

# A Bayesian framework for inter-cellular information sharing improves dscRNA-seq quantification

Avi Srivastava<sup>1,\*</sup>, Laraib Malik<sup>1</sup>, Hiral Sarkar<sup>2</sup> and Rob Patro<sup>2,\*</sup>

<sup>1</sup>Department of Computer Science, Stony Brook University, Stony Brook 11794, NY, USA and <sup>2</sup>Computer Science Department, University of Maryland, College Park 20742, MD, USA

\*To whom correspondence should be addressed.

## Abstract

**Motivation:** Droplet-based single-cell RNA-seq (dscRNA-seq) data are being generated at an unprecedented pace, and the accurate estimation of gene-level abundances for each cell is a crucial first step in most dscRNA-seq analyses. When pre-processing the raw dscRNA-seq data to generate a count matrix, care must be taken to account for the potentially large number of multi-mapping locations per read. The sparsity of dscRNA-seq data, and the strong 3' sampling bias, makes it difficult to disambiguate cases where there is no uniquely mapping read to any of the candidate target genes.

**Results:** We introduce a Bayesian framework for information sharing across cells within a sample, or across multiple modalities of data using the same sample, to improve gene quantification estimates for dscRNA-seq data. We use an anchor-based approach to connect cells with similar gene-expression patterns, and learn informative, empirical priors which we provide to *alevin*'s gene multi-mapping resolution algorithm. This improves the quantification estimates for genes with no uniquely mapping reads (i.e. when there is no unique intra-cellular information). We show our new model improves the per cell gene-level estimates and provides a principled framework for information sharing across multiple modalities. We test our method on a combination of simulated and real datasets under various setups.

**Availability and Implementation:** The information sharing model is included in *alevin* and is implemented in C++14. It is available as open-source software, under GPL v3, at <https://github.com/COMBINE-lab/alevin> as of version 1.1.0.

Contact: [asrivastava@stonybrook.edu](mailto:asrivastava@stonybrook.edu) or [rob@cs.umd.edu](mailto:rob@cs.umd.edu)

## 1 Introduction

RNA-sequencing, with subsequent gene and transcript quantification, has been an important tool for exploring genome-wide expression patterns using both bulk and single-cell experiments. With recent advancements in single-cell transcriptomic sequencing technologies, various droplet-based RNA-sequencing (dscRNA-seq) methods (Klein *et al.*, 2015; Macosko *et al.*, 2015; Zheng *et al.*, 2017) have gained popularity due to their ability to generate data with higher quantitative accuracy, sensitivity and throughput than previous approaches. These dscRNA-seq protocols have a unique output where each read is associated with a cell barcode, to facilitate separation of information between individual cells, and a unique molecular identifier (UMI) tag that allows detecting and deduplicating PCR amplified molecules. Multiple pre-processing pipelines exist that use varying algorithms and methodologies to perform cell barcode correction and whitelisting, read alignment or mapping, and UMI deduplication, and eventually provide gene quantification estimates for each cell. Some of these pipelines use complete alignment of the reads to the reference, such as *alevin* (Srivastava *et al.*, 2019), STARsolo (Dobin, 2019), Cell Ranger (Zheng *et al.*, 2017) and Hasi-T (Tan *et al.*, 2019), whereas others use lightweight mapping methods, such as *brimble* (Mehrad *et al.*, 2019). To the best of our knowledge, each method, except *alevin*, discards reads that multi-map between genes. To date, such approaches

validate accuracy by demonstrating near-perfect correlation to estimates from Cell Ranger.

In *alevin*, Srivastava *et al.* (2019) propose a novel framework for generating accurate gene-expression estimates for each cell given the read sequences from a dscRNA-seq experiment. It is shown how discarding gene multi-mapping reads, as is typically done by other existing dscRNA-seq quantification pipelines, can lead to biased and inaccurate expression estimates for certain genes and gene families. Subsequently, it is also demonstrated that *alevin* reduces this bias by providing a framework for assigning multi-mapping reads to genes rather than discarding them. Specifically, after a UMI resolution and deduplication phase (which assigns multi-mapping UMIs on the basis of parsimony), UMIs are placed into gene-level equivalence classes, associating each UMI with the set of genes to which it maps. Ambiguous reads that belong to equivalence classes with more than one gene label are probabilistically assigned using an expectation-maximization (EM) algorithm. The EM algorithm works by integrating information from reads that are confidently assigned to a single gene, either as a result of the parsimony-based UMI resolution algorithm or because this was the only gene to which the underlying read aligns. This information helps to disambiguate reads that belong to multi-gene equivalence classes, and it is shown, through various analyses, that the framework provides better gene-expression estimates than approaches that discard multi-mapping reads.

However, in situations where there is no unique evidence to disambiguate and assign a read among genes from its equivalence class with some confidence, the optimization method used by alevin uniformly divides the read count across all genes from the single equivalence class. This set of genes is then labeled as Tier 3 in the alevin output. Genes within a cell that have some unique evidence, or share equivalence class with genes that do, are labeled as Tier 2. Hence, Tier 2 genes are assigned read counts with some level of confidence by the EM algorithm. Finally, Tier 1 contains genes that have reads uniquely assigned to them at the UMI deduplication step, and hence their count can be estimated with the greatest confidence by the EM algorithm. This method of equivalence class and tier assignment is further detailed in Figure 1. In this study, we focus on genes labeled as Tier 3, and propose an approach for improving the accuracy of their quantification, instead of uniformly dividing read counts between them.

Our proposed model works by sharing information, either across closely related cells within the sample, or derived in some other fashion from the assay, such as in the case of spatial transcriptomics (ST) data. This information is integrated into the inference algorithm by introducing empirical Bayesian priors, and we show that the proposed Bayesian framework improves gene abundance estimates for Tier 3 genes under various models, based on tests using simulated and real datasets in different setups. The idea of sharing information across data modalities, using an empirical prior, has been previously considered in the context of bulk RNA-seq (Liu, 2016). Relatedly, the idea of sharing information across samples has also been applied in the context of imputation for various types of sparse genomic datasets, such as SNP genotyping and GWAS studies (Chen et al., 2016; Vander et al., 2017). However, for single-cell quantification data, most imputation methods rely on intrinsic properties of the data due the absence of an external reference and work only *post hoc* on already generated gene count matrices (Amodio et al., 2019; Asadkhan et al., 2019; Chen and Zhou, 2019; Deng et al., 2019; Esaukin et al., 2019; Gong et al., 2019; Huang et al., 2019; Li and Li, 2019; Linderman et al., 2019; Lopez et al., 2019; Miao et al., 2019; Mongkha et al., 2019; Talwar et al., 2019; Tang et al., 2019; van Dijk et al., 2019; Wagner et al., 2017; Wang et al., 2019; Zhang and Zhang, 2019). Therefore, they do not have access to either the information contained in, or the constraints imposed by, the UMI-to-gene mappings. Our approach, on the other hand, utilizes shared information directly in the quantification phase to improve UMI assignment and resolution of multi-mapping reads. Furthermore, this information is used only in the form of an empirical prior, and the resulting quantification estimates are still strictly constrained by the observed data. Hence, the likelihood of inducing globally significant false signals, as has been reported in the case of some single-cell RNA-seq imputation methods (Andrews and Hemberg, 2019), is small.

## 2 Materials and methods

### 2.1 Bayesian framework

After UMI deduplication, alevin models the read assignment problem as an optimization problem and iteratively assigns the ambiguous reads to potential candidates in a manner that maximizes (at least locally, within a cell) the joint likelihood. However, it cannot utilize the confidence information from neighboring cells, or from cells of the same type. Since a high level of sparsity is an inherent property of contemporary scRNA-seq experiments (Hicks et al., 2019), and due to the random process of capturing RNA molecules, in expectation, sampling can exhibit considerable variation across cells. Hence, we expect cells in an experiment to fall into categories of specific cell-types, and for cells of the same type to share similar expression patterns (Spear et al., 2019). However, for a specific gene, we do not expect that the molecules originating from the gene will be uniformly captured and sampled equally well across all cells of the type. Therefore, sharing confidence in the expression estimates across cells can be particularly effective in improving cell-level expression estimation. Similarly, we expect information from other assays, using either the same cells or even the same cell-type, to exhibit highly correlated gene abundances. We integrate this information using Bayesian priors by changing our optimization algorithm from an EM algorithm to a variational Bayesian optimization algorithm (Narisal et al., 2013) with an informative prior for low-information genes, i.e. genes assigned to Tier 3. This is a variant of the same collapsed variational Bayesian estimation method (VBEM) used in Salmon (Patro et al., 2017) for bulk RNA-seq abundance estimation.

Similar to Salmon's VBEM, we aim to quantify the expression, given a set of known genes  $\mathcal{G}$  and a set of gene-level equivalence classes  $\mathcal{E}$  with their associated UMIs. Each equivalence class is labeled with a set of genes and has an associated set of UMIs, such that each UMI is attributed to at least one read that multi-maps only across the set of genes in the equivalence class. Here, the set of UMIs are taken after appropriate deduplication using alevin's graph-based UMI deduplication algorithm (Srivastava et al., 2019). We use the VBEM algorithm to allow sharing quantification information across cells via the use of priors. Specifically, we define the gene-UMI count assignment matrix as  $\mathcal{X}$ , where, based on  $\mathcal{E}$ ,  $x_{ij} = 1$  if UMI  $j$  is derived from gene  $i$ . We also define the probability of generating a molecule from a particular gene according to the probability vector  $\mu$  (analogous to the nucleotide fraction in a typical bulk RNA-seq probabilistic model; Li and Dewey, 2011). Hence, we can write the probability of observing a set of deduplicated UMIs  $\mathcal{U}$  as follows:

$$\Pr[\mathcal{U}|\mathcal{X}, \mathcal{G}] = \prod_{j=1}^K \sum_{i=1}^M \Pr[\mu_i|\mu] \cdot \Pr[x_{ij}|\mu_i, x_{ij} = 1] \quad (1)$$

where  $|\mathcal{U}| = N$  is the number of total molecules in the experiment (i.e. the number of deduplicated UMIs) and  $|\mathcal{G}| = M$  is the number of genes.

In this study, we take a variational Bayesian approach to gene expression estimation. Therefore, instead of seeking the maximum-likelihood estimates, we infer (through variational approximation)

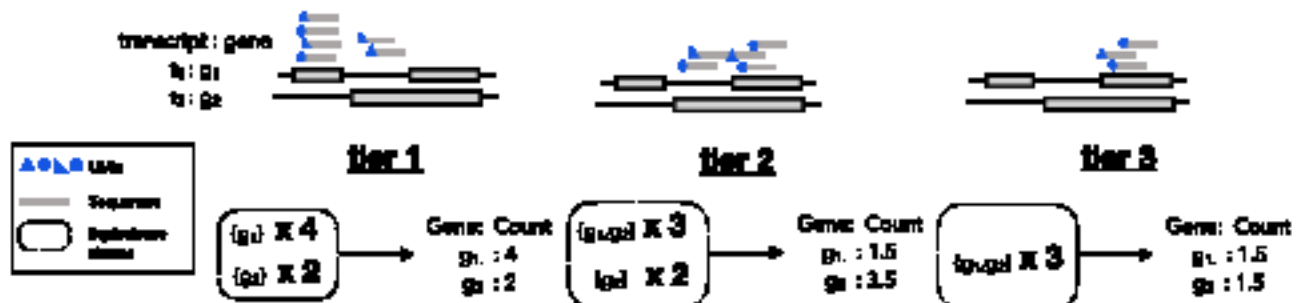


Fig. 1. Alevin categorizes the quantification estimate based on their confidence into three tiers. Here we assume two transcript  $t_1$  and  $t_2$  coming from two genes  $g_1$  and  $g_2$ , respectively. The top row depicts these tiers (from left to right) based on the mapping of its reads. Tier 1 notices two reads from gene unique equivalence classes, Tier 2 notices two reads from genes which share equivalence classes, and finally Tier 3 where no gene in the equivalence class has uniquely mapping reads.

the posterior distribution of  $\mu$ . This posterior distribution can be defined as

$$Pr[\mu|D, G] \propto \sum_{\mathcal{E}} Pr[D|\mathcal{E}, G] \cdot Pr[\mathcal{E}|\mu] \cdot Pr[\mu] \quad (2)$$

where both  $Pr[D|\mathcal{E}, G]$  and  $Pr[\mathcal{E}|\mu]$  can be estimated via a variational approach (Hessman et al., 2015). Although numerous methods for expression estimation from bulk RNA-seq data have previously adopted a variational Bayesian approach (Hessman et al., 2015; Nairki et al., 2013, 2014; Patro et al., 2017), they have all made use of uniform or uninformative priors. The novelty of our method comes from both adopting this approach in the single-cell context, and from setting the prior for  $\mu$  in an informative, data-driven, and cell-specific manner. We expect that, subject to careful selection, information in a single-cell sequencing experiment can be meaningfully shared between distinct but related cells. Note that our method aims to accurately assign reads to the genes to which they map, and does not alter the expression level of genes with zero expression in the data, as may be the case with imputation-based approaches. We explain below how information from related cells, both within a sample and across assays, can be shared, and show how this principle can be applied under various scenarios to improve gene quantification accuracy.

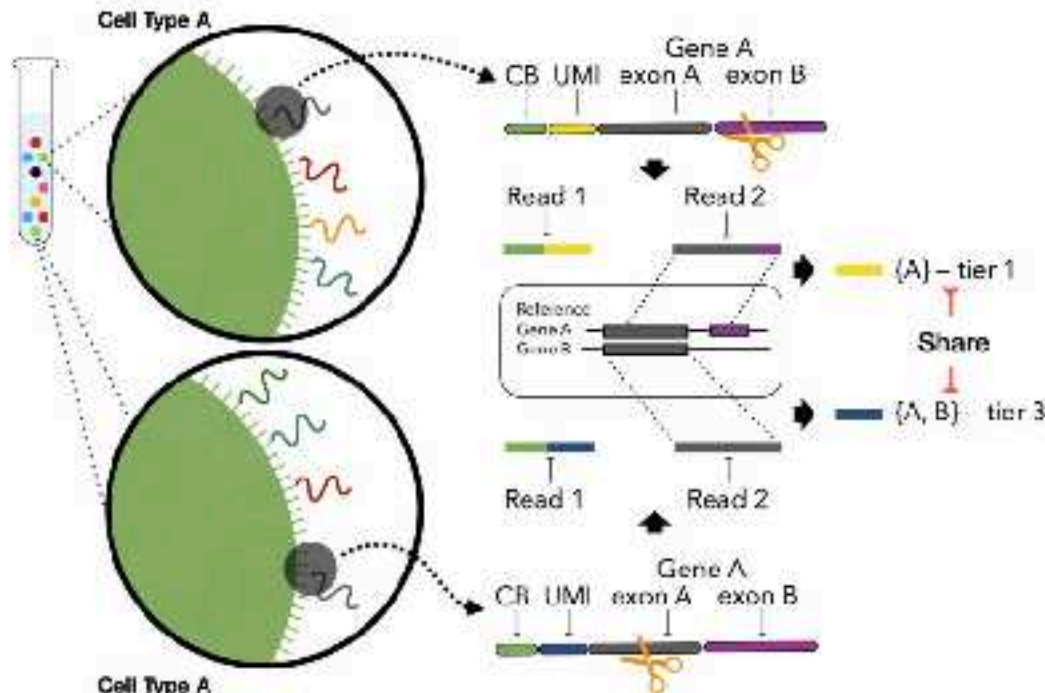
## 2.2 Anchoring to obtain informative priors

Cells of the same type within a sample share similar expression patterns (Stuart et al., 2019). However, due to both biological variability and, *ostensibly*, to the low capture rate and random sampling process in single-cell sequencing experiments, even cells of the same type do not always exhibit near-identical global gene-expression profiles. This means that a given gene from two cells of the same cell-type within a sample could have varying expression estimates, and could be assigned different tiers in individual cells by the *alevin* algorithm. Specifically, a gene may be assigned Tier 3 in one cell and Tier 1 or 2 in the other, based on the specific sequenced reads and UMIs observed, and their mapping patterns. For example, if all of

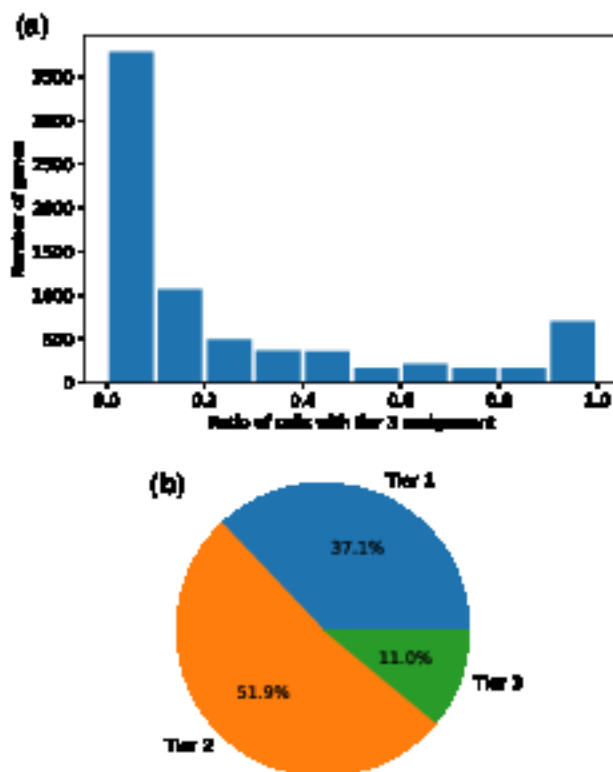
the reads arising from the gene come from an ambiguous region shared with other genes, then this gene will be assigned Tier 3. Whereas if this gene, in another cell of the same type, has sequenced reads coming from a unique region, then it will be assigned as Tier 1 (we have strong evidence of its existence in the cell). Hence, cells of the same type can potentially have different confidence levels in their gene estimates, irrespective of the associated count. This variation can be used to improve quantification of Tier 3 genes. This scenario is depicted in Figure 2, which details how this information can play an instrumental role while quantifying these genes.

To first verify that it is possible to gain information in this way, we look at the fraction of cells that assign a particular gene to Tier 3 out of the total number of cells where the gene is expressed. This is because genes will be informative only when obtained from cells where the gene has uniquely mapping reads (Tier 1) or is influenced by reads mapping uniquely to genes sharing an equivalence class (Tier 2). To do this analysis, we quantify the human PBMC 4k dataset (10x Genomics, 2017), using *alevin* supplemented with the whitelist output by Cell Ranger. This experiment contained a total of 4340 whitelisted cells. The results of this analysis, shown in Figure 3, suggest that most genes are assigned Tier 3 in <10% of the cells and, therefore, estimates from the other cells can be informative. For the 7494 genes that were assigned Tier 3 in at least one cell, 37.1% are assigned Tier 1 and 51.9% are Tier 2 in other cells where the gene is expressed. Hence, the varying degree of confidence in expression estimates across cells can be leveraged in an informative way to improve Tier 3 estimates. Note that all analyses henceforth are done using the 10x PBMC 4k dataset, except where mentioned.

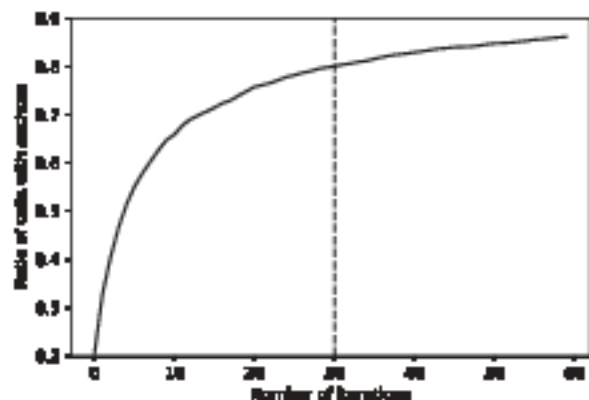
Based on these results, we can see that genes relegated to Tier 3 in a given cell frequently have unique evidence in other cells within the same sample. To take advantage of this property, the next step is finding ‘neighboring’ cells that might be useful for sharing this information to disambiguate read assignment between Tier 3 genes. We only use information from cells with a similar global expression profile within the sample. To find similar cells for sharing this evidence, we use *Scanata* (Stuart et al., 2019) cellular barcode anchoring



**Fig. 2. Motivations.** Given two cells of similar cell type A, we select two reverse transcribed cDNA molecules from Gene A, with unique Cell Barcodes (CBs) and UMIs. Since the fragmentation of a molecule happens at random, the molecule from the first cell (on the top) is fragmented from a region uniquely identifiable for Gene A while the molecule from the second cell (on the bottom) comes from sequence similar region of Gene B. Top cell thus has a high confidence, Tier 1 abundance estimate for Gene A, whereas the bottom cell has a Tier 3 estimate. Assuming the global expression profiles of these cells are similar, our proposed Bayesian model shares this information across cells to improve the quantification estimates for the second cell.



**Fig. 3.** (a) The distribution of number of genes against the fraction of cells that have Tier 3 assignment for three genes. For example, there are around 3000 genes that are assigned Tier 3 in 0.1–0.2 fraction of the total number of cells. (b) The percentage of cells assigning each tier to the genes, showing that the degree of confidence in the quantification estimate varies across cells even for a single gene. Note that both these plots are made using 7464 genes that have been assigned Tier 3 in at least one cell.



**Fig. 4.** Ratio of cells not used in each iteration of the Senseat anchoring algorithm, splitting the dataset into two random, equal sets in each iteration.

scheme that defines a framework to connect two experiments based on the similarity in the gene-expression patterns of the cells assayed in the two experiments. The algorithm works by calculating the  $\ell_2$  distance across datasets, generating two distance matrices and then defines anchors as cells that are neighbors under both distance measures. The full algorithm implemented in Senseat is more involved, and includes various scoring metrics and parameters. Although initially intended for matching cells across samples, we use this anchoring algorithm to connect cells within the sample, in order to define cell-specific priors as input for our Bayesian algorithm.

To generate cell barcode anchors, we first quantify the sample using the standard alevin algorithm (henceforth referred to as EM), and divide the quantification estimates for all cells into two equal

sets. We then run the Senseat anchoring algorithm on these sets, treating the two subsets as two separate samples. In order to identify anchors for a larger number of cells, we repeat the anchoring step multiple times, randomly dividing the quantification into two equal groups each time. We repeat the anchoring step 30 times for all experiments in this article, as we observe on the simulated data that the gain after 30 iterations is small (Fig. 4). We filter the anchors based on the score output by Senseat, using only anchors with a score  $>0.5$ . In a typical single-cell experiment, this is expected to find anchors for about 80% of the cells. The prior for a cell is then defined as the expression estimates, using the original EM-based alevin run, of the cell assigned as the anchor. However, this process can eventually assign multiple anchors for a single cell. To compensate for this, we calculate the prior by taking the average of the expression estimates from all the anchors. This prior is used to optimize the quantification estimates of Tier 3 genes with multi-mapping reads in the alevin pipeline, while keeping the prior uniform for Tiers 1 and 2 genes.

### 3 Results

#### 3.1 Improved estimates using intra-sample information

To test the hypothesis that combining the Bayesian framework with priors obtained from the anchoring procedure described in Section 2.2 can lead to improved quantification estimates for Tier 3 genes, we devised two separate experiments. We detail these two setups below, one relying on simulated data, and the other relying on experimental data with ‘equivalence class knockout (EC)’.

##### 3.1.1 Simulated data

To analyze the improvements in gene quantification estimates on simulated data, we use the empirical droplet-based single-cell RNA-seq data simulation tool MInnow (Sachse et al., 2019). MInnow models various features and protocols involved in the generation of droplet-based single-cell RNA-seq data, like PCR amplification and sequencing errors, to generate fastq files with the reads and the true cell-by-gene count matrix. We use MInnow to simulate a droplet-based single-cell RNA-seq experiment with 4340 cells and ~20 million UMIs using alevin (EM-based) quantification on the 10x PBMC dataset as input. We then compare the quantification estimates against the truth, predicted on the simulated data using alevin, with and without priors, and Cell-Range. The priors from VRM-based alevin were generated, as explained above, using the Senseat anchoring algorithm iteratively.

The results from this analysis are presented in Figure 5a, where VRM represents quantification estimates using priors and EM signifies the quantification estimates without priors. We calculate the Spearman correlation between each method and the ground truth provided by MInnow, focusing on genes assigned Tier 3 in individual cells by alevin. Although the fraction of expressed genes assigned Tier 3 in each cell is low, as shown in Figure 3b, improvement in the accuracy of the gene abundance estimates is significant across hundreds of cells and shows that using informative priors, even from within a sample, can improve quantification. The result also shows that the correlation between estimates from Cell-Range and truth is much lower. This is expected since these genes will have a high number of multi-mapping reads that will be discarded, not just when using Cell-Range but also when using other droplet-based single-cell RNA-seq quantification methods.

##### 3.1.2 Experimental data with ECs

To test that our proposed VRM method, given informative priors, can improve the accuracy of experimental data quantification, we performed an experiment that we refer to as equivalence class KO. Aleavin’s pipeline for droplet-based single-cell RNA-seq quantification has multiple phases. After the initial phase of cell barcode whitelisting and read mapping, alevin outputs an intermediate file. This file contains details of the transcript equivalence classes, including the associated cell barcodes and UMI counts. These equivalence classes are similar to the gene-level equivalence classes explained before, except that the class labels are transcripts that share UMIs after the deduplication step.

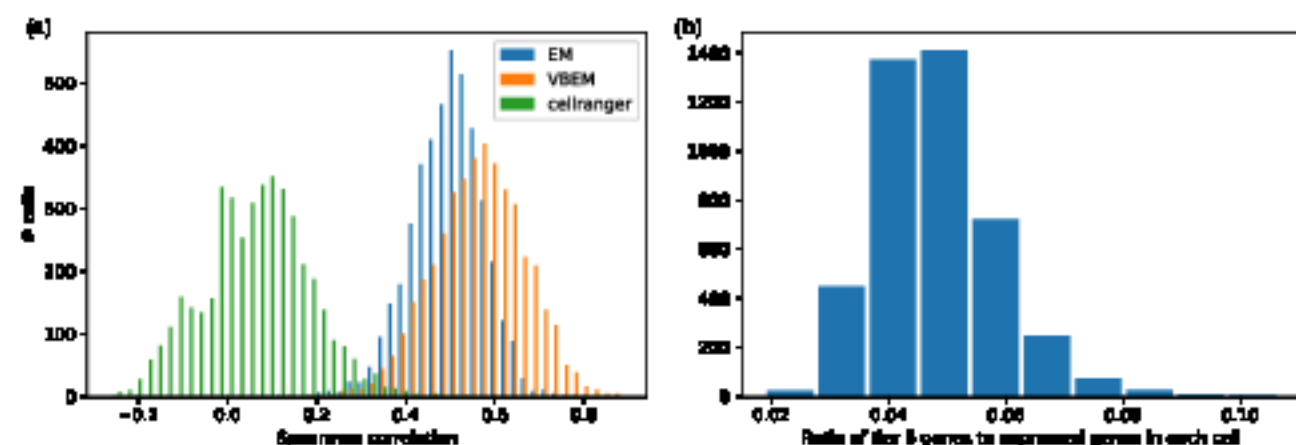


Fig. 5. (a) Comparison of the cell-wise Spearman correlation of Tier 3 genes quantified using Cell Ranger, EM, and VBEM-based aligners on simulated experiment. (b) Ratio of Tier 3 genes in each cell (shown as genes that vary in expression by the genes, leading to increased correlation with the truth)

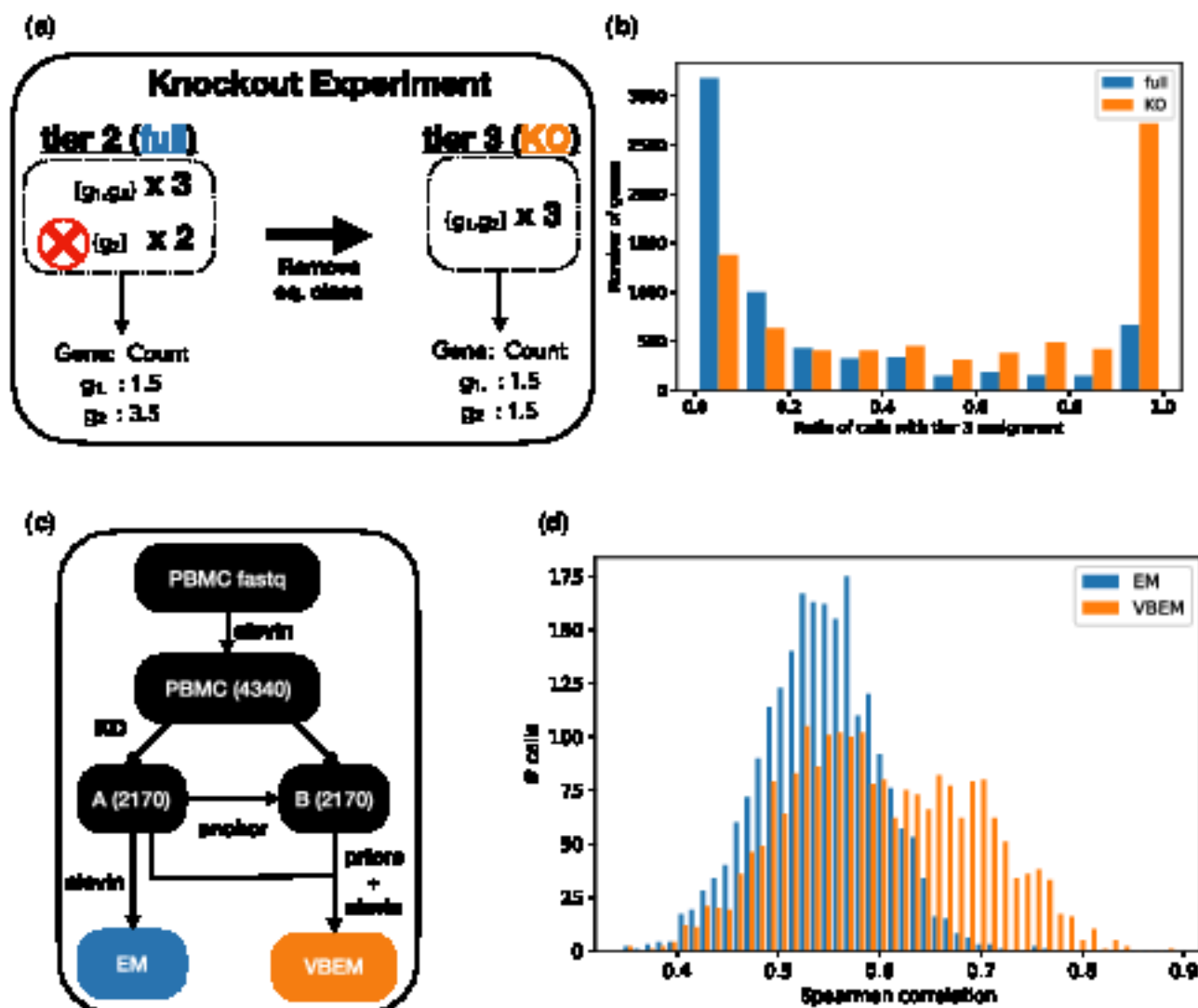


Fig. 6. (a) A toy example explaining the KO experiment. We use the equivalence classes from an aligner (EM based) run on the full PBMC dataset and remove all the transcripts using an equivalence classes to generate a KO dataset. In this example, a source such gene has a single transcript. (b) The distribution of number of genes against the fraction of cells that have Tier 3 assignment for both in the original dataset and the KO dataset, with a shift toward more Tier 3 assignments and increased multiplicity. (c) In the KO framework, to evaluate the improved by genes approach, we designed the following pipeline. We quantify the full human PBMC dataset (4340 cells) cells and randomly divide the experiment into two equal parts (A and B). We knockout using an equivalence classes in Set A (2170) cells and repeat the quantification steps to generate EM-based estimates. In parallel, we also quantify Set A KO dataset using the same pipeline. These priors are then used from Set B, without the KO and quantified in this by using the EM approach. This gives us the VBEM estimates on Set A for comparison. (d) Comparison of the cell-wise Spearman correlation for Tier 3 genes from EM-based aligners with VBEM-based, prior VBEM-based aligners on real data with KO (removal of unique equivalence classes). This shows improved estimate as under VBEM, with a higher correlation against the initial dataset, without KO.

We observe, that a majority of the genes assigned to Ties 1 and 2 after UMI deduplication are the ones associated with transcript equivalence classes of size 1 (labeled by a single transcript). In order to increase ambiguity in this data, we can remove all the transcript-unique equivalence classes from the intermediate file, with the expectation that this KO will result in a number of genes misassigning from Ties 1 to 3, as demonstrated in Figure 6a. In essence, by doing this, we are removing some of the read evidence that will eventually lead to high confidence gene abundance estimates in Ties 1 and 2. The impact of this on the distribution of Tie 3 genes in the PBMC dataset is shown in Figure 6b. This shows that the KO process results in an increased number of genes that are assigned Tie 3 across all cells. Note that the KO dataset will also have a smaller number of UMIs, because of the removal of unique equivalence classes from the intermediate file, but this will not impact our comparative analysis, as explained below. Also, note that we only knock out here equivalence classes that are transcript-unique, and that there will still be a considerable number of gene-unique equivalence classes after post-mortem UMI deduplication has taken place.

In order to ensure that we have cells with high confidence quantification estimates to provide our KO cells with an informative prior, we did not perform the KO procedure on the complete PBMC dataset. Instead, we took the alevin quantification estimates on the PBMC dataset, which has 4340 cells, and divided it into two sets containing equal numbers of cell barcodes. One of these sets, *A*, is our test set from which we KO unique equivalence classes and the other set, *B*, is used to generate priors. Iteratively using the Sensat anchoring algorithm as before, we first find anchors for set *A* in set *B*, then obtain priors from set *B* and run the alevin VREd quantification method. We also quantify the KO set *A* using the EId-based alevin method. These steps are outlined in Figure 6c. Observe that there can be a bias in the tie assignment of genes that are anchors for Tie 3 genes in the KO experiment. This is because we are removing equivalence classes only in set *A*. Hence, the ratio of anchor genes in set *B* that are assigned Ties 1 and 2 in KO may be higher than in the real dataset. This can amplify the accuracy of the VREd method in the KO experiment, but will also reflect the actual gain possible for this methodology under varying circumstances, such as in samples with higher read depth. Note that we cannot run Cell Ranger on this dataset because it utilizes the intermediate file output by alevin, which cannot be processed directly. However, we expect similar results as those observed in simulated data, since multi-mapping genes are not quantified by Cell Ranger.

In our comparison between the two methods, we find the Spearman correlation for each cell between the original, EId-based alevin estimates of the cells in set *A* and the estimates using the KO set *A* under each method. Because the original set *A* has more high confidence Ties 1 and 2 genes, we expect the estimates to be of higher accuracy. The results from this analysis are presented in Figure 6d, which shows that the cell-wise correlations of the VREd predicted abundances on the KO dataset are higher compared with the original estimates than are the EId estimates on the KO dataset. Note that these correlations are calculated for genes that are assigned Tie 3 in the KO set *A*, since those are the only genes impacted by the priors. This test shows that utilizing the anchoring procedure and extracting informative priors, combined with using a VREd-based quantification procedure, can lead to higher accuracy in abundance estimation.

It is also interesting to note that the anchoring scheme finds high scoring anchors between set *A* and set *B* for only 934 cells. The effect of this limited anchoring shows up in the correlation histogram as a bimodal distribution in the VREd correlation values, signifying that, as expected, only some of the cells—those for which we were able to find an anchor in the set *B*—have improved correlation with the original quantification estimates.

### 3.1.3 Informative sharing does not affect rare cell types

A common concern when sharing information across cells in scRNA-seq analysis is that it may contribute to loss of heterogeneity among the quantified cells (Andrew and Hemberg, 2019; Huang et al., 2019), removing not only technical ‘noise’, but also important

biological variability that leads to the detection of important features, such as rare cell types. To test the hypothesis that the proposed Bayesian framework does not ‘over-regulate’ and lose rare cell types in downstream processing, we perform the following experiment. We use the human PBMC dataset with 10k cells (10x Genomics, 2019) and quantify the cells with both the EId- and VREd-based approaches, where, for the VREd-based approach we use the same procedure of generating priors as discussed in Section 3.1.1. Next, we perform Sensat (Satija et al., 2015) based clustering on the estimates generated from both the approaches separately and compare the clusters.

In Figure 7, we show the 2D UMAP embeddings of the clustered data, colored by cell-type annotations generated using marker genes, as detailed in the Sensat pipeline (Strat et al., 2019). We observe that the clusters with relatively smaller number of cells, such as gDC, Megakaryocytes and Dendritic cell, are not lost by the Bayesian correction method. In Table 1, we show that the number of cells is almost always preserved in the most abundant cluster of each cell type across the two quantification approaches. We also observe that CD14+ Monocytes and CD8 effector cell types are divided into two subclusters when quantified with EId while they are correctly identified as one in case of VREd.

## 3.2 Improved estimates using multi-modal information

### 3.2.1 ST data

Advancements in ST have enabled scientists to relate cells with their location within a tissue. Specifically, it has been shown how combining ST with gene-expression profiling in cancer data helps understand multiple components of tumor progression and therapy outcomes (Theune et al., 2019). The 10x Genomics Visium is another interesting assay that provides higher resolution and throughput for spatial gene-expression analysis. We use the open dataset provided by 10x Genomics of the fresh frozen mouse brain tissue with 2698 spots in the tissue and process the raw reads through the alevin framework to generate a gene count matrix for each spot.

To test the Bayesian framework of alevin, we simulate 2698 cells using the gene count matrix generated by processing the mouse brain ST visium data from 10x Genomics (2019). We first run EId-based alevin on the simulated data and use the spatial 2D coordinates from the ST data to learn the prior, i.e. for each cell we use the nearest eight cells and their mean gene expression from the EId estimates to generate the prior matrix. Then, we provide alevin with the prior matrix to re-quantify the simulated data using the Bayesian method to generate VREd-based estimates. In Figure 8, we show the cell-wise Spearman correlation of Tie 3 gene estimates for both EId and VREd-based methods. We observe a global shift in the VREd quantified data, reflecting the increased accuracy obtained using informative priors from cells located spatially close together. This result is particularly interesting, as it suggests that the empirical Bayesian framework we have introduced is modular and flexible, in that the generation of an informative prior is not tied to a specific procedure (e.g. the Sensat-based anchoring). Rather, the prior can be informed by data in the same sample, by assay-specific information (nearby or differential cell clusters in spatial data as also shown by Årjō et al. (2019)) and, perhaps, even across distinct modalities (e.g. between Assay for Transposase-Accessible Chromatin using sequencing (ATAC-seq) and RNA-seq for cells assayed with both protocols in the same sample).

## 4 Discussion and conclusion

In this work, we improve upon our previously proposed alevin pipeline for quantification of droRNA-seq data. The existing alevin pipeline uses a maximum likelihood-based procedure after the UMI deduplication phase to accurately resolve multi-mapping reads, which are typically discarded by other methods. Although this approach uses unique read evidence from within a cell to optimize read assignment, it uniformly divides read counts where no unique evidence is available. The set of genes with this uniformly divided distribution are assigned Tie 3 by in the alevin output. Our

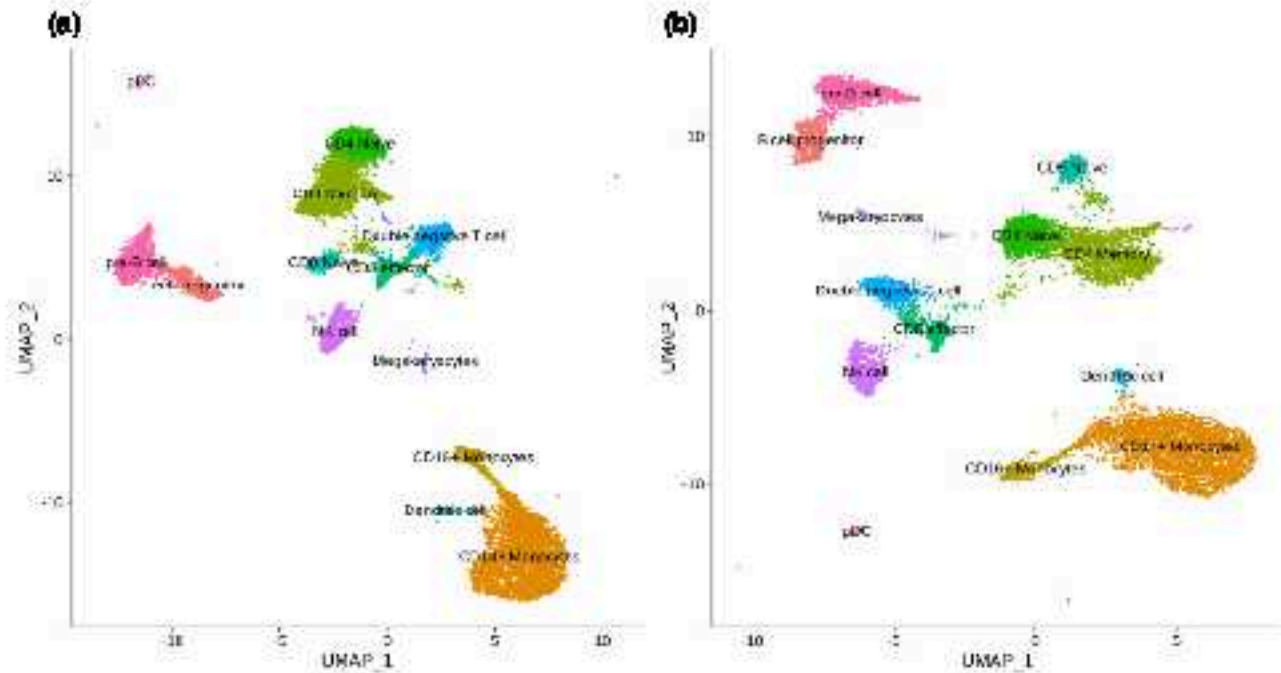


Fig. 7. We perform cell clustering using Seurat on PBMC 30k data set quantified using the EM (left) and VBEM (right) approaches in alevin, and color the cells based on their cell type as annotation generated using our star gene.

Table 1. The number of cells observed across various cell types is similar when clustering is performed on EM- and VBEM-based quantification estimates, suggesting that the information sharing approach does not eliminate meaningful heterogeneity in gene expression between cells.

Cell type/no. of cells	EM	VBEM
CD4+ monocytes	318	322
CD4 effector	222 + 144	358
CD4 naive	1015	1021
Megakaryocytes	49	49
NE cell	517	522
CD4+ monocytes	1758 + 1211	2962
pDC	68	68
CD4 naive	333	331
B-cell progenitor	455	453
Dendritic cell	74	74
CD4 memory	1428	1416
Double negative T cell	587	583
Treg-B cell	516	525

Note: The annotations are generated using Seurat's marker gene analysis.

proposed method uses a Bayesian framework to improve Tbx3 gene quantification.

This method works by sharing high confidence quantification information between cells. Information is shared only across cells that have similar gene-expression profiles (or which are spatially proximate in the case of ST data), but the exact expression estimates vary due to sparsity and uneven RNA capture in single-cell sequencing protocols. We show that, under several different experimental setups, our information-sharing framework consistently improves Tbx3 scRNA-seq quantification estimates. This approach is especially useful for highly ambiguous estimates where there is no inter-cellular unique information available to accurately quantify the genes, but where simply discarding the multi-mapping reads would lead to the loss of potentially important biological information.

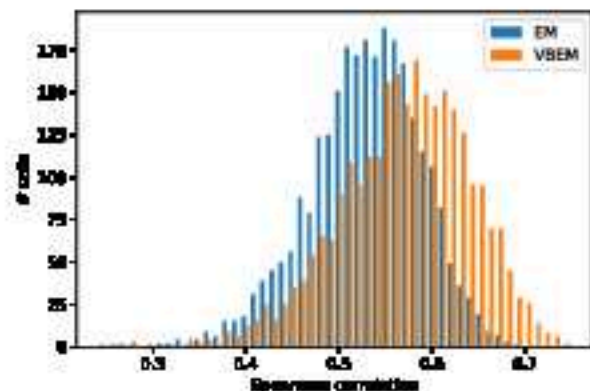


Fig. 8. Comparison of the cell-wise Spearman correlation of Tbx3 genes quantified using EM- and VBEM-based alevin on ST data.

Although we have focused on Tbx3 genes in this study, this information sharing model can be extended further to improve the UMI deduplication procedure as well, before the construction of equivalence classes. For example, instead of basing UMI deduplication on the principle of parsimony in alevin, priors can be used to drive deduplication. This can lead to improvements in abundance estimates for all genes in the reference. Similarly, with advances in single-cell sequencing protocols, this framework can be extended to incorporate priors from different technologies. For example, as we have demonstrated, spatial data can be useful for setting the prior in the proposed alevin framework. This improves accuracy by relying not on similar gene-expression profiles, but cells that are in close proximity in physical space. Further, one can imagine that other assays, like paired single-cell ATAC-seq and RNA-seq, would allow useful information sharing within the same sample but across data types and modalities. We believe this framework has the potential to open a new direction of enabling multi-modal information sharing to improve quantification of single-cell data.

## Resources

The pipeline to replicate the analysis can be found at [https://github.com/CCMIB/INB4sh/levin-pape-pipeline/tree/master/Analysis\\_of\\_levin](https://github.com/CCMIB/INB4sh/levin-pape-pipeline/tree/master/Analysis_of_levin). We used the genome20 reference for human and genome20mm10 for the mouse reference. We use Seurat version 3.0.2 and cell Ranger version 3.1 with the following commands:

1. index: cellranger mkref --genome-ref --fasta-genome-h --genes--genes.gtf --cellranger=16
2. quantification: cellranger count --id--cellranger --fastq--fastq --no-colocates=20 --no-colocates=120 --transcriptome-ref

## Acknowledgements

The authors would like to thank Michael Lové, Rabih Esfah, Tim Sauer and Charlotte Keane for fruitful discussions surrounding some of the approaches taken in this article.

## Author's contributions

A.S., I.M., E.S. and R.P. devised the methods. A.S., I.M. and R.P. implemented the methods. A.S., I.M. and R.P. devised and carried out the experiments. All authors wrote and approved the article.

## Funding

This work has been funded by National Institutes of Health [R01 HG022937], NCI [P01-173047] and CNS-176600 to R.P.].

**Conflict of interest:** R.P. is a co-founder of Ocean Genomics Inc.

## References

10X Genomics. (2017) 10x v2 human phase 4k data. <https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/phase4k>. (4 February 2022, date last accessed).

10X Genomics. (2018) 10x v3 human phase 10k data. [https://support.10xgenomics.com/single-cell-gene-expression/datasets/3.1.0/phase\\_10k\\_v3](https://support.10xgenomics.com/single-cell-gene-expression/datasets/3.1.0/phase_10k_v3). (4 February 2022, date last accessed).

10X Genomics. (2019) 10x mouse brain spatial data. [https://support.10xgenomics.com/spatial-gene-expression/datasets/1.0.0/V1\\_Adult\\_Mouse\\_Brain](https://support.10xgenomics.com/spatial-gene-expression/datasets/1.0.0/V1_Adult_Mouse_Brain). (4 February 2022, date last accessed).

Agar, T. et al. (2019) Splenic: robust estimation of a ligand spatial temporal gene expression data. *bioRxiv*, 737296.

Asanido, M. et al. (2019) Exploiting single-cell data with deep ambient-aware networks. *Nat. Methods*, 16, 1139–1137.

Andrew, T.B. and Hoshino, M. (2018) Pulse signals induced by single-cell inspection. *PLoS ONE*, 7, 1740.

Arifki Kusnata, C. et al. (2019) DeepInsight: an accurate, fast, and scalable deep neural network method to inspect single-cell RNA-seq data. *Genome Biol.*, 20, 1–14.

Chen, M. and Zhou, X. (2018) Viper: variability-preserving imputation for accurate gene expression recovery in single-cell RNA sequencing studies. *Genome Biol.*, 19, 1–15.

Chen, W.-C. et al. (2016) A combined reference panel from the 1000 Genomes and UK10K projects improved rare variant imputation in Europeans and Chinese samples. *Sci. Rep.*, 6, 35813.

Cheng, Y. et al. (2019) Feasible analysis of cell-type composition from single-cell transcriptomics using deep structure learning. *Nat. Methods*, 16, 311–314.

Debin, A. (2019) <https://github.com/alecdebin/STARsolo/blob/master/g2.7.3a>. (4 February 2022, date last accessed).

Braden, G. et al. (2019) Single-cell RNA-seq denoising using a deep count autoencoder. *Nat. Commun.*, 10, 1–4.

Gong, W. et al. (2018) DRImpute: imputing dropout events in single cell RNA sequencing data. *BMC Bioinformatics*, 19, 222.

Grassie, J. et al. (2017) Fast and accurate gene-wise inference of transcript expression from RNA-seq data. *Bioinformatics*, 32(21–22), 3229.

Heckel, S.C. et al. (2018) Missing data and technical variability in single-cell RNA-sequencing experiments. *Bioinformatics*, 33, 562–572.

Huang, M. et al. (2018) Error: gene expression recovery for single-cell RNA-sequencing. *Nat. Methods*, 15, 539–542.

Klein, A.M. et al. (2015) Simple barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell*, 160, 1187–1198.

Li, J. and Dewey, C.N. (2011) RSEM: a accurate transcript quantification from RNA-seq data with or without a reference genome. *BMC Bioinformatics*, 12, 323.

Li, W.-V. and Li, J. (2018) An accurate and robust imputation method for single-cell RNA-seq data. *Nat. Commun.*, 9, 1–9.

Linderman, G.C. et al. (2018) Zero-probability imputation of scRNA-seq data using low-rank approximation. *bioRxiv*, 397322.

Liu, F. et al. (2016) Integrative analysis with ChIP-seq elucidates the limits of transcript quantification from RNA-seq. *Genome Research*, 26, 1124–1133. [DOI: 10.1101/059174](https://doi.org/10.1101/059174).

Lopez, J. et al. (2018) Deep generative modeling for single-cell transcriptomics. *Nat. Methods*, 15, 1213–1222.

Marcenko, E.E. et al. (2015) Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell*, 161, 1121–1134.

McDonald, P. et al. (2019) Modular and efficient pre-processing of single-cell RNA-seq. *bioRxiv*, 673221.

Min, Z. et al. (2019) scRecover: Distinguishing true and false genes in single-cell RNA-seq data for imputation. *bioRxiv*, 663323.

Mongk, A. et al. (2019) McImpute: matrix completion based imputation for single cell RNA-seq data. *Front. Genet.*, 10, 9.

Narain, N. et al. (2013) TIGER: transcript isoform abundance estimation method with gap-filling of RNA-seq data by statistical Bayesian inference. *Bioinformatics*, 29, 2292–2299.

Narain, N. et al. (2014) TIGER2: sensitive and accurate estimation of transcript isoform expression with longer RNA-seq reads. *BMC Genomics*, 15, 51.

Patro, R. et al. (2017) Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods*, 14, 417–419.

Sarkar, J.E. et al. (2019) MGenow: a principled framework for rapid simulation of de novo RNA-seq data at the read level. *Bioinformatics*, 35, 2136–2144.

Saupe, J. et al. (2013) Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.*, 33, 493–492.

Schmitt, A. et al. (2019) A beta efficiently estimates accurate gene abundance from de novo RNA-seq data. *Genome Biol.*, 20, 61.

Swart, T. et al. (2019) Comprehensive integration of single-cell data. *Cell*, 177, 1825–1832.e11.

Talwar, D. et al. (2018) AutoImpute: autoencoder based imputation of single-cell RNA-seq data. *Sci. Rep.*, 8, 1–11.

Tang, W. et al. (2018) bayNorm: Bayesian gene expression recovery, imputation and normalization for single cell RNA-sequencing data. *bioRxiv*, 324326.

Thasak, K. et al. (2018) Spatially resolved transcriptomic enables detection of genetic heterogeneity in mouse brain multiplexed nucleases. *Cancer Res.*, 78, 5970–5979.

Traut, T. et al. (2019) Iso-T: an efficient and accurate approach for quantifying gene abundances from RNA-sequencing data with high rates of non-coding reads. *bioRxiv*, 520171.

van Dijk, D. et al. (2018) Recovering gene interactions from single-cell data using diffusion. *Cell*, 174, 716–729.

Vinches, F.M. et al. (2017) 10 years of GWAS discovery: biology, function, and translation. *Am. J. Hum. Genet.*, 101, 5–22.

Wagner, P. et al. (2017) K-nearest neighbor smoothing for high-throughput single-cell RNA-seq data. *bioRxiv*, 217737.

Wu, J. et al. (2019) Data denoising with tensor factorization in single-cell transcriptomics. *Nat. Methods*, 16, 873–878.

Zhang, J. and Zhang, S. (2018) FILL: an accurate single cell RNA-seq data imputation tool considering cell heterogeneity and prior expression level of droplets. *bioRxiv*, 379823.

Zhang, G.-X. et al. (2017) Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.*, 8, 14045.