

Normalizing single-cell RNA sequencing data: challenges and opportunities

Catalina A Vallejos^{1-4,10}, Davide Risso^{5,9,10}, Antonio Scialdone^{2,10}, Sandrine Dudoit^{5,6} & John C Marioni^{2,7,8}

Single-cell transcriptomics is becoming an important component of the molecular biologist's toolkit. A critical step when analyzing data generated using this technology is normalization. However, normalization is typically performed using methods developed for bulk RNA sequencing or even microarray data, and the suitability of these methods for single-cell transcriptomics has not been assessed. We here discuss commonly used normalization approaches and illustrate how these can produce misleading results. Finally, we present alternative approaches and provide recommendations for single-cell RNA sequencing users.

Single-cell RNA sequencing (scRNA-seq) has transformed the field of transcriptomics by making it possible for researchers to address fundamental questions that could not be tackled by bulk-level experiments¹. Examples include the study of tumor heterogeneity, the identification of novel cell types, and the understanding of cell fate decisions during early embryo development²⁻⁵.

Recent literature has highlighted the need for new computational methods to address the complex features, such as sparsity and technical noise, that characterize scRNA-seq data⁶⁻¹¹. Despite this, little attention has been given to normalization—a critical step in the analysis pipeline that adjusts for unwanted biological and technical effects that can mask the signal of interest. Instead, most tools developed for scRNA-seq rely on normalized expression measures obtained from methods developed for bulk RNA-seq or even microarray data analysis. However, whether these approaches are suitable for single-cell transcriptomics has not been rigorously discussed.

In this Perspective, we address normalization and focus on the most widely used strategy, global scaling, which attempts to remove cell-specific systematic biases by scaling expression measures within each cell by a constant factor. Within this framework, we illustrate that the use of bulk-based normalization methods can have serious adverse consequences for downstream analysis, such as the detection of highly variable genes before clustering. Such problems are exacerbated by the high levels of technical noise and dropout typical of scRNA-seq. We also discuss the use of extrinsic spike-in sequences (e.g., ref. 12) for normalization. To conclude, we summarize state-of-the-art methods for scRNA-seq normalization including integrated strategies, where normalization is intrinsic to a specific method, and generic tools, which provide normalized data that can be used as input to any downstream analysis pipeline.

From bulk samples to single-cell resolution

Bulk microarray and RNA-seq experiments measure gene expression levels as averages across thousands of cells. While this allows the characterization of population-level differences in overall expression, single-cell-level experiments are required to better understand the dynamics of gene expression patterns. scRNA-seq experiments can reveal heterogeneity within populations of cells. However, the additional insights offered by scRNA-seq come at the cost of more challenging experimental protocols¹³ and higher data complexity¹⁰.

A prominent feature of scRNA-seq is the sparsity of the data—i.e., the high proportion of zero read

¹MRC Biostatistics Unit, Cambridge Institute of Public Health, Cambridge, UK. ²EMBL-European Bioinformatics Institute, Wellcome Genome Campus, Cambridge, UK. ³The Alan Turing Institute, British Library, London, UK. ⁴Department of Statistical Science, University College London, London, UK. ⁵Division of Biostatistics, School of Public Health, University of California, Berkeley, Berkeley, California, USA. ⁶Department of Statistics, University of California, Berkeley, Berkeley, California, USA. ⁷Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Cambridge, UK. ⁸Wellcome Trust Sanger Institute, Wellcome Genome Campus, Cambridge, UK. ⁹Present address: Division of Biostatistics and Epidemiology, Department of Healthcare Policy and Research, Weill Cornell Medicine, New York, New York, USA. ¹⁰These authors contributed equally to this work. Correspondence should be addressed to S.D. (sandrine@stat.berkeley.edu) or J.C.M. (marioni@ebi.ac.uk).

PERSPECTIVE

counts^{7,8,14}. This ‘zero inflation’ arises for both biological reasons (e.g., subpopulations of cells or transient states where a gene is not expressed) and technical reasons (e.g., dropouts, where a gene is expressed but not detected through sequencing). Besides the dropout effect, technical noise in scRNA-seq is also reflected by high variability between technical replicates, even for genes with medium or high levels of expression⁶. Additionally, by capturing individual cells from potentially very different cell types, scRNA-seq data are highly heterogeneous, even in the absence of the technical biases discussed above. Consequently, several assumptions made when analyzing bulk RNA-seq data do not always apply in the context of scRNA-seq.

Systematic biases in scRNA-seq data sets

Data normalization strategies must capture biases that are specific to the technology of interest. For example, two-channel microarrays require normalization to account for differences in dye balance related to intensity and spatial position on the array¹⁵. By contrast, in sequencing assays, read counts must be adjusted to control for a variety of biases, including sequencing depth^{16,17}. Additionally, in any of these high-throughput assays, one needs to account for possibly more complex and putatively unknown effects, collectively known as ‘batch effects’^{18–20}.

While scRNA-seq analysis pipelines routinely include a normalization step, the sources of the systematic biases that this step captures are assay specific. To illustrate this, we focus on the Illumina sequencing platform, using a simple experimental setup where gene expression is measured in a homogeneous population of cells. The discussion below applies to whole-transcript scRNA-seq as well as to 3’ sequencing protocols and unique molecular identifier (UMI)²¹-based approaches that use barcodes to obtain molecular counts.

RNA-seq experiments are inherently stochastic, with reads being randomly sampled from a pool of amplified cDNA molecules. Typically, the quantity of interest is the expression level of each gene; the relative abundance of mRNA molecules from each gene within the overall population of mRNA molecules in each cell. There are several experimental sources of systematic biases that can affect measurements of gene expression, including gene- and cell-specific features (Fig. 1). Accordingly, we distinguish between two types of normalization—within-sample normalization, which removes gene-specific biases (e.g., GC content), and between-sample normalization, which adjusts for effects related to distributional differences in read counts between cells (e.g., sequencing depth). In this Perspective we focus on the latter type of normalization, particularly on global scaling, the most common approach in the literature.

Global-scaling normalization methods assume that the expected value of the read count for a gene in a cell is proportional to a gene-specific expression level and a cell-specific scaling factor (also known as a size factor), which is an unknown (random) variable representing nuisance effects (Box 1). Such nuisance effects include, for example, differences in reverse transcription (RT) efficiency and cell-intrinsic properties, like endogenous mRNA content. Note that, unlike endogenous mRNA content, which is fixed for a given cell, the remaining effects listed in Box 1 are random (if the same cell could be processed twice, these quantities would vary). This implies that scaling factors are inherently random. Nevertheless, most existing methods treat these scaling factors as fixed factors and/or model offsets.

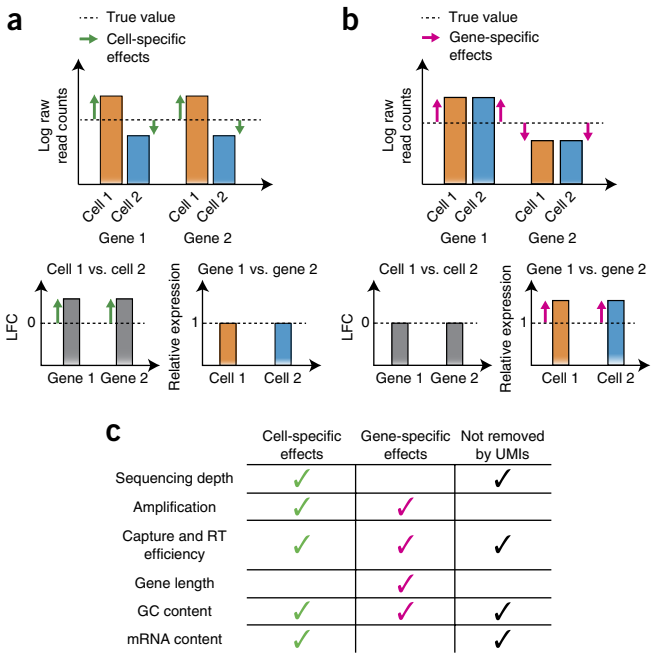


Figure 1 | Cell- and gene-specific effects in RNA-seq experiments. (a) Schematic representation of cell-specific effects. The top panel shows a pair of cells expressing two genes at the same levels. When RNA-seq is performed, cell-specific effects introduce a bias in the estimated log fold change (LFC) computed on raw read counts (bottom left panel). (b) Schematic representation of gene-specific effects. The two cells and true gene levels are the same as in a, but now gene-specific effects are shown to bias the estimation of relative gene expression (bottom right panel). In real situations, both cell- and gene-specific effects are present. (c) Table noting main cell- and/or gene-specific effects and whether these are removed by unique molecular identifiers (UMIs).

Depending on the experimental protocol, some cell-specific effects cancel out between cells. For example, if library quantification is accurate, the dilution step can remove biases related to differences in the predilution number of amplified cDNA molecules per cell. UMI-based protocols in principle remove biases related to amplification and sequencing depth, since multiple reads associated with the same UMI are collapsed into a unique count (Fig. 1c). However, this is only true if all libraries are sequenced to saturation (i.e., each uniquely tagged molecule is observed at least once). If not, some UMI-tagged cDNA molecules will be lost and, since sequencing depth randomly fluctuates between cells, systematic cell-specific differences between molecule counts can occur. Finally, since UMIs are ligated to each molecule during RT, they can neither account for differences in capture efficiency before the RT step nor for differences in cellular mRNA content.

Normalizing scRNA-seq data sets

scRNA-seq data sets are typically normalized using global-scaling normalization methods developed for bulk RNA-seq data analysis^{7,8,22}. In principle, global-scaling factors can be treated as (nuisance) model parameters and jointly estimated with other quantities of interest such as gene-specific expression levels. However, this approach is computationally intensive and necessarily tailored to a specific model (e.g., refs. 19,23).

An alternative—and widespread—approach is to compute normalized expression measures based on scaling factor estimates

BOX 1 GLOBAL-SCALING NORMALIZATION FOR scRNA-seq DATA SETS

RNA-seq experiments are inherently stochastic, with reads being randomly sampled from a pool of amplified cDNA molecules. Accordingly, let X_{ij} denote a random variable representing the read count of gene i in cell j . Typically, the parameter of interest is the expression level of each gene (see left panel), i.e., the relative abundance of mRNA molecules for a gene within the population of mRNA molecules in each cell. For the sake of simplicity, we consider here the case of a homogeneous cell population.

Intuitively, a first effect captured through the scaling factor s_j is the endogenous mRNA content n_j , the total number of mRNA molecules per cell (middle panel). Indeed, even within a homogeneous population, n_j can vary across cells. Furthermore, after cell lysis, only a fraction (F_j) of these n_j molecules, are captured and reverse transcribed into cDNA. Consequently, only $n_j \times F_j$ cDNA molecules can potentially be amplified and subsequently sequenced. Critically, the capture and reverse transcription efficiency F_j varies between cells; this introduces cell-to-cell variability that should also be handled by s_j .

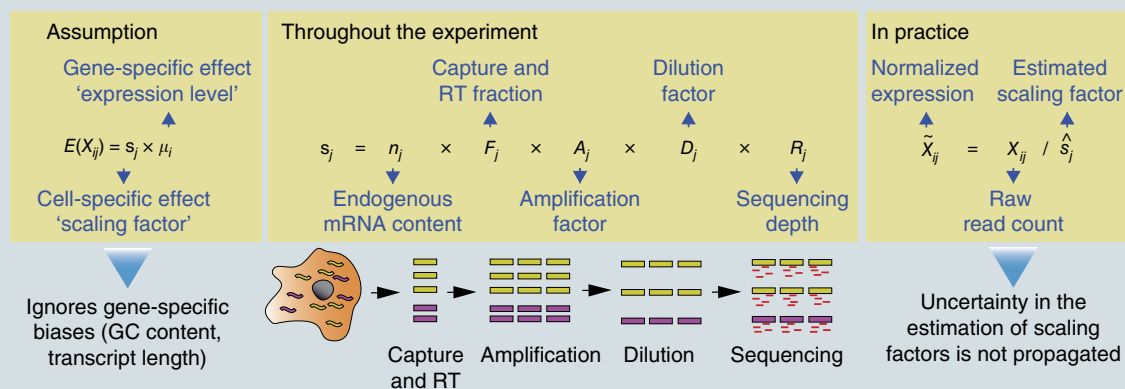
Subsequently, because of the minute amount of genetic material contained in a cell, this pool of $n_j \times F_j$ cDNA molecules must be amplified before sequencing library preparation. Variability in amplification efficiency can introduce cell- and gene-specific biases in the measurement of expression levels. We denote the cell-specific amplification factor as A_j , such that amplification leads to a pool of $n_j \times F_j \times A_j$ molecules.

Unlike microarray experiments, RNA-seq is inherently competitive, meaning that a fixed number of reads are distributed

between genes. Given this, the amplified pools are subsequently diluted by a cell-specific factor D_j , so that there are $n_j \times F_j \times A_j \times D_j$ amplified cDNA molecules to be sequenced. In principle, the dilution factor D_j can be set so that a library contains the same number of molecules from each cell by carrying out a library quantification step and setting $D_j = m / (n_j \times F_j \times A_j)$, where m is the desired number of molecules per cell. Alternatively, each cell can contribute the same volume of amplified cDNA solution to the library, such that each library will contain a different number of amplified cDNA molecules if the concentration of the solution varies between cells. In this case, $D_j = d$, where d is the proportion of amplified molecules used to prepare the sequencing library. This decision is critical for interpreting the scaling factor s_j , since it affects the number of molecules that are available for sequencing and, consequently, the scale of cell-specific read counts.

Finally, the number of sequenced reads per molecule from each cell (sequencing depth), R_j , also varies stochastically. Consequently, by considering all the above factors, we expect to observe $n_j \times F_j \times A_j \times D_j \times R_j$ reads from cell j . Hence, even within the same sequencing lane, differences in sequencing depth introduce cell-specific artifacts that will be incorporated into the global scaling factor s_j .

While the above discussion assumes a homogeneous population of cells, this interpretation of scaling factors is still valid for more realistic scenarios—with heterogeneous populations—under specific assumptions, such as that the majority of genes is not differentially expressed or that there are roughly equal numbers of upregulated and downregulated genes.



obtained during a preprocessing step (Box 1). Downstream analyses, such as clustering or differential expression, are then typically based on normalized measures (either directly or by treating the estimated scaling factors as model offsets), ignoring uncertainty related to the scaling factor estimation. While this strategy is common, there is no consensus on how to estimate the scaling factors; some popular choices are summarized below. However, all approaches share the same motivation—to bring cell-specific measures onto a common scale by standardizing a quantity of interest (e.g., total read counts per sample) across cells while assuming, for example, that most genes are not differentially expressed.

An intuitive and popular method is reads per million (RPM), which standardizes the total number of reads between cells; it is also referred to as library-size normalization and is related to RPKM²⁴ and TPM²⁵. However, these estimates can be dominated by a handful of highly expressed genes, and this can bias downstream results^{16,26}. Another possibility is to use upper quartile (UQ) normalization, which defines scaling factor estimates as proportional to the 75th percentile of the distribution of counts within each cell¹⁶. An extension of this idea (albeit outside the universe of global-scaling normalization) is full quantile (FQ) normalization, where all the quantiles of cell-specific counts are matched to

a reference distribution¹⁶. However, quantile-based normalization methods are problematic in scRNA-seq due to the high frequency of zero counts typically observed. In practice, this can lead to scaling factor estimates being set to 0 in UQ normalization, while the large number of zeros leads to ties in the gene ranking needed by FQ normalization, rendering its interpretation more difficult.

Alternative approaches have been developed in the context of bulk RNA-seq analyses. Two popular methods are DESeq²⁶ and trimmed mean of *M* values (TMM)²⁷ normalization. DESeq defines scaling factor estimates based on a pseudoreference sample, which is built with the geometric mean of gene counts across all cells. TMM trims away extreme log fold changes to normalize the counts based on the remaining set of nondifferentially expressed genes. Critically, zero inflation is problematic for DESeq, as the calculation of the pseudoreference sample is only well defined for the potentially very small set of genes with at least one read in every cell.

In the context of bulk RNA-seq, the performance of global-scaling methods was reviewed by Dillies *et al.*¹⁷, who used a variety of case studies and simulated data sets to suggest that DESeq and TMM outperform other methods. However, the performance of DESeq and TMM in the context of scRNA-seq has been given little attention.

Comparing bulk-based approaches: a case study

The use of different normalization methods can alter the results of downstream analysis. To illustrate this, we applied the three widely used normalization techniques RPM, DESeq, and TMM to a publicly available data set (Fig. 2). This data set consists of gene expression measures for 933 mouse embryonic stem cells (mESCs)²⁸. These cells were processed using a droplet-based protocol, which yielded UMI-based counts.

Overall, we observed substantial differences between the methods regarding scaling factor estimation (Fig. 2a, upper right panels). First, because of zero inflation, DESeq scaling factors are based on only 115 genes. As expected, this results in less stable estimation of the scaling factors. Moreover, we observed that—with respect to RPM (and DESeq)—TMM tends to underestimate and overestimate large and small scaling factors, respectively (this is in line with the simulation results in ref. 14). This is largely because of the sparsity of the data, with the differences between methods increasing for cells where more zero counts are observed (Fig. 2b, bottom panel).

Crucially, we observed that differences in scaling factor estimation affect gene-specific estimates of variability. This is illustrated using the squared coefficient of variation (CV^2) of the normalized expression measures per gene (Fig. 2a, lower left panels). Thus, analyses designed to uncover heterogeneity within the data are also distorted. For example, studies often start with the selection of highly variable genes (HVGs) to reduce the dimensionality of the data before clustering or other analyses. HVG selection is sensitive to the choice of normalization, and less than a third of HVGs are shared across all normalization methods (Fig. 2c).

We performed the same analyses on additional data sets, and we showed that these issues are likely general and inherent to scRNA-seq data (Supplementary Data 1). As expected, differences between scaling factors, and consequently between the lists of HVGs, are more pronounced in data sets with low sequencing depth. This is critical, as several modern experimental protocols (e.g., droplet-based methods) use shallow sequencing with fewer

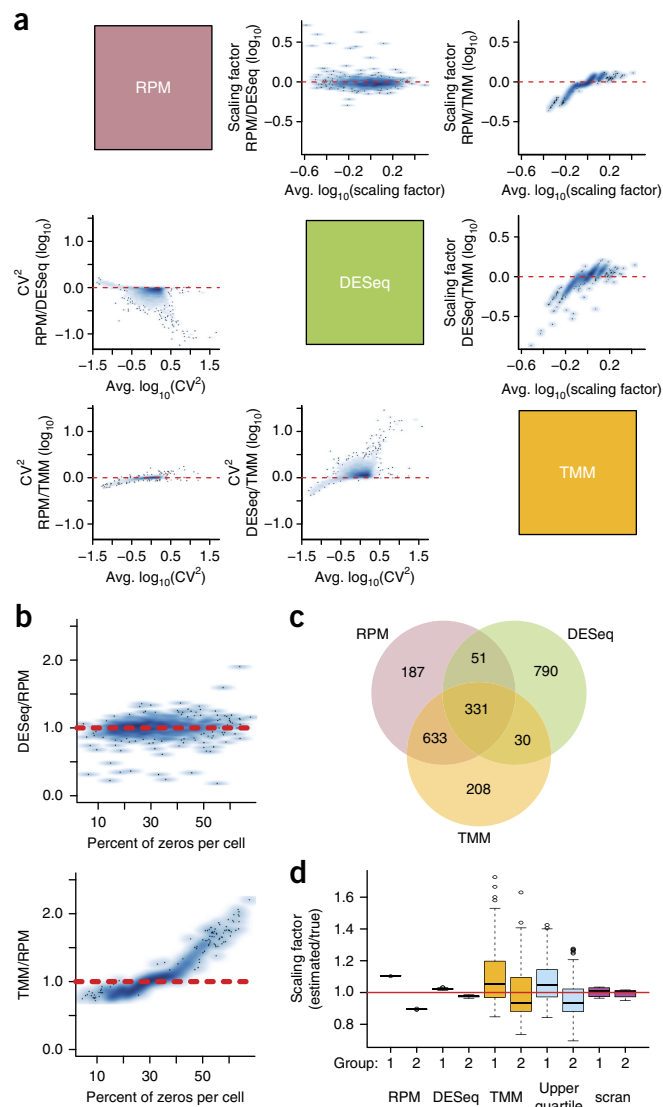


Figure 2 | Comparison of bulk-based normalization methods in real and simulated data sets. **(a)** Mean-difference plot comparing the estimated scaling factors (upper diagonal panels) and CV^2 of normalized counts (lower diagonal panels) for the data set published in ref. 28. Avg., average. **(b)** Ratio of estimated scaling factors versus proportion of zero counts (zeros) per cell for the data set²⁸. **(c)** Top 10% most variable genes identified after normalizing the data set²⁸ with three different methods. Additional data sets are analyzed in **Supplementary Data 1**. **(d)** Ratio between the estimated and the true scaling factors for the most widely used bulk-based normalization methods and a method specifically designed for scRNA-seq ('scran')¹⁴ in a simulated data set consisting of two groups of cells. See **Supplementary Data 2** for the simulation strategy and additional simulations.

than 50,000 reads per cell in order to profile a large number of cells. While shallow sequencing has been shown to allow discovery and classification of cell types in complex tissues^{29–31}, the results of more refined analyses (e.g., pseudotime ordering³²) can be distorted by differences between normalization methods.

Given the lack of ground truth in real data, we cannot determine which normalization method, if any, correctly estimates the scaling factors. To shed some light onto the relative merits of each method, we turn to simulations (**Supplementary Data 2**).

We simulated two groups of cells with varying numbers of differentially expressed genes. When the data are simulated with symmetric differential expression, all methods lead to unbiased estimates. However, with asymmetric differential expression, bulk-based methods lead to biased estimates of the scaling factors (see Fig. 2d for an example with 80% upregulated and 20% down-regulated genes in group 1; see **Supplementary Data 2** for other settings). This suggests that great care should be used if bulk-based global-scaling methods are applied to scRNA-seq data.

State-of-the-art methods

Bulk-based normalization methods are widely applied to scRNA-seq data sets, despite the problems outlined above. However, normalization methods that are specifically tailored to scRNA-seq data sets have recently been introduced. Below, we summarize state-of-the-art methods, provide practical recommendations to scRNA-seq users, and motivate the development of new methodology to address unresolved issues.

We distinguish between two different approaches. First, we consider ‘bespoke methods’ that use prenormalized expression measures (possibly using methods developed for bulk RNA-seq) and account for artifacts specific to scRNA-seq in downstream models. In the context of differential expression analyses, two examples of this approach are SCDE⁷ and MAST⁸. To attenuate the effect of technical variation in downstream analysis, SCDE introduces a two-component mixture model to capture drop-out events and events where a transcript is faithfully amplified. Alternatively, MAST uses the fraction of genes that are detectably expressed in each cell as a proxy for both technical and biological sources of variation. MAST uses a hurdle model where the expression measure of a detected gene is modeled by linear regression and the probability of detection by logistic regression.

A second strategy for normalizing scRNA-seq data sets is to use ‘generic methods’ that yield normalized expression measures that can be used as input in any subsequent analyses (e.g., refs. 32–34). A recent example of such an approach is *scrn*, which pools multiple cells in order to estimate cell-specific size factors more robustly in the presence of zero inflation and unbalanced differential expression of genes across groups of cells (Fig. 2d and ref. 14). In principle, BASiCS^{11,23} also provides a generic normalization tool, but its implementation has been coupled with specific downstream analysis.

We tested two methods motivated by scRNA-seq data, BASiCS and *scrn*, on recently published data sets; and we found that, unlike bulk-based methods, they led to very similar results in terms of scaling factor estimation and HVG selection (**Supplementary Data 3**). This likely derives from greater robustness to features of scRNA-seq data compared to bulk-based approaches. Other recent examples of normalization methods specifically designed for scRNA-seq include GRM³⁵ and SAMstr³⁶, which both rely on spike-ins, and SCnorm³⁷, which uses quantile regression to group genes with similar dependence on sequencing depth and to estimate different scaling factors for each group. However, it should be noted that GRM is not a between-sample normalization method, but rather a method to denoise gene expression levels within each cell. In addition, Qiu *et al.*³⁸ proposed the Census algorithm to convert relative RNA-seq measurements to relative transcript counts. The Census algorithm can be considered a normalization method, since it rescales TPM values by dividing them by the estimated total number of mRNA molecules.

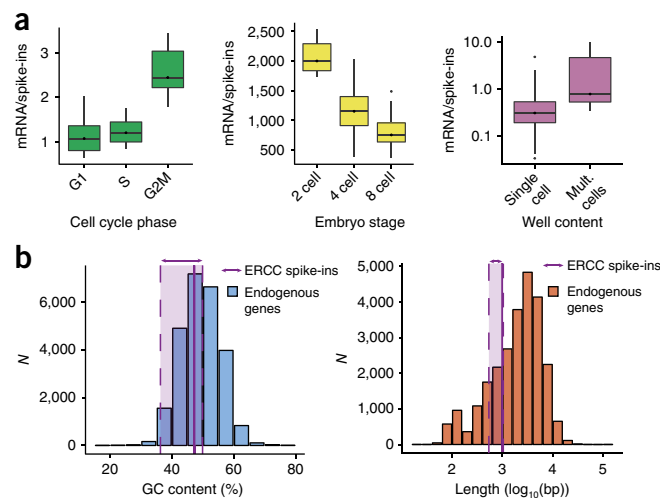


Figure 3 | ERCC spike-ins can be used to estimate mRNA content. (a) Ratio between the number of reads mapped to intrinsic genes and the number of reads mapped to ERCC spike-ins in data sets from refs. 33,41,42 (left, central and right panels, respectively). Mult., multiple. (b) Distributions of GC content (left panel) and length (right panel) for mouse genes with at least one count in one cell in the data set published in ref. 46. The purple areas show the interquartile ranges of GC content and length for ERCC spike-ins, with the medians marked by vertical purple continuous lines.

Finally, we note that, although the various global-scaling methods rely on different assumptions, these methods all fail if the number or fold change of differentially expressed genes across the cell population is too high. One strategy to alleviate this issue is to precluster the cells into smaller, more homogeneous sets (using, e.g., rank-based clustering methods, which are unaffected by global-scaling normalization). Normalization can then be performed separately for each cluster before between-cluster normalization to calculate cluster-specific offsets. This approach is used in the *scrn* method and has been shown to yield more accurate estimates of scaling factors¹⁴, as is also suggested by our simulations (Fig. 2d).

Spike-in sequences and normalization

The scaling factors introduced in **Box 1** cannot distinguish between technical biases and genuine biological differences between cells such as total mRNA content. Jiang *et al.*¹² discussed the benefits of exploiting a set of synthetic control genes—with constant expression level across all samples—to disentangle these effects in bulk RNA-seq. Extrinsic control genes have also been used in the context of scRNA-seq^{6,23,33,35,36}; spike-in sequences are added to each cell’s lysate in a theoretically constant and known amount. The most commonly used set of spike-ins is the set of 92 External RNA Control Consortium (ERCC) molecules¹². Other examples include the eight synthetic mRNAs deployed in ref. 39 and the whole-transcriptome HeLa RNA spike-in used in ref. 6. It is important to understand the utility of synthetic spike-in sequences in the context of global-scaling normalization.

One critical assumption underlying the use of spike-in sequences is that the technical effects summarized in **Box 1** equally affect the intrinsic and the extrinsic genes. If this assumption holds, additional technical scaling factors can be defined to capture these shared technical effects⁶. Thus, for any given cell,

the ratio between the scaling factor described in **Box 1** and the technical scaling factor defined above is equal to the endogenous mRNA content of the cell. As a corollary, normalization based solely on spike-in-derived scaling factors does not remove differences in endogenous mRNA content between cells, and further normalization is required to remove this effect.

This suggests that spike-in sequences can be used to obtain estimates of endogenous mRNA content per cell. At a coarse level, this is reflected in several scRNA-seq data sets (**Fig. 3a**), consistent with previously described bulk RNA-seq studies⁴⁰. Here, we look at three different data sets^{33,41,42}, for which we can stratify samples according to their expected mRNA content. The ratio of mRNA/spike-in read counts correctly indicates that mRNA content increases as mESCs progress along the cell cycle³³ and decreases across blastomeres in early mouse embryos at two-, four-, and eight-cell stages⁴¹ on account of their size difference. Analogously, in an experiment on the Fluidigm C1 instrument, wells including multiple cells are characterized by a higher mRNA content than wells where single cells are captured⁴².

However, using spike-in sequences remains challenging. In particular, calibrating the added number of spike-in molecules is nontrivial and depends on intrinsic characteristics of the studied cells, such as endogenous mRNA content. Poor calibration can invalidate the utility of the spike-ins as control genes; too many spike-ins can overwhelm signal from the intrinsic genes, while the majority of spike-in sequences can be unusable in downstream analysis if too few spike-in molecules are added¹⁹.

Additional issues arise for specific sets of spike-in sequences. In particular, for the widely used set of ERCC controls, the extreme range of concentration of spike-in molecules⁴⁰ prevents the use of the entire ERCC set in scRNA-seq. Typically, only half of the spike-in molecules are detected, and the proportion of reads mapped to the spike-in sequences may be extremely variable (**Fig. 3a**).

Moreover, potential biases in the mRNA enrichment process that are related to gene length and GC content imply that, overall, technical effects may be different for the ERCC spike-in sequences and the intrinsic genes. In fact, the ERCC set does not reflect the mammalian transcriptome in terms of gene length and GC content (**Fig. 3b**). Moreover, researchers have shown¹⁹ that ERCC spike-in signal can vary considerably between technical replicate samples. Consequently, estimates of endogenous mRNA content derived using ERCC spike-ins have large measurement uncertainties⁴⁰.

Developing a set of spike-ins specifically tailored for scRNA-seq experiments could overcome some of these limitations. Ideally, this set would closely resemble intrinsic genes in terms of the distribution of GC content, total length, and polyA tail length. Ongoing efforts to develop this set are illustrated by a recent call from the National Institute of Standard and Technology (<https://federalregister.gov/a/2015-19742>) to design an improved set of controls, which should (i) mimic endogenous RNA and (ii) not interfere with the measurement of endogenous RNA. More recently, Hardwick *et al.*⁴³ introduced ‘sequins’ (sequencing spike-ins)—a set of extrinsic spike-ins designed for bulk RNA-seq experiments.

DISCUSSION

One aim of this Perspective is to provide a straightforward understanding of the sources of variation that can be captured through global-scaling normalization in the context of scRNA-seq.

Case studies and simulated data sets indicated that a direct application of bulk RNA-seq normalization methods is not appropriate in the context of scRNA-seq, where—on account of biological heterogeneity and technical artifacts—we typically observe more heterogeneous and sparser data sets. We illustrated that the choice of the normalization method affects downstream analyses, such as HVG detection, that aim to uncover heterogeneity within the data. Although spike-in sequences carry some caveats, they can help disentangle technical artifacts from differences in endogenous mRNA content between cells. These differences occur in several contexts, such as the whole-transcriptome upregulation induced by elevated expression of the c-Myc transcription factor⁴⁴.

A variety of scRNA-seq-tailored methods that outperform bulk strategies have recently been proposed. Despite this, bulk-motivated approaches remain widely used in practice. We therefore suggest that scRNA-seq users update their analysis pipelines—matching advances in technology—to take full advantage of the rich information provided by scRNA-seq data sets. Finally, while the issue of identifying the best method for normalizing scRNA-seq data has not yet been fully resolved, many efforts are underway to develop additional robust and effective normalization techniques and to systematically assess their performance on individual data sets^{37,38,45}.

Data availability statement. We used data sets previously published in the referenced citations.

Note: Any Supplementary Information and Source Data files are available in the online version of the paper.

ACKNOWLEDGMENTS

We thank several members of the Marioni laboratory (European Molecular Biology Laboratory - European Bioinformatics Institute, EMBL-EBI; Cancer Research UK - Cambridge Institute, CRUK-CI) for support and discussions throughout the preparation of this manuscript. In particular, we are grateful to A. Lun (CRUK-CI) for constructive comments on an earlier version of the manuscript. We are also grateful to UC Berkeley collaborator J. Ngai and his group members. C.A.V., A.S., and J.C.M. acknowledge core EMBL funding. C.A.V. was supported by core MRC funding (MRC MC UP 0801/1) and by The Alan Turing Institute under the EPSRC grant no. EP/N510129/1. J.C.M. acknowledges core support from CRUK. A.S. acknowledges funding from the Wellcome Trust Strategic Award 105031/D/14/Z, “Tracing early mammalian lineage decisions by single-cell genomics.” D.R. and S.D. are supported by the US National Institutes of Health BRAIN Initiative grant no. U01 MH105979 (PI, J. Ngai).

AUTHOR CONTRIBUTIONS

C.A.V., D.R., and A.S. performed analyses. C.A.V., D.R., A.S., S.D., and J.C.M. wrote the manuscript. S.D. and J.C.M. supervised the study.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>. Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

1. Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell. *Nat. Methods* **6**, 377–382 (2009).
2. Shapiro, E., Biezuner, T. & Linnarsson, S. Single-cell sequencing-based technologies will revolutionize whole-organism science. *Nat. Rev. Genet.* **14**, 618–630 (2013).
3. Stegle, O., Teichmann, S.A. & Marioni, J.C. Computational and analytical challenges in single-cell transcriptomics. *Nat. Rev. Genet.* **16**, 133–145 (2015).
4. Saliba, A.-E., Westermann, A.J., Gorski, S.A. & Vogel, J. Single-cell RNA-seq: advances and future challenges. *Nucleic Acids Res.* **42**, 8845–8860 (2014).

5. Gawad, C., Koh, W. & Quake, S.R. Single-cell genome sequencing: current state of the science. *Nat. Rev. Genet.* **17**, 175–188 (2016).
6. Brennecke, P. *et al.* Accounting for technical noise in single-cell RNA-seq experiments. *Nat. Methods* **10**, 1093–1095 (2013).
7. Kharchenko, P.V., Silberstein, L. & Scadden, D.T. Bayesian approach to single-cell differential expression analysis. *Nat. Methods* **11**, 740–742 (2014).
8. Finak, G. *et al.* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol.* **16**, 278 (2015).
9. Pierson, E. & Yau, C. ZIFA: dimensionality reduction for zero-inflated single-cell gene expression analysis. *Genome Biol.* **16**, 241 (2015).
10. Bacher, R. & Kendziora, C. Design and computational analysis of single-cell RNA-sequencing experiments. *Genome Biol.* **17**, 63 (2016).
11. Vallejos, C.A., Richardson, S. & Marioni, J.C. Beyond comparisons of means: understanding changes in gene expression at the single-cell level. *Genome Biol.* **17**, 70 (2016).
12. Jiang, L. *et al.* Synthetic spike-in standards for RNA-seq experiments. *Genome Res.* **21**, 1543–1551 (2011).
13. Kolodziejczyk, A.A., Kim, J.K., Svensson, V., Marioni, J.C. & Teichmann, S.A. The technology and biology of single-cell RNA sequencing. *Mol. Cell* **58**, 610–620 (2015).
14. Lun, A.T., Bach, K. & Marioni, J.C. Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. *Genome Biol.* **17**, 75 (2016).
15. Smyth, G.K. & Speed, T. Normalization of cDNA microarray data. *Methods* **31**, 265–273 (2003).
16. Bullard, J.H., Purdom, E., Hansen, K.D. & Dudoit, S. Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* **11**, 94 (2010).
17. Dillies, M.-A. *et al.* A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Brief. Bioinform.* **14**, 671–683 (2013).
18. Hicks, S.C., Teng, M. & Irizarry, R.A. On the widespread and critical impact of systematic bias and batch effects in single-cell RNA-Seq data. Preprint at <http://biorxiv.org/content/early/2015/08/25/025528> (2015).
19. Risso, D., Ngai, J., Speed, T.P. & Dudoit, S. Normalization of RNA-seq data using factor analysis of control genes or samples. *Nat. Biotechnol.* **32**, 896–902 (2014).
20. Leek, J.T. svaseq: removing batch effects and other unwanted noise from sequencing data. *Nucleic Acids Res.* **42**, e161 (2014).
21. Islam, S. *et al.* Quantitative single-cell RNA-seq with unique molecular identifiers. *Nat. Methods* **11**, 163–166 (2014).
22. Grün, D. & van Oudenaarden, A. Design and analysis of single-cell sequencing experiments. *Cell* **163**, 799–810 (2015).
23. Vallejos, C.A., Marioni, J.C. & Richardson, S. BASiCS: Bayesian analysis of single-cell sequencing data. *PLoS Comput. Biol.* **11**, e1004333 (2015).
24. Mortazavi, A., Williams, B.A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat. Methods* **5**, 621–628 (2008).
25. Li, B., Ruotti, V., Stewart, R.M., Thomson, J.A. & Dewey, C.N. RNA-Seq gene expression estimation with read mapping uncertainty. *Bioinformatics* **26**, 493–500 (2010).
26. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome Biol.* **11**, R106 (2010).
27. Robinson, M.D. & Oshlack, A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol.* **11**, R25 (2010).
28. Klein, A.M. *et al.* Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell* **161**, 1187–1201 (2015).
29. Pollen, A.A. *et al.* Low-coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex. *Nat. Biotechnol.* **32**, 1053–1058 (2014).
30. Zeisel, A. *et al.* Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* **347**, 1138–1142 (2015).
31. Macosko, E.Z. *et al.* Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell* **161**, 1202–1214 (2015).
32. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
33. Büttner, F. *et al.* Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat. Biotechnol.* **33**, 155–160 (2015).
34. Haghverdi, L., Büttner, M., Wolf, F.A., Büttner, F. & Theis, F.J. Diffusion pseudotime robustly reconstructs lineage branching. *Nat. Methods* **13**, 845–848 (2016).
35. Ding, B. *et al.* Normalization and noise reduction for single cell RNA-seq experiments. *Bioinformatics* **31**, 2225–2227 (2015).
36. Katayama, S., Töhönen, V., Linnarsson, S. & Kere, J. SAMstr: statistical test for differential expression in single-cell transcriptome with spike-in normalization. *Bioinformatics* **29**, 2943–2945 (2013).
37. Bacher, R. *et al.* SCnorm: a quantile-regression based approach for robust normalization of single-cell RNA-seq data. *Nat. Methods* <http://dx.doi.org/10.1038/nmeth.4263> (2017).
38. Qiu, X. *et al.* Single-cell mRNA quantification and differential analysis with Census. *Nat. Methods* **14**, 309–315 (2017).
39. Islam, S. *et al.* Characterization of the single-cell transcriptional landscape by highly multiplex RNA-seq. *Genome Res.* **21**, 1160–1167 (2011).
40. Munro, S.A. *et al.* Assessing technical performance in differential gene expression experiments with external spike-in RNA control ratio mixtures. *Nat. Commun.* **5**, 5125 (2014).
41. Goolam, M. *et al.* Heterogeneity in Oct4 and Sox2 targets biases cell fate in 4-cell mouse embryos. *Cell* **165**, 61–74 (2016).
42. Scialdone, A. *et al.* Computational assignment of cell-cycle stage from single-cell transcriptome data. *Methods* **85**, 54–61 (2015).
43. Hardwick, S.A. *et al.* Spliced synthetic genes as internal controls in RNA sequencing experiments. *Nat. Methods* **13**, 792–798 (2016).
44. Lovén, J. *et al.* Revisiting global gene expression analysis. *Cell* **151**, 476–482 (2012).
45. Cole, M. & Risso, D. scone: Single Cell Overview of Normalized Expression data, R package version 0.99.6 (2016).
46. Kolodziejczyk, A.A. *et al.* Single cell RNA-sequencing of pluripotent states unlocks modular transcriptional variation. *Cell Stem Cell* **17**, 471–485 (2015).