

METHOD

Open Access



# Cobo1t: integrative analysis of multimodal single-cell sequencing data

Boying Gong<sup>1</sup>, Yun Zhou<sup>1</sup> and Elizabeth Purdom<sup>2\*</sup>

\*Correspondence:

[epurdom@stat.berkeley.edu](mailto:epurdom@stat.berkeley.edu)

<sup>2</sup>Department of Statistics, University of California, Berkeley, Berkeley, CA, USA

Full list of author information is available at the end of the article

## Abstract

A growing number of single-cell sequencing platforms enable joint profiling of multiple omics from the same cells. We present Cobo1t, a novel method that not only allows for analyzing the data from joint-modality platforms, but provides a coherent framework for the integration of multiple datasets measured on different modalities. We demonstrate its performance on multi-modality data of gene expression and chromatin accessibility and illustrate the integration abilities of Cobo1t by jointly analyzing this multi-modality data with single-cell RNA-seq and ATAC-seq datasets.

**Keywords:** Single cell, Multi-omics, Integration

## Background

Single-cell sequencing allows for quantifying molecular traits at the single-cell level, and there exist a wide variety of platforms that extend traditional bulk platforms, such as mRNA-seq and ATAC-seq, to the single cell. Comparison of different cellular features or *modalities* from cells from the same biological system gives the potential for a holistic understanding of the system. Most single-cell technologies require different cells as input to the platform, and therefore, there remains the challenge of linking together the biological signal from the different modalities, with several computational methods proposed to estimate the linkage between the different modalities, such as LIGER [1] and Signac (Seurat) [2, 3].

Recently, there are a growing number of platforms that allow for measuring several modalities on a single cell. CITE-seq [4] jointly sequences epitope and transcriptome; scNMT-seq [5] jointly profiles chromatin accessibility, DNA methylation, and gene expression; and sci-CAR [6], Paired-seq [7], and SNARE-seq [8] enable simultaneous measurement of transcription and chromatin accessibility (we direct readers to [9] for a comprehensive review). By directly measuring the different modalities on the same cells, these techniques greatly enhance the ability to relate the different modalities. With the emergence of joint platforms, new computational methodologies for analyzing multi-modality data have also been developed. Early methods mainly focused on CITE-seq



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

[10–12], which jointly sequences gene expressions and at most a few hundred antibodies. Recently, more methods have been proposed to enable the modeling of cells with simultaneous measurement of gene expression and higher-dimensional modalities such as chromatin accessibility, such as MOFA+ (also known as MOFA2 [13]), scMVAE [14], BABEL [15], and scMM [16].

However, single-cell datasets on a single modality are far more common and are usually of higher throughput. Indeed, it is natural that joint-modality data from the same system will be used to augment single-modality data, or vice versa. Therefore, there is a critical need for an analysis tool that is both a stand-alone application for multi-modality data as well as a tool for integration of these datasets with single-modality platforms. BABEL [15] and scMM [16], while not directly targeting this task, do allow the use of the joint-modality data to predict one single-modality dataset into another type of modality. However, neither directly integrate the data together to allow for downstream analysis of the joint set of data regardless of modality, such as cell subtype detection.

Our method **Cobolt** fills this gap by providing a coherent framework for a full integrative analysis of multi-modality and single-modality platforms. The result of **Cobolt** is a single representation of the cells irrespective of modalities, which can then be used directly by downstream analyses, such as joint clustering of cells across modalities. **Cobolt** estimates this joint representation via a novel application of Multimodal Variational Autoencoder (MVAE) [17] to a hierarchical generative model. The integration of the single-modalities is done by a transfer learning approach which harnesses the valuable information found by joint sequencing of the same cells and extends it to the cells in the single-cell platform. The end result is a single representation of all of the input cells, whether sequenced on a multi-modality platform or a single-modality platform. In this context, **Cobolt** gives an over-all integrative framework that is flexible for a wide range of modalities.

We demonstrate **Cobolt** on two use-cases. The first uses **Cobolt** to analyze only a multi-modality sequencing dataset from the SNARE-seq technology; we show that **Cobolt** provides a joint analysis that better distinguishes important facets of each modality, compared to existing methods. The second demonstrates the use of **Cobolt** to integrate multi-modality data with single-modality data collected from related biological systems, where **Cobolt** creates a joint representation that can be used for downstream analysis to provide meaningful biological insights. We show that **Cobolt** also performs better than related tools in this integrative task.

## Results

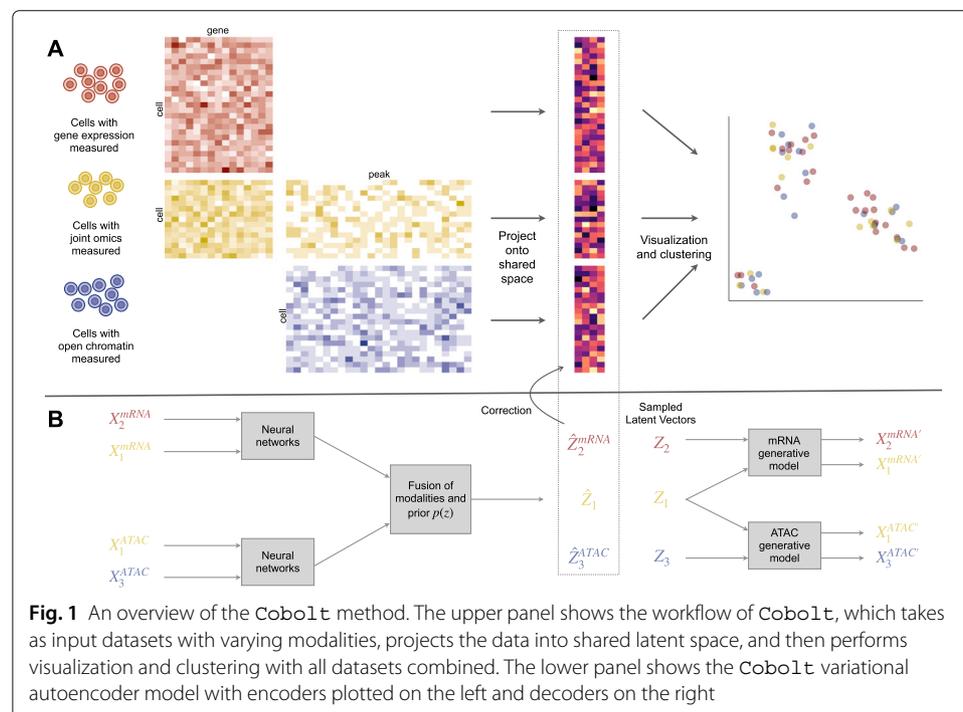
### The **Cobolt** model

We develop a novel method, **Cobolt**, that utilizes joint-modality data to enable the joint analysis of cells sequenced on separate sequencing modalities. We do this by developing a Multimodal Variational Autoencoder based on a hierarchical Bayesian generative model. We briefly describe the premise of the model using the example of two modalities: mRNA-seq and ATAC-seq (for more details in greater generality, see the “[Methods](#)” section). We assume that we have a set of cells with both mRNA-seq and ATAC-seq data collected using the joint-modality platform ( $X_1^{mRNA}$  and  $X_1^{ATAC}$ ), as well as (optionally) a set of cells with only mRNA-seq data ( $X_2^{mRNA}$ ), and a set of cells with only ATAC-seq data

( $X_3^{ATAC}$ ). Cobolt takes all of this data as input to find a representation of the cells in a shared reduced dimensionality space, regardless of modality (Fig. 1A).

Cobolt models the sequence counts from the modalities inspired by the Latent Dirichlet Allocation (LDA) [18] model. The LDA model is a popular Bayesian model for count data which has been successfully applied to genomics in areas such as clustering of single-cell RNA-seq data [19, 20], single-cell ATAC-seq analysis [21], and functional annotation [22]. Cobolt builds a hierarchical latent model to model data from different modalities and then adapts the MVAE approach to both estimate the model and allow for a transfer of learning between the joint-modality data and the single-modality data.

Cobolt assumes that there are  $K$  different types of possible categories that make up the workings of a cell. For ease of understanding, it is useful to think of these categories as biological processes of the cell, though the categories are unlikely to actually have a one-to-one mapping with biological processes. Each category will result in different distributions of features in each of the modalities—i.e., different regions of open chromatin in ATAC-seq or expression levels of genes in mRNA-seq for different categories. The features measured in a cell are then the cumulative contribution of the degree of activation of each category present in that cell. The activation level of each category is represented by the latent variable  $\theta_c$  for each cell  $c$ , which gives the relative activity of each of the  $K$  categories in the cell.  $\theta_c$  is assumed to be an intrinsic property of each cell representing the underlying biological properties of the cell, while the differences of data observed in each modality for the same cell are due to the fact that the categories active in a cell have different impacts in the modality measured (open chromatin in ATAC-seq versus gene expression in mRNA-seq). We assume  $\theta_c = \sigma(z_c)$ , where  $z_c$  is drawn from a Gaussian prior and  $\sigma$  is the soft-max transformation; this is an approximation to the standard Dirichlet prior for  $\theta_c$  that allows use of variational autoencoders to fit the model [23]. The



mean of the posterior distribution gives us an estimate of our latent variable  $z_c$  for each cell, and the posterior distribution is estimated using variational autoencoders (VAE).

In the end, **Cobolt** results in an estimate of the latent variable  $z_c$  for each cell, which is a vector that lies in a  $K$ -dimensional space. This space represents the shared biological signal of the individual cells, irregardless of modality, and can be used for the common analysis tasks of single-cell data, such as visualization and clustering of cells to find subtypes. Importantly, we can predict the latent variable  $z_c$  even when a cell does not have all modalities measured. Moreover, because of the joint-modality platforms, **Cobolt** does not require that the different modalities measure the same features in order to link the modalities together—the fact that some of the cells were sequenced on both platforms provides the link between different types of features. Therefore, ATAC-seq peaks and mRNA-seq gene expression can be directly provided as input. This is unlike methods that do not make use of the joint-modality data and require that the different modalities be summarized on the same set of features, for example by simplifying ATAC-seq peaks to a single measurement per gene.

#### **Cobolt as an analysis tool for multi-modality data**

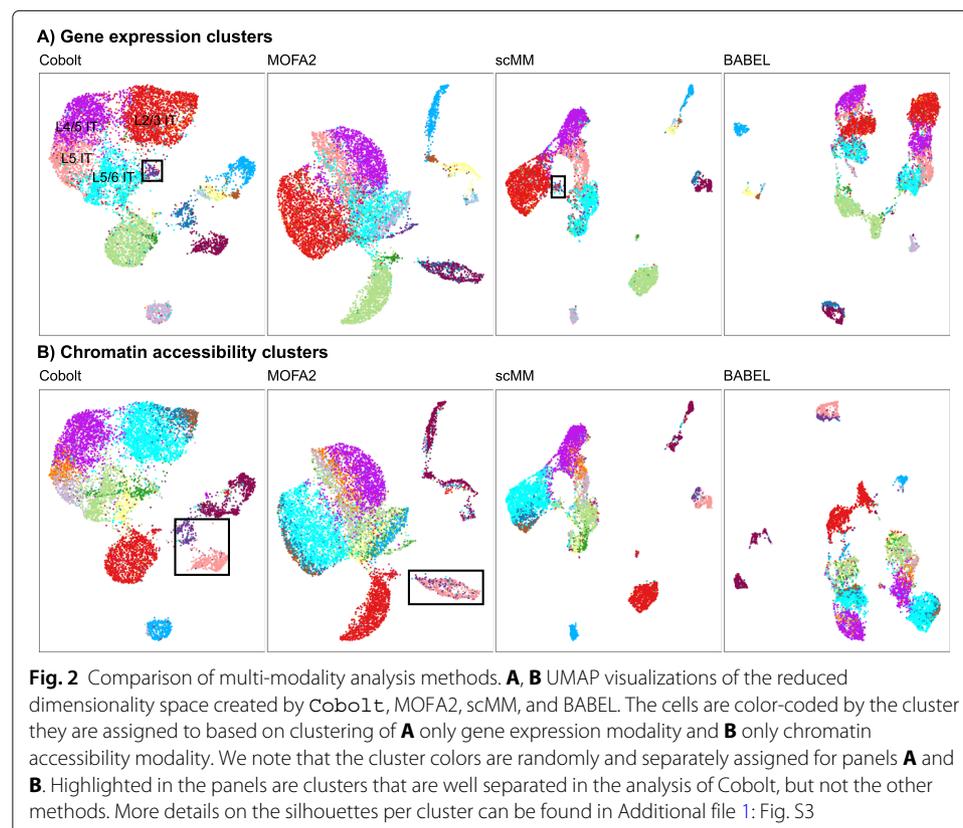
While the full power of **Cobolt** is to integrate together data from single and multiple modality datasets, in its simplest form, **Cobolt** can be used for the analysis of data from solely a multi-modality technology. We demonstrate this usage with the SNARE-seq data [8], which consists of paired transcription and chromatin accessibility sequenced on 10,309 cells of adult mouse cerebral cortices.

We first compare **Cobolt** with a simple, but common, approach for analyzing joint-modality data: the two modalities are analyzed separately and then the results are linked together. This is the strategy of [8], where the authors primarily clustered the gene expression modality to form clusters of the cells, and then performed a separate analysis on chromatin accessibility modality as a comparison. Focusing on the gene expression modality is common, since it is often assumed to have the greatest resolution in determining cell types. However, the reduced representations and clusters created on one modality may not be representative of all the underlying cell subtypes. Indeed, when we perform clustering analysis on only the gene expression modality using Seurat [3] and only the chromatin accessibility modality using *cisTopic* [21] (consistent with [8], see the “**Methods**” section), both modalities find distinct clusters that are not reflected in the other modality (Additional file 1: Fig. S1). For example, cells identified by marker genes as non-neuronal cells, such as astrocytes, oligodendrocytes, oligodendrocyte precursors, and microglial cells, are clustered into their respective cell types based on mRNA expression but are not separated based only on chromatin accessibility (Additional file 1: Fig. S2A). Similarly, a subset of layer 5/6 cells has distinct chromatin accessibility peaks but are intermingled with other layer 5/6 cells in the gene expression clusters; these differential peaks include one near the gene *Car12* which is a marker gene of the previously annotated subtype of layer 6 (L6 *Car12*, [24]) and which shows higher expression in this subset of cells (Additional file 1: Fig. S2B,C). A joint analysis with **Cobolt**, unlike the single-modality analyses, finds these subtypes detected by only one modality and not the other (Additional file 1: Fig. S1).

Next, we compare with other methods that explicitly analyze the two modalities jointly, like **Cobolt**. We consider the methods MOFA2, scMM, and BABEL. MOFA2 uses

Bayesian group factor analysis for dimensionality reduction of multi-modality datasets; BABEL trains an interoperable neural network model on the paired data that translates data from one modality to the other; and scMM [16] uses a deep generative model for joint representation learning and cross-modal generation. We apply each of these methods to the SNARE-seq data.

Since a joint analysis method should be able to reflect subtype signals captured by all modalities, we similarly evaluate the methods on how well their lower-dimensional spaces represent separate clusters identified separately in each modality, as described above. In Fig. 2, we visualize the lower-dimensional space generated by Cobolt, MOFA2, scMM, and BABEL via UMAP (Uniform Manifold Approximation and Projection [25]), where we color the cells based on the clusters found by clustering the gene expression modality data (Fig. 2A) and those based on clustering the chromatin accessibility data (Fig. 2B). We would note that BABEL does not create a single reduced-dimensionality representation for a paired cell, but rather one per modality (the two latent representations are learned jointly and are quite similar). BABEL's lower-dimensionality representation does quite poorly, separating major clusters of cells found in both modalities, such as layer 2 to 6 intratelencephalic (IT) neurons (colored red, purple, pink, and cyan in Fig. 2A). Both MOFA2 and scMM capture these large clusters, which are shared between the modalities. However, we see clusters specific to a single modality not reflected on their lower-dimensional space. For example, the highlighted gene expression cluster in Fig. 2A is practically indistinguishable in the scMM UMAP but separated in the Cobolt UMAP.



Differential mRNA expression analysis between this cluster and neighboring cells finds strong expression of known markers of layer 6 cells in this cluster (*Col24a1* [26], *Gnb4* [27], *Rxfp1* [28], *Nr4a2* [29, 30], and *Ntng2* [29], Additional file 1: Fig. S4A) as well as strong expression of *Car3* defined in [31] as a marker of a subset of layer 6 IT cells. Neighboring cells do not express these known marker genes and instead express layer 5/6 IT markers (cyan cluster) or layer 2/3 IT markers (red cluster), indicating that this cluster missed by scMM consists of a biologically meaningful subset of layer 6 IT cells. Similarly, in Fig. 2B, we highlight two clusters of cells which clearly separate in the `Cobolt` analysis and are separate clusters based on both mRNA-Seq and chromatin profiles, but are mixed together in the MOFA2 analysis. Differential mRNA expression analysis between these clusters reveal genes *Adarb2* and *Sox6* differ in expression between these groups (Additional file 1: Fig. S4B), which are known markers whose expression distinguish the CGE and Pvalb clusters, respectively [31]. Integrative analysis in the next section confirms this identification by integrating this SNARE-Seq data with annotated scRNA-Seq data and placing these cells with cells annotated as CGE and Pvalb in [31].

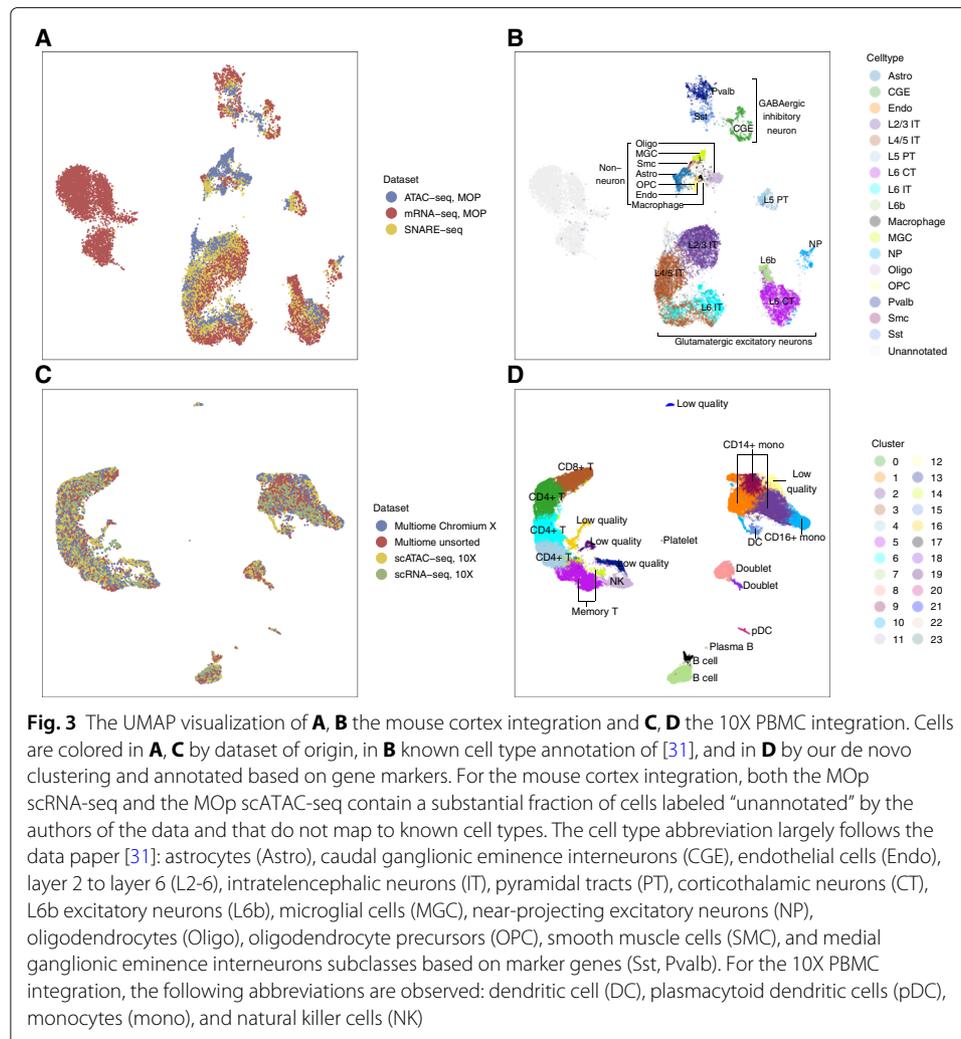
To quantitatively evaluate these observations, we calculate the average silhouette widths of the modality-specific clustering on the UMAPs generated by `Cobolt`, MOFA2, scMM, and BABEL (shown in Fig. 2C), where higher silhouette widths indicate that cells are closer to other cells in the same cluster. As expected from our observations, BABEL's representation results in extremely small silhouette widths, reflecting the many clusters separated in BABEL's representation. MOFA2 has the smallest silhouette width on chromatin accessibility clusters, supporting our observation that its joint space does not represent this modality well; similarly, scMM gives relatively small measures on the gene expression modality. `Cobolt` best represents both modalities with the highest silhouette width measure.

#### `Cobolt` for integrating multi-modality data with single-modality data

We now turn to integrating multi-modality data with single-modality data. For this use-case, we use `Cobolt` to jointly model three different datasets—the SNARE-seq of mouse cerebral cortices analyzed in the above section, together with a scRNA-seq and a scATAC-seq dataset of mouse primary motor cortex (MOp) [31]. In addition, we also demonstrate `Cobolt` for joint modeling of single-cell sequencing of human peripheral blood mononuclear cells (PBMCs): two multi-modality datasets pairing ATAC and mRNA measurements on 10,970 and 12,012 cells from different samples of the 10X Multiome platform [32, 33], combined with 23,837 cells from scRNA-Seq [34] and 9030 cells from scATAC-Seq [35].

The result in both examples is a lower-dimensional latent space that aligns the different modality data into a single representation. In Fig. 3A and B, we visualize this low-dimensional space via UMAP for the mouse cortex and human PBMCs, respectively, with cells colored by their data set of origin. We see that the cells from different datasets are well aligned regardless of their source of origin.

To consider further the biological meaning of the lower-dimensional representation, we label the MOp cells from the mouse cortex dataset in Fig. 3C by their cellular subtype as annotated in [31]. For the purposes of comparison across the modalities, we integrated some cell types in [31] into larger groupings and modified the names so as to have comparable groups (see the “Methods” section). For the SNARE-seq cells, we do not have the cell



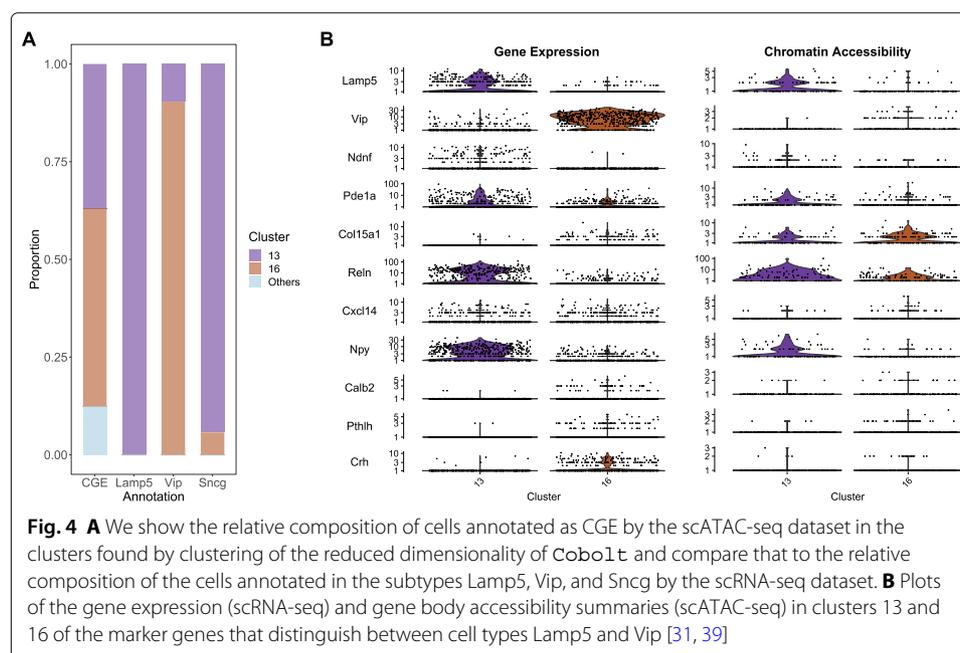
types given in [8], so we use the identifications found by our analysis of only the SNARE-seq cells (see the “Methods” section). We see that cells from the same cellular subtypes are projected closely regardless of the data source. We also see that the representation of `Cobolt` respects the larger category of cell types by grouping three major cell classes: GABAergic inhibitory neurons (CGE, Sst, Pvalb), glutamatergic excitatory neurons (IT, L5 PT, L6 CT, L6b, NP), and non-neurons.

The PBMC datasets do not have accompanying annotation, so we applied the Louvain clustering algorithm to the lower-dimensional representation from `Cobolt` to identify potential cell types. Using marker genes [3, 36–38], we classified the clusters into known subtypes expected for PBMC data (Additional file 1: Figs. S5 and S6) and we see that important cell types and functions are localized in the `Cobolt` representation (Fig. 3C).

The mouse cortex data also demonstrates the ability of our joint representation to capture subtype signals that are not shared across all of the modalities. Indeed, despite detecting mostly similar cell types, the MOp datasets profile several cell types distinct to the modality. For example, microglial cells (MGC) and smooth muscle cells (SMC) are uniquely detected in scATAC-seq. The different datasets also have different cellular compositions of their shared subtypes, where astrocytes (Astro) and oligodendrocytes (Oligo)

are much more abundant in the scATAC-seq (6.55% and 10.50%) than in the SNARE-seq (4.5% and 2.84%) and scRNA-seq (0.40% and 0.36%). As shown in Fig. 3B, cell populations unique to one dataset are grouped in the UMAP plot and are distinguishable from the other datasets/cell types. This indicates that *Cobolt* reconciles data even when one cell population is entirely absent or scarcely represented in one or more data sources, which is important in integrating datasets collected from related but slightly differing settings.

*Cobolt* also facilitates subtype identification at a finer resolution by transferring information between modalities. For example, for the mouse cortex data, Fig. 3B shows the cell types based on the published annotation. To make the annotation consistent between scRNA-Seq and scATAC-Seq, some of these cell types are the result of merging some cell types into larger categories (see the “Methods” section). One cell type, caudal ganglionic eminence interneurons (CGE, dark green in Fig. 3B), was annotated in the scATAC-seq MOP dataset as one cluster, but the scRNA-seq annotation further divided CGE cells into 3 subtypes based on marker genes—*Lamp5*, *Vip*, and *Sncg*. Our joint mapping of the cells allows us to relate the subtypes detected in scRNA-seq to scATAC-seq and provides a finer resolution breakdown of CGE in scATAC-seq. Specifically, we ran a de novo clustering of our joint mapping of all three datasets by *Cobolt* (see the “Methods” section and Additional file 1: Fig. S7). This clustering results in a cluster (cluster 13) composed of *Lamp5* and *Sncg* cells, while another (cluster 16) is mostly *Vip* cells (Fig. 4A). This subdivision is further validated by gene expression and gene activity levels in these clusters of marker genes *Lamp5* and *Vip* as well as other genes known to discriminate subtype *Lamp5/Sncg* from subtype *Vip* [39], such as *Reln* and *Npy* (Fig. 4B). We further validated the scATAC-seq clusters through de novo differential accessibility (DA) and differential expression (DE) analysis (Additional file 1: Fig. S8). We identified 94 DA genes between these two clusters. Seventy-eight of the DA genes are also found DE in the mRNA, and all of them have the same direction of fold changes in the DE and the DA analyses, i.e., the genes with lower/higher gene expression in cluster 13 compared to cluster 16 are also



**Fig. 4** **A** We show the relative composition of cells annotated as CGE by the scATAC-seq dataset in the clusters found by clustering of the reduced dimensionality of *Cobolt* and compare that to the relative composition of the cells annotated in the subtypes *Lamp5*, *Vip*, and *Sncg* by the scRNA-seq dataset. **B** Plots of the gene expression (scRNA-seq) and gene body accessibility summaries (scATAC-seq) in clusters 13 and 16 of the marker genes that distinguish between cell types *Lamp5* and *Vip* [31, 39]

less/more accessible. This shows that the joint model of `Cobolt` can help distinguish noisy cells in one dataset with additional information from other datasets or modalities.

Furthermore, `Cobolt` is robust to poor-quality cells and low-expressed genes by using a count-based model, which should naturally down-weight the influence of low-count cells. In the above analysis of the mouse cortex data, we filtered out 2.5% of MOp mRNA cells and 15.8% ATAC-seq cells due to low counts and other quality control measures following the work of [31]. Yet `Cobolt` is robust to these choices, even in the extreme case where there is no filtering on cells and only a minimal filter on genes is performed (Additional file 1: Fig. S9).

### **Comparison with existing methods**

As described in the introduction, there are few existing methods that allow full integration of multi-modality sequencing data with single-modality data. MOFA2 only analyzes joint-modality data; BABEL and scMM train a joint model on the paired joint-modality data and allow the user to apply this model to single-modality datasets to predict the other “missing” modality. Unlike `Cobolt`, these two methods do not use the single-modality data in the training of the model, nor do they provide a representation of the single-modality data in a single shared representation space regardless of modality—for example, for the purposes of joint clustering of all of the cells across modalities.

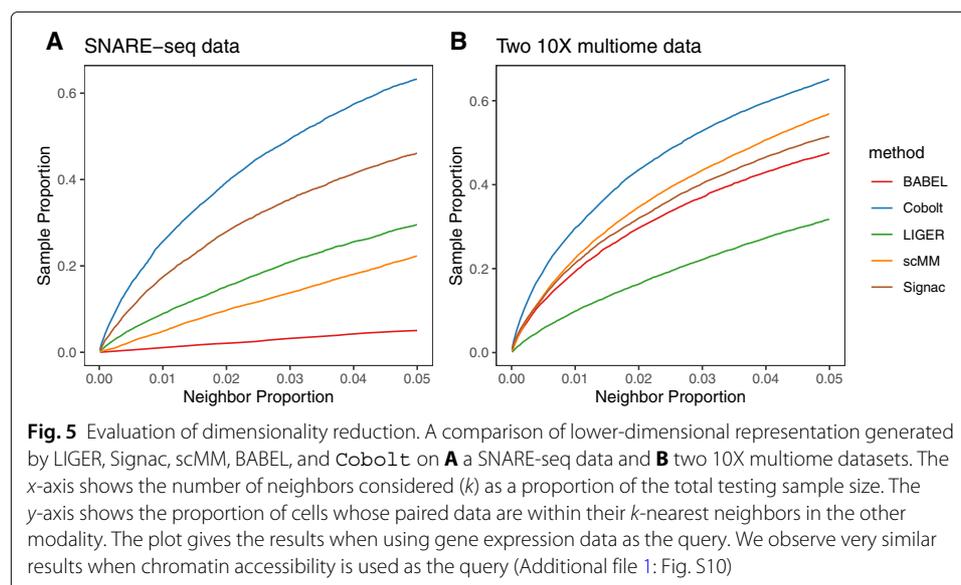
Therefore, to have additional points of comparison, we apply the LIGER and Signac methods, which are designed for integrating *unpaired* modalities (the implementation details can be found in the “Methods” section). LIGER applies an integrative nonnegative matrix factorization (iNMF) approach to project the data onto a lower-dimensional space and then builds a shared nearest neighbor graph for joint clustering. Signac implements canonical correlation analysis (CCA) for dimensionality reduction; Signac subsequently transfers cell labels by identifying mutual nearest neighbor cell pairs across modalities.

To evaluate the performance of these methods when integrating single-modality data with multiple-modality data, we return to the multi-modality datasets described above to create artificial sets of multiple-modality data and single-modality data. For the 10X Multiome data, we make use of the fact that we have two datasets from the 10X Multiome platform run on different patient samples: PBMC of a healthy male donor aged 30–35 (“Multiome Chromium X”) [32] and PBMC of a female donor aged 25 (“Multiome unsorted”) [33]. We ignore the pairing information in the Multiome unsorted data and treat the mRNA and ATAC measurements as coming from unpaired, separate sequencing experiments. For the SNARE-Seq data, we randomly assign the cells to be considered as from either the multi-modality dataset (20%) or the single-modality datasets (80%) and run each of the methods. The choice of 20% and 80% was based on the relative size of the SNARE-Seq joint-modality data to the individual MOp scATAC-Seq data and scRNA-Seq datasets and reflects the fact that single-modality datasets are much higher-throughput than paired-modality data.

We give to each of the methods the multi-modality data from the cells that remain paired and for the cells where we ignore the pairing information give the mRNA and ATAC data from those cells as if they were single-modality datasets. For LIGER and Signac, which are not designed for multiple-modality data, we hide the pairing information on all cells and treat all of the cells as if they were collected on different cells. In this way, we have a ground truth on how the cells in the single-modality datasets should

be connected to each other, and we can compare the methods by evaluating whether coordinates of the reduced dimensions  $Z$  for the pairs of cells assigned to the single-modality datasets were close together. Specifically, for each method, we evaluate for each cell in the mRNA single-modality set its coordinates  $\hat{Z}^{mRNA}$  and for a fixed number  $k$  locate its  $k$  nearest neighbors in the ATAC single-modality set based on the coordinates  $\hat{Z}^{ATAC}$ . We then calculate the percentage of mRNA single-modality cells whose paired cell in the ATAC single-modality is included in its set of nearest neighbors. The reverse analysis was done using chromatin accessibility as the query and evaluating the percentage whose nearest neighbors include their mRNA pair. Many popular clustering routines use nearest-neighbor graphs for identifying clusters, so this is a metric directly related to whether the cells assigned to the single-modality data would likely correctly cluster together across modalities, but avoids having to specify cluster parameters, especially as applied to different methods (and LIGER and Signac have their own clustering techniques specific for their methods).

As shown in Fig. 5, the `Cobolt` joint representation does a much better job for both the SNARE-Seq and 10X multiome data of assigning coordinates to the single-modality cells that place them close to their matching pair. The proportion of single-modality cells that are neighbors to their pair is much larger than any of these other methods in both datasets. Surprisingly, for the SNARE-seq data, the other methods that make use of the joint-modality data to develop their model (scMM and BABEL) do much worse than Signac and LIGER which do not have any information linking the cells together. The 10X multiome data, which has similar numbers of cells from the single-modality datasets as the multi-modality dataset, shows scMM and BABEL perform comparably to Signac, though not as well as `Cobolt`; similarly, we see improved performance of scMM and BABEL in the SNARE-Seq data when we increase the proportion of dual-modality cells to 80% and only 20% of cells being single-modality (Additional file 1: Fig. S12), but still well below the performance of `Cobolt`. This points to the power of truly integrating the high-throughput single-modality data into the analysis, particularly when there are more cells sequenced from the single-modality data, as is frequently the case.



Subtype detection is a critical aspect of single-cell data analysis, and integrating single-modality data with multiple-modality data gives the potential for higher resolution detection. Our nearest-neighbor metric is directly related to subtype detection, with a higher neighbor proportion corresponding roughly to finding larger clusters in a clustering algorithm. scMM and BABEL only provide cross-modality predictions, rather than a joint embedding of the single-modality data with the multiple-modality data. Downstream tasks such as clustering for subtype detection must be done on either the mRNA expression space or the chromatin expression space, and as we have seen in the SNARE-seq only analysis each of which can miss important features of the data. Indeed, this may be an important factor in their low performance on our nearest-neighbor analysis. LIGER and Signac do provide a joint embedding of all of the data, but do so without the use of the pairing information for the joint-modality data for training the embedding.

The previous analysis provides a comparison with a known ground-truth and the evaluation at different levels of analysis. Now, we return to the joint embedding of MOp scRNA-Seq and scATAC-Seq with the SNARE-Seq data that we considered in the previous section and consider the performance of LIGER and Signac, which performed comparatively well. We would note that LIGER and Signac have their own clustering strategies, separate from their dimensionality reduction, but we focus here on the results of their dimensionality reduction. We provided LIGER and Signac only the MOp data, as there is no clear way of including the SNARE-seq into LIGER and Signac without focusing on one of its modalities and adding extra batch correction steps. We compare the results to the clusters published with the scRNA-seq and scATAC-seq data in [31] (Additional file 1: Fig. S13). We see that *Cobolt* generates a UMAP visualization that well represents rare subpopulations and respects broader cell classes. LIGER gives greater separation between cell types but splits several subtypes into far-away islands, such as for L5 PT, MGC, and Astro. Signac adopts an asymmetric strategy of transferring scRNA-seq labels to scATAC-seq data, and as a result, Signac performs well on major cell types but poorly on under-represented subpopulations in scRNA-seq such as astrocytes (Astro), which accounts for only 0.4% of the cells in scRNA-seq but 6.55% in scATAC-seq. Furthermore, *Cobolt*, unlike LIGER and Signac, not only groups together the subtypes, but appears to also represent the three broader categories major GABAergic inhibitory neurons, glutamatergic excitatory neurons, and non-neurons.

These cluster identifications that we highlight from [31] are relatively robust, well-known cell types, representing the large structural changes in the data, for which we expect most strategies to be able to detect reasonably well. On the other hand, our nearest-neighbor analysis emphasizes the performance at a high level of resolution. Putting both of these together points to the fact that *Cobolt* provides a superior integration of the datasets across a wide spectrum of resolutions.

## Discussion

In this paper, we have shown that *Cobolt* successfully integrates multi-modal data and provides a representation that can be used for downstream analysis tasks, such as cell-type discovery. Pseudo-time estimation for reconstructing developmental order of cells [40], while not meaningful for the datasets we considered, is another important downstream application where the integrated representation of *Cobolt* allows the analysis of cells from different modalities. Future work could make use of the graphical model and

inferred parameters to establish connections between features. For example, the probability vectors generated by  $B^{(i)}$  naturally provide a reduced-dimensional space of molecular features and can potentially help in the construction of gene networks.

We would note that while we have focused on the capability of `Cobolt` to analyze data from two modalities, the underlying method can be extended to larger number of modalities and integration of different combinations of modalities, such as datasets with different pairs of modalities (see the “[Methods](#)” section). Thus, `Cobolt` provides a framework to integrate a wide range of varieties of multi-modality platforms as well as single-modality platforms.

`Cobolt` is available as a Python package at [https://github.com/epurdom/cobolt\\_manuscript](https://github.com/epurdom/cobolt_manuscript). All of the code used for the analysis is available as a github repository: [https://github.com/epurdom/cobolt\\_manuscript](https://github.com/epurdom/cobolt_manuscript).

## Conclusions

We have shown that `Cobolt` is a flexible tool for analyzing multi-modality sequencing data, whether separately or integrated with single-modality data. `Cobolt` synthesizes the varied data into a single representation, preserving meaningful biological signal in the different modalities and at different resolutions. Moreover, this latent variable space is appropriate for standard downstream analysis techniques commonly used for analyzing cells without any further specialized adjustments, allowing `Cobolt` to fit into standard analysis pipelines.

## Methods

While the most common application of joint-modality platforms consists of pairs of modalities (such as the example of mRNA-seq and ATAC-seq we described above), we will describe `Cobolt` in generality, assuming that there are  $M$  modalities.

### Modeling modality dependency

For an individual cell  $c$ , we can (potentially) observe  $M$  vectors of data  $x_c = \{x_c^{(1)}, \dots, x_c^{(M)}\}$ , each vector of dimension  $d_1, \dots, d_M$  corresponding to the number of features of each modality. We assume a Bayesian latent model, such that for each cell there is a latent variable  $z_c \in \mathbf{R}^K$  representing the biological signal of the cell, where  $z_c$  is assumed drawn from a Gaussian prior distribution. Given  $z_c$ , we assume that the data  $x_c^{(m)}$  for each modality has an independent generative process, potentially different for each modality. Specifically, we assume that the data  $x_c^{(m)}$  from each modality are conditionally independent given the common latent variable  $z_c$ . That is,

$$p(x_c^{(1)}, \dots, x_c^{(M)}, z_c) = p(z_c) \prod_{i=1}^M p(x_c^{(i)} | z_c).$$

We use  $q(z_c | x_c^{(1)}, \dots, x_c^{(M)})$  as a variational approximation of the posterior distribution  $p(z_c | x_c^{(1)}, \dots, x_c^{(M)})$ .  $q(z_c | x_c^{(1)}, \dots, x_c^{(M)})$ , the encoder, is assumed Gaussian with parameters modeled as neural networks (i.e., Variational Autoencoder, VAE [41]). This allows for estimation of the posterior distribution  $p(z_c | x_c^{(1)}, \dots, x_c^{(M)})$  and the underlying latent variable for each cell  $c$ . The posterior mean of this distribution  $\hat{z}_c$  will be our summary of the shared representation across modalities.

Importantly, this model can be estimated even when not all of the input data contains all modalities. In this case, an individual cell  $c$  contains a subset of the modalities,  $\mathcal{S}_c \subset \{1, \dots, M\}$ , and consists of data  $x_c = \{x_c^{(i)}, i \in \mathcal{S}_c\}$ . Without all modalities observed, the cell can contribute to the estimation of the model as its distribution can be explicitly written out:

$$p(x_c, z_c) = p(z_c) \prod_{i \in \mathcal{S}_c} p(x_c^{(i)} | z_c).$$

Furthermore, we can estimate latent variables for such cells by using posterior distribution of  $z_c$  when conditioning only on the observed modalities,  $q(z_c | \{x_c^{(i)}, i \in \mathcal{S}_c\})$ . Instead of using separate neural networks for  $2^M - 1$  posterior distributions of different modality combinations, we adopt a technique introduced in Multimodal Variational Autoencoder (MVAE) [17], which largely reduces the number of encoders to  $2M$  (See Additional file 2: Supplementary Methods for inference details).

As an example, if there are two modalities, mRNA-seq and ATAC-seq, and we have  $n_1$  cells with paired data from the joint modality platform,  $X_1 = (X_1^{mRNA}, X_1^{ATAC})$ ;  $n_2$  cells with only mRNA measured,  $X_2 = X_2^{mRNA}$ ; and  $n_3$  cells with only ATAC-seq measured,  $X_3 = X_3^{ATAC}$ . All  $N = n_1 + n_2 + n_3$  cells can be used in the estimation of the joint distribution of the latent variables, and estimates of the latent variables can be found as the mean of the relevant approximate posterior distributions:

$$\begin{aligned} \hat{Z}_1 &= E(Z | X_1^{mRNA}, X_1^{ATAC}) && \text{(Paired cells)} \\ \hat{Z}_2^{mRNA} &= E(Z | X_2^{mRNA}) && \text{(mRNA-seq only)} \\ \hat{Z}_3^{ATAC} &= E(Z | X_3^{ATAC}) && \text{(ATAC-seq only)} \end{aligned}$$

### Correcting for missing modalities

In practice, we find that the distributions  $q_\phi(z_c | \{x_c^{(i)}, i \in \mathcal{S}\})$  have subtle differences for different subsets  $\mathcal{S}$ , i.e., the latent variables  $\hat{Z}_1$ ,  $\hat{Z}_2^{mRNA}$ , and  $\hat{Z}_3^{ATAC}$  show distinct separations (Additional file 1: Fig. S15). One possibility could be due to platform differences between the different datasets that remain even after our batch correction. However, we also see differences in these distributions even if we only consider the joint-modality data, where we can estimate all of these posterior distributions on the same cells,  $\hat{Z}_1^{mRNA} = E(Z | X_1^{mRNA})$  or  $\hat{Z}_1^{ATAC} = E(Z | X_1^{ATAC})$  (Additional file 1: Fig. S15). Indeed, there is nothing in the optimization of the posterior distribution that requires these different posterior distributions to be the same.

While the effects are small, these subtle differences can affect downstream analyses, e.g., in clustering cells for subtype discovery. Rather than directly force these posterior distributions to match in our estimation of the model, `Cobolt` fits the model as described above (using all of the data) and then uses the paired data to train a prediction models that predict  $\hat{Z}_1$  from the modality-specific estimates  $\hat{Z}_1^{mRNA}$  and  $\hat{Z}_1^{ATAC}$ . We then apply these prediction models to  $\hat{Z}_2^{mRNA}$  and  $\hat{Z}_3^{ATAC}$  to obtain estimates  $\hat{Z}_2$  and  $\hat{Z}_3$  which are better aligned to be jointly analyzed in the same space. In practice, we find XGBoost [42] and k-nearest neighbors algorithm work equally well. We present results based on XGBoost. We would note that there is little difference in performance when we predict coordinates

into the ATAC-Seq space  $E(Z|X^{ATAC})$  or mRNA-Seq space  $E(Z|X^{mRNA})$ , rather than the joint space  $(E(Z|X_1^{mRNA}, X_1^{ATAC}))$ , see Additional file 1: Fig. S16.

### Modeling single modality of sparse counts

The choice of the generative distribution  $p_\psi(x^{(i)}|z)$  should be chosen to reflect the data and in principle can vary from modality to modality. For example, single-modality VAE models using zero-inflated negative binomial distributions (ZINB) [43] have been proposed for scRNA-seq datasets to account for sparse count data. However, we found ZINB models performed less well for technologies that measure modalities like chromatin accessibility, which results in data with sparser counts and larger feature sizes than scRNA-seq. Therefore, we develop a latent model for these types of modalities inspired by the Latent Dirichlet Allocation (LDA) [18].

Our generative model for a single modality  $i$  starts by assuming that the counts measured on a cell are the mixture of the counts from different latent categories. In the genomic setting, these categories could correspond to biological processes of the cell, for example. Each category has a corresponding distribution of feature counts. The cumulative feature counts for a cell  $c$  are then the result of combining the counts across its categories, i.e., a mixture of the categories' distributions. Specifically, each cell  $c$  has a latent probability vector  $\theta_c \in [0, 1]^K$  describing the proportion of each category that makes up cell  $c$ . Each category  $k$  has a probability vector  $\sigma(\beta_k)$  that provides the distribution of its feature counts. Here,  $\sigma$  indicates the softmax function that transforms  $\beta_k$  to a probability vector that sums to 1. The observed vector of counts  $x_c$  is a multinomial draw with probabilities  $\pi_c$ , where  $\pi_c = \sigma(B)\theta_c$ , and  $B = (\beta_1, \dots, \beta_K)$  is a matrix of the individual  $\beta_k$  vectors. To extend this model to multiple modalities, we assume a shared latent variable  $z_c$  that is common across modalities, but each modality has a different  $B^{(i)}$  that transforms the shared latent class probabilities into the feature space of the modality.

Furthermore, it is well known that there can be meaningful technical artifacts ("batch effects") between different datasets on the same modality, for example due to differences between platforms or laboratory preparations. To counter this, our model also adjusts the sampling probabilities  $\sigma(B^{(i)})\theta_c$  differently for data from different batches within the same modality  $i$ . We would note that the model can also take batch-corrected counts as input, such as are available for mRNA expression data (e.g., [44–46]), but we anticipate that for some modality types stand-alone batch correction techniques may not be as well developed. We evaluate the effect of our batch correction on the 10x Multiome data, which consists of two runs of 10x Multiome on the different patient input sampled at different times. This creates a batch effect between the multi-modality input and the single-modality input where we know the ground truth of how the single-modality data should be linked. We use the same nearest-neighbor analysis as in Fig. 5B, with and without the batch correction terms and see much improved performance using the batch correction (Additional file 1: Fig. S17). This type of quantitative nearest-neighbor analysis is not possible for the SNARE-Seq data, since we do not have two different batches of paired multi-modality data, but we visually see large improvement due to the batch correction when analyzing the single-modality datasets jointly with the multi-modality data (Additional file 1: Fig. S18).

The parameter  $\theta_c$  is the latent variable describing the contributions of each category to cell  $c$  and is shared across all modalities. In LDA models, it is typically assumed to have

a Dirichlet prior distribution. However, we use a Laplace approximation to the Dirichlet introduced in ProLDFA [23], which allows for incorporation into a VAE model. This prior assumes a latent variable  $z_c$  with a Gaussian prior and sets  $\theta_c = \sigma(z_c)$ , where  $\sigma$  is the softmax transformation. We use this approximation to the Dirichlet distribution to provide a multi-modality method appropriate for sparse sequence count data.

## Data processing and method implementation

### *SNARE-seq processing and annotation*

We downloaded the processed counts of the adult mouse cerebral cortex data (10,309 cells) [8]. We applied quality filtering that retained cells having a number of genes detected greater than 20. For genes, we used the ones detected in more than 5 cells and have a total number of counts greater than 10. For peaks, we removed the ones having nonzero counts in more than 10% of cells or less than 5 cells. We performed clustering analysis using Seurat (version 3.2.2) on the gene expression modality. The data were normalized using SCTransform function with default parameters, followed by principal component analysis (PCA) using the default 3000 variable features. Louvain algorithm was applied on the first 50 PCs with the resolution parameter equals 0.65. Cell type annotations are generated on the resulting 15 clusters using the marker genes [31]. We applied cisTopic (version 0.3.0) on the chromatin accessibility data with default parameters. Model selection was conducted based on log-likelihood using runWrapLDAModels and selectModel functions, and 30 topics are used in the results.

For the integration of SNARE-seq with the MOp data using Cobolt, we map the SNARE-seq counts to the peak set called on the MOp scATAC-seq data. Since peaks are typically called in a dataset-specific manner, the ideal integration strategy would be to redo the peak-calling with all datasets combined. However, in Additional file 1: Fig. S10, we show that our simple alternative of mapping data to peaks called on a different dataset from the same system does not result in significant performance loss for Cobolt.

### *MOp data preprocessing*

We downloaded the single-nucleus 10x v3 transcriptome dataset (90,266 cells) and the open chromatin dataset (15,731 cells, sample 171206\_3C) [31]. For mRNA-seq quality control, we filtered cells that have less than 200 genes detected or have greater than 5% mitochondrial counts. For ATAC-seq, we utilized the TSSEnrichment and blacklist region reads calculation functionalities in Signac. We subsetted cells with the blacklist ratio less than 0.1, the number of unique molecular identifiers (UMIs) greater than 50, and the TSS enrichment score greater than 2 and less than 20. A total of 88,021 and 13,249 were retained for mRNA-seq and ATAC-seq, separately. To make annotation in the two datasets consistent, we merged the layer 2/3 IT and layer 4/5 IT subclusters in ATAC-seq data. For mRNA-seq, we merged Lamp5, Vip, and Sncg into one CGE cluster. When integrating the MOp datasets with the SNARE-seq data, we used only genes detected in both the scRNA-seq and the SNARE-seq datasets.

### *10X PBMC data preprocessing*

PBMC datasets were downloaded from the 10X website ([Multiome Chromium X](#), [Multiome unsorted](#), [scRNA-seq](#), [scATAC-seq](#)). For chromatin accessibility, we mapped the scATAC-seq reads and the Multiome unsorted reads to the peaks called on the Multiome

Chromium X data. No quality filtering was applied to any of the four 10X datasets. Clusters with high mitochondrial expression are identified and annotated as low quality clusters, and clusters with the majority of cells identified by DoubletFinder [47] are annotated as doublet clusters in the downstream clustering analysis.

### **Gene activity calculation**

Gene activity matrix for chromatin accessibility is generated by counting the number of reads overlapping genes and their promoters using BEDOPS [48], where a promoter is defined as the region starting from the transcription start site (TSS) to 3000 base pairs upstream of TSS.

### **Cobolt network architecture and training**

For each modality  $i$ , the encoder takes as input the log-plus-one transformed counts. We use one fully connected layer of size 128, followed by fully connected layers for mean  $\tilde{\mu}^{(i)}$  and log-variance  $\log \tilde{\Sigma}^{(i)}$ . We tried networks with one or two hidden layers of varying sizes and found the results pretty stable. The decoders follow our probability model for sparse counts (see also Additional file 2: Supplementary methods, for details) and do not contain neural networks. We set the parameter of the Dirichlet prior to 50 divided by the number of latent variables  $K$ . The actual parameters used for the Gaussian prior are calculated using the Laplace approximation (see also Additional file 2, Supplementary methods, for details). For the ELBO objective, we set the weighting terms  $\lambda^A$  reciprocal to the number of samples available for modality combination  $\mathcal{A}$ . We set the hyperparameter weights  $\eta$  for conditional likelihood terms to 1. Adam optimizer is used, and we select a learning rate of 0.005 after tuning. We adopt a KL cost annealing schedule that linearly increases the weight of the KL term  $\gamma$  from 0 to 1 in the first 30 epochs. During training, we use a batch size of 128 and a fixed number of 100 epochs.

We note that the softmax transformation from  $z_c$  to  $\theta_c$  is not a one-to-one transformation. Therefore, we scale  $\hat{z}_c$  to mean 0 before the downstream correction, followed by clustering and visualization.

In correcting for missing modalities, we predict using XGBoost, setting the objective function to regression with squared loss, the learning rate to 0.8, and the maximum depth of a tree to 3. XGBoost is applied separately to each modality (scRNA-Seq and scATAC-Seq), resulting in our corrected estimates  $\hat{Z}_2, \hat{Z}_3$ .

The number of latent variables  $K$  was set to 10 in the SNARE-seq analysis and the method comparison analysis so as to be consistent with the default of several other methods for comparison purposes. We used  $K = 30$  for the mouse cortex integration and  $K = 10$  for the 10X PBMC integration. Estimations of the data's marginal likelihood were used to assist the selection of  $K$ .

For the mouse cortex data integration, we focused on genes and peaks that have top 30% average expression and removed the ones in the top 1%. For the SNARE-seq analysis and the 10X PBMC integration, all features were included in training. The choice of top features is less important here, which we found to have a small effect on the results.

### **Clustering and visualization of Cobolt results**

Clusters were generated on the corrected latent variables using Louvain algorithm [49]. We used the implementation of naive Louvain algorithm in FindClusters function

from the R package Seurat. All parameters, other than the resolution controlling the number of clusters identified, were set to default. The results of the clustering of the integrative analysis of the SNARE-Seq with the MOP scRNA-Seq and ATAC-Seq from the mouse cortex are given in Additional file 3: Table S1 (resolution = 0.8). The results of the clustering of the integrative analysis of the 10X PBMC data are given in Additional file 4: Table S2 (resolution = 0.8). UMAPs were generated using the `umap` function from the R package `uwot` with the number of neighbors set to 30.

#### ***MOFA2, scMM, and BABEL analysis on the SNARE-seq data***

Following the vignette of MOFA2, we used the top 2500 variable genes and `cisTopic` embeddings of the chromatin accessibility modality as input. Variable genes are selected using the `FindVariableFeatures` function from Seurat using selection method “vst”. The number of factors was set to 10. Two factors were identified as technical factors after inspecting their correlation with the total number of reads counts per cell. The UMAP representation was then generated using the rest of the factors with the number of neighbors set to 30. scMM was run with batch size equals 32, number of epochs equals 50, learning rate equals  $10^{-4}$ , number of latent dimensions equals 10, number of hidden dimension for gene expression equals 100, and number of hidden dimensions for chromatin accessibility equals 500. The parameters were chosen following the scMM paper. BABEL was run with the number of latent dimensions set to 10 and the batch size set to 256. Other parameters were chosen to follow the built-in SNARE-seq defaults.

#### ***Differential analysis***

DE analysis on gene expression and DA analysis on gene activities were performed using the Wilcoxon rank sum test followed by Bonferroni correction, implemented by the `FindMarkers` function in Seurat. Gene expression and gene activities were visualized by heatmaps using `DoHeatmap` in Seurat. DA analysis on the peaks for the SNARE-seq analysis was performed using Fisher’s exact test followed by the Benjamini–Hochberg procedure (following [8]). Peak-by-cluster matrix was normalized by size factors calculated by Monocle [50] and visualized by heatmaps (following [8]). Genes and peaks with adjusted *p*-values lower than 0.05 were called significant.

#### ***LIGER and Signac analysis on the MOp data***

LIGER and Signac (Seurat) take as input gene-level count summaries from different modalities, such as gene expression or gene body methylation/chromatin accessibility measures. The input is different from `Cobolt`, which uses the peak summaries directly without needing to summarize at the gene level. Therefore, we applied these two methods on the gene expression and gene activity matrices, where the latter is defined as the summarized chromatin accessibility counts over gene and promoter regions.

We ran LIGER (version 0.5.0) using default parameters on the filtered data. Parameter `K` for factorization is set to 30 after inspecting the plot generated by function `suggestK`. Louvain clustering was performed by setting the resolution such that 17 clusters were obtained. We ran Signac (version 1.1.0) on the same filtered data. We first performed clustering analysis on the gene expression modality and then transferred the cluster labels to the open chromatin modality. Both the gene expression matrix and gene activity matrix were normalized by running `NormalizeData` followed by `ScaleData`. For the gene

activity matrix, the `scale_factor` parameter was set to the median of UMI distribution as suggested by the vignette. Other parameters were kept as default. We then ran `FindTransferAnchors` with the default number of dimensions equals 30, which is the same as used for LIGER. Finally, we ran `TransferData` with weight reduction set to “cca” or the Latent Semantic Indexing (LSI) from analyzing the peak matrix. Results using LSI are presented in the paper as it performed better than the “cca” option.

#### ***Test training split for method comparison***

For scMM and BABEL, which do not allow the single-modality data to be used in the training of the model, we assign 20% of the cells as paired modality data that is used for the training set; the trained model was then used to generate the embedding on the rest of the cells without providing the pairing information. This provided separate estimates for the unpaired mRNA and ATAC modalities, respectively, with which we evaluated whether the paired cells were close together. For `Cobolt`, which allows the use of single-modality data for the training of the model, we assign 20% of the cells as paired modality data and the other 80% were given as single-modality data, and then we trained the model on all of the cells. For LIGER and Signac, which are not designed for multiple-modality data, we hide the pairing information on all cells and treat all of the cells as if they were collected on different cells. To be comparable with the other methods, we only evaluated their performance on the 80% of cells treated by the other methods as single-modality data.

## **Supplementary Information**

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-021-02556-z>.

**Additional file 1: Supplementary Figures.** PDF file containing Supplementary Figures referenced in the main text.

**Additional file 2: Supplementary Methods.** PDF file containing additional details regarding the `Cobolt` algorithm.

**Additional file 3: Table S1.** CSV file containing the cell clustering of the mouse cortex data integration.

**Additional file 4: Table S2.** CSV file containing the cell clustering of the 10X PBMC integration.

**Additional file 5:** Review history.

#### **Acknowledgements**

The authors would like to thank Jean-Philippe Vert for useful early discussions on the analysis of multi-modality data and Hector Roux de Bezieux for assistance in accessing the MOP data.

#### **Review history**

The review history is available as Additional file 5.

#### **Peer review information**

Barbara Cheifet was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

#### **Authors' contributions**

EP formulated the problem and supervised the work. BG developed the model and performed the analysis, with inputs from all authors. BG and YZ developed the model inference and coded the initial implementation. BG implemented the package. BG and EP drafted and revised the manuscript. YZ revised the manuscript and figures. All authors read and approved the final manuscript.

#### **Funding**

This work has been supported by the NIH grant U19MH114830.

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

An open-source package `Cobolt`, implemented in Python, is available on GitHub (<https://github.com/boyinggong/cobolt>) under the GNU General Public License v3.0; the analysis for this manuscript was done with v1.0.0 [51]. Code scripts for reproducing results presented in the paper are publicly accessible on GitHub ([https://github.com/boyinggong/cobolt\\_manuscript](https://github.com/boyinggong/cobolt_manuscript)) under the GNU General Public License v3.0 [52].

Datasets used in this paper are all publicly available. The SNARE-seq data can be downloaded from Gene Expression Omnibus with accession number [GSE126074](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE126074). The MOp scRNA-seq and scATAC-seq data can be downloaded from NeMO Archive with accession number [nemo:dat-ch1nqb7](https://www.ncbi.nlm.nih.gov/bioproject/1000000000). The 10X datasets are available on the 10x Genomics website ([Multiome Chromium X](#), [Multiome unsorted](#), [scRNA-seq](#), [scATAC-seq](#)).

## Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

The authors declare that they have no competing interests.

### Author details

<sup>1</sup>Division of Biostatistics, University of California, Berkeley, Berkeley, CA, USA. <sup>2</sup>Department of Statistics, University of California, Berkeley, Berkeley, CA, USA.

Received: 15 July 2021 Accepted: 22 November 2021

Published online: 28 December 2021

## References

- Welch JD, Kozareva V, Ferreira A, Vanderburg C, Martin C, Macosko EZ. Single-cell multi-omic integration compares and contrasts features of brain cell identity. *Cell*. 2019;177(7):1873–87.
- Stuart T, Srivastava A, Madad S, et al. Single-cell chromatin state analysis with Signac. *Nat Methods*. 2021;18:1333–41. <https://doi.org/10.1038/s41592-021-01282-5>.
- Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck III WM, Hao Y, Stoeckius M, Smibert P, Satija R. Comprehensive integration of single-cell data. *Cell*. 2019;177(7):1888–902.
- Stoeckius M, Hafemeister C, Stephenson W, Houck-Loomis B, Chattopadhyay PK, Swerdlow H, Satija R, Smibert P. Simultaneous epitope and transcriptome measurement in single cells. *Nat Methods*. 2017;14(9):865–8.
- Clark SJ, Argelaguet R, Kapourani C-A, Stubbs TM, Lee HJ, Alda-Catalinas C, Krueger F, Sanguinetti G, Kelsey G, Marioni JC, et al. scNMT-seq enables joint profiling of chromatin accessibility DNA methylation and transcription in single cells. *Nat Commun*. 2018;9(1):1–9.
- Cao J, Cusanovich DA, Ramani V, Aghamirzaie D, Pliner HA, Hill AJ, Daza RM, McFaline-Figueroa JL, Packer JS, Christiansen L, et al. Joint profiling of chromatin accessibility and gene expression in thousands of single cells. *Science*. 2018;361(6409):1380–5.
- Zhu C, Yu M, Huang H, Juric I, Abnoui A, Hu R, Lucero J, Behrens MM, Hu M, Ren B. An ultra high-throughput method for single-cell joint analysis of open chromatin and transcriptome. *Nat Struct Mol Biol*. 2019;26(11):1063–70.
- Chen S, Lake BB, Zhang K. High-throughput sequencing of the transcriptome and chromatin accessibility in the same cell. *Nat Biotechnol*. 2019;37(12):1452–7.
- Lee J, Hwang D, et al. Single-cell multiomics: technologies and data analysis methods. *Exp Mol Med*. 2020;52(9):1428–42.
- Gayoso A, Steier Z, Lopez R, et al. Joint probabilistic modeling of single-cell multi-omic data with totalVI. *Nat Methods*. 2021;18:272–82. <https://doi.org/10.1038/s41592-020-01050-x>.
- Wang X, Sun Z, Zhang Y, Xu Z, Xin H, Huang H, Duerr RH, Chen K, Ding Y, Chen W. BREM-SC: a Bayesian random effects mixture model for joint clustering single cell multi-omics data. *Nucleic Acids Res*. 2020;48(11):5814–24.
- Hao Y, Hao S, Andersen-Nissen E, Mauck WM, Zheng S, Butler A, Lee MJ, Wilk AJ, Darby C, Zager M, Hoffman P, Stoeckius M, Papalexi E, Mimitou EP, Jain J, Srivastava A, Stuart T, Fleming LM, Yeung B, Rogers AJ, McElrath JM, Blish CA, Gottardo R, Smibert P, Satija R. Integrated analysis of multimodal single-cell data. *Cell*. 2021;184(13):3573–3587.e29. <https://doi.org/10.1016/j.cell.2021.04.048>.
- Argelaguet R, Arnol D, Bredikhin D, Deloro Y, Velten B, Marioni JC, Stegle O. MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. *Genome Biol*. 2020;21:1–17.
- Zuo C, Chen L. Deep-joint-learning analysis model of single cell transcriptome and open chromatin accessibility data. *Brief Bioinforma*. 2021;22(4):287.
- Wu KE, Yost KE, Chang HY, Zou J. BABEL enables cross-modality translation between multiomic profiles at single-cell resolution. *Proc Natl Acad Sci*. 2021;118(15):e2023070118. <https://doi.org/10.1073/pnas.2023070118>.
- Minoura K, Abe K, Nam H, Nishikawa H, Shimamura T. A mixture-of-experts deep generative model for integrated analysis of single-cell multiomics data. *Cell Rep Methods*. 2021;1(5):100071. <https://doi.org/10.1016/j.crmeth.2021.100071>.
- Wu M, Goodman N. Multimodal Generative Models for Scalable Weakly-Supervised Learning. In: Bengio S, Wallach H, Larochelle H, Grauman K, Cesa-Bianchi N, Garnett R, editors. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.; 2018. <https://proceedings.neurips.cc/paper/2018/file/1102a326d5f7c9e04fc3c89d0ede88c9-Paper.pdf>.
- Blei DM, Ng AY, Jordan MI. Latent dirichlet allocation. *J Mach Learn Res*. 2003;3(Jan):993–1022.
- Yotsukura S, Nomura S, Aburatani H, Tsuda K, et al. CellTree: an R/bioconductor package to infer the hierarchical structure of cell populations from single-cell RNA-seq data. *BMC Bioinforma*. 2016;17(1):363.
- Sun Z, Wang T, Deng K, Wang X-F, Lafyatis R, Ding Y, Hu M, Chen W. DIMM-SC: a Dirichlet mixture model for clustering droplet-based single cell transcriptomic data. *Bioinformatics*. 2018;34(1):139–46.

21. González-Blas CB, Minnoye L, Papisokrati D, Aibar S, Hulselmans G, Christiaens V, Davie K, Wouters J, Aerts S. cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data. *Nature Methods*. 2019;16(5):397–400.
22. Backenroth D, He Z, Kiryluk K, Boeva V, Pethukova L, Khurana E, Christiano A, Buxbaum JD, Ionita-Laza I. FUN-LDA: a latent Dirichlet allocation model for predicting tissue-specific functional effects of noncoding variation: methods and applications. *Am J Hum Genet*. 2018;102(5):920–42.
23. Srivastava A, Sutton C. Autoencoding Variational Inference for Topic Models. In: Proceedings for the 5th International Conference on Learning Representations; 2017. <https://iclr.cc/archive/www/2017.html>. 5th International Conference on Learning Representations, ICLR 2017 ; Conference date: 24-04-2017 Through 26-04-2017.
24. Tasic B, Menon V, Nguyen TN, Kim TK, Jarsky T, Yao Z, Levi B, Gray LT, Sorensen SA, Dolbeare T, et al. Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nat Neurosci*. 2016;19(2):335–46.
25. McInnes L, Healy J, Melville J. Umap: uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426. 2018. <https://github.com/mcInnes/umap>.
26. Baker A, Kalmbach B, Morishima M, Kim J, Juavinett A, Li N, Dembrow N. Specialized subpopulations of deep-layer pyramidal neurons in the neocortex: bridging cellular properties to functional consequences. *J Neurosci*. 2018;38(24):5441–55. <https://doi.org/10.1523/JNEUROSCI.0150-18.2018>.
27. Sorensen SA, Bernard A, Menon V, Royall JJ, Glattfelder KJ, Desta T, Hirokawa K, Mortrud M, Miller JA, Zeng H, Hohmann JG, Jones AR, Lein ES. Correlated gene expression and target specificity demonstrate excitatory projection neuron diversity. *Cereb Cortex*. 2015;25(2):433–449. <https://doi.org/10.1093/cercor/bht243>.
28. Belgard TG, Marques AC, Oliver PL, Abaan HO, Sirey TM, Hoerder-Suabedissen A, García-Moreno F, Molnár Z, Margulies EH, Ponting CP. A transcriptomic atlas of mouse neocortical layers. *Neuron*. 2011;71(4):605–16. <https://doi.org/10.1016/j.neuron.2011.06.039>.
29. Zeng H, Shen EH, Hohmann JG, Oh SW, Bernard A, Royall JJ, Glattfelder KJ, Sunkin SM, Morris JA, Guillozet-Bongaarts AL, Smith KA, Ebbert AJ, Swanson B, Kuan L, Page DT, Overly CC, Lein ES, Hawrylycz MJ, Hof PR, Hyde TM, Kleinman JE, Jones AR. Large-scale cellular-resolution gene profiling in human neocortex reveals species-specific molecular signatures. *Cell*. 2012;149(2):483–96. <https://doi.org/10.1016/j.cell.2012.02.052>.
30. Fazel Darbandi S, Robinson Schwartz SE, Qi Q, Catta-Preta R, Pai ELL, Mandell JD, Everitt A, Rubin A, Krasnoff RA, Katzman S, Tastad D, Nord AS, Willsey AJ, Chen B, State MW, Sohail VS, Rubenstein JLR. Neonatal Tbr1 dosage controls cortical layer 6 connectivity. *Neuron*. 2018;100(4):831–8457. <https://doi.org/10.1016/j.neuron.2018.09.027>.
31. Yao Z, Liu H, Xie F, Fischer S, Adkins RS, Aldridge AJ, Ament SA, Bartlett A, Behrens MM, Van den Berge K, Bertagnolli D, de Bézieux HR, Biancalani T, Boeshaghi AS, Bravo HC, Casper T, Colantuoni C, Crabtree J, Creasy H, Crichton K, Crow M, Dee N, Dougherty EL, Doyle WI, Dudoit S, Fang R, Felix V, Fong O, Giglio M, Goldy J, Hawrylycz M, Herb BR, Hertzano R, Hou X, Hu Q, Kancherla J, Kroll M, Lathia K, Li YE, Lucero JD, Luo C, Mahurkar A, McMillen D, Nadaf NM, Nery JR, Nguyen TN, Niu S-Y, Ntranos V, Orvis J, Osteen JK, Pham T, Pinto-Duarte A, Poirion O, Preissl S, Purdom E, Rimorin C, Risso D, Rivkin AC, Smith K, Street K, Sulc J, Svensson V, Tieu M, Torkelson A, Tung H, Vaishnav ED, Vanderburg CR, van Velthoven C, Wang X, White OR, Huang ZJ, Kharchenko PV, Pachter L, Ngai J, Regev A, Tasic B, Welch JD, Gillis J, Macosko EZ, Ren B, Ecker JR, Zeng H, Mukamel EA. A transcriptomic and epigenomic cell atlas of the mouse primary motor cortex. *Nature*. 2021;598(7879):103–110. <https://doi.org/10.1038/s41586-021-03500-8>.
32. 10x Genomics. Pbmcs from human (multiome v1.0, chromium x), single cell multiome atac + gene expression dataset by cell ranger arc 2.0.0. 2021.
33. 10x Genomics. Pbmcs from human (no cell sorting, chromium next gem), single cell multiome atac + gene expression dataset by cell ranger arc 2.0.0. 2021.
34. 10x Genomics. Pbmcs from human (3' ht v3.1, chromium x), single cell gene expression dataset by cell ranger 6.1.0. 2021.
35. 10x Genomics. Pbmcs from human (atac v1.1, chromium x), single cell atac dataset by cell ranger atac 2.0.0. 2021.
36. Ding J, Adiconis X, Simmons SK, Kowalczyk MS, Hession CC, Marjanovic ND, Hughes TK, Wadsworth MH, Burks T, Nguyen LT, et al. Systematic comparison of single-cell and single-nucleus RNA-sequencing methods. *Nat Biotechnol*. 2020;38(6):737–746.
37. Pliner HA, Shendure J, Trapnell C. Supervised classification enables rapid annotation of cell atlases. *Nat Methods*. 2019;16(10):983–6.
38. Franzén O, Gan L-M, Björkegren JLM. PanglaoDB: a web server for exploration of mouse and human single-cell RNA sequencing data. *Database*. 2019;2019:baz046. <https://doi.org/10.1093/database/baz046>.
39. Tasic B, Yao Z, Graybuck LT, Smith KA, Nguyen TN, Bertagnolli D, Goldy J, Garren E, Economo MN, Viswanathan S, Penn O, Bakken T, Menon V, Miller J, Fong O, Hirokawa KE, Lathia K, Rimorin C, Tieu M, Larsen R, Casper T, Barkan E, Kroll M, Parry S, Shapovalova NV, Hirschstein D, Pendergraft J, Sullivan HA, Kim TK, Szafer A, Dee N, Groblewski P, Wickersham I, Cetin A, Harris JA, Levi BP, Sunkin SM, Madisen L, Daigle TL, Looger L, Bernard A, Phillips J, Lein E, Hawrylycz M, Svoboda K, Jones AR, Koch C, Zeng H. Shared and distinct transcriptomic cell types across neocortical areas. *Nature*. 2018;563(7729):72–8. <https://doi.org/10.1038/s41586-018-0654-5>.
40. Saelens W, Cannoodt R, Todorov H, Saeyns Y. A comparison of single-cell trajectory inference methods. *Nat Biotechnol*. 2019;37(5):547–54. <https://doi.org/10.1038/s41587-019-0071-9>.
41. Kingma DP, Welling M. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114. 2013.
42. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16). New York: Association for Computing Machinery; 2016. p. 785–94. <https://doi.org/10.1145/2939672.2939785>.
43. Lopez R, Regier J, Cole MB, Jordan MI, Yosef N. Deep generative modeling for single-cell transcriptomics. *Nat Methods*. 2018;15(12):1053–8.
44. Risso D, Schwartz K, Sherlock G, Dudoit S. GC-content normalization for RNA-Seq data. *BMC Bioinforma*. 2011;12(1):480. <https://doi.org/10.1186/1471-2105-12-480>.
45. Zhang Y, Parmigiani G, Johnson WE. ComBat-seq: batch effect adjustment for RNA-seq count data. *NAR Genom Bioinforma*. 2020;2(3): <https://doi.org/10.1093/nargab/lqaa078>.
46. Cole MB, Risso D, Wagner A, DeTomaso D, Ngai J, Purdom E, Dudoit S, Yosef N. Performance assessment and selection of normalization procedures for single-cell RNA-Seq. *Cell Syst*. 2019;8(4):315–3288. <https://doi.org/10.1016/j.cels.2019.03.010>.

47. McGinnis CS, Murrow LM, Gartner ZJ. DoubletFinder: doublet detection in single-cell RNA sequencing data using artificial nearest neighbors. *Cell Syst.* 2019;8(4):329–37.
48. Neph S, Kuehn MS, Reynolds AP, Haugen E, Thurman RE, Johnson AK, Rynes E, Maurano MT, Vierstra J, Thomas S, et al. BEDOPS: high-performance genomic feature operations. *Bioinformatics.* 2012;28(14):1919–20.
49. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *J Stat Mech Theory Exp.* 2008;2008(10):10008.
50. Qiu X, Hill A, Packer J, Lin D, Ma Y-A, Trapnell C. Single-cell mRNA quantification and differential analysis with Census. *Nat Methods.* 2017;14(3):309–15.
51. Gong B, Purdom E. cobolt. Github. 2021. <https://doi.org/10.5281/zenodo.5714790>.
52. Gong B. cobolt\_manuscript. Github. 2021. <https://doi.org/10.5281/zenodo.5715087>.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

