

# Splatter: simulation of single-cell RNA sequencing data

## Authors

- Luke Zappia<sup>1,2</sup> ([luke.zappia@mcri.edu.au](mailto:luke.zappia@mcri.edu.au))
- Belinda Phipson<sup>1</sup> ([belinda.phipson@mcri.edu.au](mailto:belinda.phipson@mcri.edu.au))
- Alicia Oshlack<sup>\*1,2</sup> ([alicia.oshlack@mcri.edu.au](mailto:alicia.oshlack@mcri.edu.au))

## Affiliations

1. Murdoch Childrens Research Institute, Royal Children's Hospital, 50 Flemington Rd, Parkville VIC 3052, Australia
2. School of Biosciences, The University of Melbourne, Parkville VIC 3010, Australia

\* Corresponding author

## 1 **Abstract**

2 As single-cell RNA sequencing (scRNA-seq) technologies have rapidly developed so  
3 have methods of analysis. Many of these methods have been tested and developed using  
4 simulated datasets. While this is a valid and useful approach many currently published  
5 simulations are problematic because they are not well documented, code may not be  
6 available for reproducing the simulation or their similarity to real data is not  
7 demonstrated.

8 Here we present the Splatter package for simple simulation of single-cell RNA-seq  
9 data. Splatter is a Bioconductor R package that provides a consistent, easy to use and  
10 well-documented interface for multiple scRNA-seq simulation methods. The Splatter  
11 package makes it easy to compare simulated datasets with real data to produce a  
12 realistic simulation that can be used to evaluate analysis methods. In addition we  
13 develop our own simulation, Splat, based on a gamma-poisson distribution. Splat  
14 incorporates a number of key features including high-expression outlier genes, defined  
15 library sizes, a mean-variance trend and expression-based dropout. Furthermore, Splat  
16 can simulate single populations of cells, populations with multiple cell types or  
17 differentiation paths.

## 18 **Keywords**

19 Single-cell – RNA-seq – Simulation – Software

## 20 **Background**

21 The first decade of next-generation sequencing has seen an explosion in our  
22 understanding of the genome [1]. In particular, the development of RNA sequencing  
23 (RNA-seq) has enabled unprecedented insight into the dynamics of gene expression [2].  
24 Researchers now routinely conduct experiments designed to test how gene expression  
25 is affected by various stimuli. One limitation of bulk RNA-seq experiments is that they

1 measure the average expression level of genes across the many cells in a sample.  
2 However, recent technological developments now allow the extraction and amplification  
3 of minute quantities of RNA, meaning that sequencing can now be conducted on the  
4 level of single cells [3]. The increased resolution of single-cell RNA-seq (scRNA-seq) data  
5 makes a range of new analyses possible.

6 As scRNA-seq data has become available there has been a rapid development of new  
7 bioinformatics tools attempting to unlock it's potential. Currently there are at least 100  
8 software packages that have been designed specifically for the analysis of scRNA-seq  
9 data, the majority of which have been published in peer-reviewed journals or as  
10 preprints. A table of current scRNA-seq software is available at <https://goo.gl/4wcVwn>.  
11 The focus of these packages is often different from those designed for the analysis of a  
12 bulk RNA-seq experiment. In a bulk experiment the groups of samples are known and a  
13 common task is to test for genes that are differentially expressed (DE) between two or  
14 more groups. In contrast, the groups in a single-cell experiment are usually unknown  
15 and the analysis is often more exploratory. Many of the current packages focus on  
16 assigning cells to groups based on expression profiles (clustering) before applying more  
17 traditional DE testing. This approach is taken by tools such as SC3 [4], CIDR [5] and  
18 Seurat [6] and makes sense for a sample of mature cells where it is reasonable to expect  
19 cells to have a particular type. In the developmental setting, for example, where stem  
20 cells are differentiating into mature cells, it may be more appropriate to order cells  
21 along a continuous trajectory from one cell type to another. Tools such as Monocle [7],  
22 CellTree [8] and Sincell [9] take this approach, ordering cells along a path, then looking  
23 for patterns in the changes of gene expression.

24 Existing scRNA-seq analysis packages, and any new methods that are being  
25 developed, should demonstrate two properties: 1) they can do what they claim to do,  
26 whether that is clustering, lineage tracing, differential expression testing or improved

1 performance compared to other methods and 2) that they produce some meaningful  
2 biological insight. The second criterion is specific to particular studies but it should be  
3 possible to address the first point in a more general way.

4 A common way to test the performance of an analysis method is through a  
5 simulation. Simulated data provides a known truth to test against, making it possible to  
6 see whether a method has been implemented correctly, whether the assumptions of the  
7 method are appropriate, and demonstrating it's limitations. This sort of testing is  
8 difficult with real biological data where an experiment must be specifically designed or  
9 results from an orthogonal test taken as the truth. Simulations, however, allow access to  
10 a range of metrics for assessing the performance of an analysis method. An additional  
11 advantage of assessing performance on simulated data is that many datasets, with  
12 different parameters and assumptions, can be rapidly generated at minimal cost. As  
13 such, many of the scRNA-seq analysis packages that are currently available use  
14 simulations to demonstrate their effectiveness. These simulations, however, are often  
15 not described in a reproducible or reusable way and the code to construct them may not  
16 be readily available. When code is available it can be poorly documented or written  
17 specifically for the computing environment of the authors, limiting it's reproducibility  
18 and making it difficult for other researchers to reuse. Most importantly publications do  
19 not usually devote time to demonstrating that a simulation is similar to real datasets, or  
20 in what ways it differs.

21 In this paper we present Splatter, a Bioconductor R package for reproducible and  
22 accurate simulation of single-cell RNA sequencing data. Splatter is a framework  
23 designed to provide a consistent interface to multiple simulations, enabling researchers  
24 to quickly simulate scRNA-seq count data in a reproducible fashion and make  
25 comparisons between simulations and real data. Along with the framework we develop  
26 our own simulation model, Splat, and show how it compares to previously published

1 simulations based on real datasets. We also provide a short example of how simulations  
2 can be used for assessing analysis methods.

## 3 **Results**

### 4 **The Splatter framework**

5 Splatter is designed to provide a consistent interface to multiple models for  
6 simulating scRNA-seq expression data. Currently, Splatter implements five different  
7 simulation models, each with their own assumptions (described in more detail in the  
8 following sections). The Splatter simulation process consists of two steps. The first step  
9 takes a real dataset and estimates the parameters required for the simulation. The result  
10 of this first step is a parameters object unique to each simulation model. These objects  
11 have been designed to hold the information required for the simulation and display  
12 details such as which parameters can be estimated and which have been changed from  
13 the default value. It is important that each simulation has it's own object for storing  
14 parameters as different simulations can vary greatly in the information they require. For  
15 example, some simulations only need parameters for well known statistical  
16 distributions while others require large vectors or matrices of data sampled from real  
17 datasets.

18 In the second step Splatter takes the estimated parameters, along with any  
19 additional parameters that cannot be estimated or are specified by the user, and  
20 generates a synthetic scRNA-seq dataset. In the case that there is no relevant real data to  
21 estimate parameters from, a synthetic dataset can still be generated using default  
22 parameters that can be manually modified by the user. The main result of the simulation  
23 step is a matrix of counts which is returned as an SCESet object from the scater package  
24 [10]. Briefly, the structure of this object combines cell by feature (gene) matrices for  
25 storing expression values with tables for storing metadata about cells and features

1 (further details are described in the scater documentation and the accompanying  
2 paper). This format makes it convenient to return intermediate values created during  
3 simulation as well as the final expression matrix. For example, the underlying gene  
4 means in different groups of cells are returned and could be used as a truth when  
5 evaluating differential expression testing. The documentation for each simulation  
6 describes the model, the parameters it uses and the intermediate values it returns.

7 An additional feature of Splatter is the ability to compare SCESet objects. These may  
8 be simulations with different models, different parameters, or could be real datasets  
9 from which parameters have been estimated. The comparison function takes one or  
10 more SCESet objects, combines them (keeping any cell or gene level information that is  
11 present in all of them) and produces a series of diagnostic plots comparing aspects of  
12 scRNA-seq data. The combined datasets are also returned making it easy produce  
13 further comparison plots or statistics. Alternatively, one SCESet can be specified as a  
14 reference, such as the real data used to estimate parameters, and the difference between  
15 this and the other datasets can be assessed. This approach is particularly useful for  
16 comparing how well simulations recapitulate real datasets. Examples of these  
17 comparison plots are shown in the following sections.

## 18 **Simulation models**

19 Splatter provides implementations of our own simulation model, Splat, as well as  
20 several previously published simulations. These previous simulations have either been  
21 published as R code associated with a paper or as functions in existing packages. By  
22 including them in Splatter we make them available in a single place in a more accessible  
23 way. Where only a script has been published, such as the Lun [11] and Lun 2 [12]  
24 simulations, the simulations have been re-implemented in Splatter. If the simulation is  
25 in an existing R package, for example scDD [13], we have simply written wrappers that  
26 provide consistent input and output but use the package implementation. We have

1 chosen to keep the simulations and estimation procedures as close as possible to what  
2 was originally published while keeping a consistent interface. The five different  
3 simulations currently available in Splatter are described below.

#### 4 **Simple**

5 The negative-binomial is the most common distribution used to model RNA-seq  
6 data, for example in the edgeR [14] and DESeq [15] packages, and the Simple simulation  
7 is a basic implementation of this approach. A mean expression level for each gene is  
8 simulated using a gamma distribution and the negative-binomial distribution is used to  
9 generate a count for each cell based on these means, with a fixed dispersion parameter  
10 (default = 0.1) (Additional figure 1). This simulation is primarily included as a baseline  
11 reference and is not meant to truly reproduce many of the features of scRNA-seq data.

#### 12 **Lun**

13 Published in “Pooling across cells to normalize single-cell RNA sequencing data with  
14 many zero counts” [11] the Lun simulation builds on the Simple simulation by adding a  
15 scaling factor for each cell (Additional figure 2). The cell factors are randomly sampled  
16 from a normal distribution with mean 1 and variance 0.5. The inverse- $\log_2$  transformed  
17 factors are used to adjust the gene means resulting in a matrix, where each cell has a  
18 different mean. This represents the kinds of technical effects that scaling normalisation  
19 aims to remove. The matrix of means are then used to sample counts from a negative  
20 binomial distribution, with a fixed dispersion parameter. This simulation can also model  
21 differential expression between multiple groups with fixed fold changes.

#### 22 **Lun 2**

23 In “Overcoming confounding plate effects in differential expression analyses of  
24 single-cell RNA-seq data” [12] the same authors extend the negative-binomial model  
25 from the Lun simulation. This simulation samples input parameters from real data, with

1 very little random sampling from statistical distributions. In the Lun 2 simulation the  
2 cell factors are replaced with a library size factor and an additional level of variation is  
3 added by including a batch effects factor. While the library size factor acts on individual  
4 cells the batch effects are applied to groups of cells from the same batch. This simulation  
5 is thus highly specific to the scenario when there are known batch effects present in the  
6 data, for example Fluidigm C1 plate effects. Differential expression can be added  
7 between two sets of batches and the user can choose to use a zero-inflated negative-  
8 binomial (ZINB) model. Counts are simulated using the library size and plate factor  
9 adjusted gene means and the gene-wise dispersion estimates are obtained from the  
10 data. If the ZINB model is chosen, zero inflated estimates of gene means and dispersions  
11 are used instead and an additional step randomly sets some counts to zero, based on the  
12 gene-wise proportions of zeroes observed in the data. Additional figure 3 shows the  
13 model assumptions and parameters for this simulation.

#### 14 **scDD**

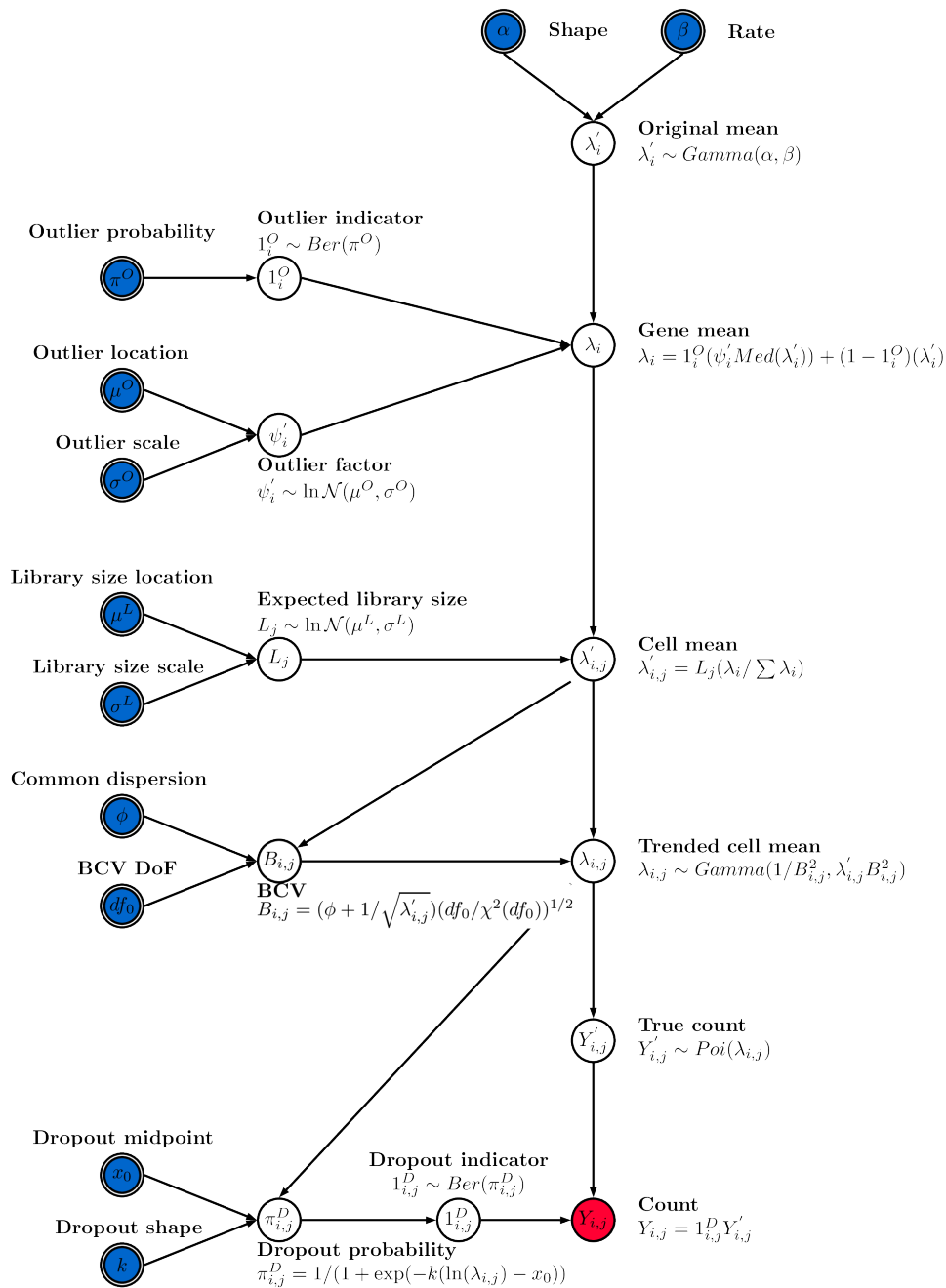
15 The scDD package aims to test for differential expression between two groups but  
16 also more complex changes such as differential distributions or differential proportions  
17 [13]. This is reflected in their simulation, which can contain a mixture of genes  
18 simulated to have different distributions, or differing proportions where the expression  
19 of the gene is multi-modal. This simulation also samples information from a real dataset.  
20 The Splatter package simply provides wrapper functions to the simulation functions in  
21 the scDD package, while capturing the necessary input and output needed to compare to  
22 other simulations. The full details of the scDD simulation are described in the scDD  
23 package vignette [16].

#### 24 **Splat**

25 We have developed the Splat simulation to capture many of the technical effects  
26 seen in real scRNA-Seq data, including high expression outlier genes, differing



1 sequencing depths between cells, trended gene-wise dispersion, and zero-inflation. Our  
2 model uses parametric distributions with hyper-parameters estimated from real data  
3 (Figure 1). The core of the Splat simulation is the gamma-Poisson hierarchical model  
4 where the mean expression level for each gene  $i$ ,  $i = 1, \dots, N$ , is simulated from a gamma  
5 distribution and the count for each cell  $j$ ,  $j = 1, \dots, M$ , is subsequently sampled from a  
6 Poisson distribution, with modifications to include expression outliers and to enforce a  
7 mean-variance trend.



1

2

3

4

5

6

7

8

9

**Figure 1: Diagram of the Splats simulation model. Input parameters are indicated with double borders and those that can be estimated from real data are shaded blue. Red shading indicates the final output. The simulation begins by generating means from a gamma distribution. Outlier expression genes are added by multiplying by a log-normal factor and the means are proportionally adjusted for each cells library size. A mean-variance trend is enforced by adjusting the means using a simulated Biological Coefficient of Variation (BCV). These final means are used to generate counts from a Poisson distribution. In the final step dropout is (optionally) simulated by randomly setting some counts to zero, based on each gene's mean expression.**

1 More specifically the Splat simulation initially samples gene means from a Gamma  
2 distribution with shape  $\alpha$  and rate  $\beta$ . While the gamma distribution is a good fit for gene  
3 means it does not always capture extreme expression levels. To counter a probability  
4 ( $\pi^0$ ) that a gene is a high expression outlier can be specified. We then add these outliers  
5 to the simulation by replacing the previously simulated mean with the median mean  
6 expression level multiplied by an inflation factor. The inflation factor is sampled from a  
7 log-normal distribution with location  $\mu^0$  and scale  $\sigma^0$ .

8 The library size (total number of counts) varies within an scRNA-seq experiment  
9 and can be very different between experiments depending on the sequencing depth. We  
10 model library size using a log-normal distribution (with location  $\mu^L$  and scale  $\sigma^L$ ) and  
11 use the simulated library sizes ( $L_j$ ) to proportionally adjust the gene means for each cell.  
12 This allows us to alter the number of counts per cell independently of the underlying  
13 gene expression levels.

14 It is known that there is a strong mean-variance trend in RNA-Seq data, where lowly  
15 expressed genes are more variable and highly expressed genes are more consistent [17].  
16 In the Splat simulation we enforce this trend by simulating the biological coefficient of  
17 variation (BCV) for each gene from a scaled inverse chi-squared distribution, where the  
18 scaling factor is a function of the gene mean. After simulating the BCV values we  
19 generate a new set of means ( $\lambda_{i,j}$ ) from a Gamma distribution with shape and rate  
20 parameters dependent on the simulated BCVs and previous gene means. We then  
21 generate a matrix of counts by sampling from a Poisson distribution, with lambda equal  
22 to  $\lambda_{i,j}$ . This process is similar to the process used by Law et al. in their simulation of bulk  
23 RNA-seq data [18].

24 One of the key features of scRNA-seq data is the high proportion of zeros [19], one  
25 cause of which is technical dropout. We use the relationship between the mean

1 expression of a gene and the proportion of zero counts in that gene to model this  
2 process, using a logistic function to produce a probability that a count should be zero.  
3 These probabilities are then used to randomly replace some of the simulated counts  
4 with zeros using a Bernoulli distribution.

5 The different steps in the Splat simulation outlined above are easily controlled and  
6 can be turned off when they are not desirable or appropriate. The final result is a matrix  
7 of observed counts  $Y_{i,j}$  where the rows are genes and the columns are cells. The full set  
8 of input parameters are shown in Table 1.

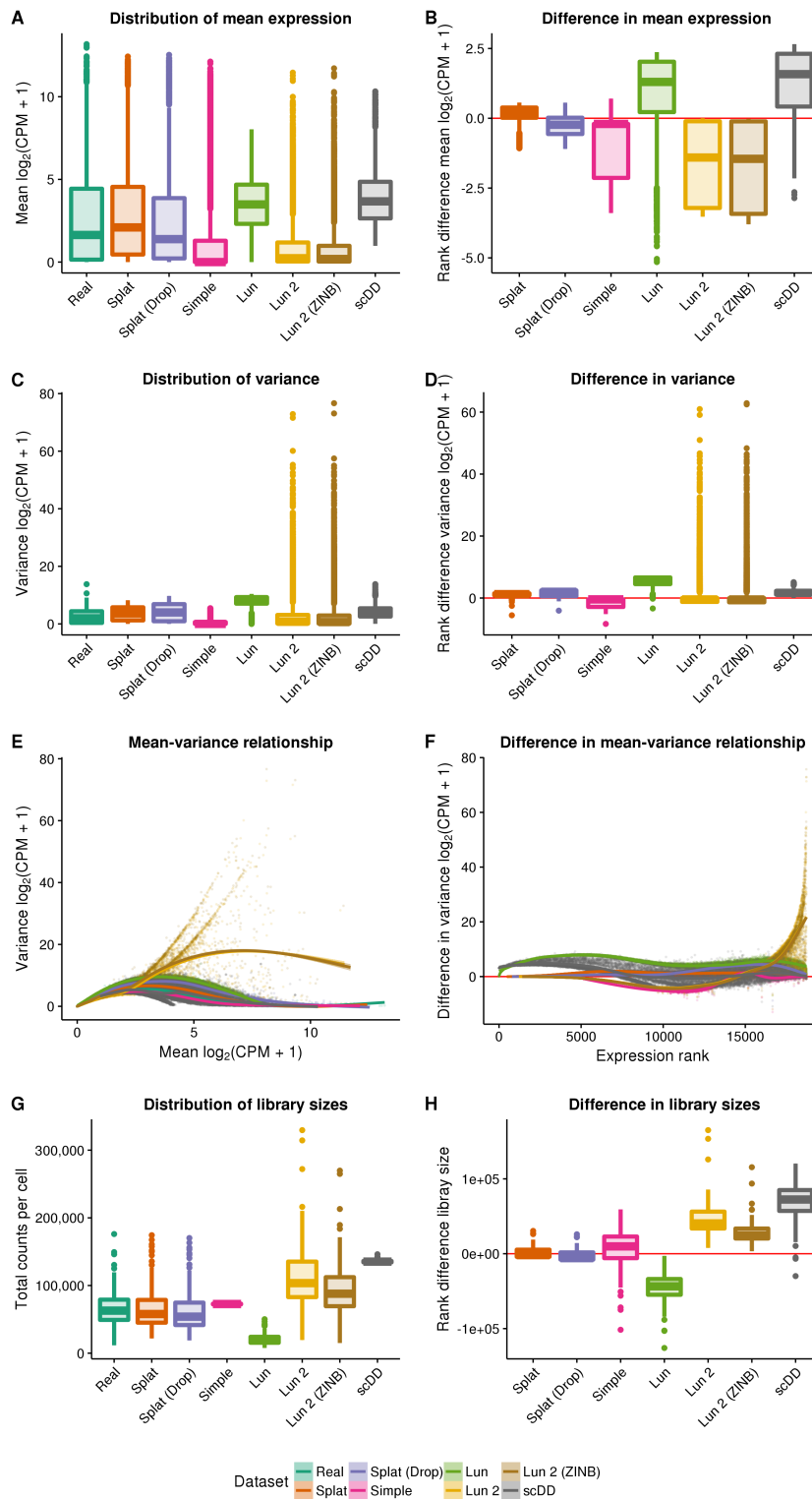
9 **Table 1: Input parameters for the Splat simulation model**

Name	Symbol	Description
Mean shape	$\alpha$	Shape parameter for the mean gene expression gamma distribution
Mean rate	$\beta$	Rate parameter for the mean gene expression gamma distribution
Library size location	$\mu^L$	Location parameter for the library size log-normal distribution
Library size scale	$\sigma^L$	Scale parameter for the library size log-normal distribution
Outlier probability	$\pi^O$	Probability that a gene is an expression outlier
Outlier location	$\mu^O$	Location parameter for the expression outlier factor log-normal distribution
Outlier scale	$\sigma^O$	Scale parameter for the expression outlier factor log-normal distribution
Common BCV	$\phi$	Common BCV dispersion across all genes
BCV degrees of freedom	$df$	Degrees of freedom for the BCV inverse chi-square distribution
Dropout midpoint	$x_0$	Midpoint for the dropout logistic function
Dropout shape	$k$	Shape of the dropout logistic function

## 10 Comparison of simulations

11 To compare the simulation models available in Splatter we estimated parameters  
12 from real datasets then generated synthetic datasets using those parameters. Both the  
13 standard and zero-inflated versions of Splat and Lun 2 simulations were included, giving  
14 a total of seven simulations. We began with the Tung dataset [20] which contains  
15 induced pluripotent stem cells from HapMap individuals.

1           No quality control of cells or filtering of genes was performed prior to estimation.  
2           We believe this presents the most challenging situation to simulate, as there are more  
3           likely to be violations of the underlying model. This scenario is also possibly the most  
4           useful as it allows any analysis method to be evaluated, from low-level filtering to  
5           complex downstream analysis. In order to reduce the computational time 200 cells were  
6           randomly sampled and used for estimation and each simulation consisted of 200 cells.  
7           Figure 2 shows some of the plots produced by Splatter to compare simulations based on  
8           the Tung dataset.

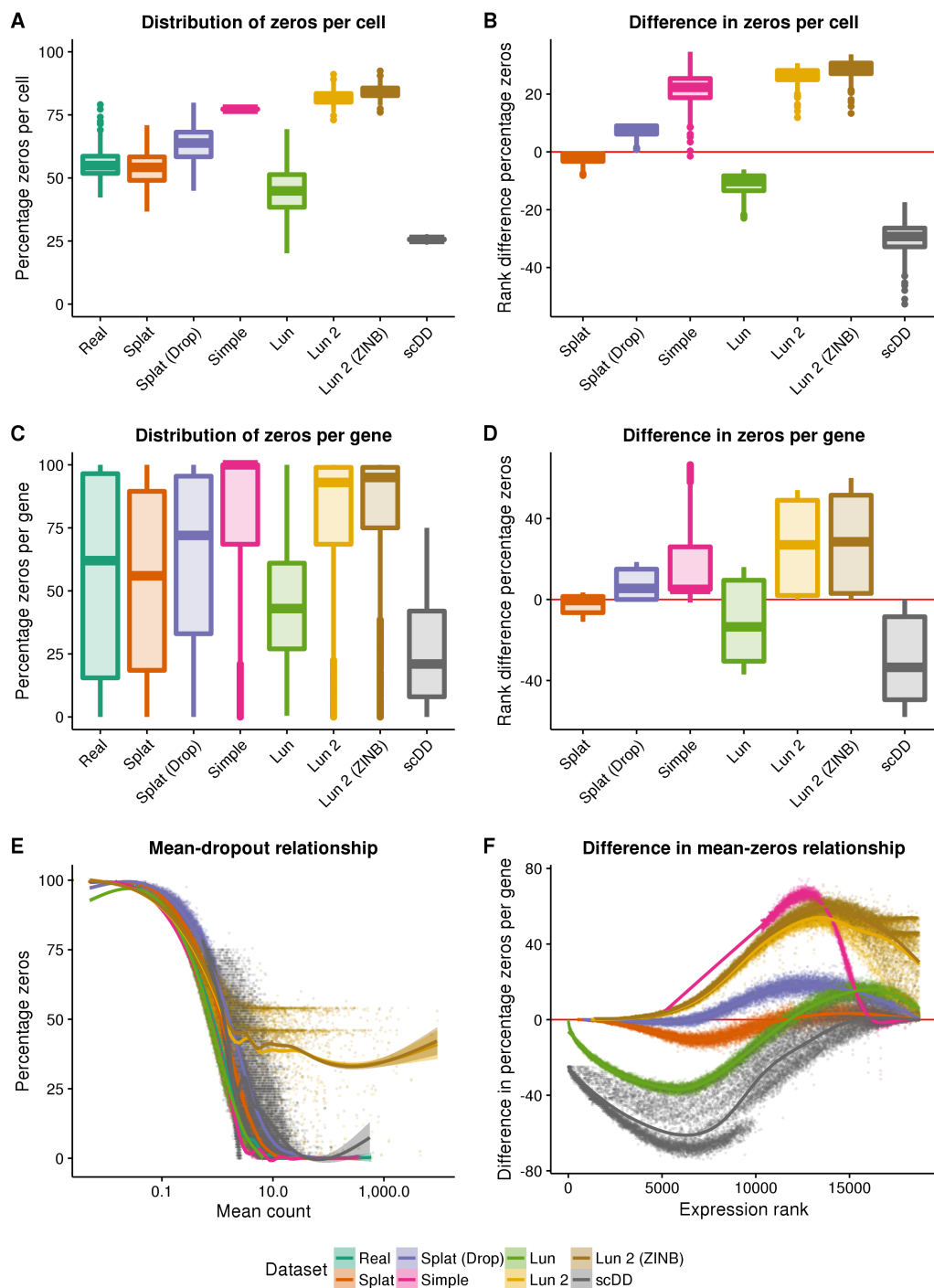


1

2 **Figure 2: Comparison of simulations based on the Tung dataset. The left column panels show the**  
 3 **distribution of mean expression (A), variance (C) and library size (G) across the real dataset and the**  
 4 **simulations as boxplots, along with a scatter plot of the mean-variance relationship (E). The right**  
 5 **column shows boxplots of the ranked differences between the real data and simulations for the**  
 6 **same statistics: mean (B), variance (D), mean-variance relationship (F) and library size (H).**

1           We compare the gene means, variances, library sizes and the mean-variance  
2 relationship. From these diagnostic plots we can evaluate how well each simulation  
3 reproduces the real dataset and in which ways it differs. One way to make comparisons  
4 is to look at the overall distributions (Figure 2, left column). Alternatively we can choose  
5 a reference (in this case the real data) and look at departures from that data (Figure 1,  
6 right column). Looking at the mean expression levels across genes we see that the Lun  
7 and scDD simulations are missing lowly expressed genes while the Simple and Lun 2  
8 simulations are skewed towards lower expression (Figure 2A , Figure 2B). The Splat  
9 simulation is a better match, likely due to the addition of high-expression outlier genes.  
10 Both versions of the Lun 2 simulation produce some extremely highly variable genes, an  
11 effect which is seen to a lesser extent in the Lun simulation. The difference in variance is  
12 reflected in the mean-variance relationship where genes from the Lun 2 simulation are  
13 much too variable at high expression levels for this dataset. Library size is another  
14 aspect in which the simulations differ from the real data. The simulations that don't  
15 contain a library size component (Simple, Lun, scDD) have different median library sizes  
16 and much smaller spreads. In this example, the Lun 2 simulations also produce some  
17 larger library sizes.

18           A key aspect of scRNA-seq data is the number of observed zeros. To properly  
19 recreate an scRNA-seq dataset a simulation must produce the correct number of zeros  
20 but also have them properly distributed across both genes and cells. In addition, a  
21 relationship between the expression level of a gene and the number of observed zeros  
22 has been established [21] and this should also be reproduced in simulations. Figure 3  
23 shows the distribution of zeros for the simulations based on the Tung dataset.



1

2 **Figure 3: Comparison of zeros in simulations based on the Tung dataset. The top row shows**  
 3 **boxplots of the distribution of zeros per cell (A) and the difference from the real data (B). The**  
 4 **distribution (C) and difference (D) in zeros per gene are shown in the middle row. The bottom row**  
 5 **shows scatter plots of the relationship between the mean expression of a gene and the percentage of**  
 6 **zeros as both the raw observations (E) and as ranked differences from the real data (F).**



1 For this dataset the Simple and Lun 2 simulations produce too many zeros across  
2 both genes and cells while the Lun and scDD simulations produce too few. Interestingly,  
3 the Splat simulation produces a better fit to this dataset when dropout is not included,  
4 suggesting that additional dropout is not present in the Tung dataset. However, this is  
5 not the case for all data and sometimes simulating additional dropout produces a better  
6 fit to the data (for example the Zeisel dataset presented below). We can also consider  
7 the relationship between expression level of a gene and the percentage of zero counts.  
8 The Lun and scDD simulations produce too few zeros at low expression levels while the  
9 Simple and Lun 2 simulations produce too many zeros at high expression levels.

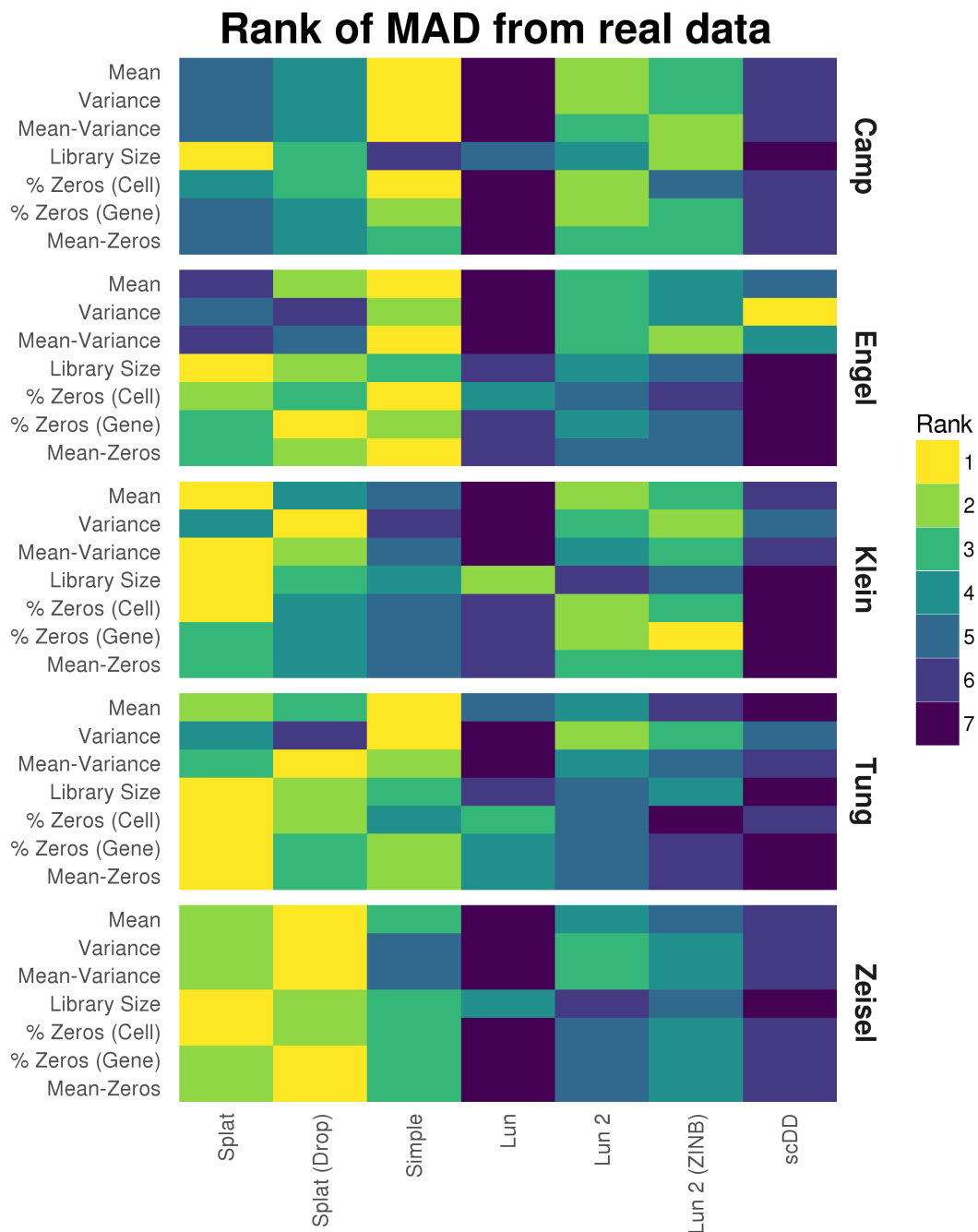
10 While the analysis presented in Figure 2 and Figure 3 allows us to visually inspect  
11 how simulations compare with a single dataset we wished to compare simulations  
12 across a variety of datasets. We performed simulations based on five different datasets  
13 (outlined in Table 2) that varied in terms of library preparation protocol, cell capture  
14 platform, species and tissue complexity. Three of the datasets used Unique Molecular  
15 Identifiers (UMIs) [22] and two were full-length protocols. Complete comparison panels  
16 for all the datasets are provided as Additional figures 4-8.

17 **Table 2: Details of real datasets**

<b>Dataset</b>	<b>Species</b>	<b>Cell type</b>	<b>Platform</b>	<b>Protocol</b>	<b>UMI</b>	<b>Number of cells</b>
Camp [23]	Human	Whole brain organoids	Fluidigm C1	SMARTer	No	597
Engel [24]	Mouse	Natural killer T cells	Flow cytometry	Modified Smart-seq2	No	203
Klein [25]	Human	K562 cells	InDrop	CEL-Seq	Yes	213
Tung [20]	Human	Induced pluripotent stem cells	Fluidigm C1	Modified SMARTer	Yes	564
Zeisel [26]	Mouse	Cortex and hippocampus cells	Fluidigm C1	STRT-Seq	Yes	3005

18 For each dataset we estimated parameters and produced a synthetic dataset as  
19 described before. We then compared simulations across metrics and datasets by

1 calculating a median absolute deviation (MAD) for each metric. For example to get a  
2 MAD for the gene expression means, the mean expression values for both the real data  
3 and the simulations were sorted and the real values were subtracted from the simulated  
4 values. The median of these absolute differences was taken as the final statistic. To  
5 compare between simulations we ranked the MADs for each metric. Figure 4  
6 summarises the ranked results for the five datasets as a heatmap with the MADs  
7 presented in Additional file 12.



1

2

**Figure 4: Comparison of simulation models based on various datasets. For each dataset**

3

**parameters were estimated and synthetic datasets generated using various simulation methods.**

4

**The Median Absolute Deviation (MAD) between each simulation and the real data was calculated for**

5

**a range of metrics and the simulations ranked. A heatmap of the ranks across the metrics and**

6

**datasets is presented here. We see that the Splat simulation (with and without dropout) performs**

7

**consistently well, with the two versions of the Lun 2 simulation also being good performers.**

1        Looking across the datasets we see that the Splat simulations are consistently good  
2 performers. On the Zeisel dataset both the zero-inflated simulations (Splat with dropout  
3 and Lun 2 ZINB) outperform their regular counterparts, suggesting that this datasets is  
4 truly zero-inflated. Interestingly the Simple simulation is the best performer on the  
5 Engel dataset and best captures the variance and mean-variance relationship in the  
6 Tung dataset. This result suggests that the additional features of the more complex  
7 simulation may be unnecessary in this case or that other models may be more  
8 appropriate. The Splat simulations were least successful on the Camp cerebral organoid  
9 dataset. The complex nature of this data (many cell types) and the full-length protocol  
10 may have contributed to this poorer performance. In this situation the semi-parametric,  
11 sampling based model of the Lun 2 simulation may have an advantage, and was the best  
12 performer. The Lun simulation is consistently amongst the worst performing but, given  
13 that this model is largely similar to the others, it is likely due to the lack of an estimation  
14 procedure for most parameters rather than significant problems with the model. The  
15 scDD simulation also often differed significantly from the real data, particularly in  
16 library size, and may benefit from the addition of fixed library sizes to the model. Most  
17 significantly we see that simulations perform differently on different datasets,  
18 emphasising the importance of trying different models and demonstrating their  
19 similarity to real datasets. The Splatter framework makes these comparisons between  
20 simulation models straightforward, making it easier for researchers to choose  
21 simulations that best reflect the data they are trying to model.

## 22        **Complex simulations with Splat**

23        The simulation models described above are sufficient for simulating a single,  
24 homogeneous population but not to reproduce the more complex situations seen in a  
25 real biological sample. For example, we might wish to simulate a population of cells  
26 from a complex tissue containing multiple mature cell types or a developmental

1 scenario where cells are transitioning between cell types. In this section we present how  
2 the Splat simulation can be extended to reproduce these complex sample types.

### 3 **Simulating groups**

4 Splat can model samples with multiple cell types by creating distinct groups of cells  
5 where several genes are differentially expressed between the different groups.

6 Previously published simulations can reproduce this situation to some degree but are  
7 often limited to fixed fold changes between only two groups. In the Splat simulation,  
8 however, modeling differential expression using a process similar to that for creating  
9 expression outliers can be used to simulate complex cell mixtures. Specifically, a  
10 proportion of genes are randomly selected to be up or down-regulated with given  
11 probabilities. Next a multiplicative factor is chosen for these genes from a log-normal  
12 distribution and applied to the underlying means. Setting the number of groups, the  
13 number of cells in each group and the proportion of genes that are differentially  
14 expressed allows flexibility in how different groups are defined. The resulting SCESet  
15 object contains information about which group each cell comes from as well as the  
16 factors applied to each gene in each group.

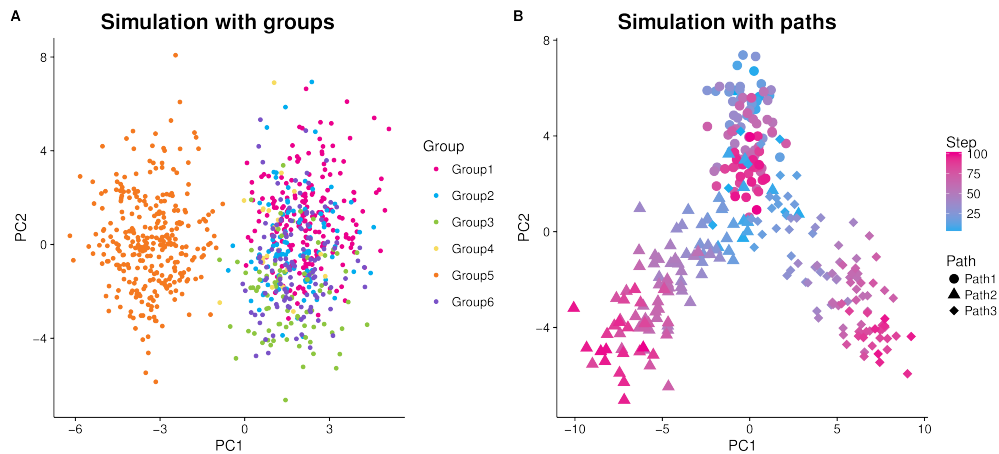
### 17 **Simulating paths**

18 An alternative situation that is often studied using scRNA-seq is cellular  
19 development and differentiation. Instead of having groups of mature cells, individual  
20 cells are somewhere on a continuous differentiation path or lineage from one cell type  
21 to another. To model this the Splat simulation uses the differential expression process  
22 described above to define expression levels of a start and end cell for each path. An  
23 interpolation process is then used to define multiple steps between the two cells types  
24 and the simulated cells are randomly assigned to one of these steps. The step each cell is  
25 allocated to is used to define the mean expression for the genes in that cell. Therefore  
26 the simulation of lineages, using Splat, is defined by the differential expression

1 parameters used to create the differences between the start and end of each path, as  
2 well as the parameters that define the path itself, the length (number of steps) and skew  
3 (whether cells are more likely to come from the start or end of the path).

4 In real data it has been observed that expression of genes can change in more  
5 complex, non-linear ways across a differentiation trajectory. For example, a gene may be  
6 lowly expressed at the beginning of a process, highly expressed in the middle and lowly  
7 expressed at the end. Splat models these kinds of changes by generating a Brownian  
8 bridge (a random walk with fixed end points) between the two end cells of a path, which  
9 is then smoothed and interpolated. This random element allows many possible patterns  
10 of expression changes over the course of a path (Additional figure 9). While non-linear  
11 changes are possible they are not the norm. Splat defines parameters that control the  
12 proportion of genes that are non-linear and how variable those genes can be.

13 Further complexity in simulating differentiation paths can be achieved by modeling  
14 lineages with multiple steps or branches. For example a stem cell that differentiates into  
15 an intermediate cell type that then changes into one of two mature cell types. These  
16 possibilities are enabled by allowing the user to set a starting point for each path.  
17 Figure 5 shows examples of using Splat to simulate datasets with multiple groups of  
18 cells or complex differentiation paths.

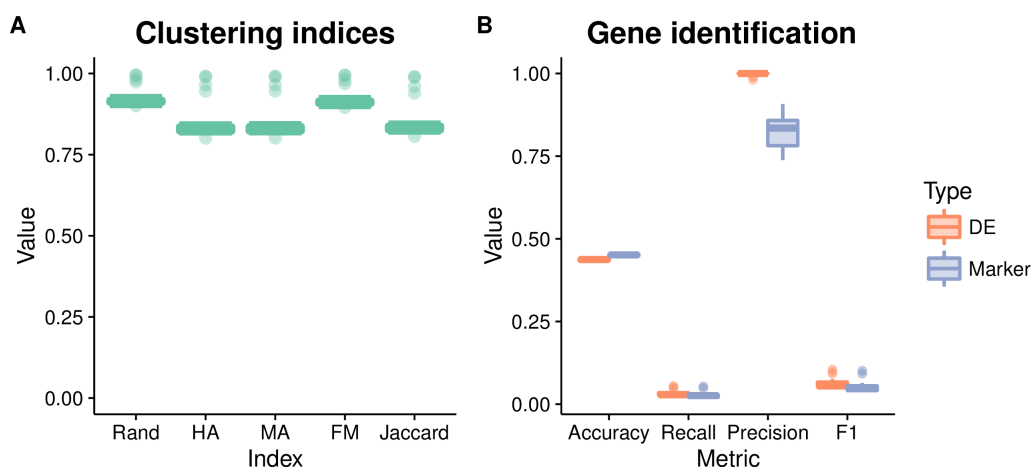


1  
2 **Figure 5: Examples of complex Splatter simulations. (A) is a PCA plot of a simulation with six**  
3 **groups with varying numbers of cells and levels of differential expression. (B) is a PCA plot of a**  
4 **simulation with differentiation paths. A progenitor cell type (circles) differentiates into an**  
5 **intermediate cell type, which becomes one of two (triangle or diamond) mature cell types. The**  
6 **coloured gradient indicates how far along a path each cell is from blue to pink.**

### 7 **Example: using Splatter simulations to evaluate a clustering method**

8 To demonstrate how the simulations available in Splatter could be used to evaluate  
9 an analysis method we present an example of evaluating a clustering method. SC3 [4] is  
10 a consensus k-means based approach available via Bioconductor. As well as assigning  
11 cells to groups, SC3 is able to detect genes that are differentially expressed between  
12 groups and marker genes that distinguish between groups. To test SC3 we estimated  
13 Splatter simulation parameters from the Tung dataset and simulated three cell populations  
14 of different sizes (200, 100 and 50 cells), with a probability of a gene being differentially  
15 expressed of 0.1. This resulted in approximately 1900 DE genes per group. We then ran  
16 SC3 with three clusters ( $k = 3$ ) and compared the results to the true groupings (Figure  
17 6A). We also assessed the detection of DE and marker genes. True DE genes were taken  
18 as genes with simulated DE in any group and true marker genes as the subset of DE  
19 genes that were DE in only a single group (Figure 6B). This procedure was repeated 20  
20 times with different random seeds to get some idea of variability and robustness of the  
21 method.

1



2

3 **Figure 6: Evaluation of SC3 results. Metrics for the evaluation of clustering (A) include the Rand**  
4 **index, Hubert and Arabie's adjusted Rand index (HA), Morey and Agresti's adjusted Rand index**  
5 **(MA), Fowlkes and Mallows index (FM) and the Jaccard index. Detection of differentially expressed**  
6 **and marker genes were evaluated (B) using accuracy, recall (true positive rate), precision and F1**  
7 **score (harmonic mean of precision and recall). All of the metrics are presented here as boxplots**  
8 **across the 20 simulations.**

9 Figure 6 shows the evaluation of SC3's clustering and gene identification on the  
10 simulated data. Five measures were used to evaluate the clustering: the Rand index  
11 (Rand), Hubert and Arabie's (HA) adjusted Rand index and Morey and Agresti's (MA)  
12 adjusted Rand index (both of which adjust for chance groupings), Fowlkes and Mallows  
13 index (FM) and the Jaccard index (Jaccard). All of these indices attempt to measure the  
14 similarity between two clusterings, in this case the clustering returned by SC3 and the  
15 true groups in the simulation. SC3 appears to identify clusters well in the majority of  
16 simulations, in some cases producing near perfect clustering. It may be interesting to  
17 examine individual cases further in order to identify when SC3 is able to perform better.  
18 Both the DE genes and marker genes identified by SC3 show a similar pattern across our  
19 classification metrics of accuracy, precision, recall and F1 score. On average  
20 approximately 3000 of the truly DE genes and 2600 of the true marker genes passed



1 SC3's automatic filtering (with additional non-DE genes). SC3 then detected around 84  
2 DE genes per simulation, along with 83 marker genes (median values). Precision (the  
3 proportion of identified genes that are true positives) is very high while recall (the  
4 proportion of true positives that were identified or true positive rate) is very low. This  
5 tells us that in this scenario SC3 is producing many false negatives but that the genes  
6 that it finds to be markers or DE are correct. This result is often desirable, particularly  
7 for marker genes.

8 While it is beyond the scope of this paper, clearly this evaluation could be extended,  
9 by including more clustering methods, more variations in simulation parameters and  
10 investigating why particular results are seen. However this data, and the code used to  
11 produce it, provides an example of how such an evaluation could be conducted using the  
12 simulations available in Splatter.

## 13 **Discussion and conclusions**

14 The recent development of single-cell RNA sequencing has spawned a plethora of  
15 analysis method and simulations can be a powerful tool for developing and evaluating  
16 them. Unfortunately, many current simulations of scRNA-seq data are poorly  
17 documented, not reproducible or fail to demonstrate similarity to real datasets. In  
18 addition, simulations created to evaluate a specific method can sometimes fall into the  
19 trap of having the same underlying assumptions as the method that they are trying to  
20 test. An independent, reproducible and flexible simulation framework is required in  
21 order for the analysis community to evaluate and develop sophisticated analysis  
22 methodologies.

23 Here we have developed Splatter, an independent framework for the reproducible  
24 simulation of scRNA-seq data. Splatter is available as an R package from Bioconductor  
25 and implements a series of simulation models. Splatter can easily estimate parameters

1 for each model from real data, generate synthetic datasets and quickly create a series of  
2 diagnostic plots comparing different simulations and datasets.

3 As part of Splatter we introduce our own simulation Splat. Splat builds on the  
4 gamma-poisson (or negative binomial) distribution commonly used to represent RNA-  
5 seq data, and adds high-expression outlier genes, library size distributions, a mean-  
6 variance trend and the option of expression-based dropout. Extensions to Splat include  
7 the simulation of more complex scenarios, such as multiple groups of cells with differing  
8 sizes and levels of differential expression, or differentiation trajectories with multiple  
9 paths and branches, and genes that change in non-linear ways.

10 We performed an evaluation of the five simulation models currently available in  
11 Splatter by comparing generated synthetic data to five published datasets. Overall Splat  
12 performed well, ranking highly on most metrics. However other simulations performed  
13 better for some metrics or better reproduced specific datasets. We found the Camp  
14 cerebral organoid dataset the most challenging to simulate, perhaps because of the  
15 complex nature of this sample with a large number of different cell types. In addition,  
16 this dataset (along with the Engel data) used a full-length protocol, which may contain  
17 additional noise compared to the UMI datasets [27].

18 One of the key features of scRNA-seq data is the high number of zero counts where  
19 no expression is observed for a particular gene in a particular cell. This can be especially  
20 challenging to simulate as not only must there be the correct number of zeros but they  
21 must be correctly distributed across genes and cells. We found that introducing dropout  
22 (in Splat) or zero-inflation (in Lun 2) was better at reproducing some, but not all,  
23 datasets. Together the results demonstrate that no simulation can accurately reproduce  
24 all scRNA-seq datasets and emphasises the variability in scRNA-seq data which arises  
25 from a complex set of biological (for example species, tissue type, cell type, treatment  
26 and cell cycle) and technical (for example platform, protocol, or processing) factors.

1 Non-parametric simulations that permute real data could potentially produce more  
2 realistic synthetic datasets but at the cost of flexibility in what can be simulated and  
3 knowledge of the underlying parameters.

4 Finally we demonstrated how Splatter could be used for the development and  
5 evaluation of analysis methods using the SC3 clustering method as an example.  
6 Splatter's flexible framework allowed us to quickly generate multiple test datasets,  
7 based on parameters from real data, and the information returned about the simulations  
8 gave us a truth to test against when evaluating the method. We found that SC3  
9 accurately clustered cells and was precise in identifying DE and marker genes.

10 The simulations available in Splatter are well documented, reproducible and  
11 independent of any particular analysis methods and Splatter's comparison functions  
12 make it easy to demonstrate how similar simulations are to real datasets. Splatter  
13 provides a framework for simulation models, makes existing scRNA-seq simulations  
14 accessible to researchers and introduces Splat, a new scRNA-seq simulation model. We  
15 hope that this framework empowers researchers to rapidly and rigorously develop new  
16 scRNA-seq analysis methods, ultimately leading to new discoveries in cell biology.

## 17 **Methods**

### 18 **Splat parameter estimation**

19 In order to easily generate a simulation that is similar to a given dataset Splatter  
20 includes functions to estimate the parameters for each simulation from real datasets.  
21 Just as with the simulations themselves the estimation procedures are based on what  
22 has been published and there is variation in how many parameters can be estimated for  
23 each model. We have given significant attention to estimating the parameters for the  
24 Splat simulation. The parameters that control the mean expression of each gene ( $\alpha$  and

1  $\beta$ ) are estimated by fitting a gamma distribution to the winsorized means of the library  
2 size normalised counts using the `fitdistrplus` package [28]. This is a basic normalisation  
3 where the counts in the original dataset are adjusted so that each cell has the same  
4 number of total counts (in this case the median across all cells) and any genes that are  
5 all zero are removed. We found that genes with extreme means affect the fit of the  
6 gamma distribution and that this effect was mitigated by winsorising to set the top and  
7 bottom 10 percent of values to the 10th and 90th percentiles respectively. Parameters  
8 for the library size distribution ( $\mu^L$  and  $\sigma^L$ ) are estimated in a similar way by fitting a  
9 log-normal distribution to the unnormalised library sizes.

10 The procedure for estimating expression outlier parameters is more complex.  
11 Taking the library size normalised counts, outliers are defined as genes where the mean  
12 expression is more than two MADs greater than the median mean gene expression level.  
13 The outlier probability  $\pi^O$  is then calculated as the proportion of genes that are outliers.  
14 Parameters for the outlier factors ( $\mu^O$  and  $\sigma^O$ ) are estimated by fitting a log-normal  
15 distribution to the ratio of the means of the outlier genes to the median mean gene  
16 expression level.

17 BCV parameters are estimated using the `estimateDispersion` function in the `edgeR`  
18 package [14]. When testing the estimation procedure on simulated datasets we  
19 observed that the `edgeR` estimate of common dispersion was inflated (Additional figure  
20 10), therefore we apply a linear correction to this value ( $\hat{\phi} = 0.1 + 0.25 \hat{\phi}_{\text{edgeR}}$ ).

21 The midpoint ( $x_0$ ) and shape ( $k$ ) parameters for the dropout function are estimated  
22 by fitting a logistic function to the relationship between the log means of the normalised  
23 counts and the proportion of samples that are zero for each gene (Additional figure 11).

## 24 Datasets

1 Each of the real datasets used in the comparison of simulations is publicly available.  
2 Raw FASTQ files for the Camp dataset were download from SRA (Accession SRP066834)  
3 and processed using a Bpipe (v0.9.9.3) [29] pipeline that examined the quality of reads  
4 using FastQC (v0.11.4), aligned the reads to the hg38 reference genome using STAR  
5 (v2.5.2a) [30] and counted reads overlapping genes in the Gencode V22 annotation  
6 using featureCounts (v1.5.0-p3) [31]. Matrices of gene by cell expression values for the  
7 Klein (Accession GSM1599500) and Zeisel (Accession GSE60361) datasets were  
8 downloaded from GEO. For the Tung dataset the matrix of molecules (UMIs) aligned to  
9 each gene available from <https://github.com/jdblischak/singleCellSeq> was used. This  
10 data is also available from GEO (Accession GSE77288). The Salmon [32] quantification  
11 files for the Engel dataset were download from the Conquer database  
12 (<http://imlspenticton.uzh.ch:3838/conquer/>) and converted to a gene by cell matrix  
13 using the tximport [33] package.

## 14 **Simulation comparison**

15 For each dataset the data file was read into R (v3.3.1) [34] and converted to a gene  
16 by cell matrix. Any genes that had zero expression in all cells or any missing values were  
17 filtered out and 200 cells were randomly selected without replacement. The parameters  
18 for each simulation were estimated from the selected cells and a synthetic dataset  
19 generated with 200 cells and the same number of genes as the real data. When  
20 estimating parameters for the Lun 2 and scDD simulations cells were randomly assigned  
21 to two groups. For the Splat and Lun 2 simulations both the regular and zero-inflated  
22 variants were used to simulate data. For the scDD simulation we set half of the genes to  
23 be “equally expressed” and the other half “equally proportioned”. The resulting seven  
24 simulations were then compared to the real data using Splatter’s comparison functions  
25 and plots showing the overall comparison produced. In order to compare simulations  
26 across the datasets summary statistics were calculated. For each of the basic metrics

1 (mean, variance, library size, zeros per gene and zeros per cell) the genes were sorted  
2 individually for each simulation and the difference from the sorted values and the real  
3 data calculated. When looking at the relationship between mean expression level and  
4 other metrics (variance, zeros per gene) genes in both the real and simulated data were  
5 sorted by mean expression and the difference between the metric of interest (eg.  
6 variance) calculated. The Median Absolute Deviation for each metric was then calculated  
7 and ranked for each dataset to give the rankings shown in Figure 4.

## 8 **Clustering evaluation**

9 Parameters for Splats simulations used in the example evaluation of SC3 were  
10 estimated from the Tung dataset. Twenty synthetic datasets were generated using these  
11 parameters with different random seeds. Each simulation had three groups of different  
12 sizes (200, 100 and 50 cells) and a probability of a gene being differentially expressed of  
13 0.1. Factors for differentially expressed genes were generated from a log-normal  
14 distribution with location parameter equal to -0.1 and scale parameter equal to 0.3. For  
15 each simulation the SC3 package was used to cluster cells with  $k = 3$  and asked to detect  
16 DE and marker genes, taking those with p-values less than 0.05. True DE genes were  
17 defined as genes where the simulated DE factor was not equal to one in one or more  
18 groups and markers genes as genes where the DE factor was not equal to one in a single  
19 group (and one in all others). Clustering metrics were calculated using the clues R  
20 package [35]. Metrics for evaluation of DE and marker genes were calculated by looking  
21 at the numbers of true negatives, true positive, false negatives and false positives.  
22 Metrics were aggregated across the 20 simulations and boxplots produced using the  
23 ggplot2 package [36].

24 Session information describing the packages used in all analysis steps is included as  
25 Additional file 13. The code and dataset files are available at  
26 <https://github.com/Oshlack/splatter-paper>.

## 1 **Declarations**

## 2 **Availability of data and materials**

3 The datasets analysed during the current study are available in from the  
4 repositories specified in the methods section or the repository for this paper,  
5 <https://github.com/Oshlack/splatter-paper>. The Splatter package is available from  
6 Bioconductor (<http://bioconductor.org/packages/splatter/>) and is being developed on  
7 Github (<https://github.com/Oshlack/splatter>).

## 8 **Competing interests**

9 The authors declare no competing interests.

## 10 **Funding**

11 Luke Zappia is supported by an Australian Government Research Training Program  
12 (RTP) Scholarship. Alicia Oshlack is supported through a National Health and Medical  
13 Research Council Career Development Fellowship APP1126157.

## 14 **Authors' contributions**

15 LZ developed the software and performed the analysis. BP contributed to the  
16 statistics and supervision. AO oversaw all aspects of the project. All authors contributed  
17 to drafting the manuscript.

## 18 **Acknowledgements**

## 19 **References**

- 20 1. Goodwin S, McPherson JD, Richard McCombie W. Coming of age: ten years of next-  
21 generation sequencing technologies. *Nat. Rev. Genet.* 2016;17:333–351.
- 22 2. Ozsolak F, Milos PM. RNA sequencing: advances, challenges and opportunities. *Nat.*  
23 *Rev. Genet.* 2011;12:87–98.
- 24 3. Tang F, Barbacioru C, Wang Y, Nordman E, Lee C, Xu N, et al. mRNA-Seq whole-  
25 transcriptome analysis of a single cell. *Nat. Methods.* 2009;6:377–382.

- 1 4. Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3:  
2 consensus clustering of single-cell RNA-seq data. *Nat. Methods*. 2017;
- 3 5. Lin P, Troup M, Ho JWK. CIDR: Ultrafast and accurate clustering through imputation  
4 for single-cell RNA-seq data. *Genome Biol*. 2017;18:59.
- 5 6. Satija R, Farrell JA, Gennert D, Schier AF, Regev A. Spatial reconstruction of single-cell  
6 gene expression data. *Nat. Biotechnol*. 2015;33:495–502.
- 7 7. Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, et al. The dynamics  
8 and regulators of cell fate decisions are revealed by pseudotemporal ordering of single  
9 cells. *Nat. Biotechnol*. 2014;32:381–386.
- 10 8. duVerle DA, Yotsukura S, Nomura S, Aburatani H, Tsuda K. CellTree: an  
11 R/bioconductor package to infer the hierarchical structure of cell populations from  
12 single-cell RNA-seq data. *BMC Bioinformatics*. 2016;17:363.
- 13 9. Juliá M, Telenti A, Rausell A. Sincell: an R/Bioconductor package for statistical  
14 assessment of cell-state hierarchies from single-cell RNA-seq. *Bioinformatics*.  
15 2015;31:3380–3382.
- 16 10. McCarthy DJ, Campbell KR, Lun ATL, Wills QF. Scater: pre-processing, quality control,  
17 normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics*. 2017;
- 18 11. Lun ATL, Bach K, Marioni JC. Pooling across cells to normalize single-cell RNA  
19 sequencing data with many zero counts. *Genome Biol*. 2016;17:1–14.
- 20 12. Lun ATL, Marioni JC. Overcoming confounding plate effects in differential expression  
21 analyses of single-cell RNA-seq data. *Biostatistics*. 2017;
- 22 13. Korthauer KD, Chu L-F, Newton MA, Li Y, Thomson J, Stewart R, et al. A statistical  
23 approach for identifying differential distributions in single-cell RNA-seq experiments.  
24 *Genome Biol*. 2016;17:222.
- 25 14. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for  
26 differential expression analysis of digital gene expression data. *Bioinformatics*.  
27 2010;26:139–140.
- 28 15. Anders S, Huber W. Differential expression analysis for sequence count data.  
29 *Genome Biol*. 2010;11:R106.
- 30 16. Korthauer K. scDD Vignette [Internet]. 2017. Available from:  
31 <https://bioconductor.org/packages/devel/bioc/vignettes/scDD/inst/doc/scDD.pdf>
- 32 17. McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-  
33 Seq experiments with respect to biological variation. *Nucleic Acids Res*. 2012;40:4288–  
34 4297.
- 35 18. Law CW, Chen Y, Shi W, Smyth GK. voom: Precision weights unlock linear model  
36 analysis tools for RNA-seq read counts. *Genome Biol*. 2014;15:R29.
- 37 19. Pierson E, Yau C. ZIFA: Dimensionality reduction for zero-inflated single-cell gene  
38 expression analysis. *Genome Biol*. 2015;16:241.



- 1 20. Tung P-Y, Blischak JD, Hsiao CJ, Knowles DA, Burnett JE, Pritchard JK, et al. Batch  
2 effects and the effective design of single-cell gene expression studies. *Sci. Rep.*  
3 2017;7:39921.
- 4 21. Andrews TS, Hemberg M. Modelling dropouts allows for unbiased identification of  
5 marker genes in scRNASeq experiments. *bioRxiv.* 2016;
- 6 22. Kivioja T, Vähärautio A, Karlsson K, Bonke M, Enge M, Linnarsson S, et al. Counting  
7 absolute numbers of molecules using unique molecular identifiers. *Nat. Methods.*  
8 2012;9:72–74.
- 9 23. Camp JG, Badsha F, Florio M, Kanton S, Gerber T, Wilsch-Bräuninger M, et al. Human  
10 cerebral organoids recapitulate gene expression programs of fetal neocortex  
11 development. *Proc. Natl. Acad. Sci. U. S. A.* 2015;112:15672–15677.
- 12 24. Engel I, Seumois G, Chavez L, Samaniego-Castruita D, White B, Chawla A, et al. Innate-  
13 like functions of natural killer T cell subsets result from highly divergent gene programs.  
14 *Nat. Immunol.* 2016;17:728–739.
- 15 25. Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, et al. Droplet  
16 barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell.*  
17 2015;161:1187–1201.
- 18 26. Zeisel A, Muñoz-Manchado AB, Codeluppi S, Lönnerberg P, La Manno G, Juréus A, et  
19 al. Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-  
20 cell RNA-seq. *Science.* 2015;347:1138–1142.
- 21 27. Phipson B, Zappia L, Oshlack A. Gene length and detection bias in single cell RNA  
22 sequencing protocols. *F1000Research.* 2017;6.
- 23 28. Delignette-Muller M, Dutang C. *fitdistrplus: An R Package for Fitting Distributions.* *J.*  
24 *Stat. Softw.* 2015;64:1–34.
- 25 29. Sadedin SP, Pope B, Oshlack A. *Bpipe: a tool for running and managing*  
26 *bioinformatics pipelines.* *Bioinformatics.* 2012;28:1525–1526.
- 27 30. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast  
28 universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.
- 29 31. Liao Y, Smyth GK, Shi W. *featureCounts: an efficient general purpose program for*  
30 *assigning sequence reads to genomic features.* *Bioinformatics.* 2014;30:923–930.
- 31 32. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-  
32 aware quantification of transcript expression. *Nat. Methods.* 2017;14:417–419.
- 33 33. Sonesson C, Love MI, Robinson MD. Differential analyses for RNA-seq: transcript-level  
34 estimates improve gene-level inferences. *F1000Research.* 2015;4:1521.
- 35 34. R Core Team. *R: A Language and Environment for Statistical Computing [Internet].*  
36 *Vienna, Austria: R Foundation for Statistical Computing; 2016. Available from:*  
37 *<https://www.R-project.org/>*
- 38 35. Chang F, Qiu W, Zamar R, Lazarus R, Wang X. *clues: An R Package for Nonparametric*  
39 *Clustering Based on Local Shrinking.* *J. Stat. Softw.* 2010;33:1–16.

1 36. Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer New York; 2010.

## 2 **Additional files**

3 **File name:** additional\_figures\_1-11.pdf

4 **File format:** PDF

5 **Title:** Additional figures

6 **Description:** Additional figures including: diagrams of other simulation models, Splatter  
7 comparison output for all datasets, example non-linear gene, dispersion estimate  
8 correction and mean-zeros fit.

9  
10 **File name:** additional12\_mads.csv

11 **File format:** CSV

12 **Title:** Table of MADs

13 **Description:** Table of the Median Absolute Deviations used to produce Figure 4 in CSV  
14 format.

15

16 **File name:** additional13\_sessionInfo.pdf

17 **File format:** PDF

18 **Title:** Session information

19 **Description:** Details of the R environment and packages used for analysis.

20