

On the widespread and critical impact of systematic bias and batch effects in single-cell RNA-Seq data

Stephanie C. Hicks^{1,2}, Mingxiang Teng^{1,2}, Rafael A. Irizarry^{1,2,*}

¹Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute,

²Department of Biostatistics, Harvard School of Public Health

*Corresponding Author

Emails:

Stephanie C. Hicks, shicks@jimmy.harvard.edu

Mingxiang Teng, mxteng@jimmy.harvard.edu

Rafael A. Irizarry, rafa@jimmy.harvard.edu

Abstract

Single-cell RNA-Sequencing (scRNA-Seq) has become the most widely used high-throughput method for transcription profiling of individual cells. Systematic errors, including batch effects, have been widely reported as a major challenge in high-throughput technologies. Surprisingly, these issues have received minimal attention in published studies based on scRNA-Seq technology. We examined data from five published studies and found that systematic errors can explain a substantial percentage of observed cell-to-cell expression variability. Specifically, we found that the proportion of genes reported as expressed explains a substantial part of observed variability and that this quantity varies systematically across experimental batches. Furthermore, we found that the implemented experimental designs confounded outcomes of interest with batch effects, a design that can bring into question some of the conclusions of these studies. Finally, we propose a simple experimental design that can ameliorate the effect of these systematic errors have on downstream results.

Single-cell RNA-Sequencing (scRNA-Seq) has become the primary tool for profiling the transcriptomes of hundreds or even thousands of individual cells in parallel. Our experience with high-throughput genomic data in general, is that well thought-out data processing pipelines are essential to produce meaningful downstream results¹⁻³. We expect the same to be true for scRNA-seq data. Here we show that while some tools developed for analyzing bulk RNA-Seq can be used for scRNA-Seq data, such as the mapping and alignment software, other steps in the processing, such as normalization, quality control and quantification, require new methods to account for the additional variability that is specific to this technology.

One of the most challenging sources of unwanted variability and systematic error in high-throughput data are what are commonly referred to as *batch effects*. Given the way that scRNA-Seq experiments are conducted, there is much room for concern regarding batch effects⁴. Specifically, batch effects occur when cells from one biological group or condition are cultured, captured and sequenced separate from cells in a second condition. Although batch information is not always included in the experimental annotations that are publicly available, one can extract surrogate variables from the raw sequencing (FASTQ) files⁵. Namely, the sequencing instrument used, the run number from the instrument and the flow cell lane. Although the sequencing is unlikely to be a major source of unwanted variability, it serves as a surrogate for other experimental procedures that very likely do have an effect, such as starting material, PCR amplification reagents/conditions, and cell cycle stage of the cells⁶⁻⁸. Here we will refer to the resulting differences induced by different groupings of these sources of variability as *batch effects*.

In a completely confounded study, it is not possible to determine if the biological condition or batch effects are driving the observed variation. In contrast, incorporating biological replicates across in the experimental design and processing the replicates across multiple batches permits observed variation to be attributed to biology or batch effects (Figure 1). To demonstrate the widespread problem of systematic bias, batch effects, and confounded experimental designs in scRNA-Seq studies, we surveyed several published data sets. We discuss the consequences of failing to consider the presence of this unwanted technical variability, and consider new strategies to minimize its impact on scRNA-Seq data.

Batch and outcomes of interest are confounded in published scRNA-Seq experiments

We examined five publicly available scRNA-Seq data sets to investigate the extent of confounding biological variation with batch effects. For each study, we downloaded the processed data available on GEO⁹ and reconstructed the study design from the sequence identifiers provided in the FASTQ files. We used the standardized Pearson contingency coefficient to assess the experimental design between processing batches and outcome of interest and found values ranging from 92.7% to 100% (perfect cofounding) (Table 1). Note that with this level of confounding it is nearly impossible to parse batch effects from biological variation.

Proportion of detected genes is a major source of technical cell-to-cell noise

Most, if not all, published studies using scRNA-Seq rely explicitly or implicitly on computing distances between the cell expression profiles. Principal component analysis is used explicitly to quantify biological or molecular distance¹⁰ or implicitly to approximate distance between individual cells. We used the processed expression data available on GEO, applied principal components analysis on the log (base 2) transformed values (adding 1 to avoid logs of 0), and computed the proportion of detected genes from the same data set with the exception of one study¹⁰. In this exception, the processed expression data available on GEO excluded most non-detected genes and the values for each gene were centered by removing the average. For this case, we computed the proportion of detected genes from the raw data. We found wide variation in the proportion of detected genes across cells: from 1% detected to 60%. Furthermore, we found strong correlation between the first principal component and the proportion of detected genes within each data set (Figure 2, Supp Fig 1).

To determine if the variability in the proportion of detected genes was biologically or technically driven we compared the variability across biological groups to the variability across batches. For most cases in which the experimental design permitted this comparison (Tables S1-S3), we found that batch explained more variability than biological group (Supp Figures 2-4). In the two studies for which batch was completely confounded with biological group (Table S4 and S5) we also observed variability across batch (Supp Figure 5). However, in these cases it is impossible to separate variability due to biology or to batch.

Batch effects lead to differences in detection rates, which lead to apparent differences between biological groups

To illustrate potential down stream effects, we examined 430 single-cells from five biological groups of interest for one of the studies in which at least one biological condition was split across two batches (Table S1). We confirmed that cells cluster by biological group (Figure 3A) as reported in the original publication. However, four of the five biological groups were confounded with batch (Table S1) and batches can also explain the clustering (Figure 3B). The one group that was not confounded with batch confirms the high level of variability explained by processing cells in different batches (Figure 3C). As expected from the previous results, different batches lead to different proportions of detected genes (Figure 3D) which we have shown to be correlated with the first principal component.

Detection rate has indirect effects on reported gene expression measurements

Using the processed scRNA-Seq data available on GEO, we computed the median of the non-zero measurements for each cell. For each study, we noticed a strong non-linear relationship between the median expression and the proportion of detected genes (Figure 4). The overall level of expression changing with the proportion of detected genes could be biologically driven, but there is a reasonable explanation of how it can be a technical artifact, which we explain using statistical notation. Let X_{ij} be the unobserved expression level for sample i and gene j . Let us consider only expressed genes ($X_{ij} > 0$). In the sequencing experiment, each expressed gene has a probability of being amplified, which means we observe a quantity proportional to $X_{ij}Z_{ij}$ where $Z_{ij} = 1$ if the gene was expressed at a high enough level to be detected and amplified and 0 otherwise. This implies that the expected amount of RNA we will obtain is a quantity proportional to

$$E(X_{ij}Z_{ij}) = \Pr(Z_{ij} = 1) E(X_{ij}|Z_{ij} = 1)$$

Because we know technical variation affects the probability detection $\Pr(Z_{ij} = 1) = p_i$, which depends on the sample i , we assume that, within a homogenous population for example, the expected level of a gene that is detected, $E(X_{ij}|Z_{ij} = 1)$, is the same across cells. Then, the total amount of RNA is proportional to the detection rate of the expressed genes p_i . Now because

experimentally we amplify to have roughly the same total amount of RNA for each sample, this implies that we have to amplify in a way that is equivalent to multiplying by $\frac{K}{p_i}$, with K the total we aspire to reach. Thus the median expression of the expressed genes will be proportional to $\frac{1}{p_i}$, which is consistent with the data (Figure 4). To confirm this, we implemented an adjustment based on these derivations. Specifically, for each study we computed Counts Per Million (CPM)

$$\text{CPM} = Y_{ij} / \left[\frac{N_i}{10^6} \right],$$

which scales the raw read counts Y_{ij} for the j^{th} gene (or transcript) in the i^{th} cell by the total number of reads N_i in the sample. Using the CPM normalized data, filtered for cells passing a data-driven filter (Supplementary Figure 6), we multiplied the i^{th} cells by an adjustment factor $\left(\frac{p_i}{\text{median}(p_1, p_2, \dots, p_I)} \right)$ motivated by derivations above. This adjustment removed much of the dependence between the median expression values and the proportion of detected genes (Supplementary Figure 7).

Finally, we also noted that the entire distribution of the non-zero genes changed with the proportion detected (Supplementary Figure 8). An important consequence of this indirect effect of experimental noise is that even after removing all genes with at least one 0, there was a strong correlation between the first principal component of this smaller subset with no 0s, and the proportion of detected genes (Supplementary Figure 9)

Experimental design solutions

There are currently no published general statistical solutions to the problem of batch effects in scRNA-Seq data. For the specific application of differential expression, a proposed solution is to account for differences in the proportion of detected genes by explicitly including it as a covariate in a linear regression model¹¹. However, given the current levels of confounding and experimental designs, this approach will not be able to distinguish biological from technical effects. Because of the nature of the experimental protocol needed to run scRNA-Seq experiments imposed by the way cells are captured and sequenced in batches, standard balanced experimental designs are not possible. An experimental design solution is to use biological replicates, namely independently repeating the experiment multiple times for each biological

condition (Figure 1). This approach allows for multiple batches of cells to be randomized across sequencing runs, flow cells and lanes as in bulk-RNA-Seq. With this design we can then model and adjust for batch effects due to systematic experimental bias.

Discussion

Batch effects and unwanted technical cell-to-cell noise remains a challenge in the analysis of scRNA-Seq data. The challenge is more complex than in previous sequencing experiments since experimental batches lead to different detection rates, which in turn lead to different transcription level estimates. In addition, detection of a gene or transcript in scRNA-Seq experiments is heavily dependent on the experimental protocol, which leads to systematic differences in the proportion of detected genes between batches of cells. The development of statistical methods that account for these systematic biases will therefore be essential in the analysis of scRNA-Seq data. Incorporating biological replicates in the experimental study design provides a solution to reducing confounding between biological condition and batch effects and will permit modeling of the technical variability that relates to processing the cells in different batches.

Acknowledgements

We thank Bradley Bernstein who provided insightful comments that greatly improved the manuscript. This research was supported by NIH R01 grants GM083084, RR021967/GM103552 and HG005220.

Competing Interests

The authors declare no competing interests.

Supplementary Material

Supplementary materials are available in a single pdf.

References

1. Akey, J.M., Biswas, S., Leek, J.T. & Storey, J.D. On the design and analysis of gene expression studies in human populations. *Nature genetics* **39**, 807-808; author reply 808-809 (2007).
2. Baggerly, K.A., Edmonson, S.R., Morris, J.S. & Coombes, K.R. High-resolution serum proteomic patterns for ovarian cancer detection. *Endocrine-related cancer* **11**, 583-584; author reply 585-587 (2004).
3. Leek, J.T. et al. Tackling the widespread and critical impact of batch effects in high-throughput data. *Nature reviews. Genetics* **11**, 733-739 (2010).
4. Stegle, O., Teichmann, S.A. & Marioni, J.C. Computational and analytical challenges in single-cell transcriptomics. *Nature reviews. Genetics* **16**, 133-145 (2015).
5. Gilad, Y. & Mizrahi-Man, O. A reanalysis of mouse ENCODE comparative gene expression data. *F1000Research* **4**, 121 (2015).
6. Shalek, A.K. et al. Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature* **510**, 363-369 (2014).
7. Brennecke, P. et al. Accounting for technical noise in single-cell RNA-seq experiments. *Nature methods* **10**, 1093-1095 (2013).
8. Buettner, F. et al. Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nature biotechnology* **33**, 155-160 (2015).
9. Edgar, R., Domrachev, M. & Lash, A.E. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic acids research* **30**, 207-210 (2002).
10. Patel, A.P. et al. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**, 1396-1401 (2014).
11. Finak, G. et al. MAST: A flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA-seq data. *bioRxiv* (2015).
12. Treutlein, B. et al. Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature* **509**, 371-375 (2014).
13. Deng, Q., Ramskold, D., Reinius, B. & Sandberg, R. Single-cell RNA-seq reveals dynamic, random monoallelic gene expression in mammalian cells. *Science* **343**, 193-196 (2014).
14. Trapnell, C. et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology* **32**, 381-386 (2014).

Study	Num ber of cells	Number of genes or transcripts	Processed data available	Confounding (%)	Correlation between PC1 and proportion of detected genes or transcripts	Ref
Patel et al. (2014)	430	5,948	TPM	98.9	0.54*	¹⁰
Treutlein et al. (2014)	198	23,745	FPKM	92.7	0.91	¹²
Deng et al. (2014)	286	22,958	RPKM	96.6	0.69-0.82**	¹³
Trapnell et al. (2014)	372	47,192	FPKM	100	0.90	¹⁴
Shalek et al. (2014)	383	27,723	TPM	100	0.93	⁶

Table 1: Description of processed single-cell RNA-Seq data sets. Column 1 shows the publications. Column 2 shows the number of cells (samples) included in the study. Column 3 shows the number of genes included in the data uploaded to the public repository. Column 4 indicates the units in which the values were reported. Column 5 shows the level of confounding between biological condition and batch effect quantified using the standardized Pearson contingency coefficient as a measure of association. The percentage ranges from 0% (no confounding) to 100% (completely confounded). Column 6 shows the Pearson correlation between the first principal component of the log transformed data and the proportion of detected genes. Column 7 provides the citation for the study.

*The available processed data was previously filtered by the authors and excluded the majority of non-detected genes, which partially explains the lower correlation (0.54 compared to 0.69-0.93).

**The biological conditions in this study were split into four groups for this analysis each with its own corresponding correlation coefficient.

The Problem of Confounding Biological Variation and Batch Effects

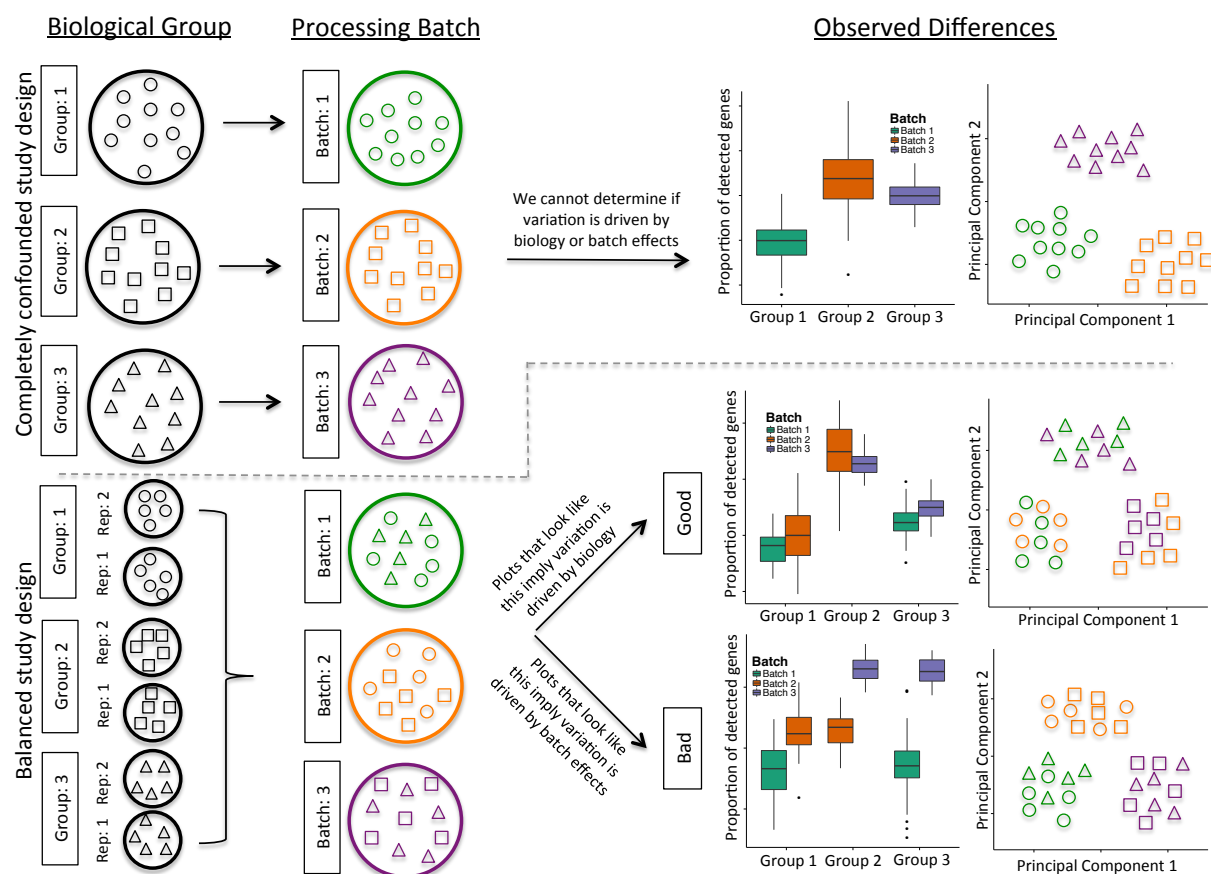


Figure 1: The problem of confounding biological variation and batch effects. The top section depicts a completely confounded study design of processing individual cells from three biological groups (represented by shapes) in three separate batches (represented by colors). In this case, we cannot determine if biology or batch effects drive the observed variation. The bottom section depicts a balanced study design consisting of multiple replicates (rep) split and processed across multiple batches. The use of multiple replicates allows observed variation be attributed to biology (cells cluster by shape) or batch effects (cells cluster by color).

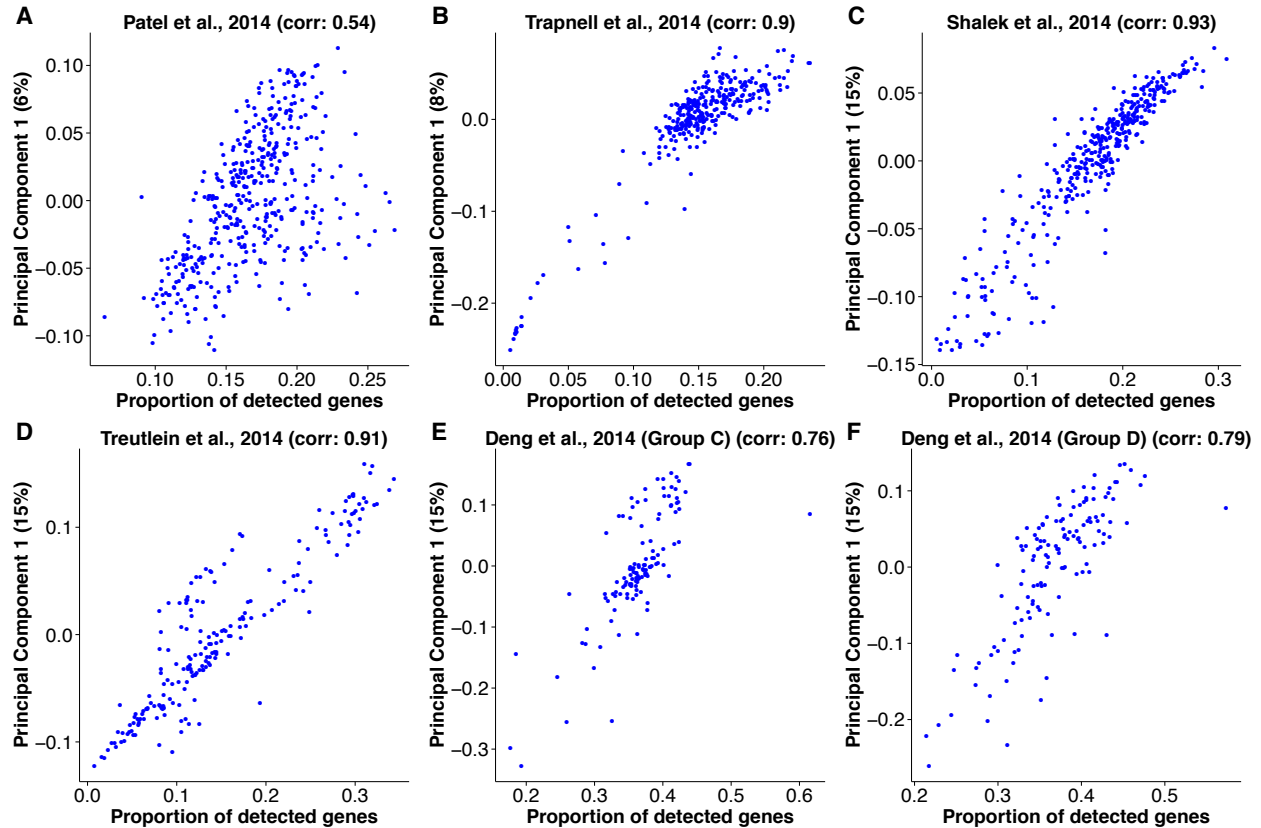


Figure 2: First principal component is strongly correlated with the proportion of detected. The principal components from the processed data available on GEO. **(A)** Patel et al. (2014). **(B)** Trapnell et al. (2014). **(C)** Shalek et al. (2014). **(D)** Treutlein et al. (2014). **(E)** The 4-cell, 8-cell, 16-cell groups (Group C) from Deng et al. (2014). **(F)** The Early, Mid and Late blastocyst groups (Group D) from Deng et al. (2014). **Note:** The proportion of detected genes was calculated using the publicly available processed data from GEO for all studies except for Patel et al. (2014). In this case, because most non-detected genes were excluded from the publicly available processed data, we computed the proportion of detected genes from the raw data.

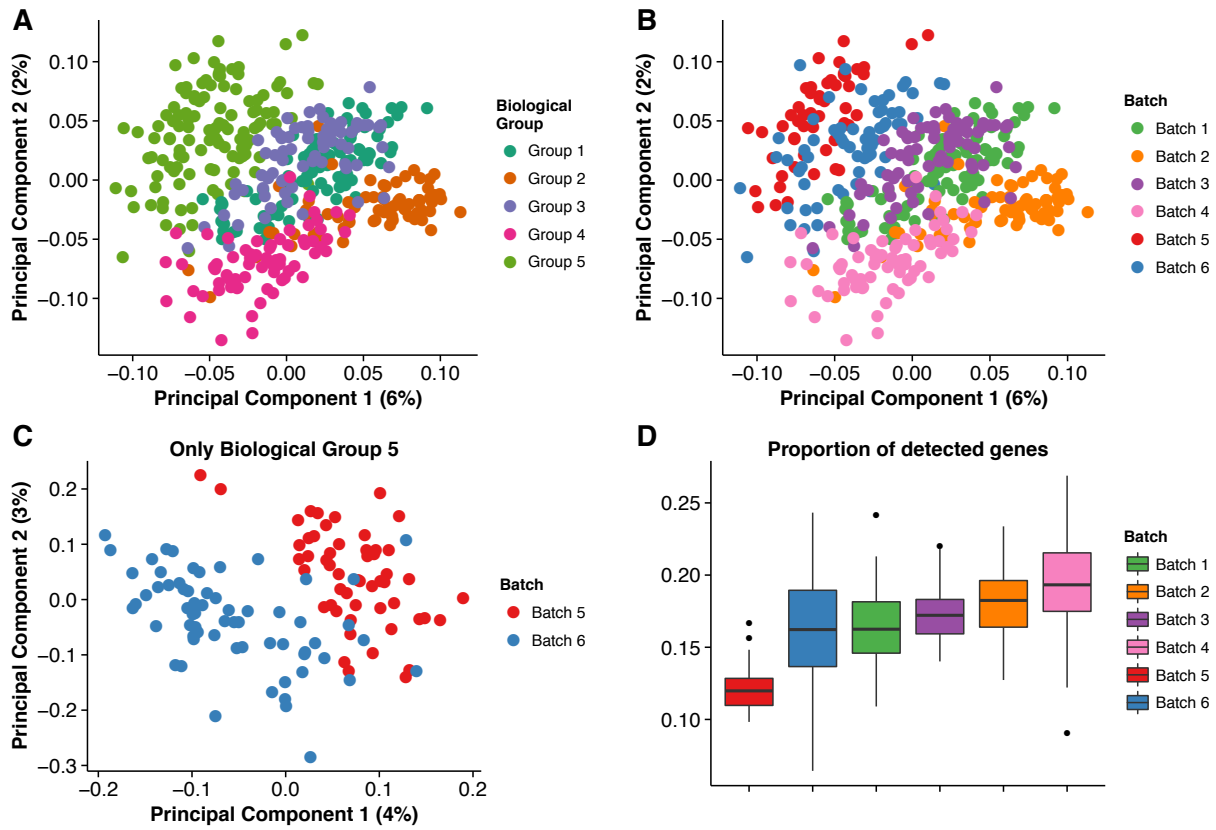


Figure 3: Illustration with public data¹⁰ of how batch effects lead to differences in detection rates, which lead to apparent differences between biological groups. **(A)** Using principal components analysis, scRNA-Seq samples cluster by biological group, but the observed biological variation across groups is confounded with **(B)** technical variation from processing the cells in different batches. **(C)** Within one group (Group 5), the cells cluster by batch. **(D)** Furthermore, individual batches of cells have different proportions of detected genes, which may be driving the observed biological variation across groups.

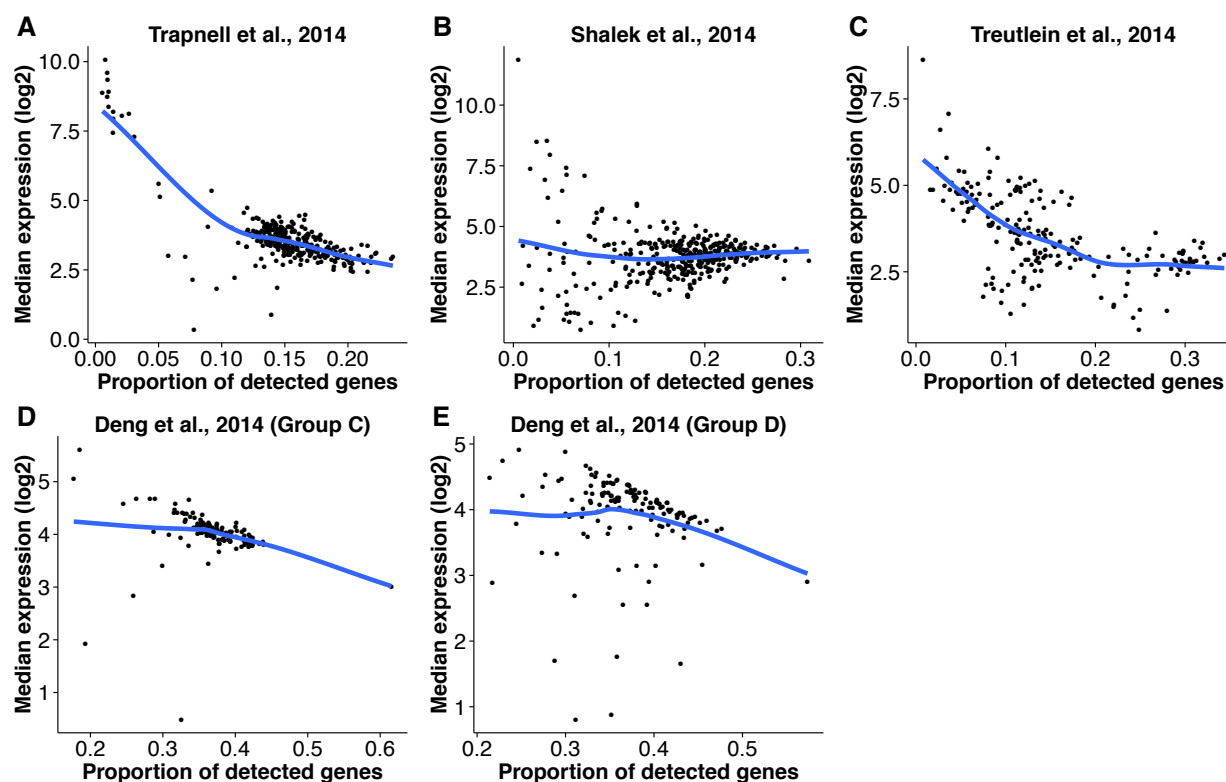


Figure 4: Non-linear relationship between the median gene expression and the proportion of detected genes using processed scRNA-Seq data available on GEO. Failure to account for differences of the proportion of detected genes between cells over-inflates the gene expression estimates of cells with a low proportion of detected genes. The blue curves were obtained by fitting a locally weighted scatter plot smoothing (loess) with a degree of 1 and span of 0.75 for all figures. Because the range of proportion of detected genes varied from study to study, the range of the x-axis differs across plots. **(A)** Trapnell et al. (2014). **(B)** Shalek et al. (2014). **(C)** Treutlein et al. (2014). **(D)** Deng et al. (2014) (Group C). **(E)** Deng et al. (2014) (Group D). **Note:** We could not include Patel et al. (2014) because of the row-standardization applied by the authors in the processed data available on GEO.