

# Distinguishing cells from empty droplets in droplet-based single-cell RNA sequencing data

Aaron T. L. Lun<sup>1,\*</sup>, Samantha Riesenfeld<sup>2,\*</sup>, Tallulah Andrews<sup>3,\*</sup>, The Phuong Dao<sup>4,\*</sup>, Tomas Gomes<sup>3,\*</sup>, participants in the 1<sup>st</sup> Human Cell Atlas Jamboree<sup>‡</sup>, John C. Marioni<sup>1,3,5,#</sup>

**1** Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Robinson Way, Cambridge CB2 0RE, United Kingdom

**2** Klarman Cell Observatory, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

**3** Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA, United Kingdom

**4** Program for Computational and Systems Biology, Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY

**5** EMBL European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, United Kingdom

\* These authors contributed equally to this work.

† Email: [aaron.lun@cruk.cam.ac.uk](mailto:aaron.lun@cruk.cam.ac.uk)

‡ The full list of participants is provided in Supplementary Table 1.

# Email: [john.marioni@cruk.cam.ac.uk](mailto:john.marioni@cruk.cam.ac.uk)

## Abstract

Droplet-based single-cell RNA sequencing protocols have dramatically increased the throughput and efficiency of single-cell transcriptomics studies. A key computational challenge when processing these data is to distinguish libraries for real cells from empty droplets. Existing methods for cell calling set a minimum threshold on the total unique molecular identifier (UMI) count for each library, which indiscriminately discards cell libraries with low UMI counts. Here, we describe a new statistical method for calling cells from droplet-based data, based on detecting significant deviations from the expression profile of the ambient solution. Using simulations, we demonstrate that our method has greater power than existing approaches for detecting cell libraries with low UMI counts, while controlling the false discovery rate among detected cells. We also apply our method to real data, where we show that the use of our method results in the retention of distinct cell types that would otherwise have been discarded.

## Introduction

Recent advances in droplet-based protocols have revolutionized the field of single-cell transcriptomics by allowing tens of thousands of cells to be profiled in a single assay [1–3]. In these technologies, individual cells are captured into aqueous droplets in a water-in-oil emulsion. Each droplet also contains a co-captured bead with primers for reverse transcription, where all primers on a single bead contain a cell barcode that is (effectively) unique to that bead. The droplets serve as isolated reaction chambers in

which cell lysis and reverse transcription are performed to obtain barcoded cDNA. This is followed by breaking of the emulsion, amplification of the cDNA and construction of a sequencing library. After sequencing, transcripts are assigned to individual droplets based on the cell barcode observed in each read sequence. This yields an expression profile for each cell, typically in the form of unique molecular identifier (UMI) counts [4] for all annotated genes. The use of droplets increases throughput by at least an order of magnitude compared to protocols based on plates [5] or conventional microfluidics [6], which is appealing for large-scale projects such as the Human Cell Atlas [7].

The complexity of the sequencing data from droplet-based technologies poses a number of interesting challenges for low-level data processing. One such challenge is the identification and removal of cell barcodes corresponding to empty droplets. An empty droplet does not contain a cell but will still contain “ambient” RNA [1], i.e., cell-free transcripts in the solution in which the cells are suspended. Ambient RNA can be actively secreted by cells or released upon cell lysis (possibly induced by the stresses of dissociation and microfluidics). The presence of ambient RNA means that many empty droplets will contain material for reverse transcription and library preparation, resulting in non-zero total UMI counts for the corresponding barcodes. However, the expression profiles for these barcodes do not originate from any individual cell and need to be removed prior to further analysis to avoid misleading biological conclusions.

Existing methods for removing empty droplets assume that droplets containing genuine cells should have more RNA, resulting in larger total UMI counts for the corresponding barcodes. Zheng *et al.* [3] remove all barcodes with total counts below 10% of the 99<sup>th</sup> percentile of the  $Y$  largest total counts, where  $Y$  is defined as the expected number of cells to be captured in the experiment. Macosko *et al.* [1] set the threshold at the knee point in the cumulative fraction of reads with respect to increasing total count. While simple, the use of a one-dimensional filter on the total UMI count is suboptimal as it discards small cells with low RNA content. Droplets containing small cells are not easily distinguishable from large empty droplets in terms of the total number of transcripts. This problem is exacerbated by variable capture and amplification efficiency across droplets, which further mixes the distributions of total counts between empty and non-empty droplets. Applying a simple threshold on the total count forces the researcher to choose between the loss of small cells or an increase in the number of artifactual “cells” composed of ambient RNA. This is especially problematic if small cells represent distinct cell types or functional states.

In this report, we propose a new method for detecting empty droplets in droplet-based single-cell RNA sequencing (scRNA-seq) data. We construct a profile of the ambient pool of RNA, and test each barcode for deviations from this profile using a Poisson-based model for the count distribution. Barcodes with significant deviations are considered to be genuine cells, thus allowing recovery of cells with low total RNA content and small total UMI counts. We combine our approach with a knee point filter to ensure that barcodes with large total counts are always retained. Using a variety of simulations, we demonstrate that our method outperforms methods based on a simple threshold on the total UMI count. We also apply our method to several real data sets where we are able to recover more cells from both existing and new cell types.

## Description of the method

### Testing for deviations from the ambient profile

To construct the profile for the ambient RNA pool, we consider a threshold  $T$  on the total UMI count. The set  $\mathcal{G}$  of all barcodes with total counts less than or equal to  $T$  are considered to represent empty droplets. The exact choice of  $T$  does not matter, as long

as (i) it is small enough so that droplets with genuine cells do not have total counts below  $T$ , and (ii) there are sufficient counts to obtain a precise estimate of the ambient profile. We set  $T = 100$  by default in our approach, motivated by examination of several real datasets (Supplementary Section 1, Supplementary Figure 1). We stress that  $T$  is not the same as the threshold used in existing methods, as barcodes with total counts greater than  $T$  are not automatically considered to be cell-containing droplets.

The ambient profile is constructed by summing counts for each gene across  $\mathcal{G}$ . Let  $y_{gb}$  be the count for gene  $g$  in barcode  $b$ . We define the ambient count for  $g$  as

$$A_g = \sum_{b \in \mathcal{G}} y_{gb} ,$$

yielding a count vector  $\mathbf{A} = (A_1, \dots, A_N)$  for all  $N$  genes. (We assume that any gene with counts of zero for all barcodes has already been filtered out, as this provides no information for distinguishing between barcodes.) We apply the Good-Turing algorithm to  $\mathbf{A}$  to obtain the posterior expectation,  $\tilde{p}_g$ , of the proportion of counts assigned to  $g$  [8], using the `goodTuringProportions` function in the `edgeR` package [9]. This ensures that genes with zero counts in the ambient pool have non-zero proportions, avoiding the possibility of obtaining likelihoods of zero in downstream calculations.

Our null hypothesis is that free-floating transcripts in solution are randomly encapsulated into the empty droplets. For a given droplet, the probability of sampling a transcript molecule for gene  $g$  is equal to  $\tilde{p}_g$ . If we condition on the total count  $t_b$  for a cell barcode  $b$ , we can model the counts for each barcode with a multinomial distribution. We define the likelihood of obtaining the counts for barcode  $b$  as

$$L_b = t_b! \prod_{g=1}^N \frac{\tilde{p}_g^{y_{gb}}}{y_{gb}!} .$$

We use a Monte Carlo approach to compute the  $p$ -value for  $b$ . We generate a new set of counts by randomly sampling from a multinomial distribution with probabilities set to  $\tilde{p}_g$  for all  $g$  and size equal to  $t_b$ . We use the above formula to calculate the likelihood for this set of counts (denoted  $L'_{bi}$ , for iteration  $i$ ), and we repeat this process for  $R$  iterations. We use the method of Phipson and Smyth [10] to define the  $p$ -value as

$$P_b = \frac{R_b + 1}{R + 1} ,$$

where  $R_b$  is the number of iterations in which  $L'_{bi} \leq L_b$ . This strategy avoids  $p$ -values of zero, which is important during multiple testing correction. See Supplementary Section 2 for a description of how these  $p$ -values are efficiently computed.

## Detecting the knee point in the log-totals

Applying a threshold on the  $p$ -value will identify barcodes that have count profiles that are significantly different from the ambient pool of RNA. We assume that this will be the case for most cell-containing droplets, as the ambient pool is formed from many (lysed) cells and is unlikely to be representative of any single cell. However, it is possible for some cell-containing droplets to have ambient-like expression profiles. This can occur if the cell population is highly homogeneous or if one cell subpopulation contributes disproportionately to the ambient pool, e.g., if it is more prone to lysis. Sequencing errors in the cell barcodes may also bias the estimates of the ambient proportions, by misassigning counts from cell-containing droplets to barcodes with low UMI totals. This may result in spurious similarities between cells and the estimated ambient profile.

To avoid incorrectly calling ambient-like cells as empty droplets, we combine our procedure with a conventional threshold on the total UMI count. We rank all barcodes in order of decreasing  $t_b$ , and consider  $\log(t_b)$  as a function  $f(\cdot)$  of the log-transformed rank, i.e.,  $\log(t_b) = f(\log r_b)$  where  $r_b$  is the rank of  $b$  in the ordered sequence of barcodes. The first “knee” point in this function corresponds to a transition between a distinct subset of barcodes with large totals and the majority of barcodes with smaller totals. This is defined at the log-rank that minimizes the signed curvature

$$\frac{f''}{(1 + f'^2)^{1.5}},$$

and represents the point at which  $f(\cdot)$  begins to drop rapidly, marking the start of the transition between large and small totals. In practice, we obtain  $f(\cdot)$  by fitting a smooth spline to  $\log(t_b)$  against the log-rank in the interval containing the knee point. The derivatives of  $f(\cdot)$  are then obtained by differentiation of the spline basis functions. This avoids multiplication of errors during numerical differentiation, which would lead to instability in the curvature values and inaccurate estimates of the knee point.

Our assumption is that any barcode with a large total count must represent a cell-containing droplet, regardless of whether its count profile resembles the ambient pool. This is based on the expectation that the distribution of the sizes of empty droplets should be unimodal, with a monotonic decreasing probability density as  $t_b$  increases past the mode. A distinct peak of large totals would not be consistent with this expected distribution. We define the upper threshold  $U$  as the  $t_b$  at the knee point and retain all barcodes with  $t_b \geq U$ , irrespective of their  $P_b$ . This guarantees recovery of any barcodes with large total counts that potentially represent cell-containing droplets. We use the knee point rather than the inflection point as the  $t_b$  at the former is larger, providing a more conservative threshold that avoids retention of empty droplets.

We stress that, despite the use of a threshold on  $t_b$ , our approach is different from existing methods due to the testing procedure. Barcodes with  $t_b$  below the knee point can still be retained if the count profile is significantly different from the ambient pool. This is not possible with existing methods that would simply discard these barcodes. Users can also set  $U$  manually if automatic detection of the knee point fails for complex  $f(\cdot)$ . Alternatively, this mechanism can be disabled completely in favour of detecting cells solely based on their  $p$ -values. This is more statistically rigorous as it avoids the selection of an *ad hoc* threshold, but may result in the failure to detect large cells.

## Correcting for multiple testing across barcodes

We correct for multiple testing by controlling the false discovery rate (FDR) using the Benjamini-Hochberg (BH) method [11]. Putative cells are defined as those barcodes that have significantly poor fits to the ambient model at a specified FDR threshold. We set the FDR threshold to 1% by default, meaning that the expected proportion of empty droplets in the set of retained barcodes is no greater than 1%.

Note that we only perform the BH correction on the  $p$ -values for barcodes that have  $t_b$  greater than  $T$ . This reduces the severity of the correction by discarding barcodes that were previously assumed to be empty droplets, thus improving detection power for barcodes with larger totals that are more likely to contain cells. In fact,  $p$ -values are not computed at all for barcodes with  $t_b \leq T$  to avoid unnecessary computational work.

Conversely, all barcodes with  $t_b \geq U$  are considered to be known true positives, regardless of how ambient-like their expression profiles are. These barcodes have their  $p$ -values set to zero during the BH correction. This approach improves power by reducing the severity of the correction in the presence of a set of known positives.

## Results

### Evaluating performance with simulated droplet-based data

We named our method “EmptyDrops” and tested it on simulated data involving cells with different RNA content (see Methods, Supplementary Figure 2). Each simulated dataset was generated from real droplet-based scRNA-seq data (Supplementary Table 1) and contained one group of large cells with high RNA content and large  $t_b$ ; one group of small cells with low RNA content and small  $t_b$ ; and a set of empty droplets with counts sampled from an ambient pool of RNA. We applied EmptyDrops at a FDR of 1% to determine the recall for each group of cells and the FDR among the detected barcodes. We also tested methods that retain all cells with total UMI counts above a threshold. The threshold was defined as the total  $U$  at the knee point, as described above; or using the quantile-based approach [3] in the CellRanger software from 10X Genomics.

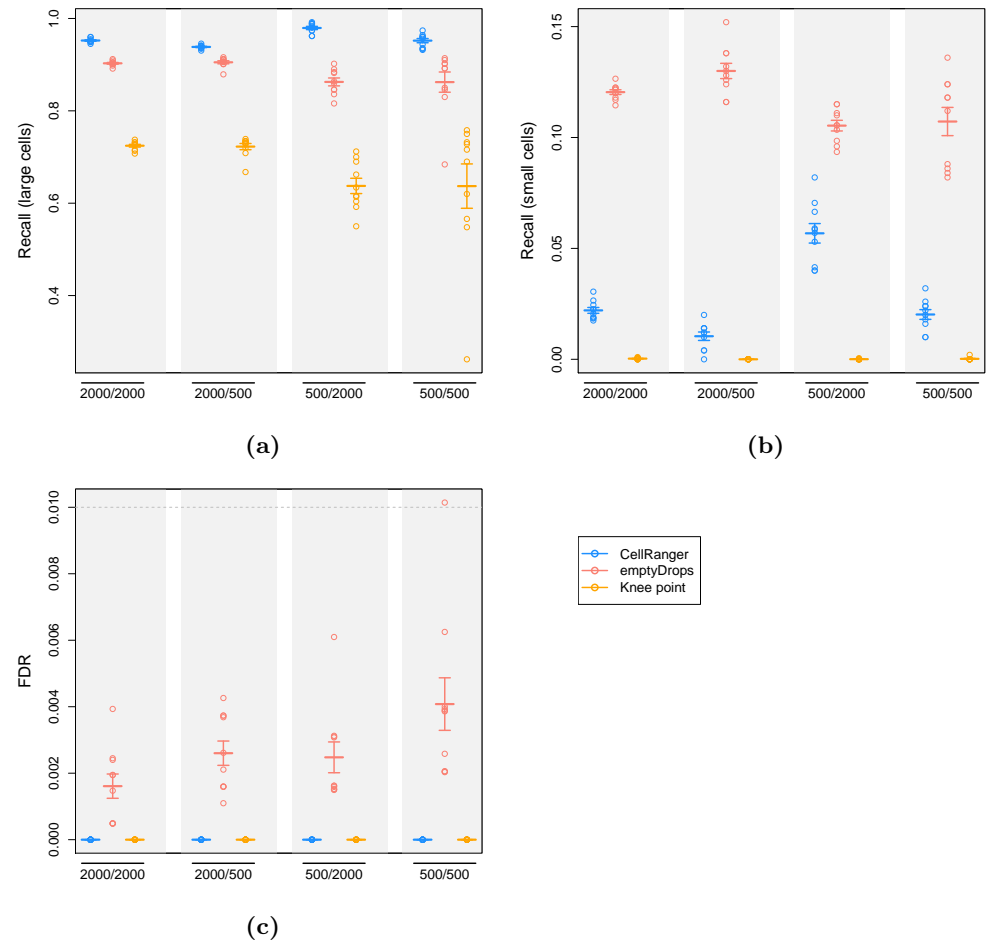
Results for a simulation based on a real dataset containing peripheral blood mononuclear cells (PBMCs) are shown in Figure 1. All methods consistently detected the large cells but recall for the group of small cells was much lower. Nonetheless, EmptyDrops was able to detect approximately 2-5-fold more small cells than the other methods in all scenarios. EmptyDrops also correctly controlled the FDR below the specified threshold of 1%. We observed similar performance in simulations based on other real datasets (Supplementary Figures 3-7). Improved detection of small cells with EmptyDrops was particularly pronounced in simulations based on the neuronal and cell line datasets, with recall of 60-90% compared to below 10% with the other methods. The performance of CellRanger was especially poor for the simulations based on the 9K neuron dataset, where over 10% of detected barcodes were false positives.

The poor performance of the threshold-based methods for small cells is expected. Barcodes corresponding to small cells with little RNA have similar total UMI counts as barcodes corresponding to large empty droplets with high levels of ambient RNA. The total UMI count cannot distinguish between these two possibilities, resulting in either reduced recall or a high false positive rate. In contrast, EmptyDrops uses the expression profile for each droplet to distinguish small cells from the ambient profile with greater power. Another benefit of EmptyDrops is its statistical rigour, allowing users to increase the stringency of the output by using a lower FDR threshold. The effect of adjusting the parameters in the quantile-based CellRanger approach is less interpretable.

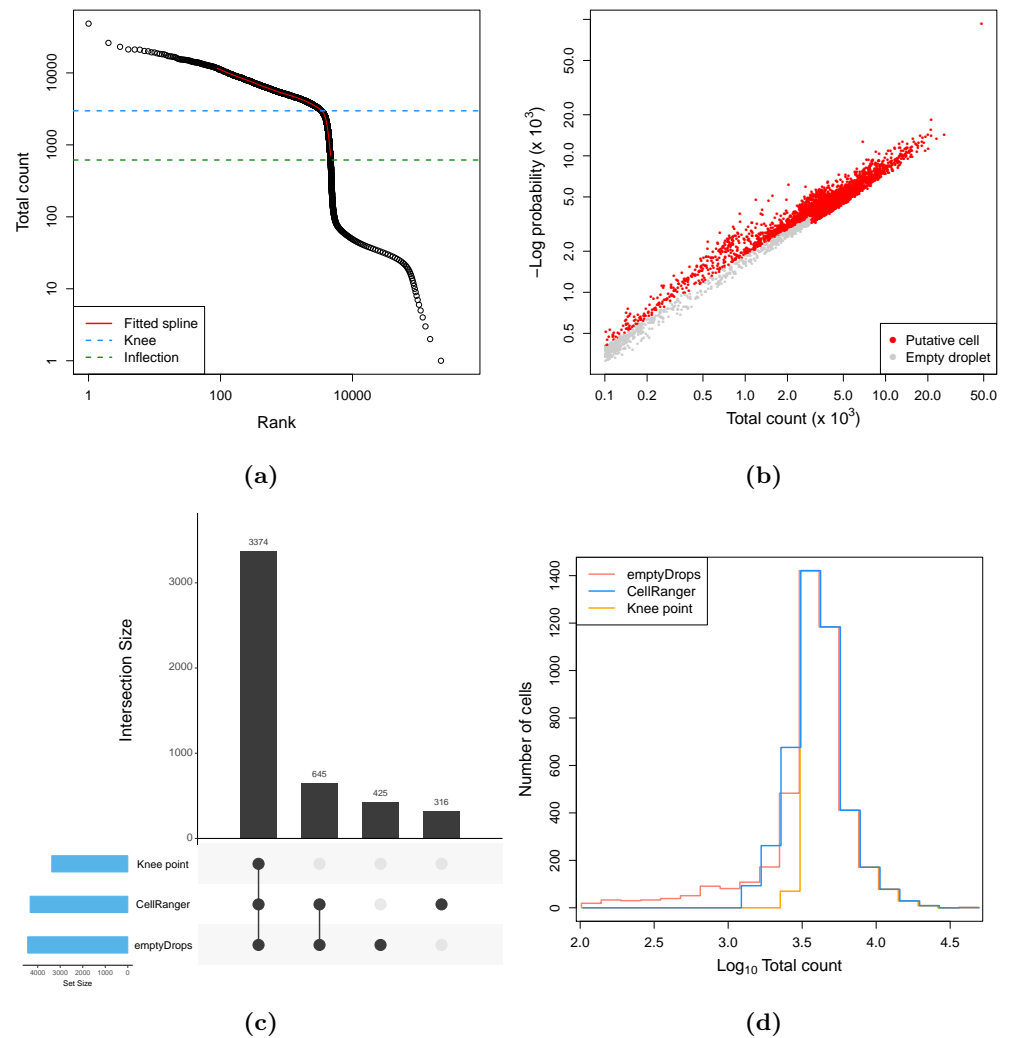
We note that CellRanger detects slightly more large cells than EmptyDrops in the PBMC-based simulation (Figure 1). We stress that this does not represent an inherent difference in performance between the two methods, and indeed, is not observed in many of the other simulations (Supplementary Figures 3-7). In fact, EmptyDrops can always be made to report a superset of the barcodes identified by CellRanger, simply by setting  $U$  to the CellRanger-identified threshold. We do not do so as the CellRanger threshold may not be appropriate, as observed for the 9K neuron simulations. It also relies on knowledge of the expected number of cells, which may not be available or accurate.

### Characterizing behaviour of EmptyDrops on real datasets

To determine how EmptyDrops behaved on real data, we applied it to detect cells in the PBMC dataset at a FDR of 1% (Figure 2). EmptyDrops identified a visually appropriate  $U$  using the knee point from the smoothed spline (Figure 2a), and detected significant barcodes as those with low probabilities under the null multinomial model (Figure 2b). Comparison of EmptyDrops to CellRanger indicated that most barcodes were detected by both methods (Figure 2c). Barcodes that were only detected by EmptyDrops had low total counts (Figure 2d), consistent with our simulation results. Conversely, CellRanger uniquely detected a number of cells with moderate total counts.



**Figure 1.** Results for the simulations based on the PBMC dataset, using three different methods for detecting cell-containing droplets. Simulation scenarios are labelled as  $G_1/G_2$  where  $G_1$  and  $G_2$  are the number of barcodes in the group of large and small cells, respectively. The recall for each group is shown as a proportion of the group size (a, b), and the FDR is calculated as the proportion of detected droplets that are empty (c). Each point represents the result of one simulation iteration, while the bar represents the mean across 10 iterations and the error bars represent the standard error of the mean. The dotted line represents the nominal FDR threshold (1%) for EmptyDrops.



**Figure 2.** Results of applying EmptyDrops and the other cell detection methods to the PBMC dataset. (a) A barcode rank plot showing the fitted spline used for knee point detection in EmptyDrops. The detected knee and inflection points are also shown. (b) The negative log-probability for each barcode in the multinomial model of EmptyDrops, plotted against the total count. Barcodes detected as putative cell-containing droplets at a FDR of 1% are labelled in red. Only barcodes with  $t_b > T$  are shown. (c) An UpSet plot [12] of the barcodes detected by each combination of methods (vertical bars). Horizontal bars represent the number of barcodes detected by each method. (d) Histogram outlines of the log-total count for barcodes detected by each method.



We observed similar results in the other tested datasets (Supplementary Figures 8-12) where EmptyDrops consistently detected the greatest number of cells. Increased retention of small cells with EmptyDrops was particularly pronounced in the neuronal datasets, where EmptyDrops was able to uniquely detect over a thousand cells. Again, a small number of barcodes were uniquely detected by CellRanger in some datasets. This is attributable to the conservativeness of the knee point threshold in EmptyDrops, which ensures that empty droplets are not inadvertently retained.

To explore the differences between methods in more detail, we generated a *t*-stochastic neighbour embedding (*t*-SNE) plot [13] of all barcodes that were detected by either method in the PBMC dataset. We observed that the CellRanger-only barcodes clustered with barcodes that were detected by both methods (Figure 3a). This suggests that the conservativeness of EmptyDrops is not a major problem, as it only results in the loss of some cells from a cluster that would have been detected irrespectively. In contrast, the EmptyDrops-only barcodes formed a number of unique clusters. One of these clusters likely contains platelets, based on the expression of platelet factor 4 (*PF4*) and pro-platelet basic protein (*PPBP*) (Figure 3b). This is not surprising as the total RNA content of a cell is often associated with its type/state, and platelets have much less RNA than other cell types [14]. Thus, EmptyDrops can capture biology associated with small cells that would have been lost with CellRanger.

A notable side-effect of retaining barcodes with low UMI totals is that a higher number of low-quality cells are also recovered. This manifests as an EmptyDrops-only cluster with high expression of mitochondrial genes and low expression of ribosomal protein genes (Figures 3c, d). Presumably, these cells were damaged during the scRNA-seq protocol, leaking cytoplasmic RNA and enriching for mitochondrial RNA [15]. EmptyDrops is technically correct in retaining the associated barcodes as even damaged cells are distinct from empty droplets. Nonetheless, such damaged cells are usually not of interest in downstream analyses, and most of them can be easily removed by applying an appropriate threshold to the proportion of mitochondrial reads.

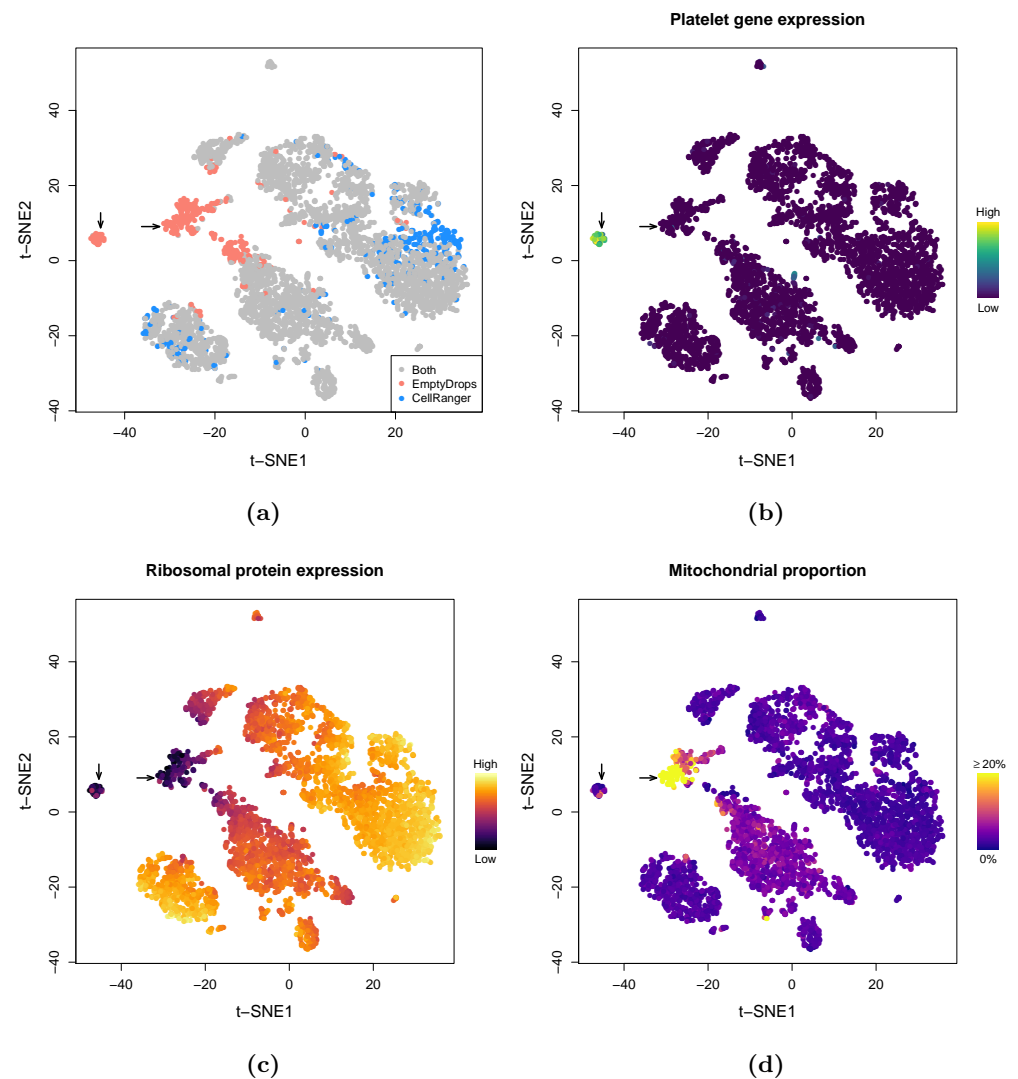
We repeated our analysis on another 10X dataset containing approximately 900 brain cells. Here, the set of cells detected by EmptyDrops was a superset of those detected by CellRanger. In particular, EmptyDrops was able to uniquely retain a cluster expressing *Gad1*, *Gad2* and *Slc6a1* (Supplementary Figure 13). This likely corresponds to a population of interneurons, which would have been almost completely lost using CellRanger. As in the PBMC dataset, we detected a cluster with high mitochondrial RNA content, probably representing damaged cells. We also detected a large population that was primarily characterized by widespread downregulation of genes such as those coding for ribosomal proteins. We hypothesize that these are cells that are so badly damaged that they have lost all cytoplasmic content, including the mitochondria.

## Discussion

Droplet-based technologies are becoming increasingly popular for high-throughput single-cell transcriptomics. However, little work has been performed to develop robust computational methods for distinguishing genuine cells from empty droplets. Here, we describe EmptyDrops, a method to detect cell-containing barcodes based on significant deviation of the expression profiles from the pool of ambient RNA. We use simulated data to demonstrate that EmptyDrops outperforms the strategy currently implemented in the CellRanger software suite. Furthermore, EmptyDrops can recover biology in real 10X data that is lost using the CellRanger strategy. Our results indicate that EmptyDrops is effective for cell detection in droplet-based scRNA-seq data.

A key assumption of our approach is that barcodes with very low UMI totals represent empty droplets. This allows us to use these barcodes to estimate the ambient





**Figure 3.** *t*-SNE plots for the PBMC dataset, constructed using the first 100 principal components of the normalized log-expression matrix. Each point represents a cell that is coloured by (a) detection with either or both EmptyDrops and CellRanger, (b) expression of platelet genes *PF4* and *PPBP*, (c) expression of ribosomal protein genes or (d) the proportion of counts assigned to mitochondrial genes. Expression in each cell was quantified as the sum of the normalized log-expression values across all genes in the relevant set. Mitochondrial proportions are capped at 20% to improve visibility. Arrows mark the putative platelet population and a population of damaged cells.

profile. However, this assumption may not be appropriate if the dataset contains a subset of cells with very low RNA content. In such cases, the estimate of the ambient expression profile will be biased, though this bias is likely to be small as few transcripts will be contributed from cells with low RNA content. Another potential source of bias may arise from sequencing errors in the cell barcode, such that transcripts from a cell-containing droplet are misassigned to an empty droplet. This effect is mitigated by the use of designed cell barcodes in the GemCode protocol, which allows for error correction based on a “whitelist” of known barcode sequences [3]. However, it may be a problem in protocols where error correction of the barcodes is not possible [1].

As we have shown, EmptyDrops is able to detect cells with low RNA content but also recovers more low-quality damaged cells in real data. Many of these cells can be removed by thresholding on the mitochondrial content or other metrics such as the proportion of ribosomal protein mRNA (e.g., if the extent of damage has stripped the cytoplasm from a cell). If this is not sufficient, manual inspection of the clustering results may be necessary to identify these cells and exclude them from further consideration. The other option is to apply a more stringent threshold on the total count, though this will also discard genuine cells with low RNA content and offset the benefits of using EmptyDrops. Nonetheless, EmptyDrops still provides an advantage over existing methods by providing a statistically rigorous framework for cell detection, without requiring any *a priori* knowledge of the expected number of cells.

We have focused exclusively on droplet-based scRNA-seq data generated using the GemCode technology from 10X Genomics. This is motivated by the widespread use of this platform as well as the availability of the unfiltered datasets (see Methods). In principle, the method can also be applied to data from other droplet-based protocols such as inDrop and Drop-seq. Cell lysis or leakage will occur in any protocol involving dissociation and microfluidics, and the formation of empty droplets containing RNA from the ambient pool is unlikely to be a phenomenon that is unique to 10X datasets.

An interesting direction for future work is whether the contribution of the ambient profile can be “subtracted” from each barcode’s expression profile, thus yielding a more accurate representation of each cell’s transcriptome. This is not straightforward as it requires an understanding of the distribution of the volume of the droplets (excluding the volume of the cell inside) to calculate the expected contribution. It is not clear whether such information can be obtained solely from the count matrix for a given dataset. Direct subtraction of the contribution from the counts is also unsatisfactory as it does not preserve the mean-variance relationship. Rather, an identity-link factor model for count data may be required, which is not trivial to implement.

Our EmptyDrops method is implemented in the DropletUtils package, available from the Bioconductor project [16]. We anticipate that it will be useful to researchers who want to extract much information as possible from their droplet-based datasets.

## Methods

### Obtaining droplet-based scRNA-seq datasets

All datasets were downloaded from the 10X Genomics website (<https://support.10xgenomics.com/single-cell-gene-expression/datasets>). Only the “raw” count matrices were used to ensure that CellRanger filtering was not already applied to the cell barcodes. Supplementary Table 1 contains a brief summary of the datasets used in our study.

## Evaluating performance with simulated data

For a given real dataset, we computed the total sum of UMI counts  $t_b$  for each barcode. We identified the inflection point in the curve of  $\log(t_b)$  against the log-rank using the `barcodeRanks` function from the `DropletUtils` package. The set of all barcodes with  $\log(t_b)$  below the inflection point was defined as the set of empty droplets  $\mathcal{G}_0$ . Counts for all  $b \in \mathcal{G}_0$  were summed together to create an ambient pool of RNA molecules. (The inflection point provides a conservative definition of empty droplets, and avoids the inclusion of cell-containing droplets in the ambient pool.) To simulate known empty droplets, we constructed expression profiles for a new set of barcodes by sampling molecules from the ambient pool without replacement. This was done such that the distribution of  $t_b$  in our set was the same as that in  $\mathcal{G}_0$ . In this manner, we recapitulated the observed number of empty droplets and their total counts in our simulations.

To obtain  $G_1$  large cells, we sampled from the set of barcodes with  $\log(t_b)$  above the inflection point. We used sampling with replacement to avoid problems in cases where  $G_1$  is greater than the number of estimated cells in the dataset. To generate  $G_2$  small cells, we sampled from the same set of barcodes and downsampled the count vector for each barcode to 10% of its original total, using the `downsampleMatrix` function from the `DropletUtils` package. This mimics the presence of small cells with low RNA content. We tested different simulation scenarios by setting  $G_1$  or  $G_2$  to 500 and 2000 cells. The various components of the simulation are visualized in Supplementary Figure 2.

We applied our EmptyDrops method to the simulated data at a FDR of 1%. The recall was defined as the proportion of known cells from each group that were successfully detected. The observed false discovery rate was defined as the proportion of detected barcodes that were known empty droplets. We repeated this evaluation using the knee point approach, where all barcodes with total counts above the knee point were retained; and with the CellRanger approach, implemented as described [3] with the expected number of cells set to  $G_1 + G_2$  (i.e., the true number of simulated cells).

We generated simulated data based on each real dataset in Supplementary Table 1. For each scenario and dataset, we repeated the simulation for 10 iterations. We used each method in each iteration and collected performance metrics across all iterations.

## Detecting cells in real data from 10X Genomics

For each dataset in Supplementary Table 1, we applied EmptyDrops to detect cells at a FDR of 1%. We also used the knee point approach, as well as the CellRanger approach where the expected number of cells was set to the reported value in Supplementary Table 1. UpSet plots were created with using the `UpSetR` package [12].

## Characterizing detected cells in real datasets

We analyzed the 10X PBMC dataset by adapting an existing workflow for scRNA-seq data analysis [17]. We performed the analysis on the union of all cells detected by either CellRanger or EmptyDrops to simplify downstream comparisons between the two methods. First, we calculated cell-specific size factors using the deconvolution method [18], with pre-clustering of cells based on the Spearman rank correlation. We divided the counts by the size factors to obtain normalized log-expression values.

We calculated the biological contribution of the variance for each gene, assuming Poisson technical noise when modelling the mean-variance trend. We performed principal components analysis on the log-expression matrix using the `irlba` package. We used the first few components as a low-rank approximation of the matrix to speed up downstream steps. The exact number of components was determined using the

`denoisePCA` function in `scrn`, which matches the sum of biological contributions across all genes to the variance explained by the chosen number of components.

We clustered cells by creating a shared nearest neighbours graph [19] and detecting communities with the Walktrap algorithm from the `igraph` package. Clusters enriched for EmptyDrops-only cells were characterised by detecting differentially expressed genes against every other cluster, using pairwise *t*-tests in the `findMarkers` function from `scrn`. A *t*-SNE plot [13] was generated using from the `Rtsne` and `scater` packages [20]. We used a perplexity of 30, though similar plots were obtained with other values.

We performed a similar analysis on the 900 brain cell dataset, with the main difference being that cells with high haemoglobin expression were removed *a priori*. These likely correspond to contaminating red blood cells that are not of interest. We also used a lower perplexity of 10 for the *t*-SNE plots as fewer cells were present.

## Implementation details

EmptyDrops is implemented as the `emptyDrops` function in the `DropletUtils` package, available from version 3.7 of the Bioconductor project (<https://bioconductor.org/packages/DropletUtils>). It is written in a combination of R and C++ and requires approximately 1-2 minutes to run on each of the tested datasets. All code for simulations and real data analysis were written in R and are available at <https://github.com/MarioniLab/EmptyDrops2017>.

## Author contributions

ATLL, SR, TA, TPD and TG developed the initial EmptyDrops algorithm and tested it on simulated data. TPD tested the algorithm on the PBMC dataset. ATLL improved the efficiency of the algorithm, incorporated it into a package, prepared new simulations for further testing, refined the real data analysis and wrote the manuscript. JCM provided guidance for the project direction.

## Funding statement

ATLL and JCM were supported by core funding from Cancer Research UK (award no. 17197 to JCM). TA was supported by a core grant to the Wellcome Sanger Institute provided by the Wellcome Trust. TG was supported by the European Union's H2020 research and innovation programme "ENLIGHT-TEN" under the Marie Skłodowska-Curie grant agreement 675395.

## Acknowledgements

We would like to thank Jonathan Griffiths for further testing of the algorithm; Elia Benito-Gutierrez for assistance with identifying neuronal markers; and Stephen Sansom for discussions on the nature of cell damage. We would also like to thank the AWS Cloud Credits for Research Program for providing computational resources during the Jamboree.

## References

1. E. Z. Macosko, A. Basu, R. Satija, J. Nemesh, K. Shekhar, M. Goldman, I. Tirosh, A. R. Bialas, N. Kamitaki, E. M. Martersteck, J. J. Trombetta, D. A.

- Weitz, J. R. Sanes, A. K. Shalek, A. Regev, and S. A. McCarroll. Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell*, 161(5):1202–1214, May 2015.
2. A. M. Klein, L. Mazutis, I. Akartuna, N. Tallapragada, A. Veres, V. Li, L. Peshkin, D. A. Weitz, and M. W. Kirschner. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell*, 161(5):1187–1201, May 2015.
3. G. X. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R. Wilson, S. B. Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu, M. T. Gregory, J. Shuga, L. Montesclaros, J. G. Underwood, D. A. Masquelier, S. Y. Nishimura, M. Schnall-Levin, P. W. Wyatt, C. M. Hindson, R. Bharadwaj, A. Wong, K. D. Ness, L. W. Beppu, H. J. Deeg, C. McFarland, K. R. Loeb, W. J. Valente, N. G. Ericson, E. A. Stevens, J. P. Radich, T. S. Mikkelsen, B. J. Hindson, and J. H. Bielas. Massively parallel digital transcriptional profiling of single cells. *Nat Commun*, 8:14049, Jan 2017.
4. S. Islam, A. Zeisel, S. Joost, G. La Manno, P. Zajac, M. Kasper, P. Lonnerberg, and S. Linnarsson. Quantitative single-cell RNA-seq with unique molecular identifiers. *Nat. Methods*, 11(2):163–166, Feb 2014.
5. S. Picelli, A. K. Bjorklund, O. R. Faridani, S. Sagasser, G. Winberg, and R. Sandberg. Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat. Methods*, 10(11):1096–1098, Nov 2013.
6. A. A. Pollen, T. J. Nowakowski, J. Shuga, X. Wang, A. A. Leyrat, J. H. Lui, N. Li, L. Szpankowski, B. Fowler, P. Chen, N. Ramalingam, G. Sun, M. Thu, M. Norris, R. Lebofsky, D. Toppani, D. W. Kemp, M. Wong, B. Clerkson, B. N. Jones, S. Wu, L. Knutsson, B. Alvarado, J. Wang, L. S. Weaver, A. P. May, R. C. Jones, M. A. Unger, A. R. Kriegstein, and J. A. West. Low-coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex. *Nat. Biotechnol.*, 32(10):1053–1058, Oct 2014.
7. A. Regev, S. Teichmann, E. S. Lander, I. Amit, C. Benoist, E. Birney, B. Bodenmiller, P. Campbell, P. Carninci, M. Clatworthy, H. Clevers, B. Deplancke, I. Dunham, J. Eberwine, R. Eils, W. Enard, A. Farmer, L. Fugger, B. Gottgens, N. Hacohen, M. Haniffa, M. Hemberg, S. K. Kim, P. Klennerman, A. Kriegstein, E. Lein, S. Linnarsson, J. Lundeberg, P. Majumder, J. Marioni, M. Merad, M. Mhlanga, M. Nawijn, M. Netea, G. Nolan, D. Pe’er, A. Philipakis, C. P. Ponting, S. R. Quake, W. Reik, O. Rozenblatt-Rosen, J. R. Sanes, R. Satija, T. Shumacher, A. K. Shalek, E. Shapiro, P. Sharma, J. Shin, O. Stegle, M. Stratton, M. J. T. Stubbington, A. van Oudenaarden, A. Wagner, F. M. Watt, J. S. Weissman, B. Wold, R. J. Xavier, and N. Yosef. The human cell atlas. *bioRxiv*, 2017.
8. William A Gale and Geoffrey Sampson. Good-Turing frequency estimation without tears. *Journal of Quantitative Linguistics*, 2(3):217–237, 1995.
9. M. D. Robinson, D. J. McCarthy, and G. K. Smyth. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140, Jan 2010.
10. B. Phipson and G. K. Smyth. Permutation P-values should never be zero: calculating exact P-values when permutations are randomly drawn. *Stat Appl Genet Mol Biol*, 9:Article39, 2010.

11. Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. B*, pages 289–300, 1995.
12. A. Lex, N. Gehlenborg, H. Strobel, R. Vuilleumot, and H. Pfister. UpSet: Visualization of Intersecting Sets. *IEEE Trans Vis Comput Graph*, 20(12):1983–1992, Dec 2014.
13. L. Van der Maaten and G. Hinton. Visualizing data using t-SNE. *J Mach Learn Res*, 9(2579-2605):85, 2008.
14. J. W. Rowley, H. Schwertz, and A. S. Weyrich. Platelet mRNA: the meaning behind the message. *Curr. Opin. Hematol.*, 19(5):385–391, Sep 2012.
15. T. Ilicic, J. K. Kim, A. A. Kolodziejczyk, F. O. Bagger, D. J. McCarthy, J. C. Marioni, and S. A. Teichmann. Classification of low quality cells from single-cell RNA-seq data. *Genome Biol.*, 17:29, Feb 2016.
16. W. Huber, V. J. Carey, R. Gentleman, S. Anders, M. Carlson, B. S. Carvalho, H. C. Bravo, S. Davis, L. Gatto, T. Girke, R. Gottardo, F. Hahne, K. D. Hansen, R. A. Irizarry, M. Lawrence, M. I. Love, J. MacDonald, V. Obenchain, A. K. Oleś, H. Pages, A. Reyes, P. Shannon, G. K. Smyth, D. Tenenbaum, L. Waldron, and M. Morgan. Orchestrating high-throughput genomic analysis with Bioconductor. *Nat. Methods*, 12(2):115–121, Feb 2015.
17. A. T. Lun, D. J. McCarthy, and J. C. Marioni. A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Res*, 5:2122, 2016.
18. A. T. Lun, K. Bach, and J. C. Marioni. Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. *Genome Biol.*, 17:75, Apr 2016.
19. C. Xu and Z. Su. Identification of cell types from single-cell transcriptomes using a novel clustering method. *Bioinformatics*, 31(12):1974–1980, Jun 2015.
20. D. J. McCarthy, K. R. Campbell, A. T. Lun, and Q. F. Wills. Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics*, 33(8):1179–1186, Apr 2017.