# REVIEWS

# Multi-omics integration in the age of million single-cell data

Zhen Miao [1,2], Benjamin D. Humphreys [3], Andrew P. McMahon[4] and Junhyong Kim [1,2] ✉

Abstract | An explosion in single-cell technologies has revealed a previously underappreciated heterogeneity of cell types and novel cell-state associations with sex, disease, development and other processes. Starting with transcriptome analyses, single-cell techniques have extended to multi-omics approaches and now enable the simultaneous measurement of data modalities and spatial cellular context. Data are now available for millions of cells, for whole-genome measurements and for multiple modalities. Although analyses of such multimodal datasets have the potential to provide new insights into biological processes that cannot be inferred with a single mode of assay, the integration of very large, complex, multimodal data into biological models and mechanisms represents a considerable challenge. An understanding of the principles of data integration and visualization methods is required to determine what methods are best applied to a particular single-cell dataset. Each class of method has advantages and pitfalls in terms of its ability to achieve various biological goals, including cell-type classification, regulatory network modelling and biological process inference. In choosing a data integration strategy, consideration must be given to whether the multi-omics data are matched (that is, measured on the same cell) or unmatched (that is, measured on different cells) and, more importantly, the overall modelling and visualization goals of the integrated analysis.

Many different technologies are now available to measure various properties of any biological system. For instance, the same physical–chemical properties can be measured with different instruments (for example, by RNA sequencing or RNA microarray) and different physical–chemical properties can be measured for the same object (for example, the protein and RNA content of a cell). In the past few years, genome-scale technologies have led to the systematic generation of very-large-scale quantitative datasets that comprise multiple measurement modalities. Although such multimodal datasets have the potential to provide unprecedented insights into biological systems, their analysis and interpretation can be complicated due to modality-specific technical problems and modelling challenges resulting from the need to draw common inference from different kinds of information.

The term 'biological data integration' has been used to describe analytic methods that combine information from multiple sources into a single biological inference. At one level, biological data integration might represent an extremely broad concept such as the integration of diverse information types, including data from Electronic Medical Records, genomic analyses, phenotypic assays and literature reviews, into a broad scientific model or hypothesis[1]. In this broad context,

the term lacks technical meaning and is not pursued further here. Rather, we focus on biological data integration in the context of integrating large-scale omics data, especially at the single-cell level[2,3]. These types of data have a high-degree of multiplexing, for example, with tens of thousands of gene measurements, leading to high-dimensional datasets. If we consider the expression level of a single gene to be one 'dimension' of our dataset, then a set of 10,000 genes would create a dataset of 10,000 dimensions. Each of these dimensions is commonly called a 'feature' of the dataset. Single-cell measurements also tend to have considerable noise and technical artefacts — a problem that is somewhat counterbalanced by the ability of new technologies to obtain measurements from thousands or even millions of cells, in a given tissue[4,5]. This large number of cell measurements alleviates some of the problems associated with high-dimensional data and high noise but creates additional challenges associated with high computational demand and biological complexity.

Despite the challenges associated with high-dimensional data, high noise and large numbers of measurements, high-throughput, single-cell omics methodologies have already provided key insights into kidney biology. For example, single-cell analyses have identified over 30 cell types along the continuous epithelial network; greater

[1]Department of Biology, University of Pennsylvania, Philadelphia, PA, USA.

[2]Graduate Group in Genomics and Computational Biology, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA.

[3]Division of Nephrology, Department of Medicine, Washington University in St. Louis, St. Louis, MO, USA.

[4]Department of Stem Cell Biology and Regenerative Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA.

✉e-mail: junhyong@ sas.upenn.edu

## Key points

- With the development of single-cell multi-omics techniques, tools and models for data integration are critically important.
- Integration problems in single-cell biology can be divided into those associated with the integration of matched and unmatched data.
- Strategies for integrating matched data include joint latent space inference, consensus of individual inferences and biological causal modelling.
- Strategies for integrating unmatched data include annotated group matching, matching with common features and aligning spaces.
- Visualization methods for integrated multimodal single-cell data are still underdeveloped.
- Future challenges include accounting for specific noise related to each modality, overcoming the need for computing efficiency and developing biologically interpretable integration strategies.

---

Assay for transposase-accessible chromatin using sequencing
(ATAC-seq). A technique that profiles the accessibility of DNA elements based on the principle that the Tn5 transposase can insert a transposon only at accessible parts of the chromosome. The insertion location is identified through DNA sequencing.

*Cis*-regulatory elements
DNA elements proximal to a gene that are required for controlling gene expression. Such elements usually include promoters and enhancers, and often contain transcription factor-binding sites.
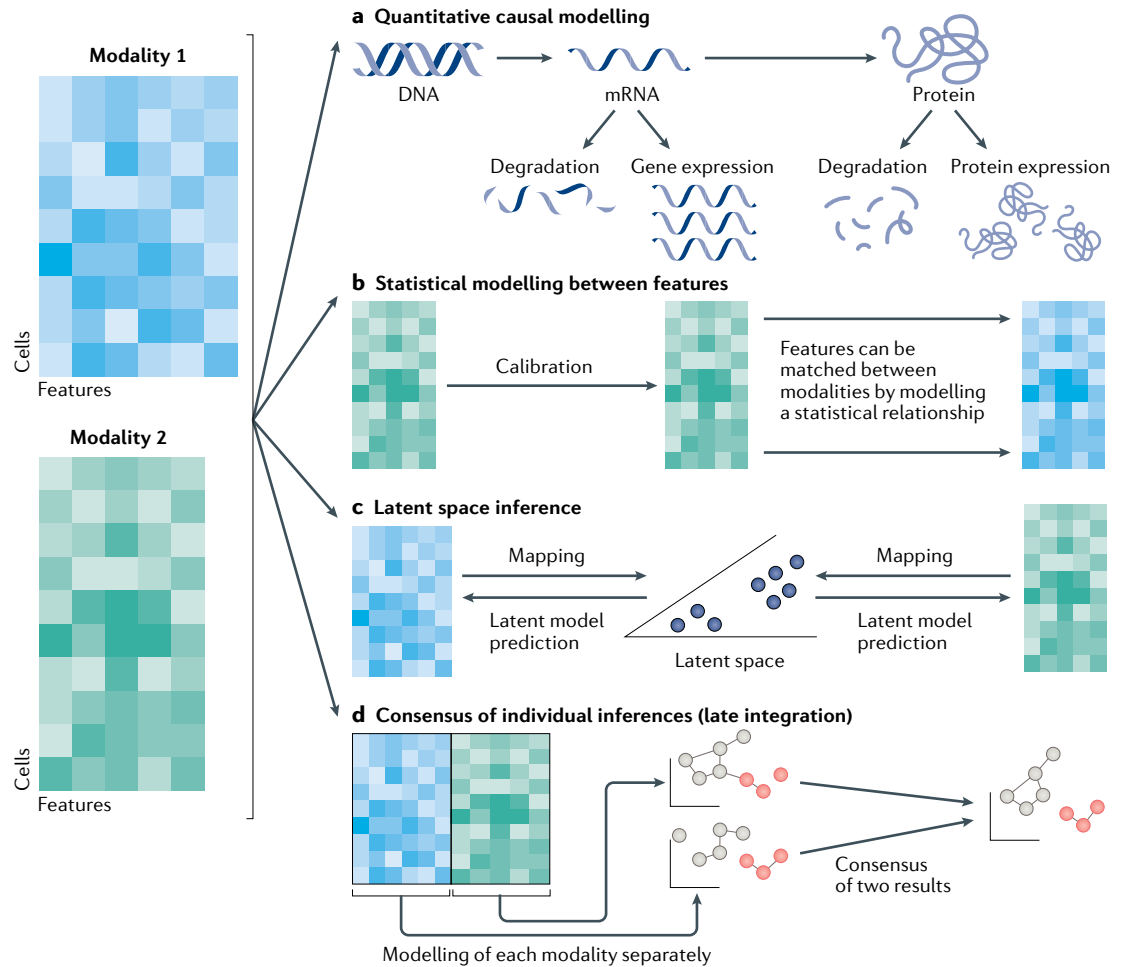
---

diversity probably exists in regional and sex-related cell states among these groupings[6,7]. As another example, time series single-nucleus RNA sequencing (snRNA-seq) has been used to identify dynamic and spatially distinct pro-inflammatory and profibrotic subsets of proximal tubule cells that fail to repair after acute kidney injury[7,8]. Single-cell data can also be combined with clinical parameters such as those regulated by the kidney, including blood pressure, blood pH, osmolarity and estimated glomerular filtration rate. One single-cell transcriptomic profiling study of human kidneys with estimated glomerular filtration rates above and below 60 ml/min/1.73 m² identified *AP1* and *NKD1* as candidate drivers of kidney fibrosis in patients with chronic kidney disease[9].

The abovementioned studies uncovered novel insights into kidney biology using single-cell transcriptomics alone. However, in the past 5 years, many single-cell measurement modalities beyond single-cell transcriptomics have been developed, including approaches to measure multiple data types in the same cell (so-called multi-omics single-cell data). More than 30 single-cell multi-omics techniques[10,11] have been developed since 2015. Although these techniques offer invaluable opportunities to interrogate the properties of cells, the integration of information from these different modalities presents an acute challenge. The high dimensionality, high noise and large number of observations underlie this challenge, in which the goal is to reconcile and make comparable distinct modalities into a coherent biological inference.

Even without explicit computational integration, combining information from different genome-scale data types can yield synergistic inferences. For example, cell-specific gene expression data can be coupled with chromatin status information in the region of an SNP variant, enabling the prioritization of causal variants for further experimental validation[12]. Multimodal data often augment independent evidence from each mode. For example, one study[13] found that a single-nucleus assay for transposase-accessible chromatin using sequencing (snATAC-seq) refined kidney cell-type clusters obtained via snRNA-seq, revealing more clusters with potential clinical relevance. In another study, use of both single-cell RNA sequencing (scRNA-seq) and snATAC-seq enabled the identification of a cell-specific regulatory network

by inferring upstream regulators from analyses of *cis*-element motifs[14]. In that study, the identification of *cis*-regulatory elements with ATAC-seq helped overcome difficulties in detecting regulatory genes, such as transcription factors, in transcriptome data because of their low abundance. This study and others exemplify that the use of multiple modes of omics information can enable combined inferences that cannot otherwise be obtained from any single mode. Thus, the integration of multimodal omics data has the potential to synthesize more knowledge than would be gained as a sum of individual measurements.

Here, we review developments in computational methods for multi-omics data integration. We first provide a general overview of the principles of data integration. Next, we take a more practical data-centric view of what methods might be applied to a particular dataset, starting with a discussion of methods for integrated analyses of multi-omics data measured on the same cell, followed by a discussion of methods for integrated analyses of multi-omics data measured on different cells. We then consider data visualization methods that can integrate different measurement modalities, and we finally discuss current and future challenges for single-cell data integration and prospects for their application to kidney biology. Throughout our Review, we focus on principles and general factors that determine the strengths and challenges of different approaches.

## Overview of single-cell data integration

As described above, available studies demonstrate that even ad hoc integration of multimodal data can yield inferences that cannot be made with a single mode of assay. Many principled computational methods are now available to aid the integration of single-cell multimodal omics data, each with different advantages and drawbacks (FIG. 1). Here, we provide an overview of the general principles of these different approaches before describing methodological details in the next section.

*Quantitative causal modelling.* The most principled form of multimodal data integration is that which considers the actual biological processes that generate the measurements (FIG. 1a). For example, chromatin states, RNA levels and protein levels represent different aspects of the single system-level molecular dynamics of a cell, where a causal relationship exists between the epigenome state, the number of RNA molecules and the number of protein molecules. An accurate quantitative systems model of the cell (BOX 1) uses associated multimodal measurements to estimate parameters with which to derive an integrated inference of the dynamic state of the cell. Some computational approaches incorporate partial systems models of the molecular dynamics of a cell. For example, the popular algorithm for RNA velocity[15] posits a differential equation model of the kinetics of transcription, splicing and degradation, and estimates the parameters of the model using exonic and intronic reads, in effect integrating the two types of read data into a single model inference. The computational tool protaccel[16] extends this kinetic model to

**Modality 1**

Cells

Features

**Modality 2**

Cells

Features

**a  Quantitative causal modelling**

DNA → mRNA → Protein

Degradation   Gene expression   Degradation   Protein expression

**b  Statistical modelling between features**

Calibration → Features can be matched between modalities by modelling a statistical relationship

**c  Latent space inference**

Mapping → ← Latent model prediction   Latent space   Mapping → ← Latent model prediction

**d  Consensus of individual inferences (late integration)**

Consensus of two results

Modelling of each modality separately

Fig. 1 | **Frameworks for the integration of single-cell multi-omics data.** Computational methods enable the integration of measured attributes (that is, features) obtained using multi-omics approaches (for example, transcriptome and protein data) from single cells. These methods can be classified into four broad categories. **a** | Integration based on quantitative causal models. For example, the rates of RNA synthesis, splicing, translation and degradation might be modelled by differential equations and single-cell multi-omics data (for example, gene and protein expression data) can be used to fit the model. **b** | Statistical modelling between features. A statistical function is used to associate data in one modality to another modality, such that the two sets of features (again, for example, gene or protein expression data) can be harmonized into one modality for downstream analyses. Such models can be calibrated from reference datasets or potentially fit to the dataset of interest. **c** | Latent space modelling. Data from different modalities are assumed to be generated from a common latent space and integrated based on the assumption that specific mapping functions are able to map the common latent space onto different modalities. The latent space can be viewed as an integrated low-dimensional embedding of the multi-omics or multimodal data, and the mapping functions can be regarded as a model of the abstract latent space to real observations. **d** | Consensus of individual inferences (late integration). Analyses (such as clustering or dimension reduction) are performed for each individual data modality, after which the results are combined to obtain common consensus outputs or complementary evidence.

include a differential equation term for proteins, allowing a model-based integration of RNA and protein data, such as can be obtained using methods that enable the simultaneous measurement of proteins and mRNAs in single cells (for example, CITE-seq[17] and REAP-seq[18]). A cell systems model-based data integration approach is ideal for the integration of multimodal data but currently impossible owing to the lack of dependable models for most dynamic molecular processes in a cell — especially of models that can predict the dynamics of small finite numbers of molecules in a single cell or in complex processes such as chromosome remodelling.

*Statistical modelling.* In the absence of a causal kinetic model, another possible integration approach is to relate different measurement modalities to each other with a statistical model (FIG. 1b). For example, a statistical relationship could be modelled between RNA levels and protein levels[19] or between the location and amount of open chromatin around a gene and its RNA levels (so-called gene activity models[20]). Therefore, one possible class of methods for the integration of different data modalities is to create a statistical model between two or more modalities such that the value from one data type can be mapped to another type. Such models could be calibrated (that is, the model parameters estimated)

## Box 1 | Computational terminology

### Model

The term 'model' is fairly generic. Here, we use the term 'model' in two different senses. In the first use, a model is a set of quantitative causal descriptions of biological processes, often abstracted to a simple form. An example would be a differential equation that describes RNA levels as a function of the rates of transcription, export and degradation. A second use of the term 'model' is to describe statistical models that relate measurements to each other; for example, a 'linear model' that relates latent space variables to observed variables as a linear mathematical function. This class of models might include more biology-motivated models, such as a gene activity model that posits a statistical relationship between the number of *cis* open chromatin regions and levels of gene expression.

### Machine learning

Machine learning (ML) is a family of computational models that tries to associate a set of input features to a set of output features. Typically, output features are discrete labels such as 'proximal tubule cells' or 'podocytes'. ML methods separate into 'supervised' methods and 'unsupervised' methods. In supervised methods, some observations of 'true' label assignment are known; for example, input features might be gene expression and true cell-type labels are available for some cells. Such ground-truth data are called 'training data'. ML methods try to tune (learn) various mathematical functions to find the association between input features and the known output features of the training data. In unsupervised methods, training data are not available and only input features are available for observations. The typical goal of unsupervised methods is to classify the input observations into groups (for example, clusters) to reveal their grouping patterns.

### Neural networks and deep learning

A neural network (NN), sometimes called artificial NN to distinguish from biological brains, is a subset of ML methodologies motivated by the modelling of a biological brain. The basic idea is to associate input features to output features using a set of mathematical functions called 'nodes'. A node generates output values as a function of all the input values. Thus, a node emulates the metaphor of a neuron integrating all the synaptic input to an axonal output firing. Multiple nodes can be applied to the input features, each of which generates values, resulting in a set of values that can be treated as input features to another set of nodes. Each set of nodes used in this manner is called a 'layer'. The complexity of the artificial NN can depend on the number of nodes in each layer, the number of layers, the input–output relationships between nodes and the type of mathematical function in each node. Deep learning is a non-technical term used to refer to the development of methods that have a very large number of nodes and layers.

### Regularization

Many statistical models can be complicated and overfit the data. For example, in the popular tSNE data visualization method, each data point has its own scale of distance, which can make pairwise relationships arbitrary. A common technique to prevent overfitting is to add some additional constraint, for example, a penalty for high model complexity, to prevent the model from being degenerate. For example, with tSNE, a constraint called 'perplexity' is introduced that constrains the observed data relationship to a certain pairwise distribution. The class of techniques used to constrain model complexity is called regularization. Regularization methods typically have a tunable parameter that controls how much regularization constraint is applied.

### Manifold

In mathematics, a manifold is a smooth connected space that locally resembles Euclidean space (that is, space where distances between points can be defined as the square root of the sum of coordinate differences). In single-cell studies, the term 'manifold' refers to the idea that ensembles of the cells may lie in a lower dimensional subset of the measurement space, which may have non-linear characteristics such as curvature and local folds.

### Loss function

In ML, loss functions are functions that need to be optimized to obtain the desired performance given the data and model. For example, in the 'least squares' regression model, the loss function is the sum of the squared error and, in Lasso regression, the loss function is the sum of the squared error with regularization of regression coefficients. The design of a loss function is key to a successful ML model.

### Encoder–decoder

A commonly used architecture in ML, where a neural net is constructed with a set of nodes that map the input to a middle layer (encoder) and another set of nodes (usually the inverse of the encoder architecture) that maps the middle layer to an output (decoder). The middle layer is typically simpler than the input, for example, with lower dimensions, and tends to encapsulate an abstract characteristic of the input dataset. The decoder then attempts to map this abstracted representation back to some observable data. In an autoencoder, the decoder tries to recapitulate the input data. If successful, the middle layer is thought to represent the essential characteristics of the input data.

### High-performance codes

In programming, there are many different ways to achieve the same computation. Some algorithms are inherently faster than others. For the same algorithm, programmes can also be written differently to speed up the execution by careful use of hardware resources. High-performance codes try to use the fastest algorithms and fine-tune the programmes for optimal speed.

---

from reference datasets or fit to the dataset of interest. When such an approach is used, in effect, all data points of one modality are converted to (mapped to) the other modality, potentially augmenting the power of the dataset. The downside of this approach is that such translation does not provide additional insights into biological processes related to each data type since this process merely converts one data type into another.

*Latent space modelling.* Converting one data type to another can be seen as constructing a (mathematical) function (often called a map) between one set of variables and another. The idea that measurements are related by functions motivates a more abstract framework for data integration. We might model data of type A, type B and type C as a mathematical function from an abstract set of states, which we call 'latent states' or 'latent space'

and the corresponding variables 'latent variables'. More concretely, the transcriptome, proteome and chromatin states might all be considered an aspect of an abstract 'latent molecular state' of a cell (FIG. 1c). That is, if the cell is in a latent state X, then mathematical functions of X will predict the number of RNA and protein molecules and the parts of chromatin that are open. Many machine learning methods (such as autoencoders; BOX 1) and statistical methods (such as factor models; see below) involve estimating a latent space, assumed to determine the observed multimodal values. This latent space, in a sense, is a representation of the integrated data because it 'explains' all of the observed data in different modalities. This concept of latent space from which the observations arise is one of the most common methods of data integration (as discussed below). Different approaches differ in the kind of mathematical functions that map

from the latent space to observations (for example, linear functions versus non-linear functions), in how they model the observed data (for example, as a probabilistic observation from the latent space), in whether they model only in the vicinity of the observed data or the entire relationship between latent and measurement spaces, and in the notion of model fit that they use. In the absence of a more mechanical or causal model, the family of latent space models encapsulates the natural idea that different types of measurements must all represent some aspect of an unknown molecular state of the cell. The main downside of such models is that the latent space typically does not have a physical or chemical interpretation, making it difficult to know what the integrated space means in terms of the actual molecular state of a cell. In addition, the same set of cells may have different latent space representations that model different hidden biological states. For example, the same set of cells might have a latent space representation of their cell cycles, another latent space representation of their circadian rhythms and yet another latent space representation of their cell-type identities. Therefore, the utility and variety of the latent space as a model of data integration depends on the goals of the biological inference.

*Late integration.* The last class of methods for data integration might be called 'late integration'[21] in the sense that this approach does not attempt to relate measurements to each other but rather attempts to use each data modality to infer a model or result unique to that data type and then attempts to integrate the output models or results (FIG. 1d). For example, we might infer gene regulatory networks from the transcriptome and from the proteome independently, and then apply an algorithm to create a consensus network. Another example might be estimating cell-type clusters in each data modality independently before applying algorithms to reconcile the clusters. The above-described study that used snATAC-seq to uncover the dynamics of transcription factor activity, which was then matched with single-cell transcriptome data to identify gene regulatory circuits involved in kidney development[14], can also be thought of as a late integration approach.

To summarize, in the best-case scenario, integrated multi-omics or multimodal analyses can help derive a causal model of cellular processes[22], for example, by using the different data modalities to fit a systems process model. Even without a causal model, analyses across modalities can generate a stronger biological inference than can be achieved with single-modality analyses. As an example, one study[23] found that correlation between chromatin accessibility and gene expression better reflects chromatin conformation than chromatin accessibility information alone. Data from different modalities can also provide independent evidence for hypothesized processes. For example, motifs in the open *cis*-chromatin regions uncovered by ATAC-seq can be used to provide additional evidence for transcriptome-based gene regulatory relationships. Approaches that convert between different data modalities or construct a common latent space can augment mutual information derived from each modality and increase the power of subsequent

inference. For example, clustering analysis on an integrated latent space might yield more stable estimates of cell types that more closely follow biological processes than analysis with single-modality inference. For exploratory analyses of diseases, integrating multiple measurement modalities might also help narrow the molecular nature of the malfunctioning processes and help determine, for example, whether a disease-related change in gene expression is caused by changes in DNA methylation or chromatin accessibility. In conclusion, the different approaches of data integration can help the resulting inference become more than the sum of its parts. Below, we take a more practical data-centric view of what methods one might apply given a particular set of data (FIG. 2).

## Integrating jointly profiled multi-omics data

The greatest challenge for single-cell measurements is recovering molecular fractions from limited amounts of material[24,25]. This problem of molecule recovery efficiency is exacerbated when attempts are made to recover different molecular compartments such as both DNA and RNA. However, simultaneous measurements from the same cell alleviates one challenge of multimodal data integration — mapping the measurement from one modality to another where each modality is measured on a different cell. Here, we refer to data with multimodal measurements on the same cell as matched data. The most popular matched multimodal technique is joint snRNA-seq and snATAC-seq, such as achieved using sci-CAR[26], SNARE-seq[27], paired-seq[28], SHARE-seq[29] and the 10X Genomics Multiome solution. Techniques are also available for joint measurement of transcriptomic and surface protein data, such as achieved using CITE-seq[17] and REAP-seq[18]. Furthermore, technology has been built to measure single-cell phenotypes along with transcriptomic data, providing an important additional dimension for single-cell profiling[30]. The technologies used for matched multi-omics have been reviewed elsewhere[10,31,32].

*Naive approaches.* A number of methods have been developed for the integration of matched multimodal data (TABLE 1). A naive approach is to transform the data in such a way that all the features (that is, the measured attributes) have homogeneous statistical characteristics. A classic approach in organismal systematic biology is to scale each feature by its variation across samples[33,34] (in our case, cells). However, this approach results in all features being considered equally important in determining cell variation, which is not biologically reasonable given their differences in functional importance. A related approach is to give each value of a feature a probabilistic score, perhaps with different models for feature sets, such that the values can have consistent probabilistic interpretation. One example of a model that uses this approach is BREM-SC, which assumes a multinomial distribution of each gene in each cell type for both RNA and protein count matrices obtained using CITE-seq. This type of model enables a probabilistic clustering of cell types[35]. We note that this approach is distinct from attempting to statistically translate

measurements of one modality onto another. These naive approaches are simple but ignore the biological context of the different modalities and instead attempt to harmonize the statistical characteristics of the different features, limiting their utility.

*Latent space approaches.* A more model-based theoretical approach is to consider each measurement, regardless of its modality, as an 'aspect' (or a 'view') of an underlying relationship between the cells. That is, we would assume the existence of a common latent space. One tool that uses this approach to dissect heterogeneity in joint transcriptome and epigenome profiling data is called single-cell aggregation and integration (scAI)[36]. To solve the problem presented by the fact that typical epigenomic information such as that obtained through scATAC-seq is often sparse with a high false-negative rate, scAI first replaces the value of a cell with a similarly weighted average of a random selection of the values of its neighbour to 'smooth over' sparse values. It then infers an underlying common latent space by assuming that the data matrix of the transcriptome and the epigenome can be approximated by a weighted linear function of the shared underlying space. An additional constraint (known as a sparseness constraint) is introduced to make the underlying space as simple as possible, along with another constraint that tries to optimize the preservation of original cell-to-cell relatedness in the underlying common space. Application of this method to joint transcriptome and epigenome data from the kidney enabled

the identification of two subpopulations with distinct open chromatin profiles but similar transcriptomes[36], indicating the need to consider both modalities in order to precisely characterize cell identities.

Latent space approaches can be thought of as integrating at the level of features (that is, early integration). Multi-omics factor analysis (MOFA) and its updated version, MOFA+, implement group factor analysis to identify shared variation across multiple modalities[37,38]. The basic models of MOFA and MOFA+ are similar to that of scAI; that is, the observed data in each modality is considered a linear weighted function of an underlying common latent space. MOFA+ adds multiple underlying latent spaces to account for group effects such as different experimental batches. The main difference with scAI is that MOFA and MOFA+ explicitly attach a probability model such that each cell's feature value is a random variable that is a function of the common latent space. Thus, while the basic mathematical structure of the model is similar to that of scAI, the way MOFA associates the model to the data is different. Although not tailored for single-cell data specifically, the utility of this tool to study a dataset with joint single-cell methylation and transcriptome profiles has been demonstrated[37].

Another tool, totalVI[39], also has similar structure to that of scAI and MOFA in that observed transcriptome and protein measurements (as achieved using CITE-seq[17]) are considered functions of a common latent space. TotalVI relates the observed data and modelled data with a machine learning model (deep neural network) that implements an encoder–decoder scheme (BOX 1). The middle layer of this encoder–decoder neural network can be interpreted as a common latent space and used as the integrated variable set to conduct downstream analyses. A potential advantage of totalVI over scAI and MOFA methods is that the neural network architecture allows more complex (non-linear) relationships between the common latent space and measured features.

*Late integration approaches.* The above methods either explicitly or implicitly aim to infer a common representation space from multi-omics data. An alternative approach involves the integration of data at the level of inferred models (that is, late integration). One such method[40], called weighted nearest neighbour analysis in Seurat V4, synthesizes a combined measure of cell-to-cell affinity from modality-specific affinity models, for example, cell-to-cell relationships calculated using RNA data and protein data. We first note that data in each modality can be used to compute neighbouring relationships of a cell; that is, we can have a neighbourhood by RNA data and a neighbourhood by protein data. Weighted nearest neighbour analysis aims to measure the informativeness of each kind of neighbourhood by assessing how well the cells in each type of neighbourhood predict the RNA or protein value of a given cell. These computations are used to synthesize a weighted average of cell-to-cell affinities from each modality. Another method, called similarity network fusion, aims to synthesize affinity relationships based on a more principled computation idea called 'message passing'[41]. In this approach, a neighbourhood relationship is first calculated for
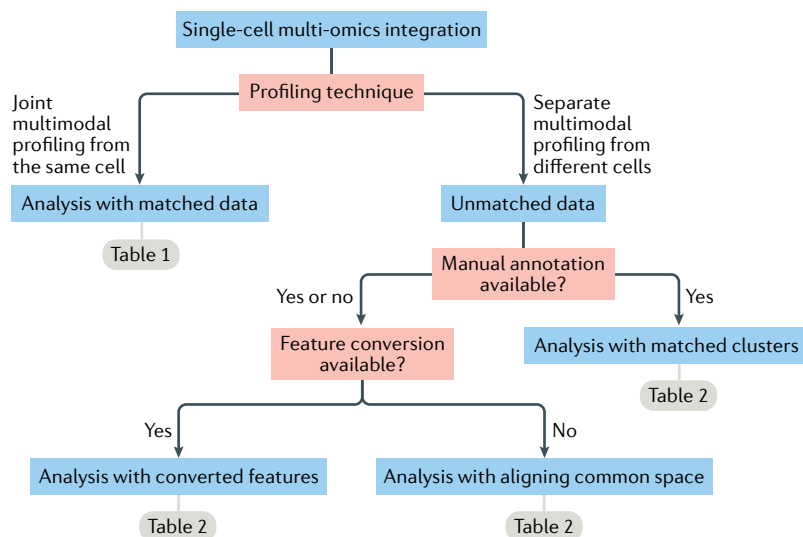


Fig. 2 | **Considerations for choosing an integration method for single-cell multi-omics analysis.** Various data integration methods can be used depending on the nature of the data and whether they are matched (that is, different modalities profiled from the same cell) or unmatched (that is, different modalities profiled from different cells). For unmatched data, analyses can be performed with matched clusters if manual annotations of cell types are available. For example, if we are only interested in the cell-type-level relationship between open chromatin and DNA methylation, we can perform clustering and cell-type annotation for each modality and integrate at the level of cell type. If manual annotations are not available or a higher resolution of integration is needed, two different strategies are available depending on whether feature conversion is possible. For data with a common feature set or converted features (for example, open chromatin to gene activity), tools developed for matching with converted features can be used. For data without common features or feature conversion, integration by aligning common spaces can be applied.

each object (that is, a cell) from the similarity (or affinity) matrix of each modality. Subsequently, the similarity matrices of each modality are 'fused' together by passing the relationship information from the set of neighbouring objects of one matrix to the other matrix, back-and-forth iteratively until they converge. This basic approach was implemented in CiteFuse[42] as a method to integrate affinity relationships from RNA and surface protein from CITE-seq.

### Integrating independent multimodal data

With current technologies, a more common problem than the integration of matched datasets is the integration of two or more independently collected datasets (that is, unmatched data), with different modalities. The emergence of comprehensive, single-modality, single-cell datasets across whole organisms[43–45] has led to an abundance of highly accessible data of this type. In general, experimental approaches for joint measurements of certain modalities are still under development or may be impossible. For example, approaches to simultaneously quantify single-cell transcriptome and whole-proteome data are extremely challenging as single-cell proteomics techniques are rapidly advancing but still lack sensitivity[46]. Single-cell lipidomics has been more successful than proteomics at quantitatively identifying molecular species[47] but we are not aware of any attempts at multimodal measurement of lipidomics data. The key problem for unmatched data is that measurements from each modality are unlikely to have cell-to-cell correspondence. That is, in measurements from one set of cells using one modality, say proteins,

and another set of cells for another modality, say the transcriptome, it is highly unlikely that there will be cells in each set that correspond exactly to the same cell state for both modalities. Thus, almost by definition, we cannot integrate information at the level of individual cells when measurements are not matched. Current integration approaches therefore attempt to match groups of cells, either at the level of distinct cell types or at the level of local ensembles (neighbouring cells). Alternatively, some methods try to statistically map one feature space to another. Here, we classify these methods into three main categories: those that match by annotated cell groups, those that match by a shared feature set and those that match without a common feature set (TABLE 2).

*Matching by annotated cell groups.* When different measurements are made on different sets of data, one coarse-grained approach to integrate those measurements is to match groups of cells (for example, clusters) between the modalities. The clusters in each modality can be associated manually if the clusters correspond to known cell types, which might have been inferred from expert knowledge (for example, through marker gene expression). If cluster label information is not available from established annotations, other features that are biologically informative can be used, such as the proximity of open chromatin to expressed genes, averaged over the ensemble of the cluster to match clusters from each modality. One study, for example, integrated scRNA and scATAC data[48] by linking open chromatin peaks of scATAC-seq cell clusters with the expression of scRNA-seq cell clusters through their proximity in the

### Table 1 | Methods for matched data analysis

| Tool | Data type | Model | Additional notes | Documentation | Ref. |
|------|-----------|-------|------------------|---------------|------|
| BREM-SC | T + P | Early integration, probabilistic modelling | This method models the observed data by multinomial distributions and assumes data from both modalities to be generated in a cluster-specific manner | https://github.com/tarot0410/BREMSC | 35 |
| scAI | T + C | Early integration, latent space modelling | scAI iteratively updates a regularized matrix factorization model to obtain an optimal common cell-loading matrix across two modalities | https://github.com/sqjin/scAI | 36 |
| MOFA+ | T + C | Early integration, latent space modelling | MOFA and MOFA+ were built on the framework of group factor analysis but extend the model to enable the integration of different data types (count versus binary) | https://github.com/bioFAM/MOFA2 | 38 |
| TotalVI | T + P | Early integration, latent space modelling | This method uses a variational autoencoder framework built on scVI. In this method, the protein measurements are modelled with a negative binomial mixture distribution to account for background reads | https://github.com/YosefLab/scvi-tools | 39 |
| CiteFuse | T + P | Late integration, latent space modelling | The similarity measurement for protein data is based on a proportionality coefficient and the similarity measurement for RNA data is constructed with the Pearson correlation | https://github.com/SydneyBioX/CiteFuse | 42 |
| Seurat 4.0 | T + P | Late integration, latent space modelling | Computes a weighted average cell affinity matrix from modality-specific affinity matrices. The weights are computed to reflect the predictive information within a cell's local neighbourhood defined within each modality | https://github.com/satijalab/seurat | 40 |

BREM-SC, Bayesian random effects mixture model-single cell; C, chromatin accessibility; MOFA, multi-omics factor analysis; P, proteome; scAI, single-cell aggregation and integration; scVI, single-cell variational inference; T, transcriptome.

Table 2 | **Methods for unmatched data analysis**

| Strategy | Tool | Data type | Feature matching | Algorithm | Additional notes | Documentation | Ref. |
|---|---|---|---|---|---|---|---|
| Group matching | Stereoscope | T + ST | R | Deconvolution | This method assumes negative binomial distributions of genes and tolerates differential gene capture efficiencies between two technologies | https://github.com/almaan/stereoscope | 53 |
| | MAESTRO | T + C | R | CCA + MNN | This method implements ChIP–seq data-based TF enrichment score calculators to define core TFs in each cell-type cluster | https://github.com/liulab-dfci/MAESTRO | 49 |
| Comon features | STvEA | MI + ET | R | MNN | This method also provides a framework to transfer cell-type annotations from one modality to another | https://github.com/CamaraLab/STvEA | 54 |
| | Clonealign | T + D | R | Variational Bayes | This method assumes correlation between DNA copy number and gene expression within the same region | https://github.com/kieranrcampbell/clonealign | 56 |
| | Seurat 3.0 | T + C | R | CCA + SNN | This method identifies anchor cells between datasets based on SNN across modalities; these anchor cells serve as a bridge for matching | https://github.com/satijalab/seurat | 57 |
| | LIGER | T + M, T + C | R | iNMF | The relative contribution of dataset-specific factors and shared factors is determined by a hyperparameter λ, which can be used to fine-tune the integration results | https://github.com/welch-lab/liger | 58 |
| Aligning spaces | MAGAN | MI + T | R | GAN | This method identifies cell-to-cell correspondence by adding a loss function defined by similarity of cell matching; such loss function requires at least some shared features between two datasets | https://github.com/KrishnaswamyLab/MAGAN | 60 |
| | MATCHER | T + C | NR | Manifold alignment | This method assumes 1D structure (pseudotime) with a pre-specified direction | https://github.com/jw156605/MATCHER | 61 |
| | MMD-MA | T + M | NR | MMD | In addition to the MMD loss, the loss function also has a distortion loss and a penalty to ensure the dimensionality and orthogonality of each projection | https://bitbucket.org/noblelab/2019_mmd_wabi/src/master/ | 62 |
| | UnionCom | T + M | NR | GUMA | The algorithm generalizes the GUMA method to achieve soft matching between datasets, enabling matching with different numbers of cells | https://github.com/caokai1073/UnionCom | 63 |
| | SCOT | T + C | NR | GWOT | A late integration method in which a similarity matrix is constructed by each modality separately, after which probabilistic transportation between datasets is achieved by GWOT | https://github.com/rsinghlab/SCOT | 64 |

C, chromatin accessibility; CCA, canonical correlation analysis; ChIP–seq, chromatin immunoprecipitation followed by sequencing; D, DNA; ET, simultaneous epitope and transcriptome; GAN, generative adversarial networks; GUMA, generalized unsupervised manifold alignment; GWOT, Gromov–Wasserstein optimal transport; iNMF, integrative non-negative matrix factorization; M, methylome; MI, multiplexed immunohistochemistry; MMD, maximum mean discrepancy; MNN, mutual nearest neighbours; NR, not required; R, required; SNN, shared nearest neighbours; ST, spatial transcriptome; T, transcriptome; TF, transcription factor.

genome, from which they inferred enhancer–promoter pairs. These enhancer–promoter pairs were consistent with prior knowledge of regulatory networks, supporting the utility of this method. Another approach, MAESTRO[49], incorporates additional information from chromatin immunoprecipitation followed by sequencing (ChIP–seq) databases to help define transcriptional regulators and match clusters based on scRNA and scATAC data.

Matching at the cell-group level is also common practice in analyses of spatial transcriptome data. Most current spatial transcriptomics technologies either lack resolution or transcriptome complexity (reviewed elsewhere[50]); however, integrating scRNA-seq with spatial data can help overcome these two limitations. For example, training of a machine learning classifier, Support Vector Machine, on highly variable genes from annotated scRNA-seq clusters enabled the classifier to identify and map major cell types from sequential fluorescence in situ hybridization (seqFISH)

data through which only 125 genes had been profiled[51]. For spatial transcriptomics data with low cellular resolution — such as that obtained using 10X Visium and Slide-seq[52] — scRNA-seq data can be used to deconvolute the spatially averaged low-resolution readout and increase resolution by estimating the frequencies of each cell type[53].

**Matching with shared feature sets.** In rare cases, measurement modalities might be different but their common molecular basis can be used to match the features. For example, STvEA[54] matches CITE-seq data with multiplexed immunohistochemistry or flow cytometry data using measurements of protein abundance as the common factor. Matching is achieved through mutual nearest neighbour (MNN) correction[55] on the two data matrices, enabling the automated annotation of multiplexed immunohistochemistry (or flow cytometry) data with labels from CITE-seq data. Given two sets of objects and a notion of distance across the datasets,

**Sequential fluorescence in situ hybridization** (seqFISH). A technique that measures mRNA quantity through sequential fluorescent probes that have combinatorially encoded information for each targeted mRNA. For example, a sequence signal, probe A then B, might encode gene X, whereas the sequence probe A then C might encode gene Y.

MNN identifies pairs of objects in the two sets that are considered each other's nearest neighbour. A classic application of MNN is in identifying homologues among gene paralogues; variations of the MNN principle have been used widely in data integration.

In the absence of a common molecular basis, measurements of one modality may be connected to features of another by some (biologically motivated) statistical model to enable joint analysis. For example, clonealign[56] assumes that an increased DNA copy number (inferred from single-cell DNA sequencing data) in cancer cells will result in increased gene expression within the corresponding region. Many scRNA–scATAC integration methods synthetically construct a 'gene activity matrix' from ATAC data, which is treated as a gene expression feature set. Multiple models have been proposed to infer gene activities from chromatin accessibility data. Seurat V3 (REF.[57]) aggregates all ATAC reads from −2 kb of the transcription start site (TSS) throughout the whole gene body to predict expression levels. MAESTRO[49] assigns weights to each peak with an exponential decay based on the distance to the TSS. The Cicero model[20] is more complex and considers read depth and distal elements that are co-accessible with the TSS. Mapping features of one modality onto another often creates systematic differences that are similar to normalization problems and batch effects. Therefore, good calibration after feature conversion is essential for matching to be successful. Calibration can begin before integration: for example, Seurat V3 and STvEA conduct normalization of both datasets before integration, whereas this step is usually skipped by other models. Seurat V3 and MAESTRO pipelines implement canonical correlation analysis to align the two datasets, which are then mapped to the same gene expression feature space; they then apply MNN correction[55] for additional alignment.

***Integration of unmatched data by latent models.*** Similar to matched data cases, data from each modality can be modelled as maps from an abstract set of common factors (latent factors). LIGER[58] uses an integrative non-negative matrix factorization[59] approach to jointly factorize multiple cell-by-feature matrices into cell-by-factor matrices and factor-by-gene matrices using a set of common factors for all matrices and another set of factors specific to each matrix. Factors here refer to hypothetical underlying (latent) features that can be thought as abstract cell states that determine observed values. Multiple modalities can be integrated through statistical modelling of features; for example, by using a quantitative measure of gene accessibility from snATAC-seq to estimate gene activity for integration with scRNA-seq. The factor loadings of each gene are usually interpreted as 'metagenes' and the magnitude of modality-specific factors is constrained and regularized (BOX 1). The factor loadings of each cell are used for clustering and cell matching. Matrix factorization methods assume that observed data are weighted linear functions of the latent decompositions but, similar to the above discussion of totalVI[39], more-complex relationships can be modelled with neural networks. MAGAN[60] implements

a type of neural network called dual generative adversarial network (dual GAN) that uses a new architecture to map two datasets from different modalities reciprocally.

Obtaining a shared feature set by mapping between modalities can be challenging or even impractical when the measurements from each cell are vastly distinct. Rather than operating on a discrete set of observed data points, another approach is to consider modelling the entire 'space' of data for each modality and map the spaces to each other. Manifold (BOX 1) alignment and related methods assume that individual cells occupy some geometric subset of the feature space of a given modality. These geometric subsets have been called a 'manifold' in the literature (with some abuse of the mathematical term). These manifolds can be thought of as smooth curved surfaces that characterize a biologically feasible set of values for a given collection of cells. Manifold alignment methods assume that a shared latent structure (a manifold) underlies each dataset and tries to learn a shared manifold among datasets to build correspondence between them; the approach is similar to linear latent variable models but with more generality.

Tools that implement manifold alignment include MATCHER[61], MMD-MA[62] and UnionCom[63]. These methods start with dimension reduction of the datasets. As an important first step, dimension reduction methods are chosen to be consistent with the model assumption and suitable to the data structure. MATCHER starts with the assumption that a one-dimensional structure exists along which all cells lie (this one-dimension can be interpreted as pseudotime). MATCHER then fits a stochastic model to infer a one-dimensional manifold structure (that is, pseudotime) for each data modality. Subsequently, a function termed monotonic warping is trained to minimize the loss function (BOX 1) matching two or more one-dimensional manifolds with a pre-specified manifold orientation. Monotonic warping refers to a function that associates two variables to each other that are strictly increasing or decreasing (that is, order preserving). Schematically, MMD-MA maps geometric relationships within each modality feature space to a common space in a way that minimizes geometric distortions between each modality space while maintaining the intraspace configuration. UnionCom embeds each modality into a distance matrix that encapsulates a low-dimensional manifold for each modality. A well-defined pairwise distance matrix is sufficient to represent the complete geometric configuration of points. Thus, two matrices in UnionCom represent the estimated geometric relationships of the cells in each measurement modality. By optimizing a notion of difference between the two geometric configurations, the configurations of two modalities are matched and probabilistic cell correspondence between the two datasets are computed. Somewhat distinct from manifold alignment, SCOT[64] uses the notion of optimal transport, which tries to define a relationship between two sets of objects, each with a number of classes (for example, cell types) and different frequency of objects in each class. The computed relationship considers both

the frequency of objects in each class and a measure of distance between the objects.

Matching different modalities by aggregation is a natural idea but tends to lose individual cell resolution. Some approaches attempt to recover individual cell resolution through initial aggregated matching and then refinement, but the results from these approaches can be highly dependent on initial conditions. Matching by applying statistical models between the features of the different modalities can provide cell-level resolution but this approach is highly dependent on the accuracy of the statistical models. Although a clear relationship exists between chromatin states and gene expression, the exact relationship, especially with respect to temporal dynamics, is unclear. Matching by latent space or manifold alignment models are somewhat more principled approaches than aggregation and refinement, but the available models are complex and their interpretation in biological terms is often unclear. In summary, the available approaches have different strengths and weaknesses and their utility is likely to be highly data and problem dependent.

### Visualization of multi-omics data

Computational visualization tools or interactive websites that allow user-friendly searches and the display of features notably promote data sharing and reuse. Two large categories of data visualization exist in the context of single-cell biology. One might be called 'unbiased' visualization and includes various dimension reduction approaches that attempt to display all data points. The other might be called 'knowledge-driven' visualization, whereby certain curated aspects of the data (for example, a focal subset of cells) are displayed. Although multiple tools have been developed to visualize scRNA-seq data, tools for the explicit visualization of single-cell multi-omics data are scarce. Below, we provide a brief overview of current methods and discuss future directions for multimodal single-cell visualization.

*Unbiased visualization.* Dimension reduction and unbiased visualization have been crucial in the interpretation of complex single-cell data. The diverse cell types and states within a single-cell dataset mean that visualizing cells as a point in a two-dimensional or three-dimensional image is useful for evaluating data qualities, cell identities, developmental trajectories and batch effects[65]. Various visualization methods have been implemented based on dimension reduction approaches, including tSNE, UMAP[66], PHATE[67] and force-directed graphs[68]. These methods extend from classic linear methods such as principal component analysis, which is based on projecting data points onto (orthogonal) directions of maximum variation, and embedding methods such as multidimensional scaling (MDS). MDS computes one set of distance relationships in the original high dimension after which points are placed in lower dimensions such that the distance relationships in the lower dimensions are as similar to those of the original dimensions as possible. Variations of MDS involve different ways to define distances or measure the distortions between high-dimensional and low-dimensional

distance relationships. The main problem faced by dimension reduction and visualization methods is that the configuration of points in a high-dimension state cannot be represented in lower dimensions without error and the methods therefore have to trade off the kinds of distortions that they allow. Typically, it is hard to uniformly spread out the distortions from smaller distances (for example, within clusters) to those from larger distances (for example, between clusters). The type of trade-off is determined by the approaches used to calculate distances and measure high-to-low distortions; typical options trade off accuracy at large distances for accuracy at smaller distances.

Methods such as tSNE, UMAP and PHATE add another twist to the dimension reduction approach by allowing inhomogeneous notions of distance or similarity. That is, a distance from point X to point Y might be different from that of point Y to point X. One interpretation of this approach is that the inhomogeneity in distances is related to curvature or (diffusion) velocity; thus, the distance of X to Y might be analogized to going uphill versus Y to X going downhill or a particular region might have high curvature and is therefore hard to traverse. Modern methods of visualization also implement non-linear notions of distance (or similarity) such that certain distances are emphasized, whereas others are de-emphasized, which often allows the resulting embeddings to highlight cluster relationships. These methods try to control the arbitrary freedom allowed by such flexibility by imposing user-defined constraints (for example, 'perplexity' in tSNE; BOX 1). We caution that the high flexibility of these methods can complicate the interpretation of data. The visualizations can also be unstable, either because the algorithms start from random initial configurations or due to the sensitivity to the addition and subtraction of points. In-depth discussion of tools for the visualization of single-cell data can be found elsewhere[69].

Unbiased visualization approaches naturally extend to multi-omics single-cell data as long as the above-described integration methods produce representation in a common space. Any of the available dimension reduction methods can be used to visualize integrated relationships within a common latent space, for instance, a shared gene expression space (by gene activity modelling) or a common layer in a neural net. For example, the multi-omics visualization arm of scAI[36], called VscAI, enables the visualization of cells, genes and (accessibility) loci by an embedding that reflects the low-dimensional latent space. However, the nature of integrative analyses suggests the need for more complex multiple views of the data. For example, we might want to see single cells laid out in the common latent space and then also see their configuration in each of the measurement modalities, in particular, with cell correspondences in each space. Although it is possible to switch views (as described elsewhere[6,44]), currently available methods do not easily show correspondence between layers. It would be desirable to have visualization systems similar to geographical information systems such as those used in landscape ecology[70], which have layers of multimodal maps.

---

**Principal component analysis**
A common dimension reduction method that aims to project the original data to a fixed smaller dimension while minimizing the squared error during data reduction. This approach can be viewed as maximizing the variance in the projected data.

**Embedding**
In mathematics, embedding is a map from one set X to another set Y, where some characteristic of X is preserved. In single-cell studies, the term embedding has been used for methods that 'place' cells in a new feature space, possibly of a lower dimension, such that notions of cell-to-cell distances are approximately preserved.

*Knowledge-driven visualization.* Single-cell data are used by researchers to derive additional biological inferences — a process that is often called downstream analysis. These downstream analyses result in the production of additional visual objects. Common examples of these visual objects include violin plots to visualize cell-type marker genes, di-graphs to visualize cellular interactions or even simple annotation overlay to visualize a focal subset of cells. Other visual devices that focus on particular knowledge-driven assumptions include displays of motif enrichment along with the expression of corresponding transcription factors[14], visualization of sequence reads along genomic tracks[71] and other associated annotation data organized by genomic coordinates[71,72]. One important approach to incorporate existing knowledge for single-cell data is to associate spatiotemporal information with single-cell visualization. Temporal trajectories have been visualized using many different pseudotime methods; for example, the RNA velocity method[15] displays estimated displacement vectors to extrapolate the 'flow' of cell differentiation states. Approaches for the visualization of single-cell data in the context of anatomical ontology (for example, KidneyCellExplorer[6]) or within detailed 3D models (for example, the NIH HuBMAP portal[73]) are under development.

*Future directions for data visualization.* Additional visualization tools and frameworks are needed to fully appreciate the complexity of multimodal data (FIG. 3). Visualization tools with greater flexibility to enable the display of multiple and coordinated views that link objects in various modalities will aid visual explorations of multimodal relationships. However, even multiple layers of data visualization will be insufficient to fully explore the biological structure of multimodal data if the visualizations are static. Complex data are best explored with interactive systems that enable dynamic modifications of views, such as the ability to re-display subsets of data or dynamically switch between different modalities. One critical consideration is the computational speed required for such interactive visualizations and analyses, especially for very large datasets (for example, those with data for $10^6$ cells[4,5]). As datasets scale to extremely large
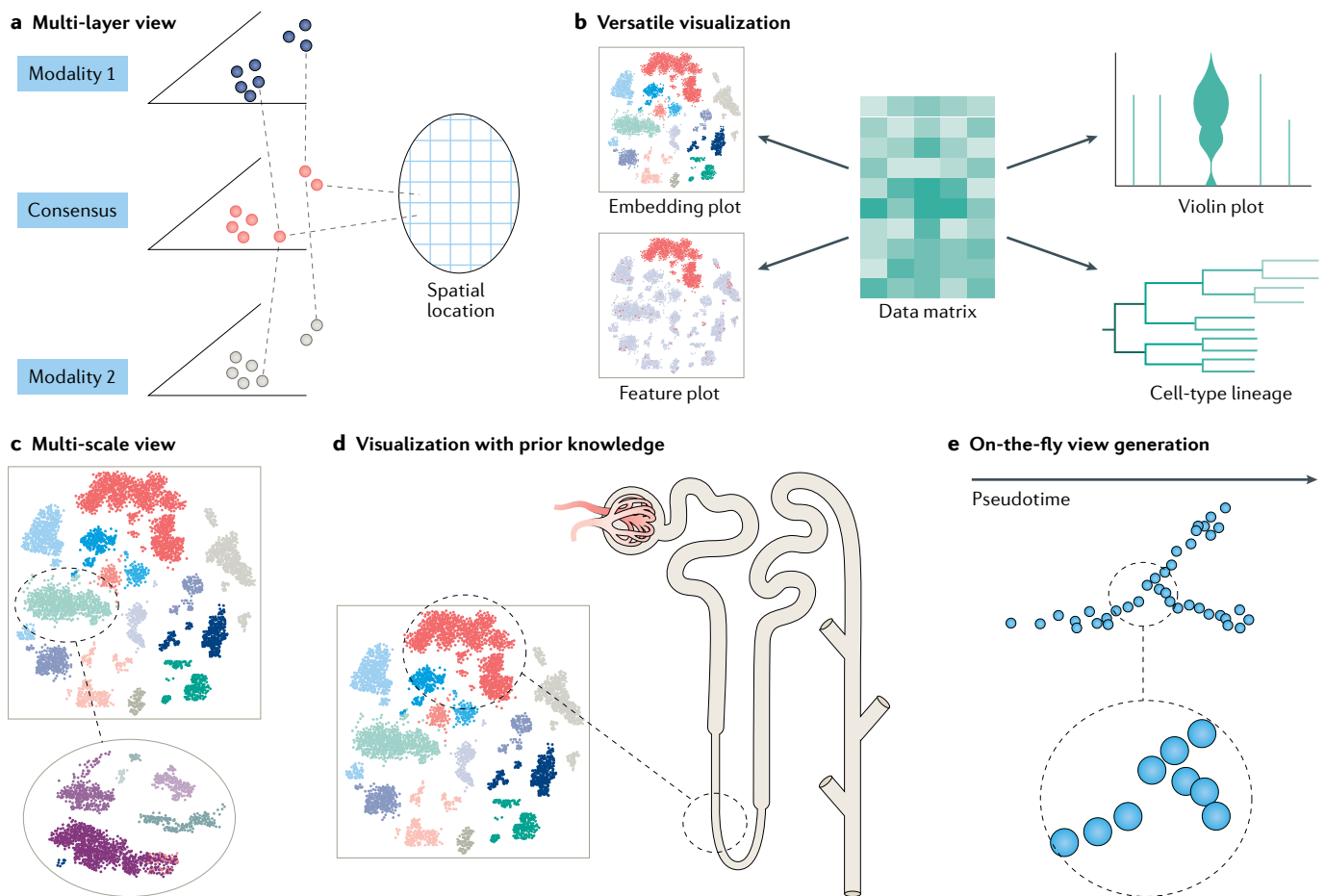


Fig. 3 | **Desired properties and functionalities of visualization tools for single-cell multi-omics.** Visualization of multi-omics data requires additional functionalities given the complex data structure; for example, the ability to switch the view between different modalities. Some other desirable features include multiple layers of data visualization based on data obtained for different modalities with mapping between each layer (ideally, the mapping between each observation and their spatial location can also be displayed as another layer of information) (part **a**); versatile and dynamic visualizations that incorporate downstream analyses or prior knowledge (part **b**); multi-scale views with multiple resolutions to assist the dissection of very large datasets (part **c**); integration of prior knowledge such as ontology and anatomy with multi-omics data to help anchor biological knowledge to the data (part **d**); and tools that enable on-the-fly or dynamic visualization of data to enable more-flexible data visualization (part **e**).

sizes, issues of where to store and compute the views (for example, in the cloud or on a client computer) become non-trivial.

Viscello[44], Cerebro[74], VscAI[36] and Giotto[75] are some of the tools that currently allow some degree of interactive multimodal single-cell data visualization. Some consortia, including ReBuilding a Kidney (RBK) and the GenitoUrinary Development Molecular Anatomy Project (GUDMAP), integrate interactive single-cell visualizers in their data archive. However, these tools are not fully interactive in the sense that they cannot recompute the visualization to an arbitrary choice of views or subsets of data.

## Challenges for single-cell multimodal data integration

The above discussion of data integration approaches for single-cell multimodal data touches only the surface of this very active area of ongoing research. All available approaches have common challenges that must be considered. These challenges can originate from the processes of data collection, data conversion or data interpretation. Here, we discuss some of the most prominent challenges to single-cell, multimodal data integration.

*Accounting for data characteristics.* It is well acknowledged that single-cell data are noisy. This noise arises from biological and technical variation. Common biological variation includes the stochastic bursting of genes[76], variation arising from circadian rhythm[77] and cell cycling[78], and variation arising from the local cell environment. The contribution of technical variations is debated but may include uneven dropouts and coverage[25,79], transcript contamination (from ambient RNA)[80] and multiplets[81]. In general, single-cell assays, especially high-throughput assays, tend to be lossy since the technologies trade off sensitivity (that is, the efficient capture of the molecules in a cell) for throughput, resulting in sparse datasets. This sparsity is a huge challenge and is typically approached by 'borrowing' local information from nearby cells, which can introduce additional biases. Multi-omics approaches have the potential to resolve some confounding factors, sparsity or noise in a single modality by 'borrowing' information in the other modality, but this integration does not always improve the prediction power achieved by a single modality[82]. Noise across multiple layers can be amplified, leading to a decrease in the signal strength[83]. An important problem in single-cell analysis is that commonly used noise models typically use off-the-shelf parametric models, such as Poisson zero inflation models and negative binomial models[84,85], whereas, in practice, single-cell noise does not seem to be well modelled by these parametric models and systematic control experiments to measure the characteristics of the noise have been rare[79,86].

Although models have been built to distinguish biological and technical variation in scRNA-seq data[87], models to account for heterogenous noise across multiple modalities still need to be developed. In some cases, the problem of heterogeneous noise is best handled by 'early integration', whereby the input datasets themselves are operated on to make a single compatible matrix (for example, by applying weights and concatenating the datasets). In other cases, the problem is best approached by 'intermediate integration' (for example, using the latent space methods approach to map the input data to theoretical common space features). In still other cases, 'late integration' might be the best approach, whereby each modality is used to infer a model, such as a gene regulatory network, and then the inferred models are combined appropriately (for example, using CiteFuse[41]). Each of these approaches have advantages and disadvantages depending on the modalities being integrated and other conditions of measurement (for example, batch effects). Earlier integration might help increase the power of ultimate downstream analysis (for example, the identification of cell clusters) both by increasing the size of the dataset and by bringing together (possibly) complementary information. Late integration can help the application of modality-specific models and methods to handle heterogeneous noise and enable the individual inferences (for example, clusters from each modality) to be combined to obtain a more robust inference.

*Data types and cell composition compatibility.* Although desirable to integrate information from all relevant sources, datasets that are to be integrated can be vastly distinct. At a simple level, gene expression profiles in scRNA-seq data are continuous variables, whereas chromatin accessibility measurements are usually binarized to indicator variables[88] (also known as dummy variables). This integration of distinct datasets requires a consistent way to match metric variables with nominal variables, which can have both technical and conceptual challenges[89].

At a more complex level, traits such as cell morphology, while having a metric representation, are difficult to statistically characterize in a meaningful manner. The emergence of machine learning methods has led to the development of approaches to integrate morphology and expression data. Fascinating insights from these studies suggest that cell morphology might predict gene expression[90], but the functional connections of such relationships are still unclear. Another more common but important challenge is that subtypes of cells that are recovered and measured with high-throughput single-cell methods can be very different for different measurement modalities. For example, immune cell populations are usually over-represented in scRNA-seq datasets, probably as a result of recovery bias, whereas snRNA-seq methodologies demonstrate bias in their recovery of different subpopulations[91]. Such differences in cell-subtype distribution can complicate data matching, especially for nearest neighbour-based methods.

*Computing millions of data points.* With the development of combinatorial indexing technologies[92] and sample multiplexing strategies[93], datasets are now available at a $10^6$ scale[4,5]. Efficient computing over such big data matrices requires different strategies to those used for smaller datasets. We note that just to compute pairwise relationships for a dataset of size $10^6$, ~$10^{12}$ computations

---

**Dropouts**
In single-cell biology, dropouts are usually the transcripts that were present in the cell but were not captured during sequencing.

**Ambient RNA**
In droplet-based single-cell RNA sequencing approaches, the measured mRNA molecules can be contaminated by mRNAs from other cells present in the suspension, for example, owing to cell rupture. These contaminating mRNAs are termed ambient RNA.

**Multiplets**
During high-throughput single-cell (or single-nucleus) isolation in droplets or similar vessels, two or more cells might be captured together creating a mixture of molecules. Computational methods have been developed to detect and remove such unwanted observations from the dataset.

need to be considered. This scale of computing is prohibitive, resulting in the use of less intensive heuristic methods. In fact, even just laying out a million points for data visualization becomes a heavy computational burden and prevents researchers from exploring different views due to the wait time involved. Future integrative analyses of single-cell data will require concerted efforts in algorithm development with incorporation of novel stochastic indexing strategies, streaming algorithms and careful heuristics, along with the development of carefully tuned high-performance codes (BOX 1). Some areas of computational biology such as phylogenetics and protein folding have long been acutely limited by computational speed and advanced algorithmics have been an inherent part of those fields. We suspect single-cell biology will soon demand similar levels of algorithm sophistication and high-performance software engineering.

*Modality mapping.* As discussed earlier, the integration of unmatched measurements is often achieved by mapping the values of one modality to those of another — a key example is the conversion of chromatin states to gene expression values. However, such conversions assume an over-simplified model between different modalities, mostly due to a lack of knowledge of whole-genome gene regulatory logic. As previously reported[29], the temporal dynamics of the open chromatin states of a cell are not at the same phase as its corresponding RNA expression; rather, gene expression lags behind the opening of its proximal chromatin. Thus, accurate mapping between the modalities requires both a precise knowledge of the mechanisms connecting the measured molecules and the temporal dynamics of the mechanisms. Similar considerations would apply when mapping between the transcriptome and the proteome or, a more complicated scenario, the connection between molecular and morphological states.

*Interpretability and validation.* Most data integration methods avoid detailed causal modelling. At the extreme end are purely data-driven machine learning methods such as autoencoders (BOX 1). For example, one autoencoder-based multi-omics data integration method[94] has been trained to create a common latent space for many different modalities. Powerful computational tools such as this can indeed integrate multiple data types, automatically and regardless of the difficulties associated with comprehensive causal modelling. In a sense, machine learning methods completely avoid the careful modelling of mechanisms and instead apply a generically complex model to a very large reference dataset to produce a well-performing model with unknown parts. Thus, interpreting the details of a machine learning method in terms of biological correspondents is difficult. More importantly, the training of complex machine learning models typically requires very large volumes of data. On the positive side, developments in high-throughput multi-omics technologies promise the availability of such training data. On the negative side, for the models to be generalizable, we need more than just replicate numbers but also large amounts of data across varying conditions such as from different cell and tissue types. Until a mechanistic model of a cell with sufficient precision to enable the integration of data under a causal model is available, both the utility and the validation of any integration method must be evaluated in terms of their application, for example, by the recovery of the identities of biologically plausible cell types.

## Conclusion and future directions

The integration of single-cell multi-omics data has been implemented in many real data analyses, revealing new biological insights. For example, multi-omics integration has identified the presence of a pro-inflammatory, 'failed repair proximal tubule cell' state in apparently healthy human kidneys[13,95]; it also facilitated the prioritization of genome-wide association study loci through the identification of methylation and gene expression changes that are likely to mediate the development of diabetic kidney disease[96] and has helped identify mechanisms of myofibroblast activation in chronic kidney disease[8].

Ideally, the process of generating data integration models and evaluating the models should itself shed light on the mechanisms of biological processes such as gene regulation. For example, cell identity is traditionally defined by the abundance of specific RNAs or proteins, but integration of these data with other omics datasets could effectively broaden the definition of cell types to other chemical–physical modalities of the cell. In addition, novel relationships across data modalities can be studied with multi-omics data integration. For instance, correlating DNA methylation with gene expression in *cis* might reveal differential functional impacts of methylation of different DNA elements (promoters or gene body). However, regardless of their utility in the modelling of biological processes, data integration often yields more or better resolved inferences than analyses of single datasets alone. For example, the addition of scATAC-seq to scRNA-seq data better distinguishes different segments of proximal tubules in the kidney[36] than does scRNA-seq data alone. Integrated data analyses can also identify underappreciated relationships that might lead to additional applications such as drug target discovery or better causal SNP inference.

Currently available computational methods have generally followed the development of the measurements themselves. The number of available methods that attempt to integrate unmatched data far outweighs the number of methods that attempt to integrate matched data simply because multi-omics measurements have only become widely available in the past couple of years. Methods for integrating cell morphologies[97], perturbations[98], spatial microenvironment[52,99] and subcellular measurements[100] (for example, of organelles) are sparse, as are the corresponding data. However, we expect that methods to integrate these data will rapidly follow the availability of such data. In addition, most current computational methods are built to integrate two modalities; however, with the development of experimental methods that jointly profile three or more modalities, more flexible computational algorithms will be required.

Among the computational tools that lag behind the analytical methods are methods for the visualization of complex multimodal data that interactively connect between different views and ancillary information. Some of the barriers to the development of these tools are the speed and capacity of the computers themselves. Approaches to enable the interactive visualization of extremely large volumes of data in the single-cell field is non-trivial and may eventually require dedicated hardware.

As discussed above, one ideal way to integrate data is in terms of a causal model between the quantitative data and the underlying molecular processes such as cell differentiation, physiology and homeostasis. Conversely, we hope that multimodal data, by providing measurements from multiple aspects of the biology of an organism, could aid the development of such causal models. The era of multi-omics single-cell biology at the scale of millions of cells is just starting and we have no doubt that the data, analytical methods and inferred models will advance our understanding of the kidney by leaps and bounds in years to come.

1. Richardson, S., Tseng, G. C. & Sun, W. Statistical methods in integrative genomics. *Annu. Rev. Stat. Appl.* **3**, 181–209 (2016).
2. Yuan, G.-C. et al. Challenges and emerging directions in single-cell analysis. *Genome Biol.* **18**, 84 (2017).
3. Eberwine, J., Sul, J.-Y., Bartfai, T. & Kim, J. The promise of single-cell sequencing. *Nat. Methods* **11**, 25–27 (2014).
4. Yao, Z. et al. A taxonomy of transcriptomic cell types across the isocortex and hippocampal formation. Preprint at *bioRxiv* https://doi.org/10.1101/2020.03.30.015214 (2020).
5. Cao, J. et al. A human cell atlas of fetal gene expression. *Science* **370**, eaba7721 (2020).
6. Ransick, A. et al. Single-cell profiling reveals sex, lineage, and regional diversity in the mouse kidney. *Dev. Cell* **51**, 399–413.e7 (2019).
   **A comprehensive kidney scRNA-seq atlas with the visualization tool Kidney Cell Explorer**.
7. Kirita, Y., Wu, H., Uchimura, K., Wilson, P. C. & Humphreys, B. D. Cell profiling of mouse acute kidney injury reveals conserved cellular responses to injury. *Proc. Natl Acad. Sci. USA* **117**, 15874–15883 (2020).
8. Kuppe, C. et al. Decoding myofibroblast origins in human kidney fibrosis. *Nature* **589**, 281–286 (2021).
9. Gerhardt, L. M. S. et al. Single-nuclear transcriptomics reveals diversity of proximal tubule cell states in a dynamic response to acute kidney injury. *Proc. Natl Acad. Sci. USA* **118**, e2026684118 (2021).
10. Ma, A., McDermaid, A., Xu, J., Chang, Y. & Ma, Q. Integrative methods and practical challenges for single-cell multi-omics. *Trends Biotechnol.* **38**, 1007–1022 (2020).
    **A comprehensive review of single-cell multi-omics technologies**.
11. Lee, J., Hyeon, D. Y. & Hwang, D. Single-cell multiomics: technologies and data analysis methods. *Exp. Mol. Med.* **52**, 1428–1442 (2020).
12. Sullivan, K. M. & Susztak, K. Unravelling the complex genetics of common kidney diseases: from variants to mechanisms. *Nat. Rev. Nephrol.* **16**, 628–640 (2020).
    **An up-to-date review on efforts to gain further understanding of kidney disease-associated genome-wide association study variants**.
13. Muto, Y. et al. Single cell transcriptional and chromatin accessibility profiling redefine cellular heterogeneity in the adult human kidney. *Nat. Commun.* **12**, 2190 (2021).
14. Miao, Z. et al. Single cell regulatory landscape of the mouse kidney highlights cellular differentiation programs and disease targets. *Nat. Commun.* **12**, 2277 (2021).
15. La Manno, G. et al. RNA velocity of single cells. *Nature* **560**, 494–498 (2018).
16. Gorin, G., Svensson, V. & Pachter, L. Protein velocity and acceleration from single-cell multiomics experiments. *Genome Biol.* **21**, 39 (2020).
17. Stoeckius, M. et al. Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* **14**, 865–868 (2017).
18. Peterson, V. M. et al. Multiplexed quantification of proteins and transcripts in single cells. *Nat. Biotechnol.* **35**, 936–939 (2017).
19. Zhou, Z., Ye, C., Wang, J. & Zhang, N. R. Surface protein imputation from single cell transcriptomes by deep neural networks. *Nat. Commun.* **11**, 651 (2020).
20. Pliner, H. A. et al. Cicero predicts cis-regulatory DNA interactions from single-cell chromatin accessibility data. *Mol. Cell* **71**, 858–871.e8 (2018).
21. Serra, A., Fratello, M., Greco, D. & Tagliaferri, R. Data integration in genomics and systems biology. in *2016 IEEE Congress on Evolutionary Computation (CEC)* 1272–1279 (IEEE, 2016).
22. Hasin, Y., Seldin, M. & Lusis, A. Multi-omics approaches to disease. *Genome Biol.* **18**, 83 (2017).
23. Liu, L. et al. Deconvolution of single-cell multi-omics layers reveals regulatory heterogeneity. *Nat. Commun.* **10**, 470 (2019).
24. Dueck, H. et al. Deep sequencing reveals cell-type-specific patterns of single-cell transcriptome variation. *Genome Biol.* **16**, 122 (2015).
25. Dueck, H. R. et al. Assessing characteristics of RNA amplification methods for single cell RNA sequencing. *BMC Genomics* **17**, 966 (2016).
26. Cao, J. et al. Joint profiling of chromatin accessibility and gene expression in thousands of single cells. *Science* **361**, 1380–1385 (2018).
27. Chen, S., Lake, B. B. & Zhang, K. High-throughput sequencing of the transcriptome and chromatin accessibility in the same cell. *Nat. Biotechnol.* **37**, 1452–1457 (2019).
28. Zhu, C. et al. An ultra high-throughput method for single-cell joint analysis of open chromatin and transcriptome. *Nat. Struct. Mol. Biol.* **26**, 1063–1070 (2019).
29. Ma, S. et al. Chromatin potential identified by shared single cell profiling of RNA and chromatin. Preprint at *bioRxiv* https://doi.org/10.1101/2020.06.17.156943 (2020).
30. Han, S. H., Choi, Y., Kim, J. & Lee, D. Photoactivated selective release of droplets from microwell arrays. *ACS Appl. Mater. Interfaces* **12**, 3936–3944 (2020).
31. Stuart, T. & Satija, R. Integrative single-cell analysis. *Nat. Rev. Genet.* **20**, 257–272 (2019).
32. Li, Y., Ma, L., Wu, D. & Chen, G. Advances in bulk and single-cell multi-omics approaches for systems biology and precision medicine. *Brief. Bioinform.* https://doi.org/10.1093/bib/bbab024 (2021).
33. Sokal, R. R. Distance as a measure of taxonomic similarity. *Syst. Biol.* **10**, 70–79 (1961).
34. Sneath, P. H. A. & Sokal, R. R. *Numerical Taxonomy: The Principles and Practice of Numerical Classification* (WF Freeman, 1973).
35. Wang, X. et al. BREM-SC: a Bayesian random effects mixture model for joint clustering single cell multi-omics data. *Nucleic Acids Res.* **48**, 5814–5824 (2020).
36. Jin, S., Zhang, L. & Nie, Q. scAI: an unsupervised approach for the integrative analysis of parallel single-cell transcriptomic and epigenomic profiles. *Genome Biol.* **21**, 25 (2020).
37. Argelaguet, R. et al. Multi-omics factor analysis — a framework for unsupervised integration of multi-omics data sets. *Mol. Syst. Biol.* **14**, e8124 (2018).
38. Argelaguet, R. et al. MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. *Genome Biol.* **21**, 111 (2020).
39. Gayoso, A. et al. Joint probabilistic modeling of single-cell multi-omic data with totalVI. *Nat. Methods* **18**, 272–282 (2021).
40. Hao, Y. et al. Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573–3587.e29 (2021).
41. Wang, B. et al. Similarity network fusion for aggregating data types on a genomic scale. *Nat. Methods* **11**, 333–337 (2014).
    **This paper introduces the similarity network fusion model, which is widely applied in multi-omics integration**.
42. Kim, H. J., Lin, Y., Geddes, T. A., Yang, J. Y. H. & Yang, P. CiteFuse enables multi-modal analysis of CITE-seq data. *Bioinformatics* **36**, 4137–4143 (2020).
43. Han, X. et al. Construction of a human cell landscape at single-cell level. *Nature* **581**, 303–309 (2020).
44. Packer, J. S. et al. A lineage-resolved molecular atlas of *C. elegans* embryogenesis at single-cell resolution. *Science* **365**, eaax1971 (2019).
    **A single-cell atlas of *Caenorhabditis elegans* with the visualization tool visCello**.
45. Cao, J. et al. The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**, 496–502 (2019).
46. Slavov, N. Single-cell protein analysis by mass spectrometry. *Curr. Opin. Chem. Biol.* **60**, 1–9 (2021).
47. Neumann, E. K., Ellis, J. F., Triplett, A. E., Rubakhin, S. S. & Sweedler, J. V. Lipid analysis of 30000 individual rodent cerebellar cells using high-resolution mass spectrometry. *Anal. Chem.* **91**, 7871–7878 (2019).
48. Zhu, Q. et al. Developmental trajectory of prehematopoietic stem cell formation from endothelium. *Blood* **136**, 845–856 (2020).
49. Wang, C. et al. Integrative analyses of single-cell transcriptome and regulome using MAESTRO. *Genome Biol.* **21**, 198 (2020).
50. Asp, M., Bergenstråhle, J. & Lundeberg, J. Spatially resolved transcriptomes—next generation tools for tissue exploration. *BioEssays* **42**, 1900221 (2020).
51. Zhu, Q., Shah, S., Dries, R., Cai, L. & Yuan, G.-C. Identification of spatially associated subpopulations by combining scRNAseq and sequential fluorescence in situ hybridization data. *Nat. Biotechnol.* **36**, 1183–1190 (2018).
52. Rodriques, S. G. et al. Slide-seq: a scalable technology for measuring genome-wide expression at high spatial resolution. *Science* **363**, 1463–1467 (2019).
53. Andersson, A. et al. Single-cell and spatial transcriptomics enables probabilistic inference of cell type topography. *Commun. Biol.* **3**, 565 (2020).
54. Govek, K. W. et al. Single-cell transcriptomic analysis of mIHC images via antigen mapping. *Sci. Adv.* **7**, eabc5464 (2021).
55. Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427 (2018).
    **This paper introduces the MNN method that became popular in single-cell biology with multiple applications**.
56. Campbell, K. R. et al. clonealign: statistical integration of independent single-cell RNA and DNA sequencing data from human cancers. *Genome Biol.* **20**, 54 (2019).
57. Stuart, T. et al. Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902.e21 (2019).
58. Welch, J. D. et al. Single-cell multi-omic integration compares and contrasts features of brain cell identity. *Cell* **177**, 1873–1887.e17 (2019).
59. Yang, Z. & Michailidis, G. A non-negative matrix factorization method for detecting modules in heterogeneous omics multi-modal data. *Bioinformatics* **32**, 1–8 (2016).
60. Amodio, M. & Krishnaswamy, S. MAGAN: aligning biological manifolds. *Proc. Machine Learn. Res.* **80**, 215–223 (2018).
61. Welch, J. D., Hartemink, A. J. & Prins, J. F. MATCHER: manifold alignment reveals correspondence between single cell transcriptome and epigenome dynamics. *Genome Biol.* **18**, 138 (2017).
62. Liu, J., Huang, Y., Singh, R., Vert, J.-P. & Noble, W. S. in *19th International Workshop on Algorithms in Bioinformatics (WABI 2019)* (eds Huber, K. T. & Gusfield, D.) Vol. 143 10:1–10:13 (Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2019).
63. Cao, K., Bai, X., Hong, Y. & Wan, L. Unsupervised topological alignment for single-cell multi-omics integration. *Bioinformatics* **36**, i48–i56 (2020).

64. Demetci, P., Santorella, R., Sandstede, B., Noble, W. S. & Singh, R. Gromov-Wasserstein optimal transport to align single-cell multi-omics data. Preprint at *bioRxiv* https://doi.org/10.1101/2020.04.28.066787 (2020).

65. Li, X. et al. Deep learning enables accurate clustering with batch effect removal in single-cell RNA-seq analysis. *Nat. Commun.* **11**, 2338 (2020).

66. McInnes, L., Healy, J. & Melville, J. UMAP: uniform manifold approximation and projection for dimension reduction. Preprint at *arxiv* https://arxiv.org/abs/1803.00385 (2020).

67. Moon, K. R. et al. Visualizing structure and transitions in high-dimensional biological data. *Nat. Biotechnol.* **37**, 1482–1492 (2019).

68. Costa, F., Grün, D. & Backofen, R. GraphDDP: a graph-embedding approach to detect differentiation pathways in single-cell-data using prior class knowledge. *Nat. Commun.* **9**, 3685 (2018).

69. Wu, Y. & Zhang, K. Tools for the analysis of high-dimensional single-cell RNA sequencing data. *Nat. Rev. Nephrol.* **16**, 408–421 (2020).
    **A comprehensive review of scRNA-seq data analysis pipelines and computational tools**.

70. Steiniger, S. & Hay, G. J. Free and open source geographic information tools for landscape ecology. *Ecol. Inform.* **4**, 183–195 (2009).

71. Raney, B. J. et al. Track data hubs enable visualization of user-defined genome-wide annotations on the UCSC genome browser. *Bioinformatics* **30**, 1003–1005 (2014).

72. Ou, J. & Zhu, L. J. trackViewer: a bioconductor package for interactive and integrative visualization of multi-omics data. *Nat. Methods* **16**, 453–454 (2019).

73. Snyder, M. P. et al. The human body at cellular resolution: the NIH human biomolecular atlas program. *Nature* **574**, 187–192 (2019).

74. Hillje, R., Pelicci, P. G. & Luzi, L. Cerebro: interactive visualization of scRNA-seq data. *Bioinformatics* **36**, 2311–2313 (2020).

75. Dries, R. et al. Giotto: a toolbox for integrative analysis and visualization of spatial expression data. *Genome Biol.* **22**, 78 (2021).

76. Larsson, A. J. M. et al. Genomic encoding of transcriptional burst kinetics. *Nature* **565**, 251–254 (2019).

77. Chakrabarti, S. et al. Hidden heterogeneity and circadian-controlled cell fate inferred from single cell lineages. *Nat. Commun.* **9**, 5372 (2018).

78. Zhong, L. et al. Single cell transcriptomics identifies a unique adipose lineage cell population that regulates bone marrow environment. *eLife* **9**, e54695 (2020).

79. Lahens, N. F. et al. IVT-seq reveals extreme bias in RNA sequencing. *Genome Biol.* **15**, R86 (2014).

80. Marquina-Sanchez, B. et al. Single-cell RNA-seq with spike-in cells enables accurate quantification of cell-specific drug effects in pancreatic islets. *Genome Biol.* **21**, 106 (2020).

81. Xi, N. M. & Li, J. J. Benchmarking computational doublet-detection methods for single-cell RNA sequencing data. *Cell Syst.* **12**, 176–194.e6 (2021).

82. Franzosa, E. A. et al. Gut microbiome structure and metabolic activity in inflammatory bowel disease. *Nat. Microbiol.* **4**, 293–305 (2019).

83. Tini, G., Marchetti, L., Priami, C. & Scott-Boyer, M.-P. Multi-omics integration — a comparison of unsupervised clustering methodologies. *Brief. Bioinform.* **20**, 1269–1279 (2019).

84. Pierson, E. & Yau, C. ZIFA: dimensionality reduction for zero-inflated single-cell gene expression analysis. *Genome Biol.* **16**, 241 (2015).

85. Kharchenko, P. V., Silberstein, L. & Scadden, D. T. Bayesian approach to single-cell differential expression analysis. *Nat. Methods* **11**, 740–742 (2014).

86. Marinov, G. K. et al. From single-cell to cell-pool transcriptomes: stochasticity in gene expression and RNA splicing. *Genome Res.* **24**, 496–510 (2014).

87. Zhang, L. & Nie, Q. scMC learns biological variation through the alignment of multiple single-cell genomics datasets. *Genome Biol.* **22**, 10 (2021).

88. Fang, R. et al. Comprehensive analysis of single cell ATAC-seq data with SnapATAC. *Nat. Commun.* **12**, 1337 (2021).

89. Velleman, P. F. & Wilkinson, L. Nominal, ordinal, interval, and ratio typologies are misleading. *Am. Stat.* **47**, 65–72 (1993).

90. He, B. et al. Integrating spatial gene expression and breast tumour morphology via deep learning. *Nat. Biomed. Eng.* **4**, 827–834 (2020).

91. Wu, H., Kirita, Y., Donnelly, E. L. & Humphreys, B. D. Advantages of single-nucleus over single-cell RNA sequencing of adult kidney: rare cell types and novel cell states revealed in fibrosis. *J. Am. Soc. Nephrol.* **30**, 23–32 (2019).

92. Cao, J. et al. Comprehensive single-cell transcriptional profiling of a multicellular organism. *Science* **357**, 661–667 (2017).

93. McGinnis, C. S. et al. MULTI-seq: sample multiplexing for single-cell RNA sequencing using lipid-tagged indices. *Nat. Methods* **16**, 619–626 (2019).

94. Yang, K. D. et al. Multi-domain translation between single-cell imaging and sequencing data using autoencoders. *Nat. Commun.* **12**, 31 (2021).

95. Dhillon, P. et al. The nuclear receptor ESRRA protects from kidney disease by coupling metabolism and differentiation. *Cell Metab.* **33**, 379–394.e8 (2021).

96. Sheng, X. et al. Systematic integrated analysis of genetic and epigenetic variation in diabetic kidney disease. *Proc. Natl Acad. Sci. USA* **117**, 29013–29024 (2020).

97. Wu, P.-H. et al. Single-cell morphology encodes metastatic potential. *Sci. Adv.* **6**, eaaw6938 (2020).

98. Dixit, A. et al. Perturb-seq: dissecting molecular circuits with scalable single cell RNA profiling of pooled genetic screens. *Cell* **167**, 1853–1866.e17 (2016).

99. Lindström, N. O. et al. Spatial transcriptional mapping of the human nephrogenic program. Preprint at *bioRxiv* https://doi.org/10.1101/2020.04.27.060749 (2020).

100. Khaladkar, M. et al. Subcellular RNA sequencing reveals broad presence of cytoplasmic intron-sequence retaining transcripts in mouse and rat neurons. *PLoS ONE* **8**, e76194 (2013).
     **The first subcellular RNA sequencing method**.

### RELATED LINKS

**GenitoUrinary Development Molecular Anatomy Project:** https://www.gudmap.org/
**HuBMAP portal:** https://portal.hubmapconsortium.org/
**KidneyCellExplorer:** https://cello.shinyapps.io/kidneycellexplorer/
**ReBuilding a Kidney:** https://www.rebuildingakidney.org/