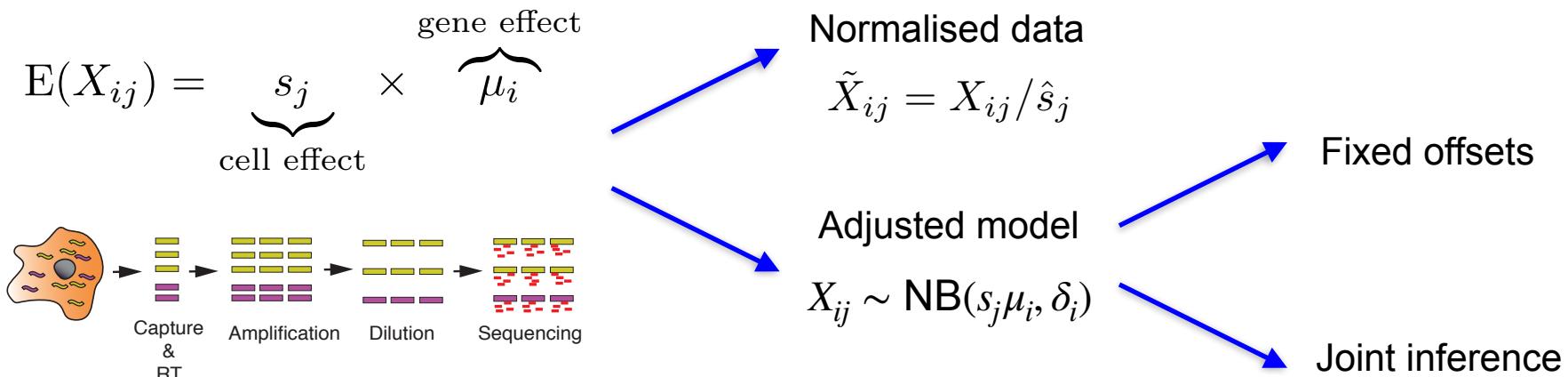


@CataVallejosM

Questions

- Q1.** How do you define normalization? Why is it important?
- Q2.** How do you normalize? (This can be a discussion of your method)
- Q3.** How do you demonstrate success? How do you know if you've normalized properly?
- Q4.** Where does your method break? Or where do you see challenges with your method or others?
- Q5.** What is your suggestion for how we spend the second day of the workshop to collaborate on computational problems? Are there particular tasks or challenges that you think would be good to address?

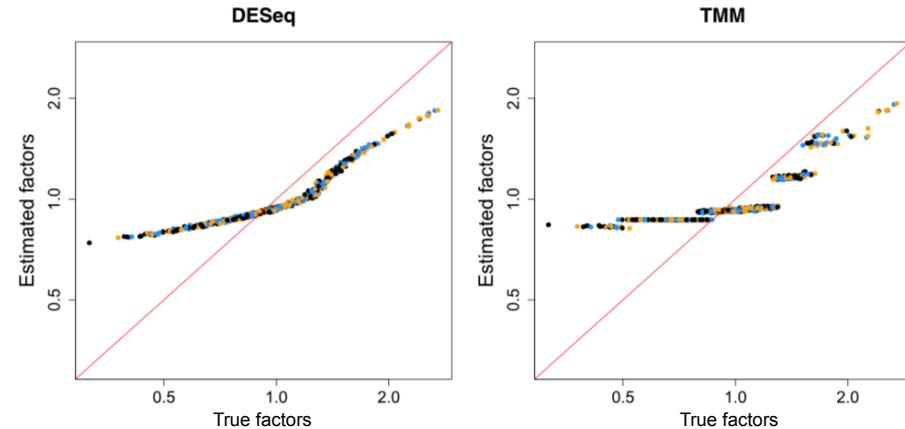
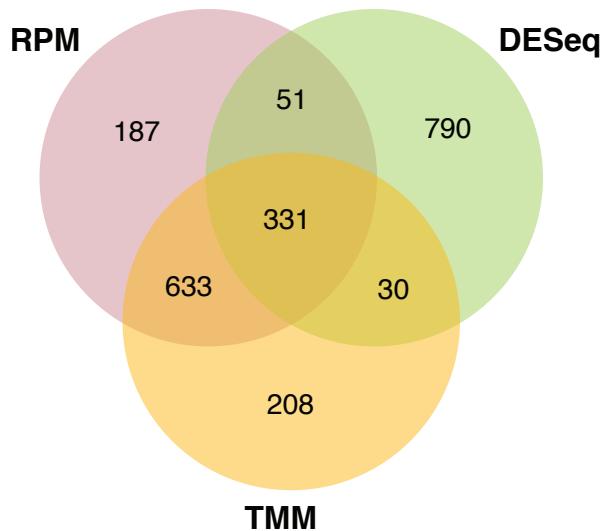
Normalising scRNAseq data — global scaling



Note: some of these are removed when using UMIs

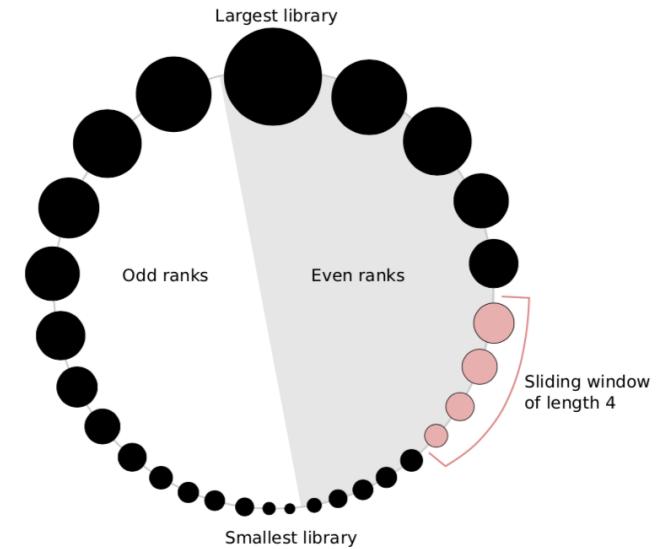
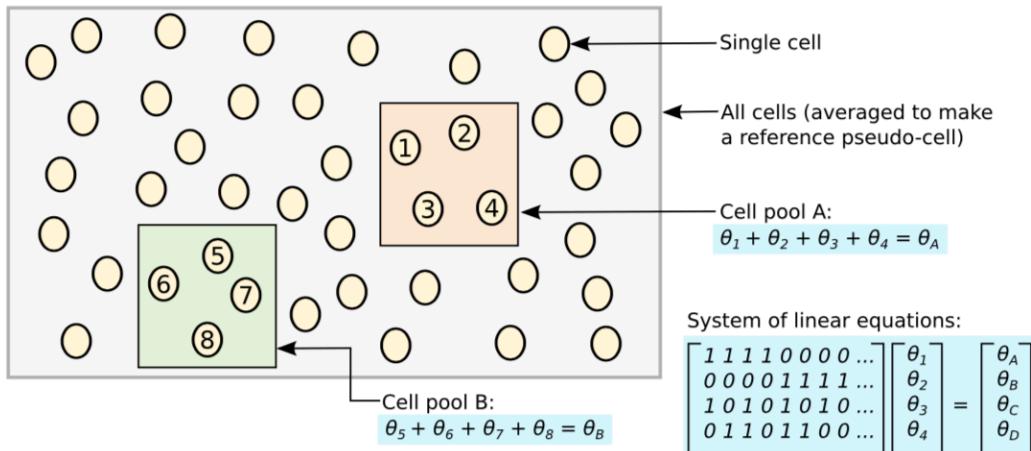
Normalising scRNASeq data — bulk methods?

Top 10% highly variable genes



Vallejos, Risso, Scialdone *et al* (2017) *Nature Methods*
Lun *et al* (2016) *Genome Biology*

scran: pooling across cells to normalise scRNAseq data

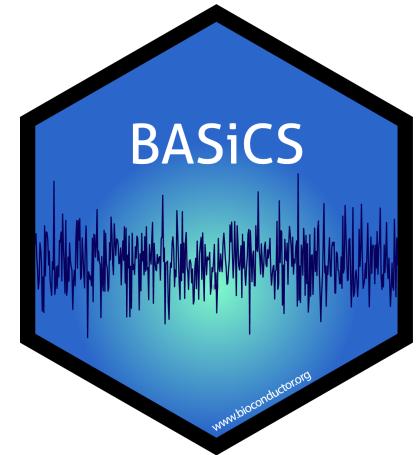


Lun et al (2016) *Genome Biology*

BASiCS: Bayesian Analysis of Single Cell Sequencing data

BASiCS is an integrated Bayesian statistical framework that **simultaneously** performs:

- Data normalisation
- Technical noise quantification
- **Supervised** downstream analyses
 - Detection of highly/lowlly variable genes **within** a population
 - Differential mean and variability testing **between** populations



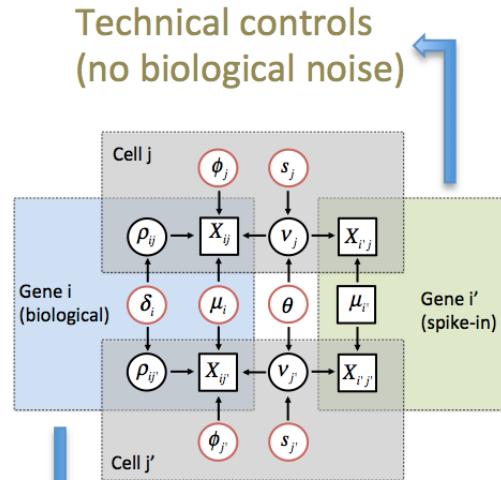
Integrated approach to propagate statistical uncertainty

Vallejos, Marioni and Richardson (2015) *PLoS Computational Biology*

Vallejos, Richardson and Marioni (2016) *Genome Biology*

Eling et al (2018) *Cell Systems*

BASiCS: Bayesian Analysis of Single Cell Sequencing data



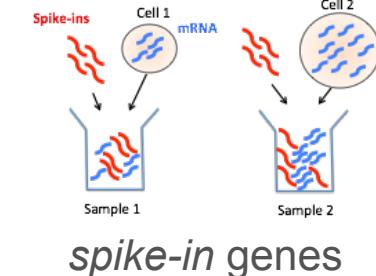
Hierarchical Bayesian model to capture **technical** and **biological** over-dispersion

$$X_{ij} | \mu_i, \phi_j, \nu_j, \rho_{ij} \sim \text{Poisson}(\phi_j \nu_j \mu_i \rho_{ij})$$

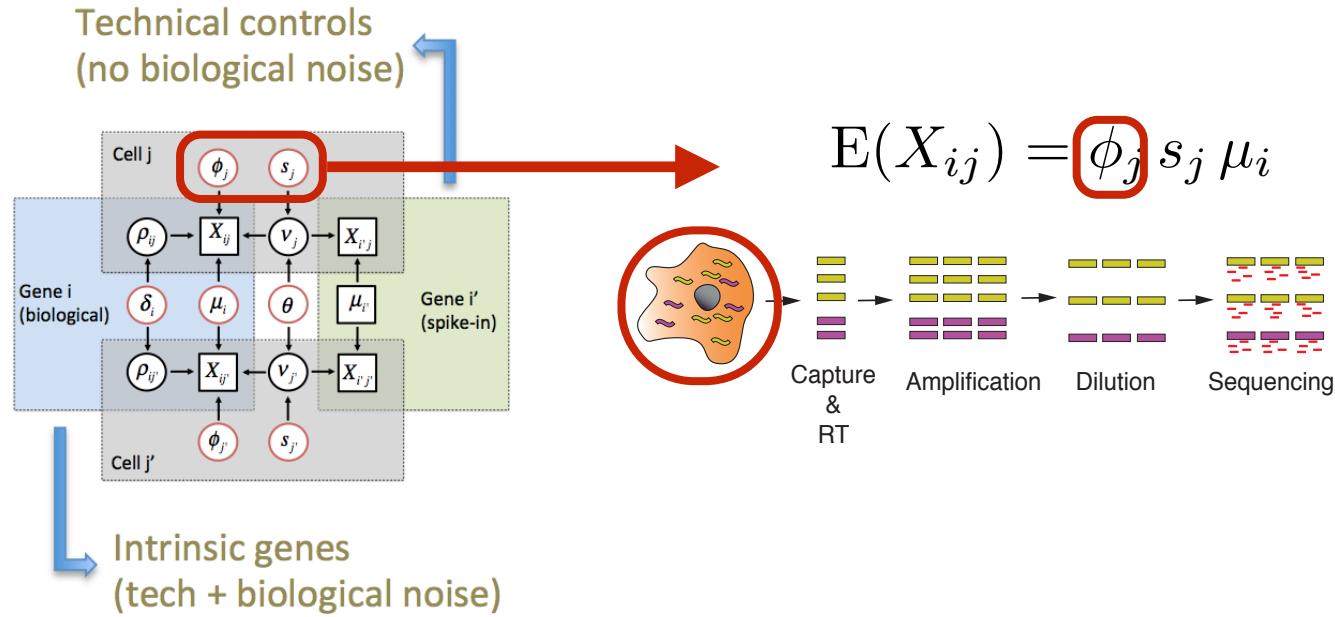
$$\nu_j | s_j, \theta \sim \text{Gamma}(s_j^{-1}, (s_j \theta)^{-1})$$

$$\rho_{ij} | \delta_i \sim \text{Gamma}(\delta_i^{-1}, \delta_i^{-1})$$

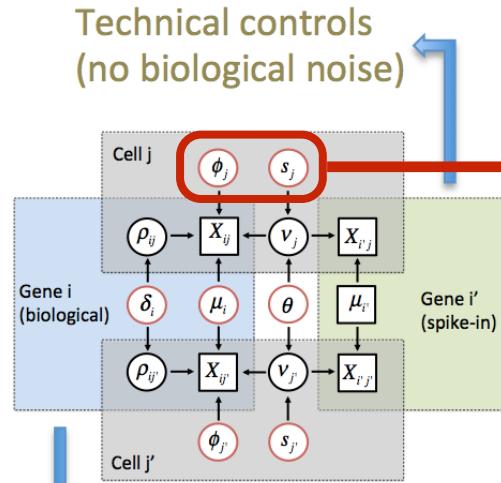
Requires pre-filtered data



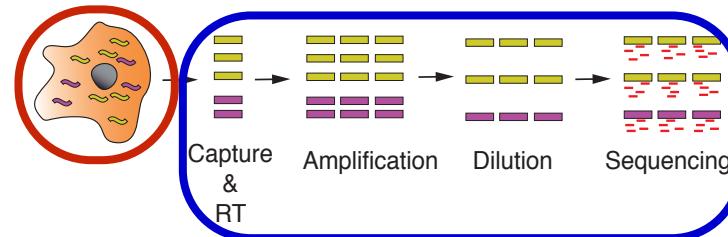
BASiCS: Bayesian Analysis of Single Cell Sequencing data



BASiCS: Bayesian Analysis of Single Cell Sequencing data



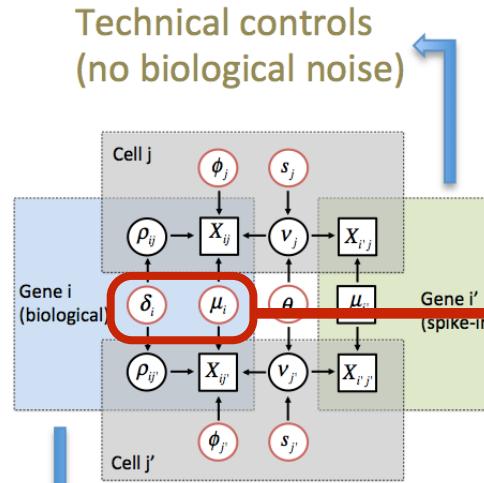
$$E(X_{ij}) = \phi_j s_j \mu_i$$



Intrinsic genes
(tech + biological noise)

One set of scaling factors if spike-ins are not available

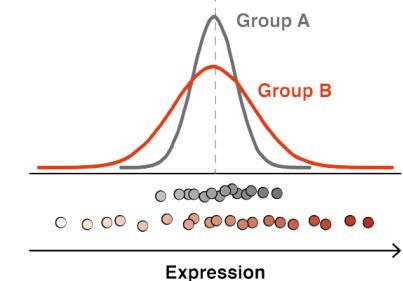
BASiCS: Bayesian Analysis of Single Cell Sequencing data



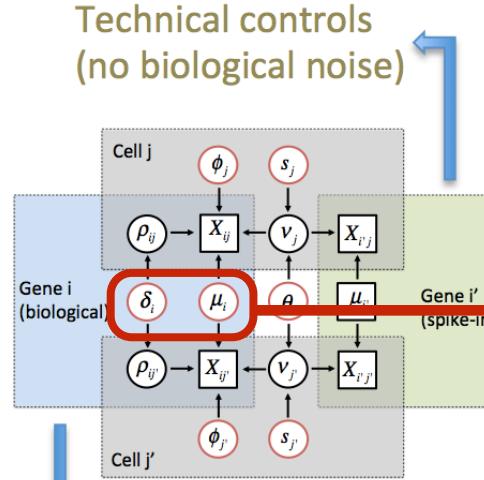
Intrinsic genes
(tech + biological noise)

Gene-specific parameters:

- mean expression μ_i
- biological over-dispersion δ_i



BASiCS: Bayesian Analysis of Single Cell Sequencing data

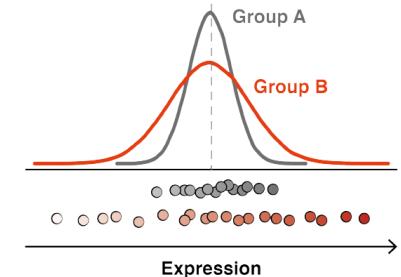
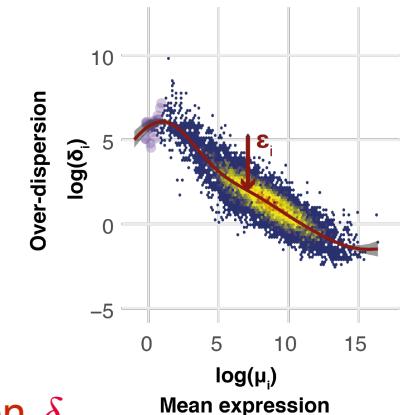


Technical controls
(no biological noise)

Intrinsic genes
(tech + biological noise)

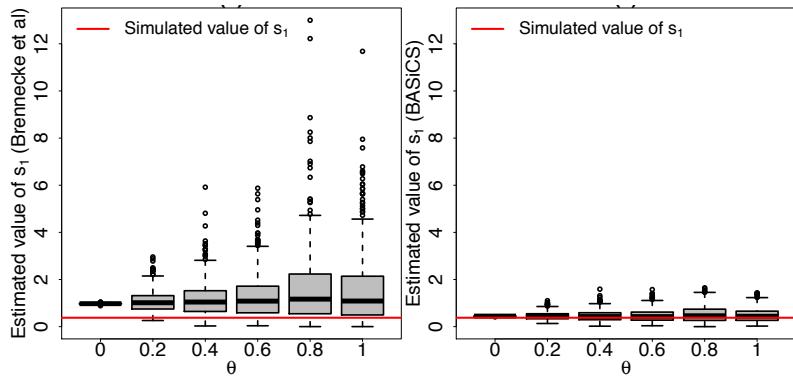
Gene-specific parameters:

- mean expression μ_i
- biological over-dispersion δ_i

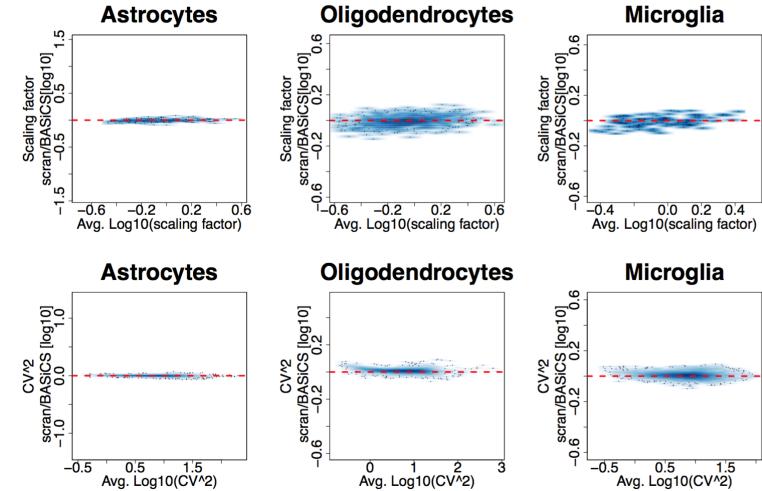


Does it work?

BASiCS and DESeq¹ — synthetic data

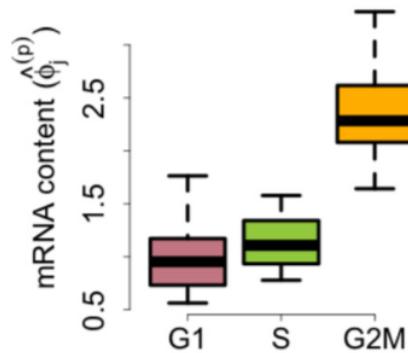


BASiCS and scran² — Zeisel et al data³

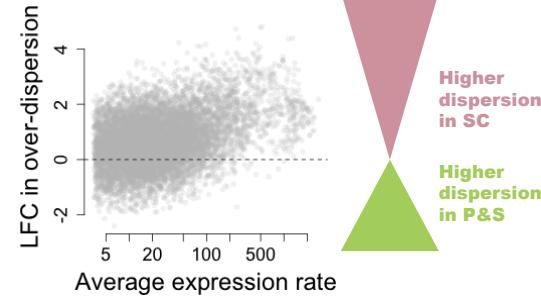
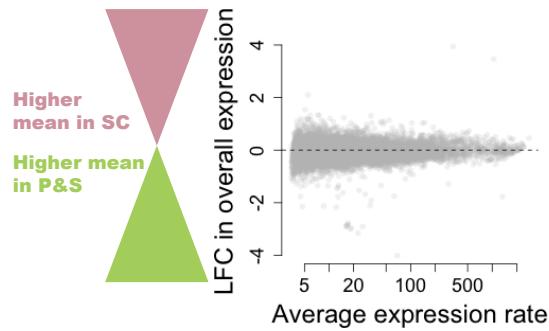


Does it work?

mRNA content across cell cycle



A pool and split experiment



Does it work?

Predictive checks

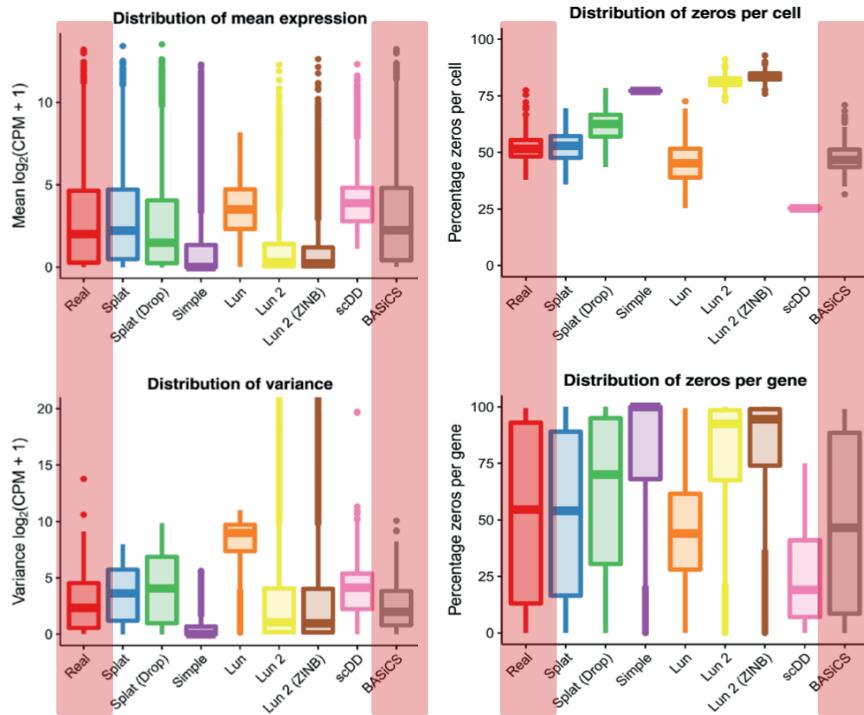
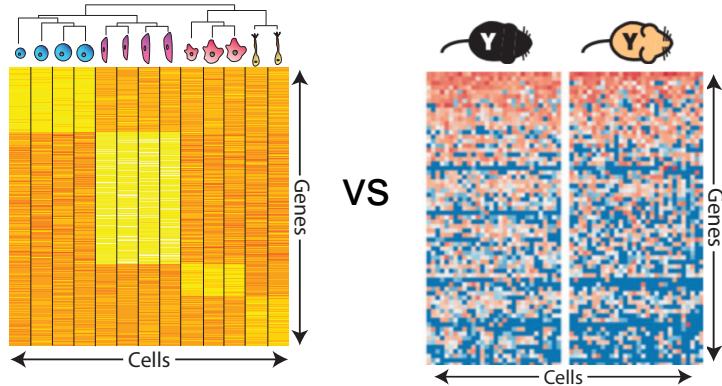


Image adapted from Zappia et al (2017) *Genome Biology*

Dealing with structural heterogeneity — scran & BASiCS



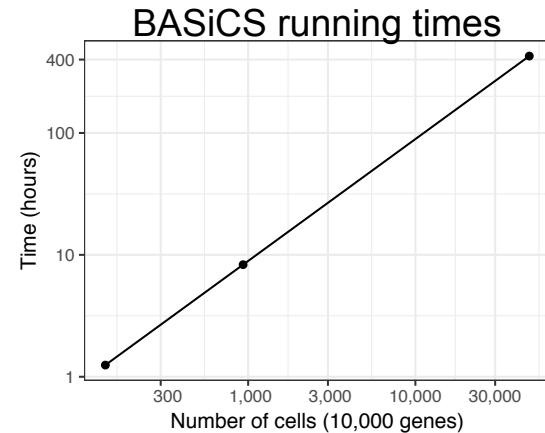
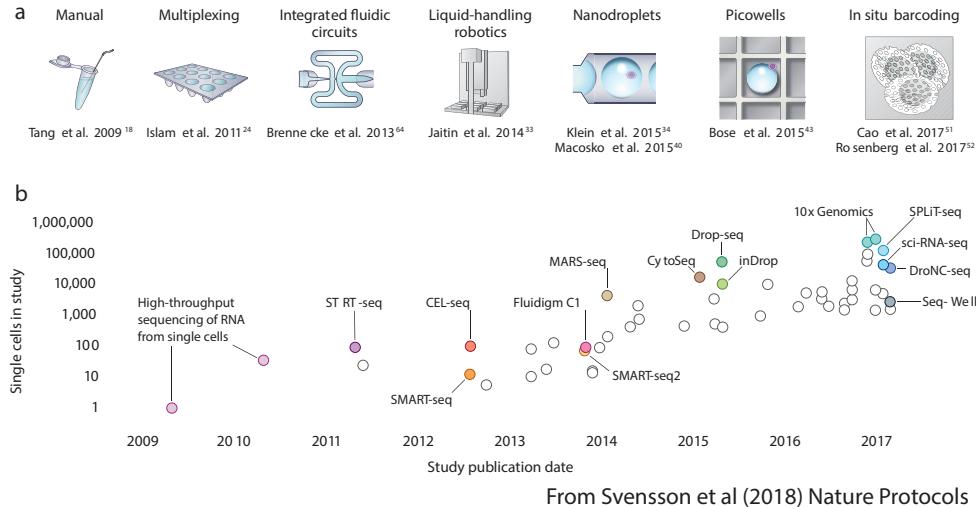
From Martinez-Jimenez,
Eling et al (2017) Science

If structural heterogeneity (strong DE) exists:

- *Pre-cluster* cells
- Normalise cells within each cluster
- Normalise between clusters



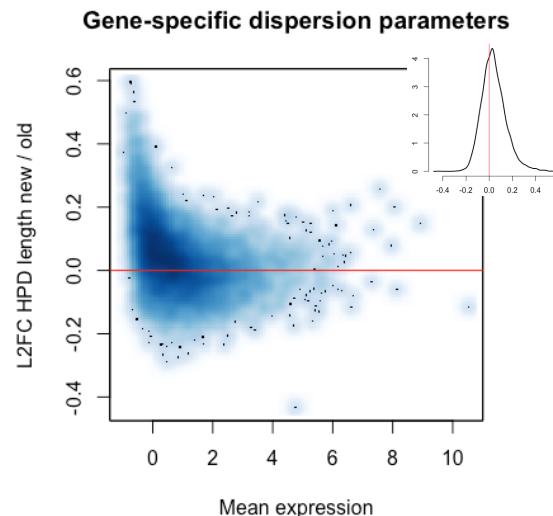
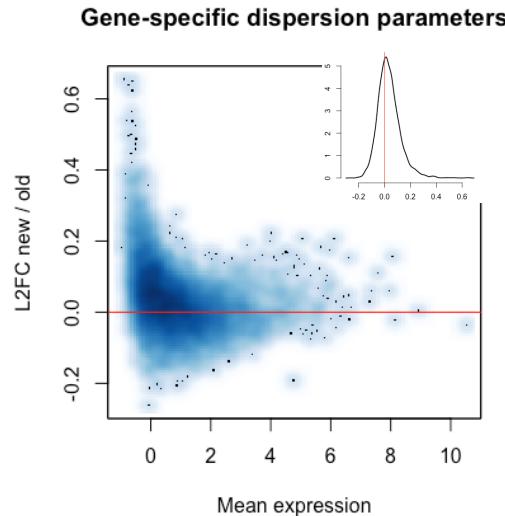
Scalability - ongoing work



@AlanBOCallaghan

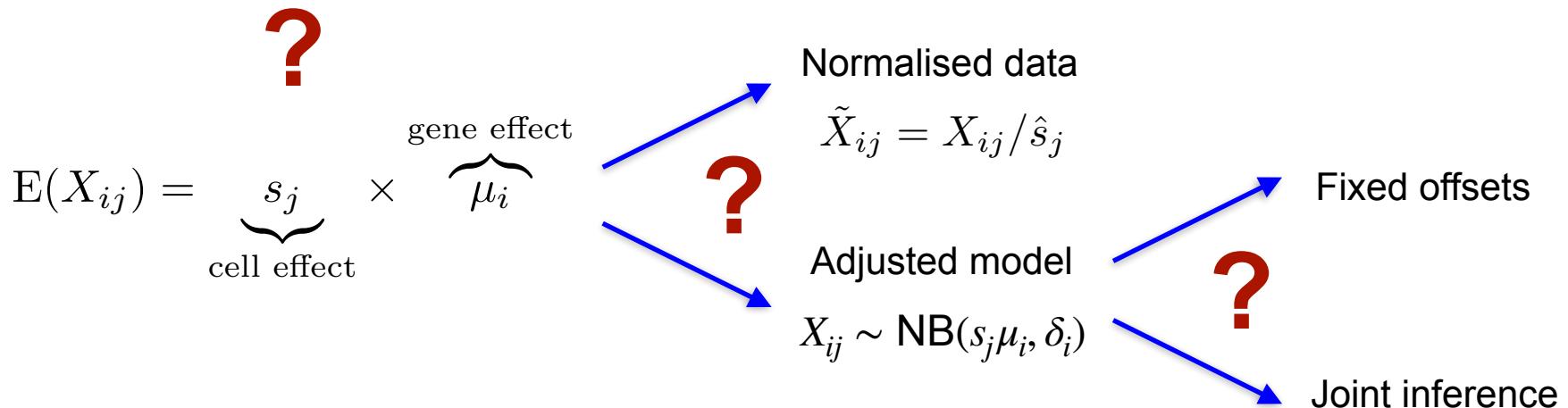
Integrated approach to propagate statistical uncertainty ...

... what if we run BASiCS with scran size factors as fixed offsets?



Note: posterior uncertainty about the normalisation parameters was low

For tomorrow?



Acknowledgements



Vallejos Group

[Alan O'Callaghan](#)

Andreas Kapourani
Andrew Papanastasiou
Ava Khamseh
Christos Maniatis
Evgenii Lobzaeb
Fraser McPhie
Nathan Constantine-Cooke
Sylvia Richardson



[HelmholtzZentrum münchen](#)
German Research Center for Environmental Health

Berkeley
UNIVERSITY OF CALIFORNIA



John Marioni
Nils Eling
Arianne Richard
Aaron Lun

Antonio Scialdone

Sandrine Dudoit
Davide Risso



The
Alan Turing
Institute



THE UNIVERSITY
of EDINBURGH