

MarkerCount: A stable, count-based cell type identifier for single cell RNA-Seq experiments

HanByeol Kim

Dankook University

Joongho Lee

Dankook University

Keunsoo Kang

Dankook University

Seokhyun Yoon (✉ syoon@dku.edu)

Dankook University <https://orcid.org/0000-0002-0464-2233>

Research Article

Keywords: Downstream Analysis, Gene Expression, Batch Effect, Cluster-basis

Posted Date: April 20th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-418249/v2>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

MarkerCount: A stable, count-based cell type identifier for single cell RNA-Seq experiments

HanByeol Kim¹, Joongho Lee¹, Keunsoo Kang² and and Seokhyun Yoon^{3,*}

¹ Dept. of Computer Science, College of SW Convergence, Dankook University, Yongin-si, Korea

² Dept. of Microbiology, College of Natural Sciences, Dankook University, Cheonan-si, Korea

³ Dept. of Electronics and Electrical Eng., College of Engineering, Dankook University, Yongin-si Korea

* Corresponding: syoon@dku.edu

Abstract: Cell type identification is a key step to downstream analysis of single cell RNA-seq experiments. Indispensible information for this is gene expression, which is used to cluster cells, train the model and set rejection thresholds. Problem is they are subject to batch effect arising from different platforms and preprocessing. We present MarkerCount, which uses the number of markers expressed regardless of their expression level to initially identify cell types and, then, reassign cell type in cluster-basis. MarkerCount works both in reference and marker-based mode, where the latter utilizes only the existing lists of markers, while the former required pre-annotated dataset to train the model. The performance was evaluated and compared with the existing identifiers, both marker and reference-based, that can be customized with publicly available datasets and marker DB. The results show that MarkerCount provides a stable performance when comparing with other reference-based and marker-based cell type identifiers.

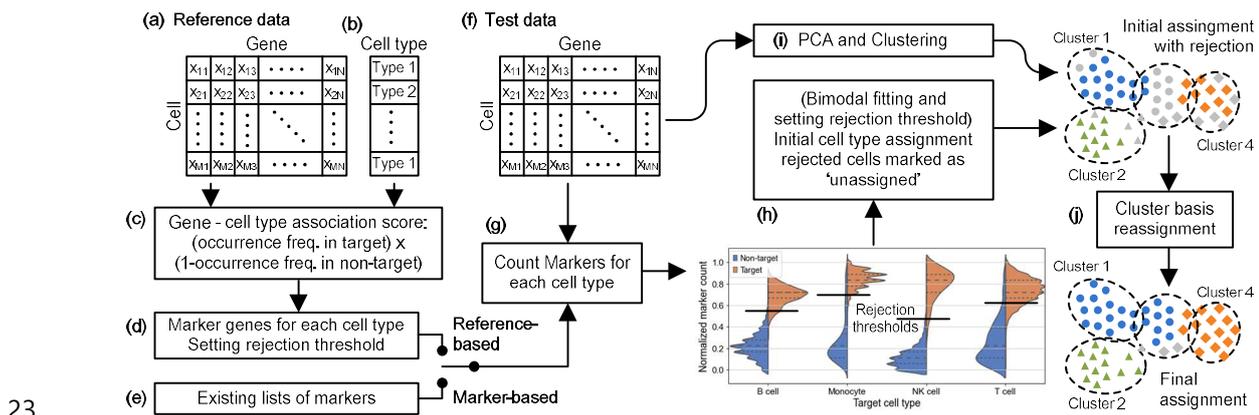
Keyword: cell type identification, single-cell RNA-Seq., marker-based identification, reference-based identification, cell type markers

Background

Single-cell RNA-Seq technology [1-3] enabled transcriptomic analysis of micro environment in heterogeneous tissues [4], such as tumor micro environment, where many cells of different types co-exist. A key step to its analysis is the cell type identification, for which there are many manual annotation tools, such as Seurat[5], SC3[6], SCANPY[7] and Alona[8], by which cell type annotations were made publicly available for further research. Automatic annotation pipelines were also developed, one of which is marker-based approach. This approach utilizes lists of known markers to identify cell types using gene expression profiles obtained from single cell RNA-Seq experiments. Garnett [9], SCINA[10], scSorter[11] and CellAssign[12] fall into this class, for which, there are several utilizable databases for cell type markers such as Panglao DB[13] and CellMarker DB[14]. Another approach is reference-based methods that utilize the existing annotation to obtain cell-type profiles to be used for identification, e.g., SingleR[15], scPred[16], scmap[17], CaSTLe[18] and CHETAH[19]. Various machine

1 learning techniques, such as classification and clustering, were used as key enabling tools for such types
2 of cell type identification. They are mostly work well, if all the cell types in the test data is already
3 known, i.e., they are in the references or their markers are available. As also noticed in [20], however,
4 one of the important functions of cell type identification is to precisely identify unknown cell clusters
5 too, which can be implemented by suitably setting the rejection thresholds. Although most of the cell
6 type identifiers provide such functions, they look not working so well as shown in [19]. Other things to
7 note are the measurement noise due to insufficient number of RNA molecules per cell and the batch
8 effect caused by different platforms and preprocessing. Most of the cell type identifiers, especially the
9 reference-based ones, utilizes directly the genes expression profiles to train the identification models.
10 Gene expression profiles, however, is subject to noise and batch effect, which may limit the performance.
11 To overcome such problems, we present MarkerCount, a count-base cell type identifier that support
12 both marker-based and reference-based identification. The overall procedure is shown in Figure 1 and
13 the detailed description can be found in the method section. Briefly speaking, the data processing
14 pipeline consists of two steps, (1) selecting markers for given reference and (2) identification of cell type
15 utilizing “marker counts” and cluster-wise cell type correction, where the second step can be slightly
16 modified to make it marker-based identification utilizing the existing cell type markers, e.g., those in
17 [13] and [14]. The key is to suitably set rejection threshold to determine ‘unknown’ cell type, for example,
18 unknown tumor cells, that is not in the set of target cell types. Most of the identifiers, either marker
19 based or reference based, showed reasonable performance when there are no unknowns. However, a
20 good identifier must be able to successfully identify such cell clusters for further study on their genomic
21 characteristics not identified so far.

22



24 **Figure 1** An overview of MarkerCount data processing. MarkerCount operates in both reference-based
25 and marker-based mode. In the former, it requires reference data (a) along with cell type annotation (b)
26 to find makers of each cell type (d), for which gene-cell type association score (c) is used. In marker-
27 based mode, it uses the existing markers (e). In test phase, it take cell type markers and the test data (f)
28 as input. Following data processing includes marker count (g), initial cell type assignment and rejection
29 of unclear cells (h), PCA and clustering (i) and the cluster basis reassignment for final identification (j).

1 Results

2 For evaluation, we ran MarkerCount in both reference-based and marker-based mode and
3 compared separately with the existing reference-based identifiers [15, 17-19] and the marker-based ones
4 [9-11] with the cell type markers in [14], respectively. To this end, we used 8 single cell RNA-Seq data
5 for which manual annotation is available. Among them, two are the Pancreas data reported in [21]
6 (Pancreas 2K) and [22] (Pancreas 8K), two are peripheral blood [23] (CBMC 8K) and [24] (PBMC 68K)
7 consisting of various immune related cells, two are lung data used in [25] and [26], which can be
8 downloaded from the human cell atlas (<https://data.humancellatlas.org/>) and the remaining two are
9 from the tissue with tumor [27] (Melanoma 5K) and [28] (Head and Neck 6K). A summary on these
10 datasets is shown in Supplemental Table 1. There are other software packages for cell type identification,
11 such as Moana [29], DigitalCellSorter[30], Cell-BLAST[31], scMatch[32], ACTINN[33], SCSA[34] and
12 mores in [20]. However, some of them do not provide functions to build custom models either utilizing
13 list of markers or reference data with annotation. Therefore, we consider only those with which one can
14 build customized model with gene expression patterns and cell type annotation for reference-based
15 methods and with list of existing markers for marker-based methods.

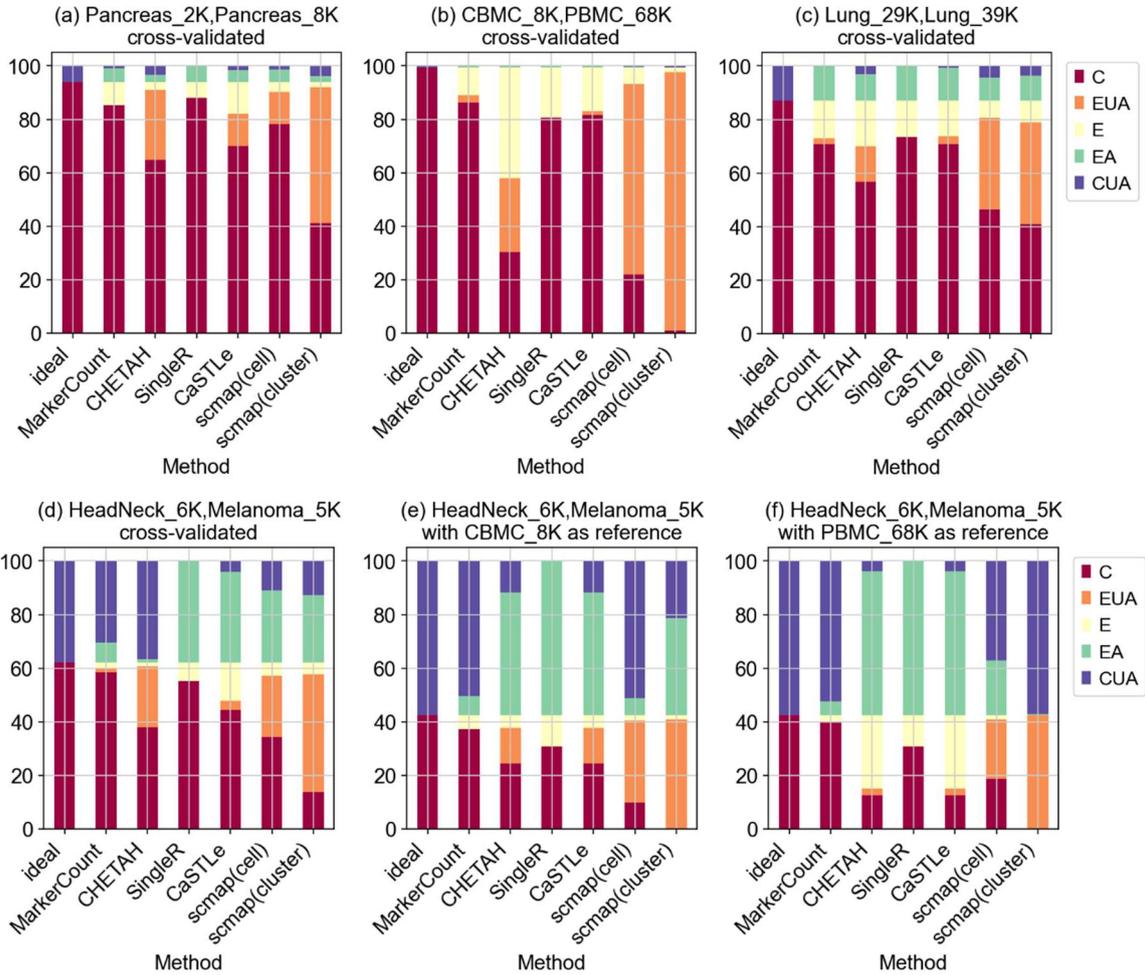
16 In the reference-based identification, we considered four pairs of data sets (Pancreas 2K and Pancreas
17 8K), (CBMC 8K and PBMC 68K), (Lung 29K and Lung 39K) (Melanoma 5K, HeadNeck 6K) separately.
18 For each of these pairs, we performed cross-validation, i.e., circularly used one as test data and others
19 as reference to identify marker genes and to set rejection threshold. To show the stable performance,
20 we also tested tumor data using CBMC 8K or PBMC 68K as reference to find markers for various
21 immune related cells including T cells, B cells, monocytes, dendritic cells and so on. Note that no tumor
22 cells were found in CBMC 8K and PBMC 68K, while there are many tumor cells in Melanoma and
23 HeadNeck data. Ideally, those tumor cells and the cell types not annotated in the reference must be
24 determined as 'unknown' or 'unassigned'.

25 Figure 2 summarizes the results, where, similar to those in [19], we considered 5 criteria, that is,
26 correct(C), error(E), erroneously assigned (EA), erroneously unassigned (EUA) and correctly
27 unassigned (CUA). These are defined according to (1) whether any valid label assigned by the identifier
28 or not (i.e., marked as unknown, unclear or unassigned), (2) whether the assigned label is in the
29 reference or not and (3) if they are in the reference cell type, whether the predicted label is the same as
30 the assigned label or not. The detailed definitions of the five measures are described in the method
31 section.

32 Ideally, only C and CUA should exist, which was depicted in the first bar in all figures. Imperfect
33 identification causes E, EA and/or EUA, each of which may have different impact according to the
34 purpose of the analysis. Although too much EUA is undesired, E and EA might have worse impact than
35 EUA in most cases as one can manually characterize clusters predicted as unknown or unassigned.
36 Roughly speaking, it is desirable that the predicted results is as close to the ideal case (the first bar) as

1 possible. In pancreas, blood and lung data, there is very small portion of cells of which their type is not
 2 in the references while, the melanoma and head-neck data contain a large portion of tumor cells that
 3 must be identified ideally as unknown (or unassigned).

4



6

7 **Figure 2** Comparisons of the reference-based cell type identification performances of MarkerCount,
 8 CHETAH, SingleR, CaSTLe, scmap (cell) and scmap (cluster). (a) Pancreas data, (b) Peripheral Blood
 9 (immune system) and (c) Tissues with tumor. C: Correct, EUA: Erroneously unassigned, E: Error, EA:
 10 Erroneously assigned, CUA: Correctly unassigned. The first bar in each figure shows the ideal
 11 performance, where CUA account for tumor cells and those cell types not in the reference. Ideally, they
 12 must be predicted as unknown (or unassigned).

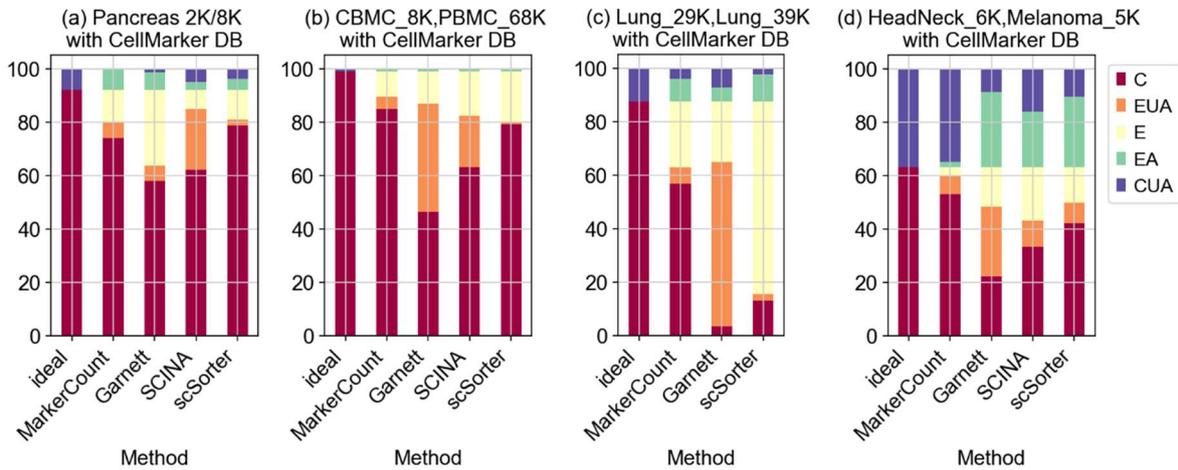
13

14 For the pancreas and lung data, SingleR has shown to be the best in terms of the correct prediction
 15 and MarkerCount, scmap(cell), CaSTLe, CHETAH, scmap (cluster) follow. However, if E and EA are
 16 desired as small as possible, scmap(cell), scmap(cluster) or CHETAH look better than SingleR,
 17 MarkerCount and CaSTLe, even though scmap(cluster) and CHETAH unassigned cell type for large
 18 portion of cells. In blood data, MarkerCount showed the best performance in terms of correct prediction.
 19 Although scmap(cell) and scmap(cluster) showed very small error, their EUAs are too large and the

1 correct precision were very small. The prediction accuracy in terms of CUA is highlighted in Figure 1
 2 (d) to (f) for tumor datasets with different references, where, in (e) and (f), we used CBMC 8K and
 3 PBMC 68K as references, respectively. In all 3 results, MarkerCount showed the closest pattern to ideal
 4 case with small portion of E and EA. Other cell type identifiers, such as CHETAH, scmap(cell) and
 5 scmap(cluster), showed smaller E than MarkerCount in (e) and (f). But, they erroneously assigned most
 6 of tumor cells to other types of normal cell type. Overall, the MarkerCount showed better performance
 7 than others for the blood and tumor datasets. Although its performance for pancreas and lung was not
 8 the best, it performed reasonably good.

9 Unlike other reference-based methods, MarkerCount provides also marker-based identification with
 10 a slight modification. Therefore, we also compared MarkerCount with the existing marker-based
 11 methods, including Garnett [9], SCINA [10] and scSorter [11] (We couldn't run CellAssign [12] due to
 12 installation error). Marker-based approach requires only lists of cell markers for which we used
 13 CellMarker DB [14]. We extracted cell markers for pancreas, lung, blood and peripheral blood, where
 14 the last two were used for CBMC 8K, PBMC 68K and the two tumor tissues, Melanoma 5K and
 15 HeadNeck 6K. The results are summarized in Figure 2 for the four pairs of datasets, pancreas,
 16 peripheral blood (CBMC and PBMC), lung and tumor tissues. Similar to the comparison of reference-
 17 based cell type identifiers, we measured 5 criteria, C, EUA, E, EA and CUA. For pancreas dataset,
 18 scSorter performed slightly better than MarkerCount. However, for all other datasets, MarkerCount
 19 showed the best performance.

20



21
 22 **Figure 3** Comparisons of the marker-based cell type identification performances of MarkerCount,
 23 scSorter, Garnett, and SCINA. (a) Pancreas data, (b) Peripheral Blood, (c) Lung and (d) Tissues with
 24 tumor, using CellMarker DB. C: Correct, EUA: Erroneously unassigned, E: Error, EA: Erroneously
 25 assigned, CUA: Correctly unassigned. The first bar in each figure shows the ideal performance, where
 26 CUA account for tumor cells and those cell types not in the reference cell type. Ideally, they must be
 27 predicted as unknown (or unassigned).

28

1 Discussion

2 Marker-based approach and reference-based approach have their own advantage and disadvantage.
3 In both approaches, one of key components is clustering and, for a good identification, each cell types
4 should be well clustered, which is not always the case. Although many clustering algorithms, such as
5 partitioning-based, distribution-based or graph-based, are available, performance difference with
6 different clustering seems not critical. Rather, handling unclear clusters of highly differentiable cells, as
7 in monocytes, macrophage and dendritic cells, seems to be more important issue.

8 Reference-based approaches select their own markers internally and look more robust if the
9 reference annotation is sufficiently reliable. However, they can be subject to batch effect when working
10 with datasets from different platforms and preprocessing. Therefore, using the binary indication of
11 gene expression, as in MarkerCount, makes it more robust to the batch effect.

12 Another point to consider in reference-based approach is the uneven number of cell type population.
13 In Supplemental Figure 1, we plotted Sankey diagram between the manual annotation and the
14 predicted one for blood and tumor datasets, respectively. The figure shows that the cell types having
15 sufficient number of cells shows better performance than those with smaller number of cells. Having
16 large number of cells, they can be well characterized to obtain the best set of markers so that one can
17 better identify those cell types while the identification model for those cell types with small number of
18 cells tends to be overfitted to the data. Doubtlessly, we need more cells to improve the identification
19 performance in most of the reference-based identifiers. The problem is that uneven cell population is
20 inherent. One may combine two or more dataset into one reference data to train the model to improve
21 the performance. However, it may not be always true because of the batch effect. Most of the cell type
22 identifiers utilizes clustering, anyhow, and the identification model is fitted in cluster basis. Problem
23 arises if there are two or more clusters far apart for one cell type due to batch or tissue specific effect.
24 Another issue to consider is the cell types that are highly differentiable, such as monocytes,
25 macrophages and dendritic cells. As shown in Supplemental Figure 1 and 2, some of these cell types
26 were hardly recovered and interchangeably detected. In many cases, they are not well clustered with
27 sufficient separation as shown in Supplemental Figure 2.

28 Compared to reference-based approach, the marker-based approach does not require reference
29 annotation, which is a big advantage compared to reference-based approach. However, it depends
30 highly on the selection of markers and the best set of markers should be selected taking the specific
31 identification procedures into account. One reason to high dependency on the set of markers seems to
32 be the uneven number of cell markers. In CellMarker DB, several subtypes of dendritic cells have
33 hundreds of marker genes while others, e.g., natural killer cells and B cells, have only around ten
34 markers. Such uneven number of markers may degrade the identification performance and this is why
35 MarkerCount (marker-based mode) reselect markers if too many markers is provided in marker DB
36 (see method section). Depending on the specific procedures and algorithm of cell type identifiers, the

1 best set of markers will also be different. Literally, 'cell marker' means that it is expressed only in a
2 specific cell type and not expressed at all elsewhere. However, it seems not be always hold since the
3 definition of cell type can be quite different according to the analytical purposes and applications. For
4 example, one can required only identify T cells as a whole in some applications, while one may require
5 identification of specific subtypes, such as cytotoxic T cells, memory T cells, helper T cells and
6 regulatory T cells. Depending on what specific set of cell types are required, different approach and/or
7 different tools must be used. Supplemental Figure 3 and 4 show the Sankey diagram and the UMAP
8 plots similar to those in the reference-based identification experiments. In Supplemental Figure 3 and
9 4, we see that specific subtypes are interchangeably identified to other subtypes within the same
10 broader cell types. Having markers of both broad and specific cell type, simultaneous identification of
11 broad cell types and specific cell types may cause unexpected results as shown in Supplemental Figure
12 3. As it is also risky to identify specific cell type directly using only a few markers, a better option is to
13 use hierarchical identification, i.e., identify broad cell type first using many markers and then determine
14 their subtypes using a few specific markers. What specific hierarchy should be used depends on specific
15 tissues to analyze.

16 **Conclusion**

18 Single-cell RNA-Sequencing is certainly a useful tool to analyze heterogeneous tissues, for example,
19 to study tumor micro-environment and cell-differentiation and proliferation, for which cell type
20 identification is a first step to the downstream analysis and automatic identification of cell type may
21 speed up those researches by replacing time-consuming manual annotations with fast and reliable
22 annotation techniques. For medical application, identifying different tumor types and immune cells
23 that co-exist in a tissue may give us a new insight into the personalized cancer therapy. Although the
24 application of single-cell RNA-Seq technology to these studies is still in its early stage, many of the
25 clinical factors and measurements can be obtained by using one single-cell RNA-Seq experiment on the
26 target tissue. And, to this end, more data should be filed with appropriate annotations and cross-
27 checked with different annotation tools.

28 **Methods**

30 The operation of MarkerCount can be divided into two steps, (1) Finding markers for each cell types
31 using reference data and (2) Use markers to identify cell types in cluster basis.

32 **Finding markers**

33 Using the prior annotation provided with the gene expression matrix, reference cell types are identified
34 first. Let G be the set of genes and x_{ij} be the genes expression of the i th cell and j th gene. We first
35 convert them to binary indicator $b_{ij} \in \{0,1\}$ representing whether the gene is expressed in that cell or

1 not, i.e., $b_{ij} = 1$ if $x_{ij} > 0$ or 0 otherwise. We then compute marker score $s_{m,j}$ for a cell type m as

$$2 \quad s_{m,j} = \left(\frac{1}{|C_m|} \sum_{i \in C_m} b_{ij} \right) \left(1 - \frac{1}{|\bar{C}_m|} \sum_{i \in \bar{C}_m} b_{ij} \right)^n \quad (1)$$

3 where C_m is the set of cells annotated as of type m and \bar{C}_m is the complementary set of C_m . $|C_m|$ is
4 the cardinality of C_m . $\frac{1}{|C_m|} \sum_{i \in C_m} b_{ij}$ is simply the occurrence frequency of the j th gene in the m th cell
5 type. n is an hyperparameter relatively weighing the second term. We set it 2 in all the experiments
6 performed in this work. Then, for each cell type m , we sort $s_{m,j}$ in descending order and select first
7 N genes as its markers. N is also an hyper-parameter and we fixed it to 18 for all the cell types in this
8 work. Note that the selection of marker genes is done separately for each cell type and it is possible for
9 two or more cell types to share some common markers. In addition, before comparing the score in (1),
10 one can narrow down the candidate genes by enforcing the condition, $\frac{1}{|C_m|} \sum_{i \in C_m} b_{ij} \geq f_{th}$ for some
11 occurrence frequency threshold f_{th} which we set 0.9.

12 **Marker count**

13 In test phase, for given binary expression profile of the i th cell, b_{ij} , we obtain the normalized marker
14 count $y_{i,m} = \frac{1}{N} \sum_{j \in M_m} b_{ij}$ where M_m is the set of markers of the m th cell type. The cell type is initially
15 determined by taking the maximum of $y_{i,m}$, i.e., $m_i^* = \operatorname{argmax}_m y_{i,m}$, which is accepted if $\max_m y_m$ is
16 greater than or equal to the cell type specific rejection threshold, t_m . Otherwise, the cell is marked as
17 ‘unassigned’.

18 **Obtaining the rejection threshold**

19 The count threshold t_m can be obtain in various ways. In this work, we considered two approach, i.e.,
20 parametric and non-parametric. In non-parametric approach, we directly obtain from y_m 's. Consider
21 two set of cells for given rejection threshold t , $C_m(t)$ and $C_m^*(t)$. The former is the set of cells (in the
22 reference data) of which the true (manually annotated) cell type is m and its normalized marker count
23 $y_{i,m} \geq t$, while the latter is the set of cells that are decided to be the m th cell type according to $y_{i,m}$
24 satisfying (1) $y_{i,m} > y_{i,m'}$ for all other cell types m' and (2) $y_{i,m} \geq t$. Then, the false positive rate can
25 be defined as

$$26 \quad FPR(t) = 1 - \frac{|C_m(t) \cap C_m^*(t)|}{|C_m(t)|} \quad (2)$$

27 where the second term is the true positive rate. The objective is to find t such that $FPR(t) \approx p$ for
28 given target FPR , p . In non-parametric approach, one can find t by sorting $y_{i,m}$ in descending order
29 and find minimum $y_{i,m}$ such that $FPR(t) \leq p$. In parametric approach, we use univariate Gaussian
30 mixture model with two components, i.e., $y_{i,m} \sim N(\mu_1, \sigma_1^2)$ for $m \in C_m$ and $y_{i,m'} \sim N(\mu_0, \sigma_0^2)$ for $m' \in$
31 \bar{C}_m . Denoting the relative size of C_m and \bar{C}_m as π_0 and π_1 , respectively, $FPR(t)$ can be defined as

$$FPR(t) = \frac{F_0(t)}{F_1(t)+F_0(t)} \text{ with } F_k(t) \equiv \int_t^\infty \pi_k N(y; \mu_k, \sigma_k^2) dy \quad (3)$$

The parametric approach especially useful for marker-based operation of MarkerCount, where reference data is not available. Given test data only, one can resort to the bimodal fitting, for which expectation maximization (EM) algorithm can be used to find π_k, μ_k, σ_k^2 for $k = 0, 1$. Note that SCINA also used bimodal fitting. But, it applied to gene expression level, while MarkerCount applies it to the normalized marker counts $y_{i,m}$ to determine the rejection threshold. In reference-based mode, the rejection thresholds for all the cell types are computed in the training phase using the reference data using either (2) or (3), while, in marker-based mode, it can be obtained by applying bimodal fitting to the test data. Supplemental Figure 5 shows comparisons of the two distributions of C_m and \bar{C}_m for some cell types in pancreas, CBMC 8K and PBMC68K.

Cluster basis cell type reassignment

Although marker count based cell type identification works reasonably well, one can improve the identification accuracy, specifically the recall, by employing cluster basis reassignment. The procedure is as follows. We first perform clustering. While some clusters have all its cells assigned by a cell type, others may not be fully covered or not covered at all. The latter might be unknown cells that does not exist in the reference cell type. However, the former may be supposed to be a specific cell type that was not fully covered due to high rejection threshold. We can reassign cell types to those unassigned cells in a partially covered region in cluster by cluster fashion. To this end, we first identify cell types that are partially occupying the cluster. With their centroids, covariance matrices and size, we reassign a cell type to those unassigned by comparing their distances from the centroids as follows.

$$d(\mathbf{z}, \boldsymbol{\mu}_m; \boldsymbol{\Sigma}_m, \pi_m) = \frac{1}{\pi_m} (\mathbf{z} - \boldsymbol{\mu}_m)^T (\boldsymbol{\Sigma}_m + \rho \mathbf{I})^{-1} (\mathbf{z} - \boldsymbol{\mu}_m) \quad (4)$$

where \mathbf{z} is a dimension-reduced gene expression profile of a cell, $\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m, \pi_m$ are the centroid (mean), covariance matrix and the relative size of the m th cell type that reside in the cluster and ρ is a regularization constant, which we set 0.1 of the average variances. By comparing the weighted Mahalanobis distance in (4) of an unassigned cell from the centroids for the cell types in that cluster, we decide its cell type to the closest one. For the dimension reduction, we applied principal component analysis (PCA) as typically used in the preprocessing for clustering. The reassignment is performed if the portion of the unassigned cells is below a certain threshold, say 0.8. Although more sophisticated methods can be devised, we used this rather simple heuristic approach.

Reference-based and marker-based operation of MarkerCount

Using the functions described above, the operation of MarkerCount can now be described more succinctly. In the reference-based operation, MarkerCount uses the reference data and its annotation to find the best markers for each cell type in the reference. Using these markers, it computes the

1 normalized marker counts and decide initial prediction of cell types for all the cells in the reference
2 data. Using the initial prediction and the annotated cell type, the rejection thresholds are obtained. The
3 outputs of the training phase consist of the sets of markers and the rejection threshold for each cell type,
4 which are then used for cell type identification for the test data. In the test phase, we first assign cell
5 types to those cells whose normalized marker count is above the rejection threshold. Then clustering is
6 performed for the test data and the distance-based reassignment is applied in cluster by cluster basis.
7 In marker-based mode, on the other hand, we first compute the normalized marker counts for all the
8 cells and all the cell types utilizing the existing markers. For each cell type, bimodal fitting is applied to
9 determine the rejection threshold for each cell type. One possible problem in the marker-based mode
10 is the uneven number of markers. In this case, cell types with a few markers tend to get higher priority
11 than those with much larger number of markers. This does not happen in the reference-based mode as
12 we select the same or at least similar number of markers for each cell type. To handle this problem, we
13 applied the following tricks.

14 **Resolving uneven number of markers**

15 One possible solution is to reselect markers for those cell types with large number of markers. That is,
16 in the first round of cell type assignment, we assign cell types with higher rejection threshold and use
17 the prediction in the first around to obtain the score in (1) for all the markers and finally reselect markers
18 that have higher scores. This is done only for those cell types that has the number of markers larger
19 than the desired numbers. This approach, however, solve the problem only partially since there exist
20 cell types having only a few or sometimes only one specific markers. To resolve this problem, we
21 applied penalty weight per cell type according to their number of markers as $w = (1 + 3e^{-(n-1)/2})^{-1}$,
22 where n is the number of markers. The penalty weight w is multiplied to the normalized marker
23 count before obtaining the rejection threshold and making cell type assignment. Although these tricks
24 were devised for the case when the existing marker DB is used, the best way to avoid such problem in
25 MarkerCount is to use the same or similar number of markers for each cell type, in which case, the
26 penalty weight will be the same for all the cell types so that it does not affect the MarkerCount operation.

27 **Dataset and cell type renaming**

28 To properly evaluate and compare performances, we collected 6 single cell RNA-Seq data with manual
29 annotation available. Those data were used in [19] and partially in other works. The description of the
30 data is summarized in Supplemental Table 1. In all data, we did not take into count those cells annotated
31 as 'unknown' or 'unclear' in the performance evaluation even though we performed identification
32 anyway.

33 Two tumor data and CBMC 8K used rather broad cell types, such as T cells, NK cells and dendritic
34 cells, while PBMC 68K used specific types, such as CD8+/CD45RA+ Naive Cytotoxic, CD4+ T Helper2,

1 CD4+/CD45RA+/CD25- Naive T, CD8+ Cytotoxic T, CD4+/CD45RO+ Memory and CD4+/CD25 T Reg.
2 Therefore, in the comparison of reference-based identifiers, we renamed PBMC cell types to their
3 corresponding broad cell type, e.g., Naive Cytotoxic, T Helper2, Naive T, Cytotoxic T, Memory T and
4 T Reg were all mapped to one broad type, T cell, to train the models for all the reference-based identifier.
5 In marker-based, the two marker DB provide markers for specific cell type along with those for broad
6 cell type. For example, CellMarker DB provides 87 T cell markers along with several specific markers
7 for CD4+, CD8+, CD4+ memory T cells, and so on. Therefore, we used the markers as is for identification
8 of specific cell type. The problem was that, when identifying specific cell types, the performances were
9 unacceptable for all the identifiers we examined, as shown in Supplemental Figure 3. Hence, in the
10 comparison of marker-based identifiers, we renamed after identification, using the predicted specific
11 cell types. As a matter of fact, immune related cells are highly differentiable so that they do not form
12 clearly separable clusters as shown in Supplemental Figure 2 and 4. Moreover, it is questionable that
13 the manual annotation can be used as ground truth since it is typically done in cluster basis, which
14 limits the resolution of cell type clustering.

15 **Performance criteria**

16 Five performance measures, C, E, EA, CUA and EUA, are defined as follows:

- 17 A. C is the portion of cells that has an assigned label, one from the reference cell type, and the
18 predicted label is the same as the original assignment.
- 19 B. E is the portion of cells that has an assigned label existing in the reference, while the predicted
20 label is different from the original one.
- 21 C. EA is those cells that has an assigned label but that is not in the reference and the identifier
22 predicted as one of the reference cell type
- 23 D. CUA is those that has an assigned label but that is not in the reference, including tumor cells, and
24 the identifier predicted as unknown (or unassigned).
- 25 E. EUA is those that has an assigned label in the reference and the identifier predicted as unknown
26 (or unassigned).

27 In general, there exists tradeoff between the two. Although it is also informative to show precision
28 versus recall, we used these more specific measures to give better insight into the performance of cell
29 type identifiers.

30 **Cross validation for the reference-based identification**

31 To show the effectiveness of the reference-based cell type identification, we used cross validation,
32 where, with M datasets, one is selected for test while others are used as reference to find cell type
33 markers and this is repeated by circularly shifting their role in the experiment. This was done for the
34 results in Figure 1 (a) to (d). For the results in Figure 1 (e) and (f), the reference and the test data were
35 fixed.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

Ethics approval and consent to participate

Not applicable

Consent for publication

Not applicable

Availability of data and materials

All datasets are publicly available and their download site are summarized in Supplemental Table 1. The python code and the example in Jupyter notebook were deposited to Github and can be obtained at <https://github.com/combio-dku/MarkerCount/tree/master> (Project name: MarkerCount, license: GPL 3.0, Operating system(s): Platform independent, Programming language: python 3, other requirement: None).

Competing Interest

The authors declare that they have no competing interest.

Funding

This work was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (NRF-2020R1F1A1066320).

Authors' contribution

HK and SY devised the concept and key idea. SY and HK developed the python code. HK and JL performed experiments. SY guided experiments and data analysis. KK gave interpretation and feedback on the results. All authors wrote, read, and approved the final manuscript.

References

1. Tang, F. C., Barbacioru, C., Wang, Y. Z., Nordman, E., Lee, C., Xu, N. L., et al. (2009). mRNA-Seq whole-transcriptome analysis of a single cell. *Nat. Methods* 6, 377–386. doi: 10.1038/nmeth.1315
2. Picelli, S., Bjorklund, A. K., Faridani, O. R., Sagasser, S., Winberg, G., and Sandberg, R. (2013). Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat. Methods* 10, 1096–1098. doi: 10.1038/nmeth.2639
3. Kolodziejczyk, A. A., Kim, J. K., Svensson, V., Marioni, J. C., and Teichmann, S. A. (2015). The technology and biology of single-cell RNA sequencing. *Mol. Cell* 58, 610–620. doi: 10.1016/j.molcel.2015.04.005

- 1 4. Chung, W., Eum, H. H., Lee, H. O., Lee, K. M., Lee, H. B., Kim, K. T., et al. (2017). Single-cell RNA-
2 seq enables comprehensive tumour and immune cell profiling in primary breast cancer. *Nat.*
3 *Commun.* 8:15081. doi: 10.1038/ncomms15081
- 4 5. Satija, Rahul, et al. Spatial reconstruction of single-cell gene expression data. *Nature biotechnology*
5 2015;33:495-502.
- 6 6. Kiselev, Vladimir Yu, et al. SC3: consensus clustering of single-cell RNA-seq data. *Nature methods*
7 2017;14:483-486.
- 8 7. Wolf, F. Alexander, Philipp Angerer, and Fabian J. Theis. SCANPY: large-scale single-cell gene
9 expression data analysis. *Genome biology* 2018;19:1-5.
- 10 8. Franzén, Oscar, and Johan LM Björkegren. alona: a web server for single-cell RNA-seq analysis.
11 *Bioinformatics* 2020;36:3910-3912.
- 12 9. Pliner, Hannah A., Jay Shendure, and Cole Trapnell. "Supervised classification enables rapid
13 annotation of cell atlases." *Nature methods* 2019;16:983-986.
- 14 10. Zhang, Ze, et al. SCINA: a semi-supervised subtyping algorithm of single cells and bulk samples.
15 *Genes* 2019;10:531.
- 16 11. Guo, Hongyu, and Jun Li. scSorter: assigning cells to known cell types according to marker genes.
17 *Genome biology* 2021;22: 1-18.
- 18 12. Campbell, K. R., S. P. Shah, and A. W. Zhang. Assigning scRNA-seq data to known and de novo
19 cell types using CellAssign. 2019;2019
- 20 13. Franzén, Oscar, Li-Ming Gan, and Johan LM Björkegren. PanglaoDB: a web server for exploration
21 of mouse and human single-cell RNA sequencing data. *Database* 2019;2019
- 22 14. Zhang, Xinxin, et al. CellMarker: a manually curated resource of cell markers in human and mouse.
23 *Nucleic acids research* 2019;47:D721-D728.
- 24 15. Aran, Dvir, et al. Reference-based analysis of lung single-cell sequencing reveals a transitional
25 profibrotic macrophage. *Nature immunology* 2019;20:163-172.
- 26 16. Alquicira-Hernandez, Jose, et al. scPred: accurate supervised method for cell-type classification
27 from single-cell RNA-seq data. *Genome biology* 2019;20:1-17.
- 28 17. Kiselev, Vladimir Yu, Andrew Yiu, and Martin Hemberg. scmap: projection of single-cell RNA-seq
29 data across data sets. *Nature methods* 2018;15:359-362.
- 30 18. Lieberman, Yuval, Lior Rokach, and Tal Shay. CaSTLe—classification of single cells by transfer
31 learning: harnessing the power of publicly available single cell RNA sequencing experiments to
32 annotate new experiments. *PloS one* 2018;13:e0205499.
- 33 19. de Kanter, Jurrian K., et al. CHETAH: a selective, hierarchical cell type identification method for
34 single-cell RNA sequencing. *Nucleic acids research* 2019;47:e95-e95.
- 35 20. Abdelaal, T., Michielsen, L., Cats, D. *et al.* A comparison of automatic cell identification methods for
36 single-cell RNA sequencing data. *Genome Biol* **20**, 194 (2019)
- 37 21. Muraro, Mauro J., et al. A single-cell transcriptome atlas of the human pancreas. *Cell systems*
38 2016;3:385-394.

- 1 22. Baron, Maayan, et al. A single-cell transcriptomic map of the human and mouse pancreas reveals
2 inter-and intra-cell population structure. *Cell systems* 2016;3:346-360.
- 3 23. Stoeckius, Marlon, et al. Simultaneous epitope and transcriptome measurement in single cells.
4 *Nature methods* 2017;14:865-868.
- 5 24. Zheng, Grace XY, et al. Massively parallel digital transcriptional profiling of single cells. *Nature*
6 *communications* 2017;8:1-12.
- 7 25. Madisson, E., Wilbrey-Clark, A., Miragaia, R.J. *et al.* scRNA-seq assessment of the human lung,
8 spleen, and esophagus tissue stability after cold preservation. *Genome Biol* **21**, 1 (2020)
- 9 26. Reyfman PA, Walter JM, Joshi N, et al. Single-Cell Transcriptomic Analysis of Human Lung
10 Provides Insights into the Pathobiology of Pulmonary Fibrosis. *Am J Respir Crit Care Med.*
11 2019;199(12):1517–1536.
- 12 27. Tirosh, Itay, et al. Dissecting the multicellular ecosystem of metastatic melanoma by single-cell
13 RNA-seq. *Science* 2016;352:189-196.
- 14 28. Puram, Sidharth V., et al. Single-cell transcriptomic analysis of primary and metastatic tumor
15 ecosystems in head and neck cancer. *Cell* 2017;171:1611-1624.
- 16 29. Wagner, Florian, and Itai Yanai. Moana: A robust and scalable cell type classification framework
17 for single-cell RNA-Seq data. *BioRxiv* 2018;456129.
- 18 30. Domanskyi, Sergii, et al. Polled Digital Cell Sorter (p-DCS): Automatic identification of
19 hematological cell types from single cell RNA-sequencing clusters. *BMC bioinformatics* 2019;20:1-
20 16.
- 21 31. Cao, Zhi-Jie, et al. Searching large-scale scRNA-seq databases via unbiased cell embedding with
22 Cell BLAST. *Nature communications* 2020; 11:1-13.
- 23 32. Hou R, Denisenko E and Forrest ARR, scMatch: a single-cell gene expression profile annotation
24 tool using reference datasets, *Bioinformatics* 2019; 35(22): 4688–4695
- 25 33. Ma F and Pellegrini M, ACTINN: automated identification of cell types in single cell RNA
26 sequencing, *Bioinformatics* 2020; 36(2):533–538
- 27 34. Cao Y, Wnag X and Peng G, SCSA: A Cell Type Annotation Tool for Single-Cell RNA-seq Data,
28 *Front. Genet.* 2020; 11:490
29

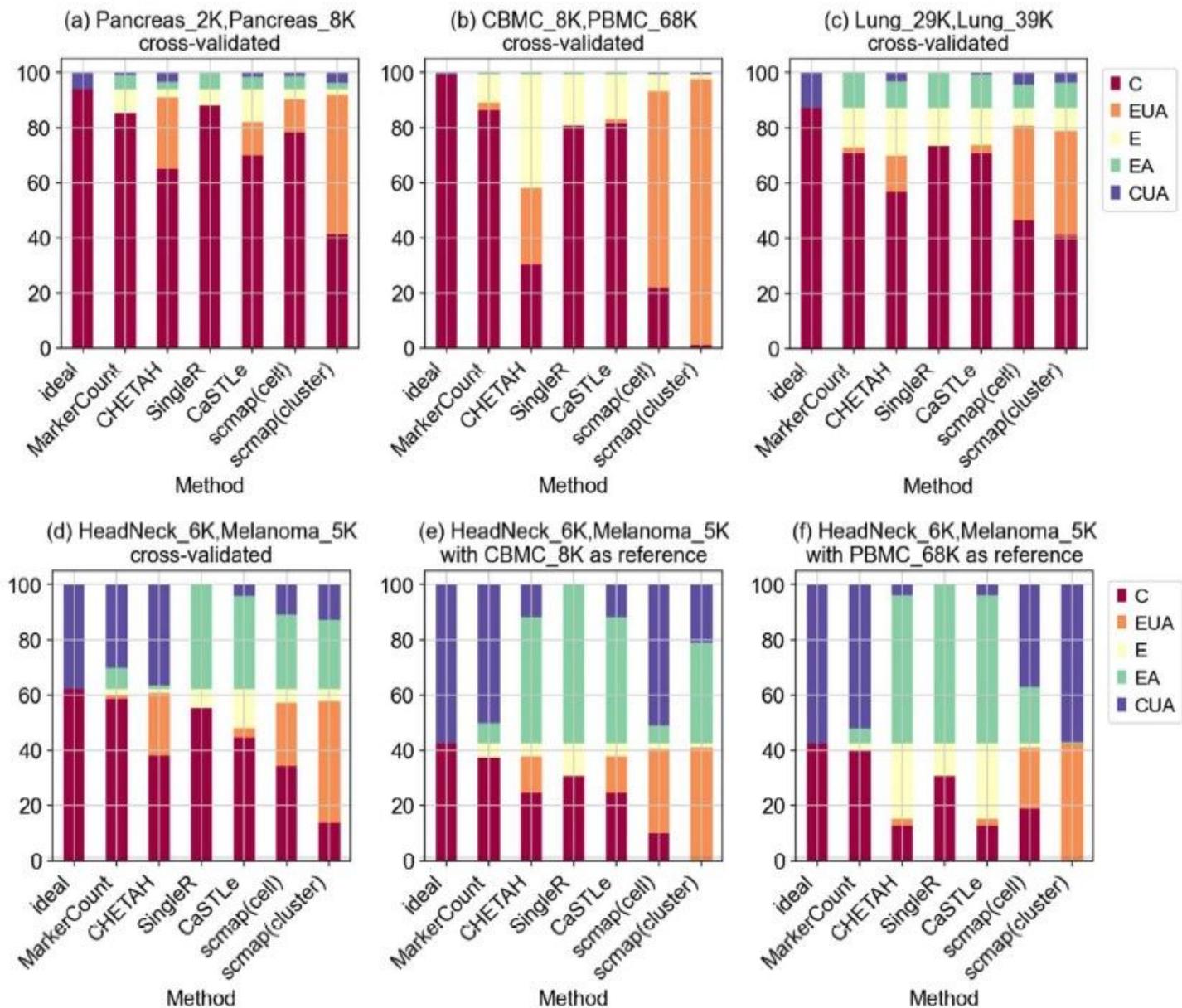


Figure 2

Comparisons of the reference-based cell type identification performances of MarkerCount, CHETAH, SingleR, CaSTLe, scmap (cell) and scmap (cluster). (a) Pancreas data, (b) Peripheral Blood (immune system) and (c) Tissues with tumor. C: Correct, EUA: Erroneously unassigned, E: Error, EA: Erroneously assigned, CUA: Correctly unassigned. The first bar in each figure shows the ideal performance, where CUA account for tumor cells and those cell types not in the reference. Ideally, they must be predicted as unknown (or unassigned).

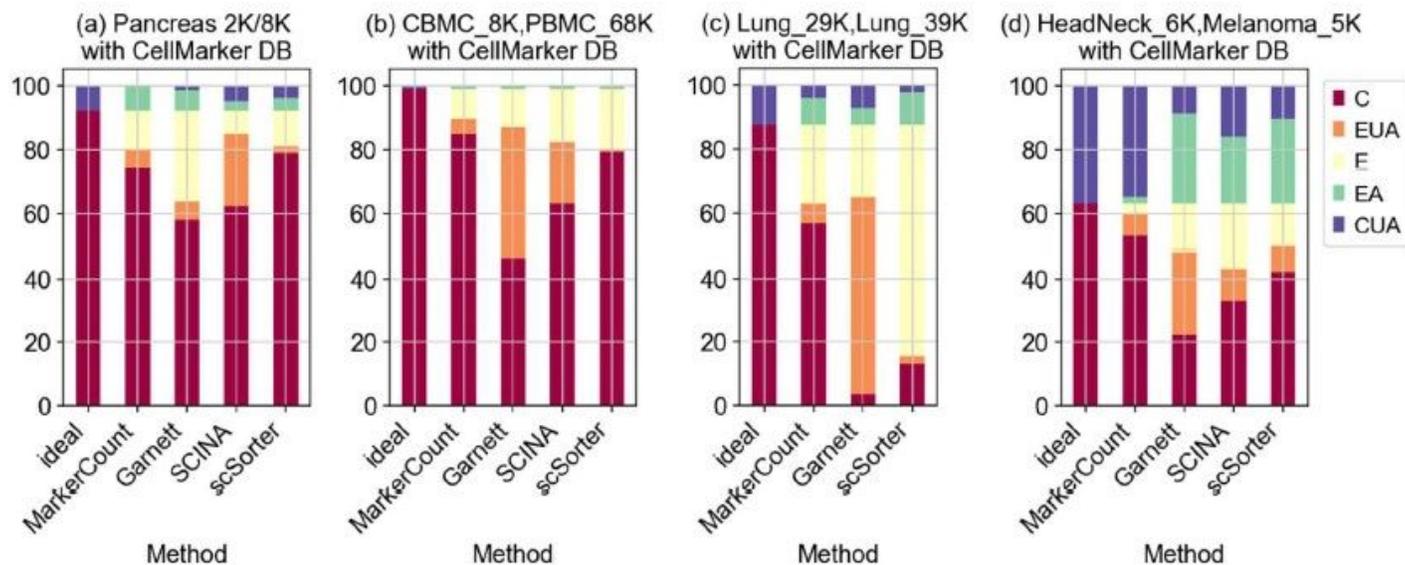


Figure 3

Comparisons of the marker-based cell type identification performances of MarkerCount, scSorter, Garnett, and SCINA. (a) Pancreas data, (b) Peripheral Blood, (c) Lung and (d) Tissues with tumor, using CellMarker DB. C: Correct, EUA: Erroneously unassigned, E: Error, EA: Erroneously assigned, CUA: Correctly unassigned. The first bar in each figure shows the ideal performance, where CUA account for tumor cells and those cell types not in the reference cell type. Ideally, they must be predicted as unknown (or unassigned).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementalMaterials.pdf](#)