

## **Tutorial: guidelines for annotating single-cell transcriptomic maps using automated and manual methods**

**Zoe A. Clarke<sup>1,2,\*</sup>, Tallulah S. Andrews<sup>2,3,4\*</sup>, Jawairia Atif<sup>3,4\*</sup>, Delaram Pouyabahar<sup>1,2,\*</sup>,  
Brendan T. Innes<sup>1,2</sup>, Sonya A. MacParland<sup>3,4,5+</sup>, Gary D. Bader<sup>1,2,6,7+</sup>**

1 - Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada

2 - The Donnelly Centre, University of Toronto, Toronto, Ontario, Canada

3 - Ajmera Transplant Centre, Toronto General Hospital Research Institute, Toronto, Ontario, Canada

4 - Department of Immunology, University of Toronto, Toronto, Ontario Canada

5 - Department of Laboratory Medicine and Pathobiology, University of Toronto, Toronto, Ontario, Canada

6 - Department of Computer Science, University of Toronto, Toronto, Ontario, Canada

7 - Lunenfeld-Tanenbaum Research Institute, Toronto, Ontario, Canada

\*equal contribution

+corresponding: [s.macparland@utoronto.ca](mailto:s.macparland@utoronto.ca), [gary.bader@utoronto.ca](mailto:gary.bader@utoronto.ca)

# Abstract

Single-cell transcriptomics can profile thousands of cells in a single experiment and identify novel cell types, states and dynamics in a wide range of tissues and organisms. Standard experimental protocols and analysis workflows have been developed to create single-cell transcriptomic maps from tissues. This Tutorial focuses on how to interpret these data to identify cell types, states and other biologically relevant patterns with the objective of creating an annotated map of cells. We recommend a three-step workflow including automatic cell annotation (wherever possible), manual cell annotation and verification. Frequently encountered challenges are discussed, as well as strategies to address them. Guiding principles and specific recommendations for software tools and resources that can be used for each step are covered and an R notebook is included to help run the recommended workflow. Basic familiarity with computer software is assumed and basic knowledge of programming (e.g. in the R language) is recommended.

# Introduction

Single-cell genomics enables the molecular profiling of thousands of cells in a single experiment<sup>1–3</sup> to create comprehensive maps of cellular heterogeneity in multicellular systems<sup>4,5</sup>. In particular, single-cell RNA sequencing (scRNA-seq) and single-nuclei RNA sequencing (snRNA-seq) can be used to measure single-cell transcriptomes and map novel cell types<sup>6</sup>, states<sup>7</sup> and dynamics<sup>8</sup> in a wide range of tissues and organisms.

Single-cell transcriptomics data are often presented as a two-dimensional “map” organizing cells based on the similarity of their gene expression profiles. Data visualized in this way naturally identifies groups (or “clusters”) of highly similar cells, as well as gradients and other transcript-based patterns. Such artifacts must be interpreted and annotated to define cell types and states to support biological discovery (Figure 1). Standard experimental protocols and analysis workflows detail how to create single-cell transcriptomic maps from tissues<sup>9–12</sup>. Briefly, tissues are dissociated into single cells and profiled using a single-cell transcriptomic technology. Computational analysis is then used to perform quality control filtering on the results (e.g. removing low-quality cells), quantify the expression of each mapped gene in each cell<sup>13</sup>, identify groups of similar cells using a clustering algorithm<sup>14–18</sup>, and visualize all cells in two dimensions using techniques such as tSNE<sup>19</sup> or UMAP<sup>20</sup> to produce an unannotated “single-cell map” image (Box 1)<sup>21</sup>. To interpret this map biologically, it is necessary to determine which cell types or cell states are represented by clusters or other patterns (e.g. gradients) observed in the data. These interpretations can then be labeled on the map, which helps place them in a conceptual framework useful for better understanding tissue biology. This Tutorial offers a guide to the map interpretation and labeling process, starting from clustered data and resulting in a completely annotated single-cell map (Figure 1). The general workflow for annotating cells in scRNA-seq data has three major steps: automatic annotation, manual annotation, and verification (Figure 2). First, automatic annotation uses a predefined set of “marker genes” (i.e. genes that are specifically expressed in a known cell type) or reference single-cell data (i.e. an existing expertly annotated single-cell map) to identify and label individual cells or cell clusters by matching their gene expression patterns (signatures) to those of known cell types. A second major step is manual annotation, which involves studying genes and gene functions specific to each cell cluster or pattern to verify automatic cell annotations and identify novel cell types and states. Finally, verification can confirm the identity and function of select cell types using independent methods, such as new validation experiments.

## Step 1: Automatic cell annotation

Automatic cell annotation is an efficient way to label cells or cell clusters using a computer algorithm and an appropriate set of prior biological knowledge. The general principle is to identify a gene expression signal (pattern, signature) in a single cell or cell cluster that matches a characteristic gene expression signature of a known cell type or state; the cell or cluster is then assigned the respective label. Labels often have an associated confidence score.

There are two major automatic cell annotation approaches. One is to use known marker genes for each of the cell types that are likely to be found in the sample to be annotated (referred

to as “marker-based automatic annotation”). In this case, known relationships between marker genes and cell types are obtained from databases, such as SCSig<sup>22</sup>, PanglaoDB<sup>23</sup>, and CellMarker<sup>24</sup>, or manually from the literature. Then cells or clusters are labeled according to the marker genes they characteristically express. The second approach is to compare single-cell RNA-seq data to be annotated (the ‘query’ data set) to an existing, similar, expertly annotated scRNA-seq data set (the ‘reference’ data set), and transfer the label from a reference cell or cluster to a sufficiently similar one in the query (referred to as “reference-based automatic annotation”). Reference single-cell data are obtained from sources such as Gene Expression Omnibus (GEO)<sup>25</sup>, the Single Cell Expression Atlas<sup>26</sup> or cell atlas projects<sup>27,28</sup>.

Automatic cell annotation methods can be applied to individual cells (either before or after clustering) or to clusters of cells, which occurs only after clustering the cells. In the case of annotating clusters, the gene expression profile for each cluster is determined by averaging the expression profiles of all cells within the cluster. Annotating individual cells is ideal, as this reduces the chance of missing important differences between cells. However, some scRNA-seq experimental data are based on low numbers of transcript reads per cell, so there may be insufficient data for cell-based annotation to function correctly, making clustered data sets easier to work with. Annotating clusters is faster, as there are fewer clusters than cells to process; it can also be more accurate than the single-cell approach, considering it is based on more reliable expression level estimates averaged across all cells in a cluster. However, not all cells can be easily grouped into clusters, especially for dynamic systems like developing tissues<sup>29</sup> or tissues that contain gene expression gradients<sup>30,31</sup>.

A major challenge with automatic cell annotation is that many cell types do not have well-characterized gene expression signatures, resulting in incomplete or inaccurate labeling for some cells. Automated methods typically work better for major cell types and may not be able to effectively distinguish subtypes. Thus, automatic cell annotation is useful to quickly identify known cell types and highlight unknown cell types for further exploration. The main caveats and recommendations for automatic cell annotation are summarised in Table 1.

## Marker-based automatic annotation

Marker-based automatic annotation labels cells or cell clusters based on the characteristic expression of known marker genes. To be successful, the marker gene or gene set (a collection of marker genes) should be specifically and consistently expressed in a given cell, cluster or class of cells (e.g. immune cells). Markers are readily available for well-characterized organisms and cell types (e.g. human PBMC samples<sup>32</sup>). Marker-based automatic annotation works well once a relevant and sufficiently large set of marker genes is collected<sup>33</sup>.

To label individual cells, one of the most reliable marker-based annotation tools is Semi-supervised Category Identification and Assignment (SCINA)<sup>34</sup>. SCINA assumes each marker follows a bimodal gene expression distribution, where one peak corresponds to cells from the associated cell type and the other peak contains the rest of the cells in the experiment. A cell of a particular type is assumed to have expression in the upper part of this distribution for all the markers of that cell type, consequently requiring markers provided as input to SCINA to be specific to only one cell type. AUCell<sup>35</sup> is another good marker-based labeling method that classifies

individual cells or clusters. AUCell ranks the genes in each cell by decreasing expression value, and cells are labeled according to their most active (highly expressed) marker gene sets. AUCell works best with cell types that have a sufficiently large set of marker genes such that multiple markers are detected in each cell. It has the advantage of scoring a whole set of marker genes at once, which may increase sensitivity over methods that examine each marker gene independently.

To label whole clusters, Gene Set Variation Analysis<sup>36</sup> (GSVA) has been benchmarked to be fast and reliable<sup>37</sup>. GSVA works similarly to AUCell - given a database of marker gene sets, it identifies sets that are enriched in the gene expression profile of a cluster. The GSVA software has a practical advantage that it can annotate all clusters in one operation.

Marker-based automatic cell annotation methods often have the advantage that they will only assign labels to cells associated with known markers and other cells will remain unlabeled<sup>33</sup>. However, this depends on the specific tool and the parameters used; see Table 2 and Supplementary table 1 for details on which tools have the option to leave cells unlabelled. A disadvantage of these tools is that markers are not easily accessible for all cell types.

## Reference-based automatic cell annotation

Reference-based cell annotation is based on the concept of “guilt-by-association”, whereby a cell or cluster label in the reference data is transferred to an unlabeled cell or cluster in the query data with a similar gene expression profile. Consequently, this approach is only possible if high-quality and relevant annotated reference single-cell data are available. Studying the original clustering and annotation steps performed on the reference data can help determine its quality, and ensure that errors in the reference will not be propagated to new data. Tissue-specific reference data can be obtained from public databases (e.g. the Gene Expression Omnibus<sup>25</sup> or the Expression Atlas<sup>26</sup>) or large cell atlas projects (e.g. the Human Cell Atlas<sup>27</sup>, the Tabula Muris or Mouse Cell Atlas<sup>5</sup>, or others<sup>4,28,38–40</sup>), although the required associated cell annotations are not always easily available. These atlases typically contain hundreds of thousands of cells and dozens of different annotated cell types.

scmap<sup>41</sup> is one of the best performing tools for reference-based automatic cell or cluster annotation, in terms of both accuracy of assigned labels and avoiding incorrect labeling of novel cell types<sup>33</sup>. Other tools for reference-based automatic annotation include SingleCellNet<sup>42</sup> and SingleR<sup>43</sup>. SingleCellNet has high accuracy when all cell types are well represented in the reference data but with low accuracy if the reference data are incomplete or represent a poor match<sup>33</sup>. The main advantage to SingleR is that a reasonable, general reference data set is included with the tool, but this may not perform as well as a reference specifically matched to the query data set. An alternative to using specific software packages for reference-based cell annotation is to train a machine learning tool, such as a support vector machine (SVM)<sup>44</sup> or random forest classifier<sup>45</sup>, on selected reference data. This model can then be applied to classify cells or clusters as specific cell types in novel data. These methods can outperform any of the prepackaged automatic-annotation software tools<sup>33</sup> but require substantial computational expertise to use.

Another approach to reference-based cell annotation is to integrate a query data set with a reference data set using an integration algorithm, enabling clusters to be identified that span

both data sets (Box 2). The reference labels can then be transferred to within-cluster query data cells. This approach supports the identification of novel cell types, distinct cell types, and gradients in cell state, but can be computationally expensive to run and additional problems, such as over-integration, may be encountered.

## Refining automatic annotation

Benchmarking studies show variable performance of automatic annotation tools, depending on the data set and distinctiveness of the gene expression profiles of the cell types to be annotated<sup>33,37</sup>. For instance, distinguishing T-cells from B-cells is relatively straightforward, but automatic tools sometimes cannot accurately distinguish CD8+ cytotoxic T cells from natural killer (NK) cells (Figure 3). Thus, we recommend applying multiple, complementary annotation tools with multiple available marker gene databases to a single data set.

When applying multiple cell annotation methods to a data set, cells or clusters will acquire multiple, sometimes conflicting, cell-type labels. A set of annotations on a cell or cluster can easily be resolved to a single label if all labels are in agreement. If conflicts exist, most tools provide label confidence scores that can be used to identify a single high-scoring label. However, confidence scores are not standardized between tools, so they are often not comparable. Conflicts can also be resolved via a majority-rule approach, which selects the most frequent label (Figure 4), or percent agreement between methods. If no label can be confidently decided, the cells or cluster must be manually annotated.

Conflicting annotations within a cluster may reflect important information about that cluster, such as whether it contains cell subtypes. However, if subtypes cannot be clearly defined, a more general cell-type annotation may be more appropriate. For example, if a cluster is annotated as regulatory T-cells, naive T-cells and helper T-cells by different methods, it may be most appropriate to assign the general label of “T-cells”. In this case, the original clustering parameters should be altered to better capture cell subtypes (see section “The impact of experimental and analysis parameters on annotation”).

If the conflicting annotations are not subtypes of the same cell type, then the cluster may represent an intermediate cell state or gene-expression gradient. As many automatic annotation tools assume discrete cell types, they often assign clusters or cells within a larger gradient to a well-defined endpoint. However, gradients often contain cells of various phenotypes, so multiple methods may assign the same cell to different ends of the gradient. Recommendations for handling gradients are discussed in the section “Annotating cell states and gradients”. Alternatively, a conflicting label on a cell could indicate that the cell is actually a doublet; a scenario in which two or more cells of different types are captured by the same cell-barcode. This case can be detected using doublet-finding methods<sup>46–48</sup>.

Most automatic annotation tools are designed to annotate individual cells (Table 2, Supplementary Table 1). Advantages of this approach are the ability to identify insufficiently resolved cell types and cellular gradients, as well as the independence to choose clustering resolution, feature selection, and dimensionality reduction parameters. Interestingly, the resulting annotations can be used to inform these analysis choices. For example, cell annotation could help optimize the clustering process to result in one cluster per cell type.

Finally, a cluster may have a novel cell identity that is absent from the reference data. Often, this results in widely varying results from automatic annotation methods or insufficient confidence for any tool to assign any label. In such situations, manual annotation must be performed.

## Step 2: Expert manual cell annotation

Although automated cell annotation methods are convenient and systematic, they require an appropriate reference database and do not always result in high confidence annotations. When these methods result in lower confidence, conflicting or absent cell labels, expert manual annotation is required. In manual cell annotation, cells are manually examined for clues to their function using a variety of resources, following the same principles as marker-based automatic annotation. Manual annotation usually operates at the cluster level for convenience, but rare cells can be individually examined. Expert manual annotation is typically regarded as the gold-standard method for annotating cells; however it is slow, labour-intensive, and can be subjective.

If automatic annotation has not been performed, marker-based annotation should first be manually applied. Usually, each known marker gene is individually visualized on the 2D projected data map (Box 1) to create a "gene expression overlay" plot (Figure 5). The entire list of markers may also be simultaneously visualized across clusters as a heat map (Extended data figure 1) or dot plot (Figure 6). A dot plot is more informative than a heat map, as it can communicate mean detected gene expression levels and the proportion of cells in a cluster in which each gene is detected, whereas a heat map only typically describes average gene expression levels per cluster. If many marker genes for a known cell type are highly expressed across cells in a cluster, this is often sufficient support for it to be labeled as that cell type. Easy-to-use software like the free Loupe Cell Browser for 10x Genomics scRNA-seq data supports this visualization and analysis process. Challenges in this approach are that well-known markers are often too few in number to completely annotate a scRNA-seq data set, and some well-known markers may not be as specific within a scRNA-seq data set as expected. Often, additional markers must be manually found via searching the literature and mining existing single-cell transcriptomic data for gene expression signatures related to the query data set. Master transcription factors that drive cell fate<sup>49,50</sup> often make better gene expression markers than cell-surface proteins that are commonly used to classify cell populations with methods like flow cytometry<sup>5,39</sup>, presumably because mRNA and protein levels may not be strongly correlated<sup>51</sup>. Furthermore, there may not be any single distinguishing gene expression marker; in which case, multiple genes must be used together to distinguish a cell type from others in the data.

The ideal primary source for cell-defining genes is a single-cell atlas from a relevant organism, organ and disease context. In the absence of this, gene expression markers can be collected from bulk RNA-seq data from purified cell populations isolated from the same tissue source<sup>52</sup>. Given that protein expression may correlate with mRNA expression, protein expression markers can be gathered and used as potential gene expression markers<sup>53</sup> from published evidence of staining patterns within the tissue (i.e. using immunohistochemistry or immunofluorescence), flow cytometry and western blots. Integrating markers from independent sources can be challenging due to conflicts between lists. For instance, PanglaoDB<sup>23</sup> contains 220 markers for B-cells and CellMarker<sup>24</sup> contains 1426 markers, yet only 66 are shared. If

species-specific data are scarce, then data can be transferred by orthology from model organisms (Box 3) or other models (e.g. *in vitro* cell culture or organoids).

Ideally, each cluster will uniquely express the markers of one cell type. However, in some instances a cluster may not express markers of any known cell type; conversely, it may express markers of more than one cell type. Clusters that express markers of more than one cell type may represent doublets. Typically, such clusters will be very small compared to the clusters of true single cells, and they may express more genes than single cells. There are a variety of doublet detection tools that can help determine if a cluster is composed of doublets<sup>46–48</sup>. If a cluster does not express markers of any known cell type, it may contain poor quality cells or represent a novel cell type.

Once cell-type information from known markers is exhausted, cells that have not been confidently annotated must be manually examined, cluster by cluster. Potential novel markers are identified by computing differential expression between a cluster and all other cells<sup>54–56</sup> (Figure 6, Extended data figure 1). All marker genes are then manually researched to find functional information that may help identify the cell type of the cluster they are associated with. Pathway enrichment analysis should also be applied to each cluster to identify cluster-specific pathways using standard workflows<sup>57</sup> and tools, like GSVA<sup>36</sup> or ssGSEA<sup>58</sup>. Pathway enrichment analysis simultaneously scores multiple functionally related genes for gene expression activity within a cluster at once, and can be more sensitive than individual gene-based analysis.

Some cells may be challenging to annotate, including novel cell types, which can be described based on the function of genes they express. Furthermore, it can be particularly difficult to differentiate between tissue-resident cells (e.g. tissue-resident macrophages) and non-tissue-resident cells (e.g. monocytes circulating in the blood) of the same overall type. One approach to identify tissue-resident cells is to modify the experimental design to remove passenger cells from the tissue in question with a perfusion step. However, the number and types of cells removed by flushing will depend on the specific tissue and protocol. In situations where this flushing is not possible, one may profile peripheral blood mononuclear cells (PBMCs) from the same individual and then subtract those cell signatures from the tissue map using the cell annotation methods mentioned above.

Ultimately, when annotating cell types, it is prudent to use standard nomenclature, such as from the Cell Ontology (CL), which is a hierarchically-organized controlled vocabulary of cell types and subtypes<sup>59–61</sup>. This enables maps to be more easily integrated across studies.

## Annotating cell states and gradients

When analyzing and characterizing novel cell types, it is important to determine whether they represent a stable cell type or contain multiple cell states. The definitions of cell type and state are not yet standardized, but a stable cell type may be expected to have homogeneous gene expression across a cluster and be compact in a 2D projection plot (Box 1), whereas cell gradients appear as a spread-out string of cells and cell states (e.g. cell cycle state) (Figure 6). Expression gradients indicate continuous differences that are present in the cell population, which could represent states like the cell cycle<sup>62</sup>, immune activation<sup>63</sup>, spatial patterning<sup>64</sup> or transient developmental stages<sup>65,66</sup>. Care must be taken to distinguish biologically meaningful cell states and experimental batch effects, which can manifest in a similar way (Box 2).



Annotating the intermediate stages of a gradient is often difficult, as these regions rarely express unique marker genes. It is often easier to label the ends of a gradient and then characterize intermediate stages using the order of specific genes that mark these ends as increasing or decreasing across the gradient. Extracting the cells in the gradient and performing PCA on them is often a useful visualization for gradients, as it preserves the large-scale distances between cells (Figure 6). There are currently no automated gradient annotation methods, so they must be manually annotated, making use of known structure and cell-type transitions relevant to the particular experiment<sup>7,59</sup>.

Similarly, homogenous or similar cell states or cell types are often difficult to annotate because they share many of the same marker genes (Figure 3). For instance, when annotating T-cells within a tissue sample, it is common for all the T-cell subtypes to exhibit common T-cell markers; the subtype-specific markers are hidden within or below the general T-cell signal. In this case, it is often useful to subcluster the population, or to test specifically each subpopulation against the other related clusters to identify the subtype-specific markers. Very fine distinctions between highly similar cell types may not be visible transcriptionally and may only be visible in other genomic layers, such as chromatin state (ATAC-seq, DNA methylation).

## Step 3: Annotation verification

The above tools and approaches can provide confident cell-type labels for scRNA-seq data. However, due to the various challenges discussed above, it is important to confirm cell annotation labels using independent methods, such as statistical methods<sup>67</sup> or by consulting an expert. Furthermore, as mRNA measurements only partially define cell type and function, important conclusions about novel cell types must be experimentally validated.

As an example, cell-type labels of tissue-resident immune cells can be refined using T-cell receptor (TCR)<sup>68</sup> and B-cell receptor (BCR)<sup>69</sup> clonotyping, to examine the transcriptional signature of T- and B-cells as stratified by the TCRs and BCRs that they express. For instance, mucosal-associated invariant T-cells express the marker genes *SLCA4A10* and *KLRB1*,<sup>70</sup> which can be identified in a scRNA-seq experiment, as well as the known semi-invariant TCRs that are found in MAIT cells (*TRAV1-2/TRAJ12/20/3*), which can be revealed by TCR clonotyping. In addition, identifying the B-cell receptor repertoire within single-cell data sets enables annotation of naive vs. mature B-cells. Naive B-cells express both IgM and IgD heavy chains, whereas mature B-cells, which have undergone antibody class switching by V(D)J recombination, express IgG, IgA or IgE heavy chains. Other traditional methods to increase cell annotation confidence include *in vitro* functional assays such as cytokine secretion, proliferative capacity, and cytotoxic potential measures, imaging experiments (such as fluorescence in situ hybridization (FISH)<sup>71</sup> of source tissue samples)<sup>72</sup>, and single-cell qPCR to verify the co-expression of a novel combination of marker genes in a larger number of samples<sup>73,74</sup>. Complementary single-cell genomic methods are also useful, such as Cellular Indexing of Transcriptomes and Epitopes by sequencing (CITE-seq)<sup>51</sup>, which simultaneously immunophenotypes cell surface proteins and measures scRNA-seq, single-cell ATAC-seq, which maps chromatin state, and spatial transcriptomics, which combines cell imaging and scRNA-seq to measure spatial transcript patterns and cell morphology in one experiment<sup>75</sup>.

In the context of tumour biology, mutations are important to distinguish cancer and normal cell types. Genetic alterations like single nucleotide variants (SNV) and copy number variants (CNV) can be detected in single-cell data using tools developed for bulk RNA-seq data<sup>76</sup>, despite challenges with sequencing coverage. CNV inference methods identify consistently up or down expression values relative to a reference across large numbers of genomically contiguous genes to call amplification or deletion events, respectively. HoneyBADGER<sup>77</sup> and CaSpER<sup>78</sup> methods predict CNVs using other cells in the data as a reference, whereas InferCNV<sup>79,80</sup> uses a given set of normal cells.

## Experimental design considerations

### The impact of experimental and analysis parameters on annotation

Cell-type annotation quality is affected by many data analysis pipeline parameters, such as data filtering and data quality settings, and the selected clustering resolution. Quality control filtering often involves removing cells from the data set where the marker genes are highly enriched in mitochondrial, heat shock, or other stress-response genes<sup>81–83</sup>, but this must be balanced to retain important biological signals that should be kept and annotated<sup>84</sup>.

Choosing an appropriate clustering resolution is critical for accurate cell annotation. If the clustering resolution is too low, rare cell types may be merged with larger clusters or related cell subtypes may be merged with each other. If the clustering resolution is too high, a single cell type may be split across multiple clusters with few unique markers that are often a result of experimental noise rather than distinguishing biological function. To identify rare cell types, it may be necessary to use a feature selection tool that specifically identifies markers of rare cell types (e.g. GiniClust<sup>85</sup>) before clustering the cells. However, this can lead to over clustering of data sets that do not contain rare cell types. If cell-based annotation identifies multiple cell types within a cluster, then increasing the clustering resolution or subsetting the cluster and rerunning the clustering on the resulting smaller group of cells to create a zoomed-in map can help isolate these unique cell groups. Tools like scClustViz<sup>86</sup>, Seurat<sup>21</sup> and clustree<sup>87</sup> help select an appropriate clustering resolution.

Cell-based automatic annotation tools are often useful to choose an appropriate clustering resolution, as the results are independent of the clustering pipeline. Thus, the clustering parameters can be tuned to identify clusters that correctly segregate cells annotated to different types. Alternatively, the presence of cluster-specific differentially expressed genes can be used to tune the clustering parameters either by gradually increasing or decreasing the resolution until the maximum number of clusters that still exhibit unique differentially expressed genes is identified.

In some cases, the original gene expression matrix generated by droplet-based technologies can be contaminated by cell-free, or ‘ambient’, mRNA within the cell suspension. Frequently, the ambient RNA is derived from one or more cell types that are sensitive to the tissue dissociation or cell-handling steps of the scRNA-seq experiment and break apart from the rest.

As a result, markers of the contaminating cell types may be spread to all other cell clusters, which will interfere with marker determination. It is possible to estimate and correct the background contamination using methods such as SoupX<sup>88</sup>, which looks for non-specific expression of cell-type markers, or CellBender<sup>89</sup>, which uses machine learning to learn and correct cell expression profiles. However, care is needed to avoid over- or under-correcting the data.

## Workflow recommendations

The preferred starting approach for transcriptomic cell map annotation depends on the level of computational skills of the annotator. We recommend starting with automatic annotation because it is fast and reproducible, thus efficient for large data sets with many samples. Automated methods require programming, database and data science skills to operate (mainly using R or python programming languages). A little programming knowledge goes a long way, as many recommended software packages are well documented and easy to use with basic programming knowledge. We strongly recommend anyone working regularly with single-cell genomics data to learn programming. R programming is a good language to start with, due to its prevalent use in single-cell genomics and ease of use. This recommendation may change, as an increasing number of point-and-click tools are being developed that package automated methods into easy-to-use workflows<sup>90–92</sup>. A second recommendation is to use a powerful computer with lots of memory (e.g. 128GB RAM), as current analysis and visualization methods load all data in memory for processing.

If needed, the map can be completely annotated manually by investigating gene expression patterns of cells and associated gene functions using point-and-click software (e.g. Loupe Browser, GSEA, Cerebro<sup>93</sup>) without programming skills. This process is easier for those knowledgeable about the biology and markers of the cells in the sample, and is sufficient for many projects, but is time consuming, especially if it must be repeated for multiple analysis parameters, such as different cluster resolutions. Even if automatic cell annotation is used, some manual annotation is usually needed due to the incompleteness of known marker databases and single-cell atlases that automated methods depend on.

Manual annotation should begin by identifying major well-known cell types, indicated by clearly defined, discrete cell clusters, as these are easiest to work with. Each cluster presents its own challenge; in the same experiment, one cluster could be annotated easily using a single well-known marker and another may require iterations of data preprocessing pipeline optimization to accurately annotate. More challenging aspects of the map, such as cell subtypes, gradients, highly homogeneous data or poorly defined clusters, can then be progressively annotated. Sometimes, it is useful to split the data into broad subsets (e.g. immune, endothelial, tumour) and apply our recommended workflow on each one. Consider each label on the map a scientific claim that a functionally distinct biological entity (e.g. cell type) exists, and that this must be supported by evidence.

In addition, not all tools are applicable to all data sets; it is important to consider the availability of reliable known markers, high quality reference data sets, or if there is sufficient diversity in a sample to detect differentially expressed genes before applying methods that rely on that information. Utilizing an approach based on known marker genes when there aren't any repeatedly reproduced markers for the cell type of interest may lead to false conclusions.

Likewise, annotating cells or cell clusters with a poor quality or inaccurately labeled reference data set will likely lead to the propagation of incorrect cell type identifications. As a final example, calculating differentially expressed genes on a largely homogenous sample will typically result in a list of markers that are false positives, or genes subject primarily to experimental or technical noise that are unrelated to the actual biology.

## Concluding remarks

Although the field of single-cell genomics is rapidly advancing and new technologies are being developed that will improve our ability to interpret, annotate and validate single-cell maps, we expect the overall workflow described here to remain valid over time. We expect major improvements in automatic annotation due to rapidly growing reference atlases, improvements to resources like the Cell Ontology<sup>59</sup> and improved data set integration algorithms. These methods will also need to scale up to much larger data set sizes with millions of cells<sup>94,95</sup>. New experimental technologies are being developed to measure more molecular details about each cell, including multi-omics technologies (e.g. mRNA, ATAC-seq<sup>96</sup>, methylation<sup>97</sup>, surface proteins<sup>98</sup>) that can measure multiple types of information about individual cells, and these are expected to greatly improve our ability to understand multicellular systems. For instance, epigenetic information will help define stem cell subtypes that are not detectable using transcriptomics<sup>99</sup>. Data sets acquired from millions of cells across hundreds of patients will create computational challenges related to efficient analysis and annotation, requiring analysis to be performed on high-performance computing or cloud computing systems. In addition, meta-analyses across many single-cell maps will more clearly identify cell-type markers (e.g. macrophage or endothelial) across tissues and states (e.g. inflammation). We also expect the focus of map interpretation to gradually shift to comparisons across disease, age, or other conditions, as the number of samples per study increases.

# Boxes

## Box 1: Visualizing single-cell data in a 2D projection

A scRNA-seq data set is typically visualized as a two-dimensional (2D) scatter plot where cells (points) with similar transcriptomes are placed near each other. This 2D representation is projected from a higher dimensional space where each cell is described by the expression of thousands of genes, each of which is considered a separate dimension. The three most popular projection methods used for scRNA-seq data are t-Distributed Stochastic Neighbor Embedding (t-SNE)<sup>19</sup>, Uniform Manifold Approximation and Projection (UMAP)<sup>20</sup> and Principal Component Analysis (PCA)<sup>100</sup>.

t-SNE (Figure 6c) is a non-linear projection that preserves local groups of similar cells, while equalizing the density of cells within each group<sup>101</sup>. The scale of a “local group” is controlled by the “perplexity” parameter, with higher values creating larger local groups<sup>102</sup>. t-SNE effectively visualizes distinct robust clusters, making it easy to observe discrete cell types; however, global relationships between cell types are not maintained, and thus cluster-to-cluster relationships cannot be inferred and may be misleading. Cell subtypes can be combined into one large cluster or split into distinct plot regions depending on the perplexity<sup>101</sup>.

UMAP (Extended data figure 1) is a non-linear projection method that differentiates discrete cell clusters<sup>20</sup>. UMAP is typically regarded as better for visualizing global relationships and gradients than t-SNE, although these differences are likely due to default parameters<sup>103</sup>. UMAP is often less computationally intensive to run than t-SNE<sup>104</sup>.

PCA (Figure 6b) performs a linear transformation of normalized and scaled scRNA-seq data, to identify independent principal components (PCs) that capture major axes of variation in the data, which could represent biological factors, like cell types and states, or technical factors<sup>105</sup>. PCs are ranked in decreasing order of variance and typically the first two PCs are used to visualize the data, but more can be considered to detect more subtle expression patterns between cells<sup>100</sup>. PCA can be useful for visualizing cell gradients and states.

Although these methods visually group similar cells and help visualize clusters, they do not define clusters and, therefore, are not clustering algorithms. Cell clustering algorithm output is typically visualized as colours on the plot and these colours may or may not correspond to patterns observed in the 2D plot.

## Box 2: Correcting confounding factors

ScRNA-seq data contains a mix of biological (e.g. cell types, states, age, sex, and disease condition) as well as technical (e.g. batch effects) factors<sup>106</sup>. It is important to correct for undesired (i.e. “confounding”) factors while maintaining biological signals of interest. Confounding factors can either be regressed out of the data, or adjusted for when integrating data from multiple samples (Figure 7). Batch effects can be identified when cells from different batches form distinct stripes within groups or completely separate groups in a 2D visualization (Box 1). Harmony, mnnCorrect, Seurat v3, and LIGER are among the top-performing scRNA-seq integration or batch-correction tools<sup>107</sup>.

Harmony<sup>108</sup> iteratively merges data sets represented by top PCs, which are then used to cluster cells. Each cell is iteratively adjusted based on an estimated correction vector to shift it closer to the centre of its cluster until convergence. Mutual nearest neighbour (MNN) approaches, such as mnnCorrect/FastMNN<sup>109</sup> or Seurat v3<sup>21</sup>, identify the most similar cells (MNNs), called “anchors”, across data sets that are used to estimate and correct the cell-type-specific batch effects. LIGER<sup>110</sup> identifies shared (common biology) and unique (biological or technical) factors between data sets using non-negative matrix factorization. LIGER is recommended when specific cell types appear to be present in some of the data sets and missing in others<sup>107</sup>. Integration methods can suffer from overcorrection, where different cell types are merged, or undercorrection, when resulting clusters contain cells from only one input data set. Multiple integration methods may need to be evaluated to find a balance that best represents the data.

### Box 3: Cell annotation across species

Sometimes the best reference single-cell map to use for cell annotation is from a different organism. To use such a reference for cell annotation, genes from the query organism must be mapped by orthology to genes from the reference, using databases such as Ensembl<sup>111</sup> or EggNOG<sup>112</sup>, or tools such as OrthoFinder<sup>113</sup>, before being input to data integration or marker-based annotation methods. Typically, one-to-one orthologs are used,<sup>114,115</sup> which better ensures the conservation of function, although it is possible to use one-to-many and many-to-many relationships to increase gene coverage. The latter can be accomplished by grouping paralogs to create artificial gene groupings called ‘meta-genes’<sup>116</sup>. If homologous genes are unavailable, genes from each species can be aggregated into pathways or “biological process activities” (BPAs), which are compared across species to improve sensitivity of cross-species mapping<sup>117</sup>. If the query species genome has not been sequenced, RNA transcripts can be assembled *de novo* from the entire pool of RNA-seq reads from all cells, which are then used to quantify gene expression and identify orthologs. Evolutionarily close reference species should be chosen; otherwise, integration may not be able to map similar cell types for annotation transfer.

# Figures

## Figure 1: An annotated single-cell transcriptomic map

A completely annotated single-cell transcriptomic map of the human liver from data generated using scRNA-seq applied to five human liver samples (8,444 cells) reported in MacParland et al.<sup>115</sup> and visualized using a t-SNE plot.

## Figure 2: Cell annotation workflow

The recommended cell annotation process is composed of three major steps: automatic cell annotation, manual cell annotation, and verification. The scRNA-seq data typically enter the workflow as a clustered gene-by-cell matrix, which is visualized using a dimensionality reduction method. An automatic cell annotation method is used to annotate cells either by comparison of the data with annotated reference data (e.g. a single-cell atlas) or using known marker genes indicative of a specific cell type. Manual annotation confirms or provides further detail for annotated cells or clusters, or identifies the cell type of unlabeled clusters. Cell type can be manually inferred using a combination of marker genes, pathway analysis, and differentially expressed genes with known functional information. Cell annotations are often verified using independent sources, such as new validation experiments or comparison to complementary data, such as spatial transcriptomics data.

## Figure 3: Automatic annotation results depend on marker genes used

Peripheral blood mononuclear cells (PBMCs) from a 10x Genomics 3' sequencing protocol<sup>118</sup> (68,579 cells) were automatically annotated with SCINA. The user provides lists of marker genes associated with all expected cell types, and SCINA assigns cell-type labels to individual cells based on the expression levels of the marker genes. SCINA was provided with (a) the top 20 marker genes from each previously annotated scRNA-seq cell cluster along with their associated cell type<sup>118</sup>, and (b) Literature-derived PBMC cell-type markers gathered by Diaz-Mejia et al.<sup>37</sup>. (c) Sankey plot coloured by the cell-type labels found in (a) compares the assigned cell-type labels between those from SCINA annotations (a) and (b). Changes in label assignment demonstrate the variability of automatic annotation based on the marker genes chosen by the user to represent each cell type or subtype. The PBMC data set is available from [https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/fresh\\_68k\\_pbmc\\_donor\\_a](https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/fresh_68k_pbmc_donor_a).

## Figure 4: Refining cluster labels from automatic annotation

Raw scmap-cluster<sup>41</sup> annotations provide every cluster with a unique label based on the identity assigned to the majority of cells in each cluster. These labels can be refined by aggregating identical labels, and label confidence can be assessed by viewing the proportion of cell identities in each cluster ("majority rule"). (a,c) are t-SNE maps of the liver transcriptomic map constructed from the data set reported in MacParland et al. 2018<sup>115</sup> (8,444 cells). (a) The cells are coloured by the clusters identified by Seurat using default parameters. Each cluster was assigned a cell-type label by scmap-cluster. Two clusters (labeled 4 and 10) are identified as different immune cells with labels of various levels of support, further described in (b). (b) scmap-cluster identifies

cluster labels as well as individual cell labels. Cluster 4 has approximately 83% of cells identified as inflammatory macrophages. Cluster 10 consists almost entirely of non-inflammatory macrophages with some cells unassigned. (c) Consensus cluster-level annotations are determined after aggregating identical cluster labels (e.g. clusters 1 and 3 from (a)) and using a majority rule (selecting a cluster label from the majority of cells assigned to a cell type).<sup>115</sup>

### Figure 5: Visualizing well-known markers

Human liver marker genes are visualized as gene expression data overlaid on a t-SNE plot. The various markers represent: (a) immune cells and immune cell subtypes that can be easily identified by well-known markers, and (b) hepatocytes and inferred hepatocyte subtypes that are difficult to identify because the markers are not as well characterized and may vary across individual samples. Gene expression data are reported in MacParland et al.<sup>115</sup>

### Figure 6: How to identify and visualize a cell-type gradient

Gradients across cell types can be identified by: (a) finding marker genes that are expressed at varying levels across clusters, (b) observing clear gradients of marker gene expression between cell types in PCA gene expression overlay plots, and (c) identifying closely related cell types that are split across clusters in a t-SNE plot. (a) Dot plots of marker genes distinguish clusters of scRNA-seq data from young and old mouse brains reported in Ximerakis et al.<sup>119</sup> (37,069 cells). (b,c) These data can be visualized in two dimensions with gene expression of *Milt* overlaid across cells; this occurs in a gradient across neuronal-restricted precursors (NRP), immature neurons (ImmN), neuroendocrine cells (NendC), and mature neurons (mNEUR).

### Figure 7: Batch correction

Three peripheral blood mononuclear cell (PBMC) samples were assayed with the 10X platform using different library construction protocols: 5' (green, 7,726 cells), 3' V1 (orange, 5,419 cells) and 3' V2 (blue, 7,726 cells). (a, b) UMAP diagrams showing clusters annotated by the SCINA using PBMC markers collected by Diaz-Mejia et al.<sup>37</sup> (c, d, e, f) Bar graphs indicating the proportions of cells per cluster. The left column (a, c, e) shows the data merged without batch correction, and the right column (b, d, f) shows the data integrated using the Harmony batch correction method. Before Harmony, cells group by experimental protocol, clusters rarely contain cells from multiple experiments, and multiple clusters of the same cell type exist (e.g. there is one B-cell cluster for each experimental protocol, circled in e). After Harmony, cells group by cell type, clusters contain cells from various protocols, and fewer clusters share a cell identity.

### Extended Data Figure 1: Heat map and UMAP visualizations

An extension of Figure 6 incorporating: (a) the visualization of marker genes for identified cell types as a heat map, and (b) a *Milt11* expression overlaid on a UMAP plot. *Milt11* is expressed at various levels across clusters, suggesting a cell-type gradient across neuronal-restricted precursors (NRP), immature neurons (ImmN), neuroendocrine cells (NendC), and mature neurons (mNEUR). Both plots are generated from scRNA-seq data from young and old mouse brains<sup>119</sup>.



**Data availability:**

The data used to generate this Tutorial are openly available at the following sources:

- The sequence data used to generate Figures 1 and 4 are available from MacParland et al.<sup>115</sup> through the NCBI GEO accession GSE115469. The analysed data from which the map was directly created can also be accessed interactively as the R package HumanLiver from <https://github.com/BaderLab/HumanLiver>.
- The sequence data used to generate Figures 3 and 7 are available from 10X Genomics and can be downloaded from: <https://support.10xgenomics.com/single-cell-gene-expression/datasets>.
- The sequence data used to generate Figure 6 are available through the NCBI GEO accession GSE129788, as reported by Ximerakis et al.<sup>119</sup>. The analysed data can be accessed interactively at <http://shiny.baderlab.org/AgingMouseBrain/>.
- The human bulk RNA-seq data used to generate the reference data set in the accompanying R code (<https://github.com/BaderLab/CellAnnotationTutorial> and <https://codeocean.com/capsule/d67541eb-43f8-4cae-a258-5ef0069e5301/>) are available from the Database of Immune Cell Expression (DICE) and can be downloaded in R through the package “celldex”<sup>43</sup> by the command `DatabasImmuneCellExpressionData()`.
- The query data set used in the accompanying R code is available from 10X Genomics and can be downloaded from: [https://cf.10xgenomics.com/samples/cell-exp/1.1.0/pbmc3k/pbmc3k\\_filtered\\_gene\\_bc\\_matrices.tar.gz](https://cf.10xgenomics.com/samples/cell-exp/1.1.0/pbmc3k/pbmc3k_filtered_gene_bc_matrices.tar.gz)
- The collection of PBMC marker genes used in the accompanying R code is available from Diaz-Mejia JJ et al.<sup>37</sup> with read data from NCBI Sequence Read Archive (SRA) accession number SRX1723926. The supplementary data from the Diaz-Mejia paper can be accessed from: <https://zenodo.org/record/3369934/#.X3CGN5NKjGI>.

**Code availability:**

An R script that implements the main workflow described in this proposal is available at <https://github.com/BaderLab/CellAnnotationTutorial>, and a version-controlled copy is available through Code Ocean at <https://codeocean.com/capsule/d67541eb-43f8-4cae-a258-5ef0069e5301/>.

**Acknowledgements:**

We acknowledge Shirley Hui, Daniel Stueckmann and Ronald Xie for their help in reviewing the tutorial.

**Funding:**

This project has been made possible in part by grant number CZF2019-002429 from the Chan Zuckerberg Initiative DAF, an advised fund of Silicon Valley Community Foundation.

**Competing Interests:**

No competing interests have been identified.

**Author Contributions:**

Z.A.C., T.S.A., J.A., D.P., and S.A.M. initially wrote different sections of the document, and collaboratively joined them together. Z.A.C., T.S.A., J.A., D.P., and B.T.I. designed the figures. Z.A.C. and G.D.B. organized the process and undertook major edits. All authors read and approved the final manuscript.

**References:**

1. Sasagawa, Y., Hayashi, T. & Nikaido, I. Strategies for Converting RNA to Amplifiable cDNA for Single-Cell RNA Sequencing Methods. *Adv. Exp. Med. Biol.* 1129, 1–17 (2019).
2. Klein, A. M. et al. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell* 161, 1187–1201 (2015).
3. Macosko, E. Z. et al. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* 161, 1202–1214 (2015).
4. Han, X. et al. Construction of a human cell landscape at single-cell level. *Nature* 581, 303–309 (2020).
5. Tabula Muris Consortium et al. Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. *Nature* 562, 367–372 (2018).
6. Grün, D. et al. Single-cell messenger RNA sequencing reveals rare intestinal cell types. *Nature* 525, 251–255 (2015).
7. Xia, B. & Yanai, I. A periodic table of cell types. *Development* 146, (2019).
8. Schiebinger, G. et al. Optimal-Transport Analysis of Single-Cell Gene Expression Identifies Developmental Trajectories in Reprogramming. *Cell* 176, 928-943.e22 (2019).
9. Ziegenhain, C. et al. Comparative Analysis of Single-Cell RNA Sequencing Methods. *Mol.*

- Cell 65, 631-643.e4 (2017).
10. Lafzi, A., Moutinho, C., Picelli, S. & Heyn, H. Tutorial: guidelines for the experimental design of single-cell RNA sequencing studies. *Nat. Protoc.* 13, 2742–2757 (2018).
  11. Hwang, B., Lee, J. H. & Bang, D. Single-cell RNA sequencing technologies and bioinformatics pipelines. *Exp. Mol. Med.* 50, 96 (2018).
  12. Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol. Syst. Biol.* 15, e8746 (2019).
  13. Henry, G. H., Mathews, J. A. & Malladi, V. S. BICF Cellranger count Analysis Workflow. Zenodo (2019). doi:10.5281/zenodo.3373749
  14. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* 36, 411–420 (2018).
  15. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* 19, 15 (2018).
  16. Duò, A., Robinson, M. D. & Soneson, C. A systematic performance evaluation of clustering methods for single-cell RNA-seq data. [version 2; peer review: 2 approved]. *F1000Res.* 7, 1141 (2018).
  17. Kiselev, V. Y., Andrews, T. S. & Hemberg, M. Challenges in unsupervised clustering of single-cell RNA-seq data. *Nat. Rev. Genet.* 20, 273–282 (2019).
  18. Menon, V. Clustering single cells: a review of approaches on high-and low-depth single-cell RNA-seq data. *Brief. Funct. Genomics* 17, 240–245 (2018).
  19. Maaten, L. van der & Hinton, G. Visualizing Data using t-SNE. *Journal of Machine Learning Research* (2008).
  20. Becht, E. et al. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* 37, 38–44 (2018).
  21. Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* 177, 1888-1902.e21

- (2019).
22. GSEA | MSigDB | SCSig collection: Signatures of Single Cell Identities. at  
<[https://www.gsea-msigdb.org/gsea/msigdb/supplementary\\_genesets.jsp#SCSig](https://www.gsea-msigdb.org/gsea/msigdb/supplementary_genesets.jsp#SCSig)>
  23. Franzén, O., Gan, L.-M. & Björkegren, J. L. M. PanglaoDB: a web server for exploration of mouse and human single-cell RNA sequencing data. *Database (Oxford)* 2019, (2019).
  24. Zhang, X. et al. CellMarker: a manually curated resource of cell markers in human and mouse. *Nucleic Acids Res.* 47, D721–D728 (2019).
  25. Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res.* 30, 207–210 (2002).
  26. Papatheodorou, I. et al. Expression Atlas update: from tissues to single cells. *Nucleic Acids Res.* 48, D77–D83 (2020).
  27. Regev, A. et al. The human cell atlas. *elife* 6, (2017).
  28. HuBMAP Consortium. The human body at cellular resolution: the NIH Human Biomolecular Atlas Program. *Nature* 574, 187–192 (2019).
  29. Yuzwa, S. A. et al. Developmental Emergence of Adult Neural Stem Cells as Revealed by Single-Cell Transcriptional Profiling. *Cell Rep.* 21, 3970–3986 (2017).
  30. Kurial, S. N. T. & Willenbring, H. Transcriptomic traces of adult human liver progenitor cells. *Hepatology* 71, 1504–1507 (2020).
  31. Stanley, G., Gokce, O., Malenka, R. C., Südhof, T. C. & Quake, S. R. Continuous and discrete neuron types of the adult murine striatum. *Neuron* 105, 688-699.e8 (2020).
  32. Satpathy, A. Curated, multi-omic, ML-driven single-cell atlas for characterizing the human immune system across disease states. *The Journal of Immunology* 204, 159.11-159.11 (2020).
  33. Abdelaal, T. et al. A comparison of automatic cell identification methods for single-cell RNA sequencing data. *Genome Biol.* 20, 194 (2019).

34. Zhang, Z. et al. SCINA: A Semi-Supervised Subtyping Algorithm of Single Cells and Bulk Samples. *Genes (Basel)* 10, (2019).
35. Aibar, S. et al. SCENIC: single-cell regulatory network inference and clustering. *Nat. Methods* 14, 1083–1086 (2017).
36. Hänzelmann, S., Castelo, R. & Guinney, J. GSVA: gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics* 14, 7 (2013).
37. Diaz-Mejia, J. J. et al. Evaluation of methods to assign cell type labels to cell clusters from single-cell RNA-sequencing data. [version 3; peer review: 2 approved, 1 approved with reservations]. *F1000Res.* 8, (2019).
38. Cao, J. et al. Comprehensive single-cell transcriptional profiling of a multicellular organism. *Science* 357, 661–667 (2017).
39. Han, X. et al. Mapping the Mouse Cell Atlas by Microwell-Seq. *Cell* 172, 1091-1107.e17 (2018).
40. Regev, A. et al. The Human Cell Atlas White Paper. Apollo - University of Cambridge Repository (2019). doi:10.17863/cam.40032
41. Kiselev, V. Y., Yiu, A. & Hemberg, M. scmap: projection of single-cell RNA-seq data across data sets. *Nat. Methods* 15, 359–362 (2018).
42. Tan, Y. & Cahan, P. SingleCellNet: A Computational Tool to Classify Single Cell RNA-Seq Data Across Platforms and Across Species. *Cell Syst.* 9, 207-213.e2 (2019).
43. Aran, D. et al. Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. *Nat. Immunol.* 20, 163–172 (2019).
44. Cortes, C. & Vapnik, V. Support-vector networks. *Mach. Learn.* 20, 273–297 (1995).
45. Breiman, L. Random forests. *Machine learning* (2001).
46. Wolock, S. L., Lopez, R. & Klein, A. M. Scrublet: Computational Identification of Cell Doublets in Single-Cell Transcriptomic Data. *Cell Syst.* 8, 281-291.e9 (2019).
47. Bais, A. S. & Kostka, D. scds: computational annotation of doublets in single-cell RNA

- sequencing data. *Bioinformatics* 36, 1150–1158 (2020).
48. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet Detection in Single-Cell RNA Sequencing Data Using Artificial Nearest Neighbors. *Cell Syst.* 8, 329–337.e4 (2019).
  49. Lambert, S. A. et al. The human transcription factors. *Cell* 172, 650–665 (2018).
  50. Niwa, H. The principles that govern transcription factor network functions in stem cells. *Development* 145, (2018).
  51. Stoeckius, M. et al. Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* 14, 865–868 (2017).
  52. Clark, J. Z. et al. Representation and relative abundance of cell-type selective markers in whole-kidney RNA-Seq data. *Kidney Int.* 95, 787–796 (2019).
  53. Uhlén, M. et al. Proteomics. Tissue-based map of the human proteome. *Science* 347, 1260419 (2015).
  54. Dal Molin, A., Baruzzo, G. & Di Camillo, B. Single-Cell RNA-Sequencing: Assessment of Differential Expression Analysis Methods. *Front. Genet.* 8, 62 (2017).
  55. Adossa, N. A., Schauser, L., Gregersen, V. G. & Elo, L. L. Feature extraction approach in single-cell gene expression profiling for cell-type marker identification. *BioRxiv* (2019). doi:10.1101/686659
  56. Soneson, C. & Robinson, M. D. Bias, robustness and scalability in single-cell differential expression analysis. *Nat. Methods* 15, 255–261 (2018).
  57. Reimand, J. et al. Pathway enrichment analysis and visualization of omics data using g:Profiler, GSEA, Cytoscape and EnrichmentMap. *Nat. Protoc.* 14, 482–517 (2019).
  58. Barbie, D. A. et al. Systematic RNA interference reveals that oncogenic KRAS-driven cancers require TBK1. *Nature* 462, 108–112 (2009).
  59. Diehl, A. D. et al. The Cell Ontology 2016: enhanced content, modularization, and ontology interoperability. *J. Biomed. Semantics* 7, 44 (2016).

60. Meehan, T. F. et al. Logical development of the cell ontology. *BMC Bioinformatics* 12, 6 (2011).
61. Aeevermann, B. D. et al. Cell type discovery using single-cell transcriptomics: implications for ontological representation. *Hum. Mol. Genet.* 27, R40–R47 (2018).
62. Hsiao, C. J. et al. Characterizing and inferring quantitative cell cycle phase in single-cell RNA-seq data analysis. *Genome Res.* 30, 611–621 (2020).
63. Azizi, E. et al. Single-Cell Map of Diverse Immune Phenotypes in the Breast Tumor Microenvironment. *Cell* 174, 1293-1308.e36 (2018).
64. Adler, M., Korem Kohanim, Y., Tendler, A., Mayo, A. & Alon, U. Continuum of Gene-Expression Profiles Provides Spatial Division of Labor within a Differentiated Cell Type. *Cell Syst.* 8, 43-52.e5 (2019).
65. Liu, S. & Trapnell, C. Single-cell transcriptome sequencing: recent advances and remaining challenges. [version 1; peer review: 2 approved]. *F1000Res.* 5, (2016).
66. Schumacher, L. J. Neural crest migration with continuous cell states. *J. Theor. Biol.* 481, 84–90 (2019).
67. Chung, N. C. Statistical significance of cluster membership for unsupervised evaluation of cell identities. *Bioinformatics* 36, 3107–3114 (2020).
68. Rosati, E. et al. Overview of methodologies for T-cell receptor repertoire analysis. *BMC Biotechnol.* 17, 61 (2017).
69. Setliff, I. et al. High-Throughput Mapping of B Cell Receptor Sequences to Antigen Specificity. *Cell* 179, 1636-1646.e15 (2019).
70. Park, D. et al. Differences in the molecular signatures of mucosal-associated invariant T cells and conventional T cells. *Sci. Rep.* 9, 7094 (2019).
71. Moter, A. & Göbel, U. B. Fluorescence in situ hybridization (FISH) for direct visualization of microorganisms. *J. Microbiol. Methods* 41, 85–112 (2000).
72. Ren, X. et al. Reconstruction of cell spatial organization from single-cell RNA sequencing

- data based on ligand-receptor mediated self-assembly. *Cell Res.* 30, 763–778 (2020).
73. Porter, J. R., Telford, W. G. & Batchelor, E. Single-cell Gene Expression Profiling Using FACS and qPCR with Internal Standards. *J. Vis. Exp.* (2017). doi:10.3791/55219
  74. Wu, A. R. et al. Quantitative assessment of single-cell RNA-sequencing methods. *Nat. Methods* 11, 41–46 (2014).
  75. Chen, K. H., Boettiger, A. N., Moffitt, J. R., Wang, S. & Zhuang, X. RNA imaging. Spatially resolved, highly multiplexed RNA profiling in single cells. *Science* 348, aaa6090 (2015).
  76. Liu, F. et al. Systematic comparative analysis of single-nucleotide variant detection methods from single-cell RNA sequencing data. *Genome Biol.* 20, 242 (2019).
  77. Fan, J. et al. Linking transcriptional and genetic tumor heterogeneity through allele analysis of single-cell RNA-seq data. *Genome Res.* 28, 1217–1227 (2018).
  78. Serin Harmanci, A., Harmanci, A. O. & Zhou, X. CaSpER identifies and visualizes CNV events by integrative analysis of single-cell or bulk RNA-sequencing data. *Nat. Commun.* 11, 89 (2020).
  79. Tirosh, I. et al. Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. *Science* 352, 189–196 (2016).
  80. Tickle, T., Gc Ti, Brown, M. & Haas, B. inferCNV of the Trinity CTAT Project. (2019).
  81. AlJanahi, A. A., Danielsen, M. & Dunbar, C. E. An Introduction to the Analysis of Single-Cell RNA-Sequencing Data. *Mol. Ther. Methods Clin. Dev.* 10, 189–196 (2018).
  82. van den Brink, S. C. et al. Single-cell sequencing reveals dissociation-induced gene expression in tissue subpopulations. *Nat. Methods* 14, 935–936 (2017).
  83. Zhao, Q. et al. A mitochondrial specific stress response in mammalian cells. *EMBO J.* 21, 4411–4419 (2002).
  84. Guantes, R. et al. Global variability in gene expression and alternative splicing is modulated by mitochondrial content. *Genome Res.* 25, 633–644 (2015).
  85. Jiang, L., Chen, H., Pinello, L. & Yuan, G.-C. GiniClust: detecting rare cell types from



- single-cell gene expression data with Gini index. *Genome Biol.* 17, 144 (2016).
86. Innes, B. T. & Bader, G. D. scClustViz - Single-cell RNAseq cluster assessment and visualization. [version 2; peer review: 2 approved]. *F1000Res.* 7, (2018).
  87. Zappia, L. & Oshlack, A. Clustering trees: a visualization for evaluating clusterings at multiple resolutions. *Gigascience* 7, (2018).
  88. Young, M. D. & Behjati, S. SoupX removes ambient RNA contamination from droplet based single cell RNA sequencing data. *BioRxiv* (2018). doi:10.1101/303727
  89. Fleming, S. J., Marioni, J. C. & Babadi, M. CellBender remove-background: a deep generative model for unsupervised removal of background noise from scRNA-seq datasets. *BioRxiv* (2019). doi:10.1101/791699
  90. Mohanraj, S. et al. Crescent: cancer single cell expression toolkit. *Nucleic Acids Res.* 48, W372–W379 (2020).
  91. David, F. P. A., Litovchenko, M., Deplancke, B. & Gardeux, V. ASAP 2020 update: an open, scalable and interactive web-based portal for (single-cell) omics analyses. *Nucleic Acids Res.* 48, W403–W414 (2020).
  92. Franzén, O. & Björkegren, J. L. M. alona: a web server for single-cell RNA-seq analysis. *Bioinformatics* 36, 3910–3912 (2020).
  93. Hillje, R., Pelicci, P. G. & Luzi, L. Cerebro: interactive visualization of scRNA-seq data. *Bioinformatics* 36, 2311–2313 (2020).
  94. Zhang, A. W. et al. Probabilistic cell-type assignment of single-cell RNA-seq for tumor microenvironment profiling. *Nat. Methods* 16, 1007–1015 (2019).
  95. Miao, Z. et al. Putative cell type discovery from single-cell gene expression data. *Nat. Methods* 17, 621–628 (2020).
  96. Buenrostro, J. D., Wu, B., Chang, H. Y. & Greenleaf, W. J. ATAC-seq: A Method for Assaying Chromatin Accessibility Genome-Wide. *Curr. Protoc. Mol. Biol.* 109, 21.29.1-21.29.9 (2015).

97. Angermueller, C. et al. Parallel single-cell sequencing links transcriptional and epigenetic heterogeneity. *Nat. Methods* 13, 229–232 (2016).
98. Baron, M. & Yanai, I. New skin for the old RNA-Seq ceremony: the age of single-cell multi-omics. *Genome Biol.* 18, 159 (2017).
99. Guilhamon, P. et al. Chromatin blueprint of glioblastoma stem cells reveals common drug candidates for distinct subtypes. *BioRxiv* (2018). doi:10.1101/370726
100. Wold, S., Esbensen, K. & Geladi, P. Principal component analysis. *Chemometrics and Intelligent Laboratory Systems* 2, 37–52 (1987).
101. Kobak, D. & Berens, P. The art of using t-SNE for single-cell transcriptomics. *Nat. Commun.* 10, 5416 (2019).
102. Halladin-Dąbrowska, A., Kania, A. & Kopeć, D. The t-SNE Algorithm as a Tool to Improve the Quality of Reference Data Used in Accurate Mapping of Heterogeneous Non-Forest Vegetation. *Remote Sens (Basel)* 12, 39 (2019).
103. Kobak, D. & Linderman, G. C. UMAP does not preserve global structure any better than t-SNE when using the same initialization. *BioRxiv* (2019). doi:10.1101/2019.12.19.877522
104. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv* (2018).
105. Ringnér, M. What is principal component analysis? *Nat. Biotechnol.* 26, 303–304 (2008).
106. Hicks, S. C., Teng, M. & Irizarry, R. A. On the widespread and critical impact of systematic bias and batch effects in single-cell RNA-Seq data. *BioRxiv* (2015). doi:10.1101/025528
107. Tran, H. T. N. et al. A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol.* 21, 12 (2020).
108. Korsunsky, I. et al. Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat. Methods* 16, 1289–1296 (2019).
109. 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).
110. Welch, J. D. et al. Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell Identity. *Cell* 177, 1873-1887.e17 (2019).
  111. Clamp, M. et al. Ensembl 2002: accommodating comparative genomics. *Nucleic Acids Res.* 31, 38–42 (2003).
  112. Huerta-Cepas, J. et al. eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.* 47, D309–D314 (2019).
  113. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biol.* 20, 238 (2019).
  114. Hodge, R. D. et al. Conserved cell types with divergent features between human and mouse cortex. *BioRxiv* (2018). doi:10.1101/384826
  115. MacParland, S. A. et al. Single cell RNA sequencing of human liver reveals distinct intrahepatic macrophage populations. *Nat. Commun.* 9, 4383 (2018).
  116. Geirsdottir, L. et al. Cross-Species Single-Cell Analysis Reveals Divergence of the Primate Microglia Program. *Cell* 181, 746 (2020).
  117. Ding, H., Blair, A., Yang, Y. & Stuart, J. M. Biological process activity transformation of single cell gene expression for cross-species alignment. *Nat. Commun.* 10, 4899 (2019).
  118. Zheng, G. X. Y. et al. Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* 8, 14049 (2017).
  119. Ximerakis, M. et al. Single-cell transcriptomic profiling of the aging mouse brain. *Nat. Neurosci.* 22, 1696–1708 (2019).
  120. Pedregosa, F. et al. Scikit-learn: Machine Learning in Python. *J Mach Learn Res* 12, 2825–2830 (2011).
  121. Van de Sande, B. et al. A scalable SCENIC workflow for single-cell gene regulatory network analysis. *Nat. Protoc.* 15, 2247–2276 (2020).

122. Subramanian, A., Kuehn, H., Gould, J., Tamayo, P. & Mesirov, J. P. GSEA-P: a desktop application for Gene Set Enrichment Analysis. *Bioinformatics* 23, 3251–3253 (2007).
123. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* 33, 495–502 (2015).

# Tables

Table 1: Comparison of the caveats and recommendations for different approaches to cell annotation

Stage of Analysis	Aspect of analysis	Potential Caveats	Recommendation
	<b>All methods</b>	<b>Fast, but not effective for poorly characterized cells.</b>	<b>Use manual annotation for poorly characterised cells.</b>
	Annotating clusters	May miss important differences between cells.	Use automatic annotation of clusters to get a general idea of cell type and then refine labels manually. Also use multiple cluster-based methods and compare results.
	Annotating individual cells	Ideal, but requires high reads per cell.	Experiments with low reads per cell require cluster-based annotation.
	Marker-based annotation methods	Marker genes not easily accessible for all cell types; may result in conflicting or absent cell labels.	Requires expert knowledge to curate more extensive marker lists.
	Reference-based annotation methods	Perform poorly with incomplete or poorly-matched reference data, which may result in conflicting or absent cell labels.	Use well-matched reference data or marker-based methods if such data are unavailable.
		Often requires batch correction, which may reduce the accuracy of results.	Analyse the reference data for strong biological signals. Use a good experimental protocol that will prevail over batch effects.
		Mistakes in reference data get carried over to results.	Analyse reference data for potential errors before using.
	Comparing results from different automatic annotation methods	Results may not agree with each other.	Compare confidence scores of respective labels, and consider label agreement (majority rule); resolve conflicts using manual annotation.
			Consider the possibility of cell subtypes, new cell types, or gradients and cell states.
Expert manual cell annotation	<b>All methods</b>	<b>Slow, labour-intensive.</b>	<b>Whenever possible, begin with automatic annotation to determine general cell labels.</b>
		<b>Subjective.</b>	<b>Work with an expert; consider multiple cell-type conclusions.</b>
	Marker-based annotation	Cell types not distinguishable by a single marker.	Use multiple markers for each cell type.

		Known markers not distinguishing cell types.	Curate larger lists of markers from the literature, additional experiments or experts.
		Conflicting marker gene sets between sources.	Select a marker gene set that best represents the biological signal being looked for in the data (e.g. if looking for cell subtypes, use more extensive gene sets than what is used for general cell-type annotation).

Table 2: Summary of referenced annotation tools

Tool	Type	Language	Resolution	Approach	Allows "None"	Notes
singleCell Net <sup>42</sup>	Reference - based	R	Single cells	Relative-expression gene-pairs + Random Forest	Yes, but rarely does so even when it should <sup>33</sup>	10X-100X slower than other methods. High accuracy.
scmap-cluster <sup>41</sup>	Reference - based	R	Single cells	Consistent correlations	Yes	Fastest method available. Balances false-positives and false-negatives.. Includes web-interface for use with a large pre-built reference or custom reference set.
scmap-cell <sup>41</sup>	Reference - based	R	Single cells	Approximate nearest neighbours	Yes	Assigns individual cells to nearest neighbour cells in reference; allows mapping cell trajectories. Fast and scalable.
singleR <sup>43</sup>	Reference - based	R	Single cells	Hierarchical clustering, Spearman correlations	No	Includes a large marker reference. Does not scale to data sets of 10,000 cells or more. Includes web-interface with marker database
Scikit-learn <sup>120</sup>	Reference - based	python	Multiple possible	k-nearest neighbours (KNN), Support vector machine (SVM), Random forest (RF), Nearest mean classifier (NMC), Linear discriminant analysis (LDA)	(optional)	Expertise required for correct design and appropriate training of classifier while avoiding over-training.
AUCell <sup>121</sup>	Marker-based	R	Single cells	Area Under the Curve to estimate marker gene set enrichment	Yes	Due to low detection rates at the level of single cells, it requires many markers for every cell-type.
SCINA <sup>34</sup>	Marker-based	R	Single cells	Expectation-Maximization, Gaussian mixture model	(optional)	Simultaneously clusters and annotates cells. Robust to the inclusion of incorrect marker genes.
GSEA/ GSVA <sup>36,122</sup>	Marker-based	R/Java	Clusters of cells	Enrichment test,	Yes	Marker gene lists must be reformatted in GMT format.

						Markers must all be differentially expressed in the same direction in the cluster.
Harmony <sup>108</sup>	Integration (Box 2)	R	Single cells	Iterative clustering and adjustment	Yes	Integrates only lower-dimensional projection of the data. Seamlessly integrated into Seurat pipeline. May over-correct data.
Seurat-CCA <sup>123</sup>	Integration (Box 2)	R	Single cells	MNN-anchors + Canonical Correlation Analysis	Yes	Accuracy depends on the accuracy of MNN-anchors, which are automatically-identified corresponding cells across data sets.
mnnCorrect <sup>109</sup>	Integration (Box 2)	R	Single cells	MNN-pairs + SVD	Yes	Accuracy depends on the accuracy of MNN-pairs (cells matched between data sets). Referred to in Box 2.
LIGER <sup>110</sup>	Integration (Box 2)	R	Single cells	Non-negative matrix factorization	Yes	Allows interpretation of data-set-specific and shared factors of variation. Referred to in Box 2.

Figure 1

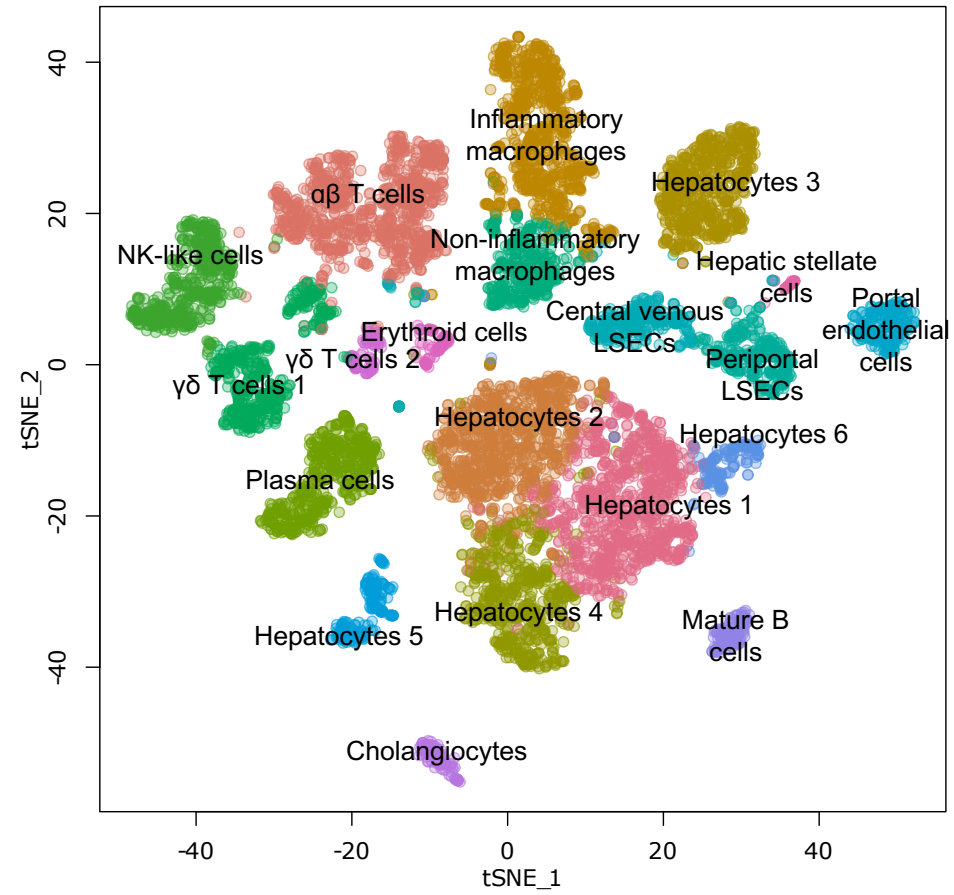




Figure 2

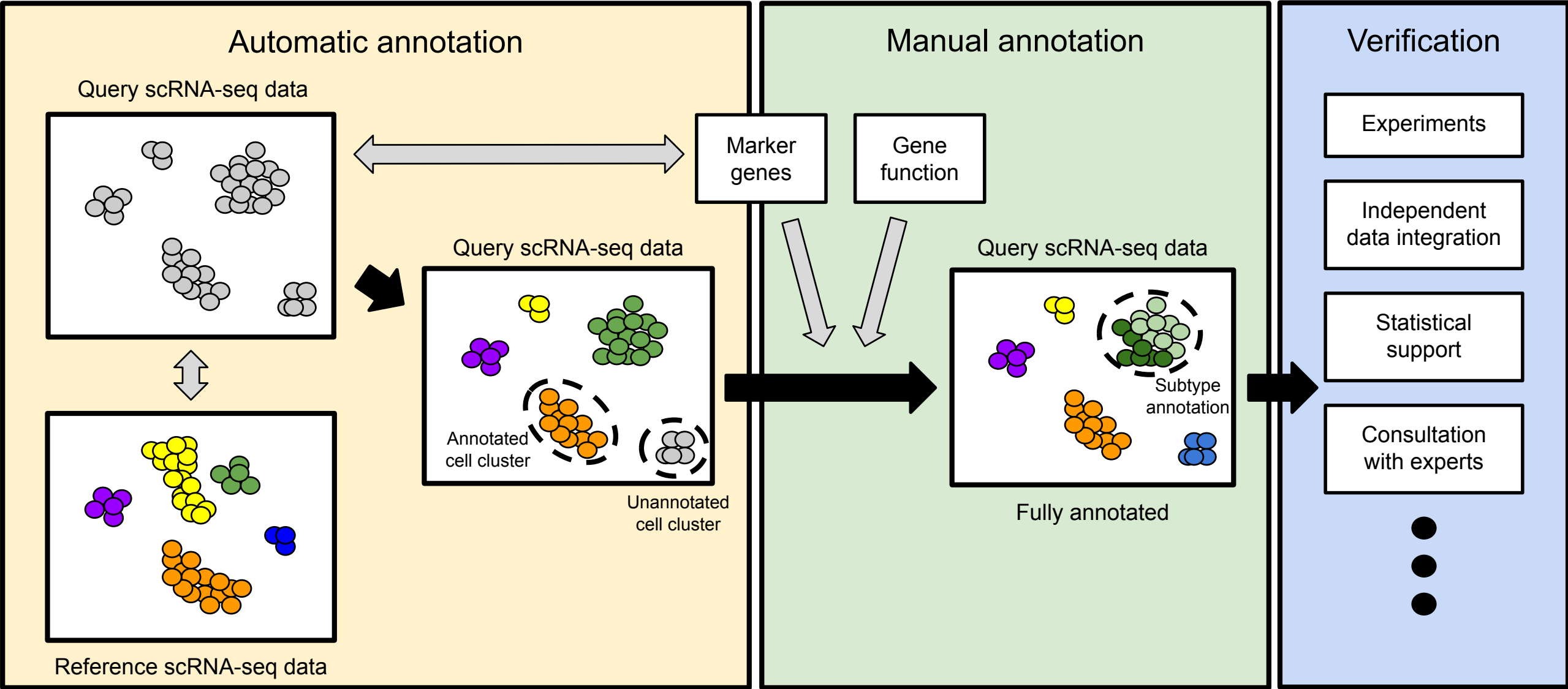
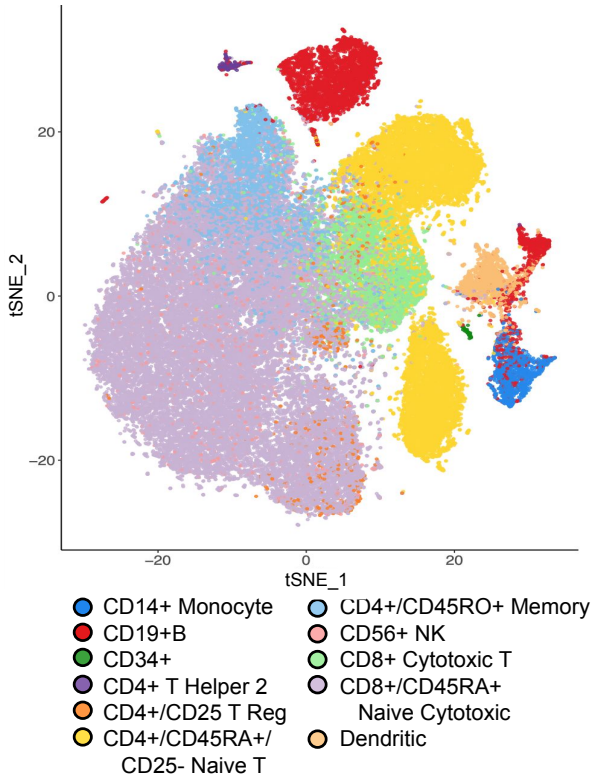


Figure 3

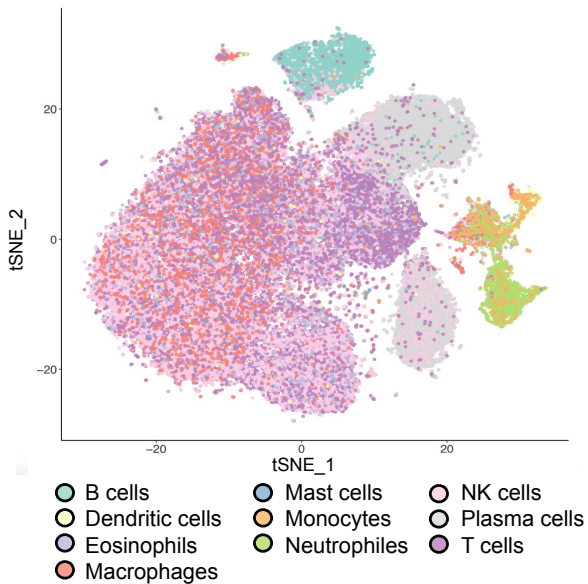
a

SCINA-based cell annotations  
from scRNA-seq  
cluster-derived markers



b

SCINA annotations  
with literature-derived  
PBMC markers



c

Cell annotations from  
SCINA scRNA-seq (a)

Cell annotations from  
SCINA literature (b)

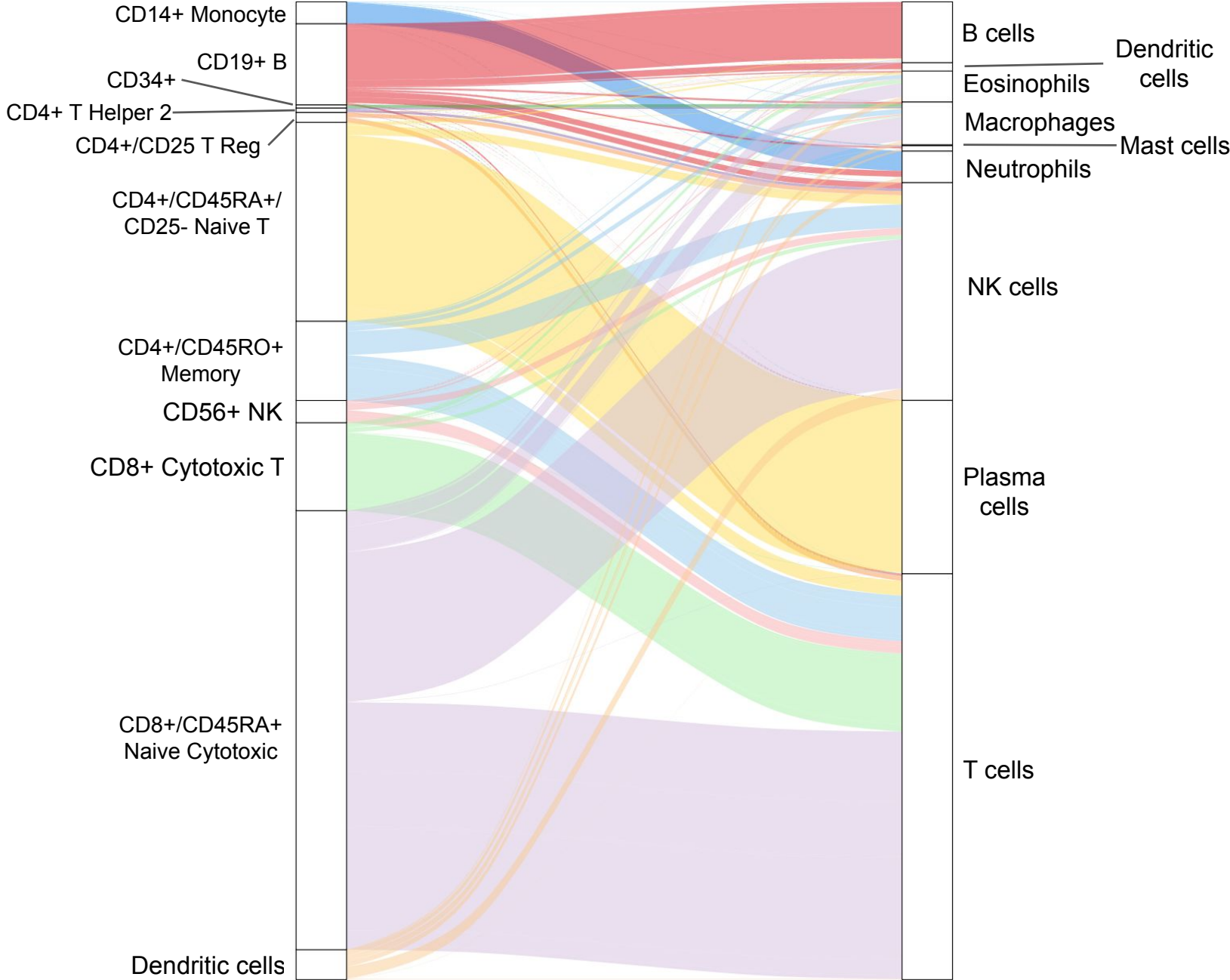


Figure 4

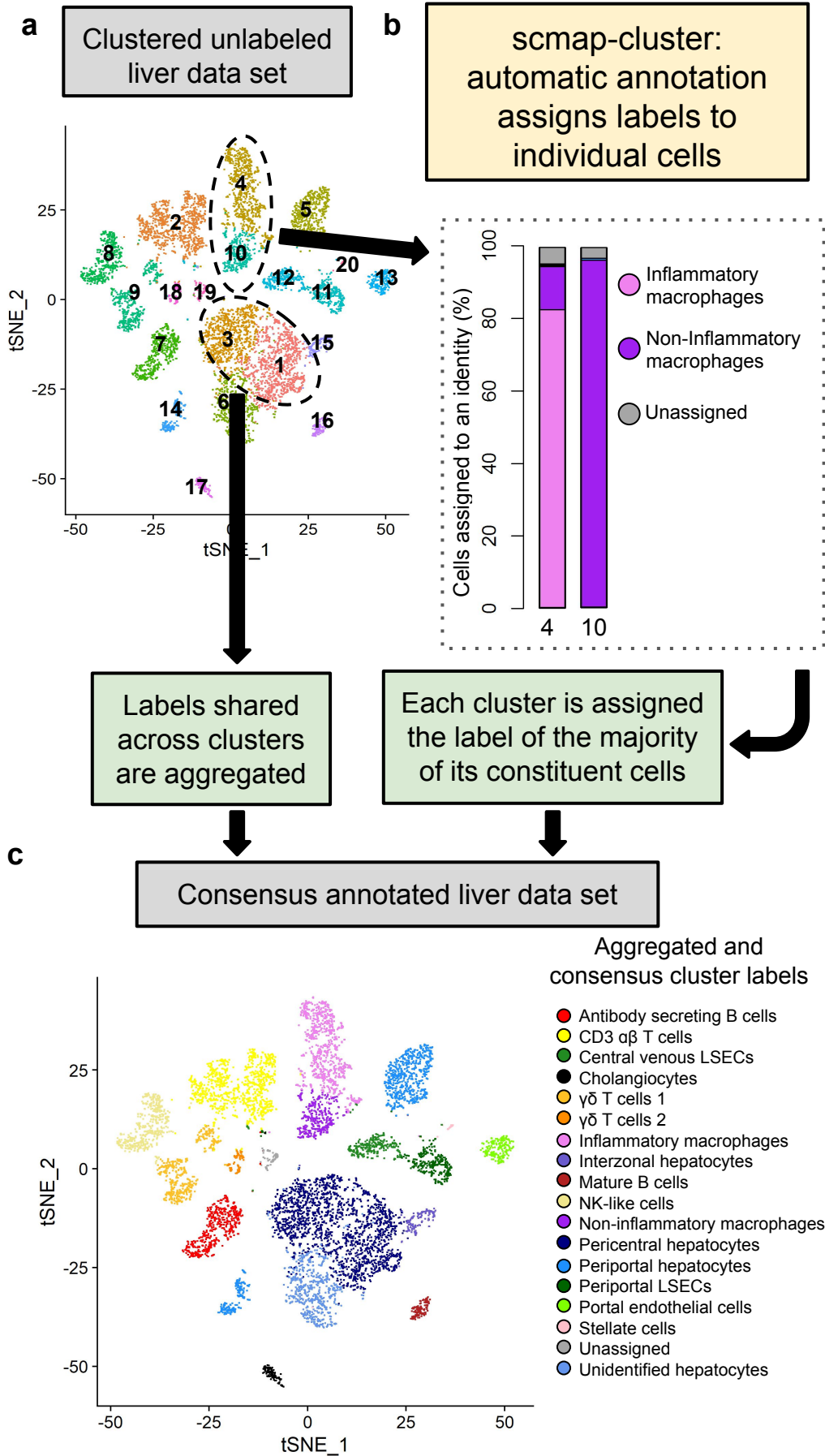
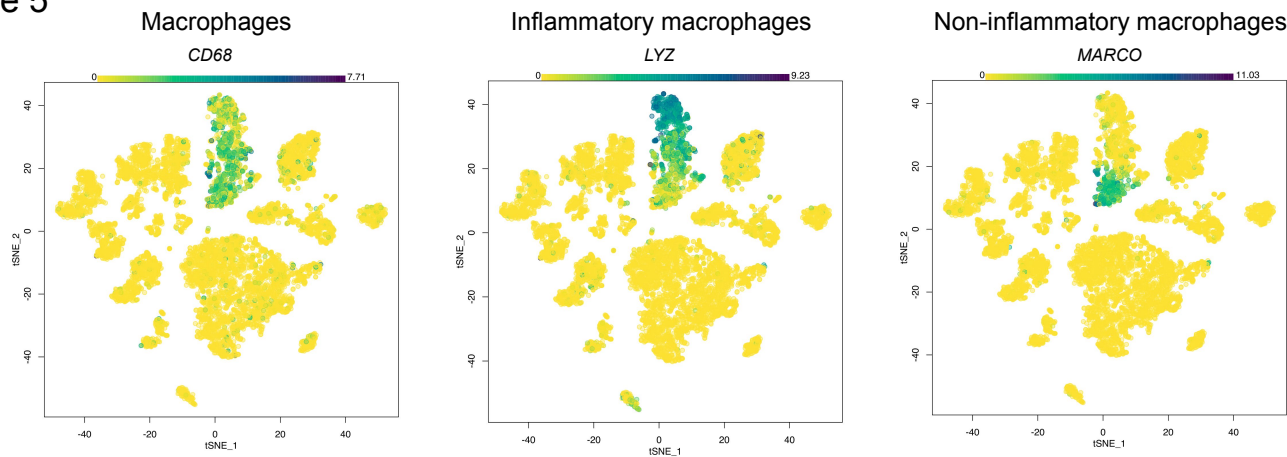


Figure 5

**a**



**b**

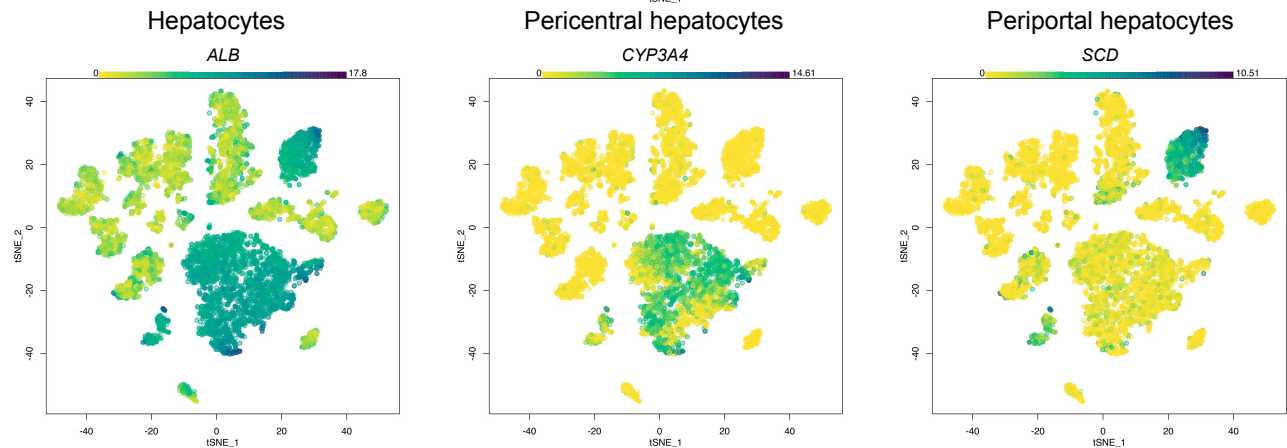
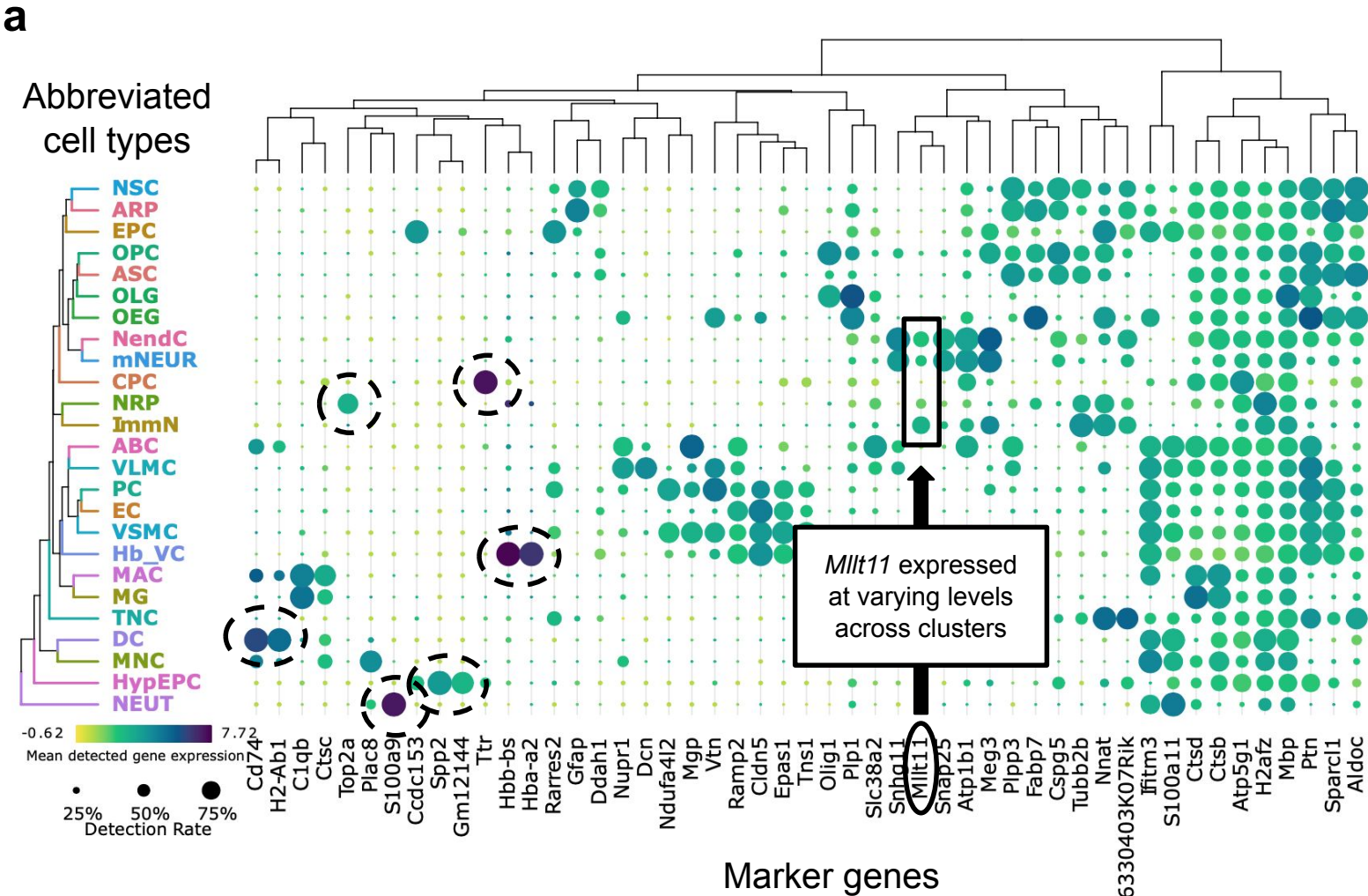




Figure 6



○ = clear marker genes

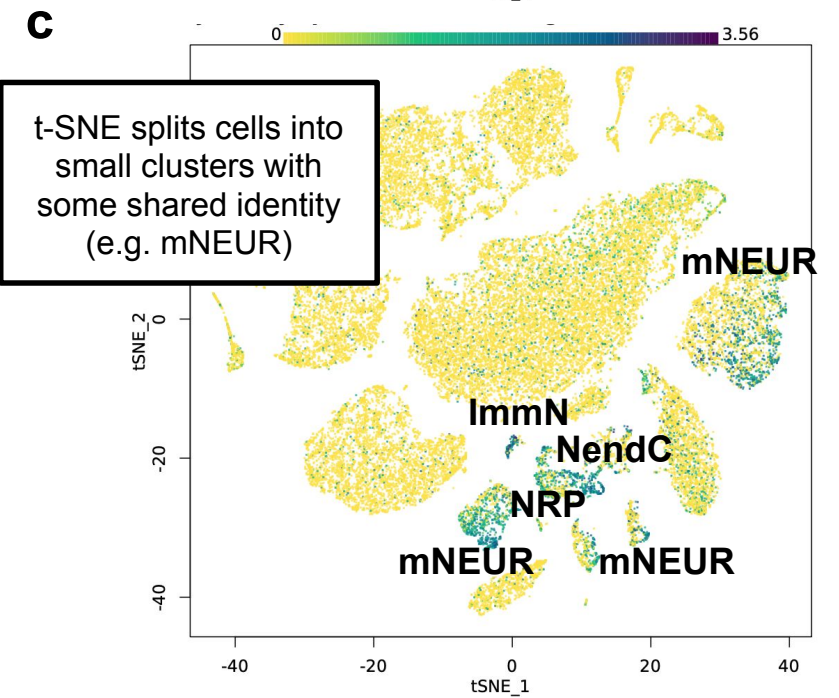
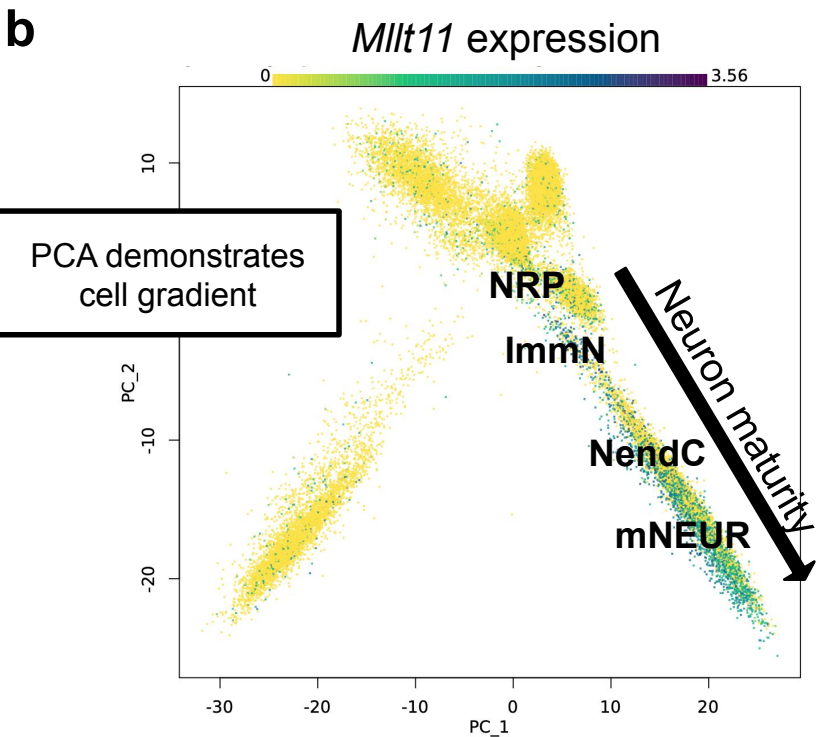


Figure 7

Cells sequenced  
with 3 different  
technologies



Batch correction  
applied across  
data sets

