1
2
# clustifyr: An R package for automated single-cell RNA sequencing cluster classification

3

4   Rui Fu[1], Austin E. Gillen[1], Ryan M. Sheridan[1], Chengzhe Tian[2], Michelle Daya[3], Yue Hao[4], Jay
5   R. Hesselberth[1,5], Kent A. Riemondy[1*]

6

7   [1] RNA Bioscience Initiative, University of Colorado School of Medicine, Aurora, CO, 80045,
8   USA. [2] Department of Biochemistry, University of Colorado Boulder, Boulder, CO, 80303, USA. [3]
9   Biomedical Informatics & Personalized Medicine, University of Colorado Anschutz Medical
10  Campus, Aurora, CO, 80045, USA. [4] Bioinformatics Research Center, North Carolina State
11  University, Raleigh, NC, 27695, USA. [5] Department of Biochemistry and Molecular Genetics,
12  University of Colorado School of Medicine, Aurora, CO, 80045, USA. [*] Corresponding author.
13  Contact: kent.riemondy@cuanschutz.edu

14
15
## ABSTRACT

16  **Background**: In single-cell RNA sequencing (scRNA-seq) analysis, assignment of likely cell
17  types remains a time-consuming, error-prone, and biased process. Current packages for identity
18  assignment use limited types of reference data, and often have rigid data structure
19  requirements. As such, a more flexible tool, capable of handling multiple types of reference data
20  and data structures, would be beneficial.

21

22  **Findings:** To address difficulties in cluster identity assignment, we developed the clustifyr R
23  package. The package leverages external datasets, including gene expression profiles from
24  scRNA-seq, bulk RNA-seq, microarray expression data, and/or signature gene lists, to assign
25  likely cell types. We benchmark various parameters of a correlation-based approach, and also
26  implement a variety of gene list enrichment methods. By providing tools for exploratory data
27  analysis, we demonstrate the feasibility of a simple and effective data-driven approach for cell
28  type assignment in scRNA-seq cell clusters.

29

30  **Conclusions:** clustifyr is a lightweight and effective cell type assignment tool developed for
31  compatibility with various scRNA-seq analysis workflows. clustifyr is publicly available at
32  https://github.com/rnabioco/clustifyr

33
## KEYWORDS

34  Single-cell RNA sequencing, cell type classification, gene expression profile, R package

35

## INTRODUCTION

Single-cell mRNA sequencing promises to deliver improved understanding of cellular mechanisms, cell heterogeneity within tissue, and developmental transitions[1–5]. A key challenge in scRNA-seq data analysis is the identification of cell types from single-cell transcriptomes. Manual inspection of the expression patterns from a small number of marker genes is still standard practice, which is cumbersome and frequently inaccurate. Unfortunately, current implementations of scRNA-seq suffer from several limitations[3,6,7] that further compound the problem of cell type identification.  One, only RNA levels are measured, which may not correlate with cell surface marker or gene expression signatures identified through other experimental techniques. Two, due to the low capture rate of RNAs, low expressing genes may face detection problems regardless of sequencing depth. Many previously established markers of disease or developmental processes suffer from this issue, such as transcription factors. On the data analysis front, over or under- clustering may generate cluster markers that are uninformative for cell type labeling. In addition, cluster markers that are unrecognizable to an investigator may indicate potentially interesting unexpected cell types, but can be very intimidating to interpret.

For these reasons, many investigators struggle to integrate scRNA-seq into their studies due to the challenges of confidently identifying previously characterized or novel cell populations. Formalized data-driven approaches for assigning cell type labels to clusters will greatly aid researchers in interrogating scRNA-seq experiments. Currently, multiple cell type assignment packages exist but they are specifically tailored towards input types or workflows[8–10].

We developed the R package clustifyr, a lightweight and flexible tool that leverages a wide range of prior knowledge of cell types to pinpoint target cells of interest or assign general cell identities to difficult-to-annotate clusters. Here, we demonstrate its applications with transcriptomic information of external datasets and/or signature gene profiles, to explore and quantify likely cell types. The clustifyr package is built with compatibility and ease-of-use in mind to support other popular scRNA-seq tools and formats.

## METHODS

### Extracting information from existing R objects

For clustifyr, query data and reference data can take the form of raw or normalized expression matrices and corresponding metadata tables. To better integrate with standard workflows that involve S3/S4 R objects, methods for clustifyr are written to directly recognize Seurat[11] or SingleCellExperiment[12] objects, retrieve the required information, and reinsert classification results back into an output object (**Fig. 1A**). A more general wrapper is also included for compatibility with other common data structures, and can be easily extended to new object types.

This approach also has the added benefit of forgoing certain calculations such as variable gene selection or clustering, which may already be stored within input objects. clustifyr is designed to perform per-cluster or per-cell classification after previous steps of analysis by other informatics tools. Therefore, it relies on, and is agnostic to, common external packages for cell clustering and variable feature selection. It has been tested against scRNA-seq data analyzed by Seurat[11] and Bioconductor SingleCellExperiment (SCE)[12]. We envision it to be compatible with all scRNA-seq processing, clustering, and marker gene discovery workflows. Simple and non-package-dependent functions for k-means clustering and selection of high variance genes are implemented as alternatives.

### Measuring correlation and comparing gene lists

To assess similarity between query and reference cell types, Spearman, Pearson, Kendall, and Cosine correlation calculations are implemented in clustifyr. Multiple methods are implemented to assess cell identity based on curated gene lists including hypergeometric tests, Jaccard Index, GSEA via the fgsea R package[13], mean percentage of cells that express marker genes, and marker scoring based on mean per-cell Spearman ranked correlation.

### Benchmarking

clustifyr was tested against scmap v1.8.0[8], SingleR v1.0.1[9], and Seurat v3.1.1[11]. scRNA-seq Tabula Muris data was downloaded from https://tabula-muris.ds.czbiohub.org/ as seuratV2 objects. Human pancreas data was downloaded from https://hemberg-lab.github.io/scRNA.seq.datasets/ as SCE objects. In all instances, to mimic the usage case of

98   clustifyr, clustering and dimension reduction projections are acquired from available metadata,

99   in lieu of new analysis.

100   An R script was modified to benchmark clustifyr following the approach and data sets of

101   scRNAseq_Benchmark[14], using M3Drop[15] variable gene selection for every test. R code

102   used for benchmarking, and preprocessing of other datasets, in the form of matrices and tables,

103   are documented in R scripts available in the clustifyr GitHub repository.

104

# FINDINGS

106

107   Prior knowledge of cells types should facilitate cell identity assignment in scRNA-seq analysis.

108   However, in practice, differences between flow cytometry, microarray data, bulk RNA-seq, and

109   the implementations of scRNA-seq, including but not limited to Dropseq, Microwell-seq, 10X

110   genomics 3' end seq, and 5' end seq, make cross-platform comparisons difficult. We therefore

111   set out to build a flexible framework that could compare single-cell transcriptomes across

112   different experimental methods.

113

114   Using clustifyr, which adopts correlation-based methods to find reference transcriptomes with

115   the highest similarity to query cluster expression profiles, peripheral blood mononuclear cell

116   (PBMC) clusters are correctly labeled using either bulk-RNA seq references generated from the

117   ImmGen database[9,16], processed microarray data of purified cell types[17], or previously

118   annotated scRNA-seq results[11] (**Fig. 1B**). We reached similarly satisfactory results in scRNA-

119   seq brain transcriptome data from mouse and human samples, as detailed by

120   scRNAseq_Benchmark[14] (F1-score of 1 for all 4 identity mapping pairs, on 3 main cell types,

121   data not shown).

122

123   To assess the performance of clustifyr, we used the Tabula Muris dataset[5], which contains

124   data generated from 12 matching tissues using both 10x 3' end seq ("drop") and SmartSeq2

125   ("facs") platforms. Using references built from "facs" Seurat objects, we attempted to assign cell

126   type identities to clusters in "drop" Seurat objects. In benchmarking results, clustifyr is

127   comparably accurate versus other automated classification packages (**Fig. 1C**). Cross-platform

128   comparisons are inherently more difficult, and the approach used by clustifyr is aimed at being

129   platform- and normalization-agnostic. Mean runtime, including both reference building and test

130    data classification, in Tabular Muris classifications was ~ 1 second if the required variable gene

131    list is extracted from the query Seurat object (**Fig. 1D**). Alternatively, variable genes can be

132    recalculated by other methods such as M3Drop[15], to reach similar results.

133

134    We further benchmarked clustifyr against a suite of comparable datasets, PBMCbench[18],

135    generated from 2 PBMC samples using multiple scRNA-seq methods. Notably, for each

136    reference dataset cross-referenced to other samples, clustifyr achieved a median F1-score of

137    above 0.94 using Spearman ranked correlation (**Fig. 2A**). Other correlation methods are on par

138    or slightly worse at cross-platform classifications, which is expected based on the nature of

139    ranked vs unranked methods. We therefore selected Spearman as the default method in

140    clustifyr, with other methods also available, as well as a wrapper function to find consensus

141    identities across available correlation methods (see **Fig. 3B**).

142

143    For scRNA-seq reference data, matrices are built by averaging per-cell expression data for each

144    cluster, to generate a transcriptomic snapshot similar to bulk RNA-seq or microarray data. An

145    additional argument to subcluster the reference dataset clusters is also available, to generate

146    more than one expression profile per reference cell type. The number of subclusters for each

147    reference cell type is dependent on the number of cells in the cluster (n), and the sub-clustering

148    power argument (x), following the formula $n^x$[9]. This approach does not improve classification

149    in the PBMCbench data (**Fig. 2B**), however. We envision its utility would greatly depend on the

150    granularity of the clustering in the reference dataset.

151

152    We also tested a general reference set built from the Mouse Cell Atlas[19], and found

153    classification of the Tabula Muris data to be of high accuracy (**Fig. 2C**). Therefore, clustifyr is

154    useful in identity-mapping across different techniques, or simple exploratory analysis using

155    generalized pre-built references. As expected, with further downsampling of the number of cells

156    in each query cluster, we observe decreased accuracy. Yet, even at 15 cells per tested cluster,

157    clustifyr still performed well, with a further increase in speed. Based on these results, we set the

158    default parameters in clustifyr to exclude classification of clusters containing less than 10 cells.

159

160    Recognition of missing reference cell types, so as to avoid misclassification, is another point of

161    great interest in the field. From general usage of clustifyr, we find using a minimum correlation

162    cutoff of 0.5 or 0.4 is generally satisfactory. Alternatively, the cutoff threshold can be determined

163    heuristically using 0.8 * highest correlation coefficient among the clusters. One example is

164   shown in **Fig. 2D**, using benchmark data modified by the SciBet package[20]. Megakaryocytes

165   were removed from reference data, and labeled as "neg.cells" for ground truth in test data.

166   clustifyr analysis found the "neg.cells" to be dissimilar to all available reference cell types, and

167   hence left as "unassigned" under the default minimum threshold cutoff. Next, we applied

168   clustifyr to a series of increasingly challenging datasets from the scRNAseq_Benchmark[14]

169   unseen population rejection test. Without the corresponding cell type references, 57.5% of T

170   cells were rejected and unassigned. When only CD4+ references were removed, 28.2% of test

171   CD4+ T cells were rejected and unassigned. clustifyr was unable to reject CD4+/CD45RO+

172   memory T cells, mislabeling them as CD4+/CD25 T Reg instead when the exact reference was

173   unavailable. However, these misclassifications are also observed with other classification tools

174   benchmarked in the scRNAseq_Benchmark study[14].

175

176   As the core function of clustifyr is ranked correlation, feature selection to focus on highly

177   variable genes is critical. In **Fig. 2E**, we compare correlation coefficients using all detected

178   genes (>10,000), feature selection by the package M3Drop, variable genes selected by Seurat

179   VST (default takes top 2,000), and using 1,000 genes with highest variance in the reference

180   data. As seen, a basic level of feature selection is sufficient to classify the pancreatic cells. In

181   the case of other cell type mixtures, especially ones without complete knowledge of the

182   expected cell types, clustering and feature selection will be of greater importance. clustifyr does

183   not provide novel clustering or feature selection methods on its own, but instead is built to

184   maintain flexibility to incorporate methods from other, and future, packages. We view these

185   questions as fast-moving fields[21,22], and hope to benefit from new advances, while keeping

186   the general clustifyr framework intact.

187

188   Reliable and high-quality full transcriptome datasets are often not available for many cell types

189   and therefore biologists must use a short list of marker genes established from literature to

190   identify cell types. To replace the inefficient experience of plotting the expression of a handful of

191   key marker genes and manually assigning cell types, clustifyr also implements quick methods of

192   gene list enrichment analysis. Using ranked and unranked lists, respectively, clustifyr can

193   correctly annotate PBMC and pancreas scRNA-seq clusters (data not shown). We tested the

194   gene list functionality of clustifyr against the same test of 12 Tabula Muris reference and test

195   pairs, as described above for the ranked correlation approach. With automated marker gene

196   selection, ~85% of clusters were classified correctly (clustifyr_lists in **Fig. 1C**).  In real world use

197   cases, we expect the marker gene lists to be more carefully tailored, and hence perform better.

198   In **Fig. 3A**, we compare the various calculated metrics of clustifyr, using ranked correlation on

199   variable genes or a list of 5 previously established markers, and observe a consensus result

200   identifying the alpha, beta, and delta cell clusters correctly. To combine all analysis, a function

201   assesses consensus results across multiple classification methods within clustifyr and plots

202   consensus cell types (**Fig. 3B**).


203   # CONCLUSIONS

204   We present a flexible and lightweight R package for cluster identity assignment. The tool

205   bridges various forms of prior knowledge and scRNA-seq analysis. Reference sources can

206   include scRNA-seq data with cell types assigned (or average expression per cell type, which

207   can be stored at much smaller file sizes), sorted bulk RNA-seq, microarray data, and ranked or

208   unranked gene lists. clustifyr, with minimal package dependencies, is compatible with a number

209   of standard analysis workflows such as Seurat or Bioconductor, without requiring the user to

210   perform the error-prone process of converting to a new scRNA-seq data structure, and can be

211   easily extended to incorporate other data storage object types. Benchmarking reveals the

212   package performs well in mapping cluster identity across different scRNA-seq platforms and

213   experimental types.

214

215   On the user end, clustifyr is built with simple out-of-the-box wrapper functions, sensible defaults,

216   yet also extensive options for more experienced users. Instead of building an additional single-

217   cell-specific data structure, or requiring specific scRNA-seq pipeline packages, it simply handles

218   basic data.frames (tables) and matrices (**Fig. 1A**). Input query data and reference data are

219   intentionally kept in expression matrix form for maximum flexibility, ease-of-use, and ease-of-

220   interpretation. Also, by operating on predefined clusters, clustifyr has high scalability and

221   minimal resource requirements on large datasets. Using per-cluster expression averages results

222   in rapid classification. However, cell-type annotation accuracy is therefore heavily reliant on

223   appropriate selection of the number of clusters. Users are therefore encouraged to explore cell

224   type annotations derived from multiple clustering settings. Additionally, assigning cell types

225   using discrete clusters may not be appropriate for datasets with continuous cellular transitions

226   such as developmental processes, which are more suited to trajectory inference analysis

227   methods. As an alternative, clustifyr also supports per-cell annotation, however the runtime is

228   greatly increased and the accuracy of the cell type classifications are decreased due to the

229 sparsity of scRNA-seq datasets, and requires a consensus aggregation step across multiple
230 cells to obtain reliable cell type annotations.
231
232 To further improve the user experience, clustifyr provides easy-to-extend implementations to
233 identify and extract data from established scRNA-seq object formats, such as Seurat[11],
234 SingleCellExperiment[12], URD[4], and CellDataSet (Monocle)[23]. Available in flexible wrapper
235 functions, both reference building and new classification can be directly achieved through
236 scRNA-seq objects at hand, without going through format conversions or manual extraction.
237 The wrappers can also be expanded to other single cell RNA-seq object types, including the
238 HDF5-backed loom objects, as well as other data types generated by CITE-seq and similar
239 experiments[24]. Tutorials are documented online to help users integrate clustifyr into their
240 workflows with these and other bioinformatics software.

# AVAILABILITY

242 clustifyr is submitted for review as a Bioconductor package and is licensed under the MIT
243 license. Up-to-date source code, tutorials, and prebuilt references are available at
244 https://github.com/rnabioco/clustifyr. Data used in examples and prebuilt references can also be
245 found at https://github.com/rnabioco/clustifyrdata.

# ABBREVIATIONS

247 PBMC: peripheral blood mononuclear cell; scRNA-seq: single-cell RNA sequencing; SCE:
248 SingleCellExperiment.

# COMPETING INTERESTS

250 The author(s) declare that they have no competing interests.

# FUNDING

256

## AUTHOR'S CONTRIBUTIONS

258   **Software** RF, AEG, RMS, CT, MD, YH, JRH, KAR; **Conceptualization** RF, AEG, KAR;

259   **Writing – original draft** RF**; Writing – review & editing** RF, AEG, RMS, CT, MD, YH, JRH,

260   KAR; **Supervision** JRH;

## FIGURE LEGENDS

262

263   **FIGURE 1. clustifyr uses many types of expression data for cluster identity assignment.**

264   A) Schematic of input data types supported by clustifyr. B) UMAP projections of PBMC cells

265   colored by known cell type (Ground truth cell types) or cell types assigned by clustifyr using

266   reference transcriptome data from microarray, sorted bulk RNA-seq, and scRNA-seq

267   experiments. C) Accuracy of classifications generated by clustifyr or existing methods using the

268   Tabula Muris to benchmark cell type classifications across sequencing platforms. D) Run-time

269   of clustifyr or existing methods on the Tabula Muris cross-platform classification.

270

271   **FIGURE 2. Parameter considerations for clustifyr.** A) Comparison of accuracy of different

272   correlation methods for classifying across platforms using the PBMCbench dataset. B)  An

273   assessment of the accuracy of using single or multiple averaged profiles as reference cell types

274   was conducted using the PBMCbench test set. The number of reference expression profiles to

275   generate for each cell type is determined by the number of cells in the cluster (n), and the sub-

276   clustering power argument (x), with the formula $n^x$.  C) Accuracy and performance were

277   assessed with decreasing number of query cluster cell numbers using the PBMCbench test. D)

278   Heatmap showing correlation coefficients between query cell types and the reference cell types.

279   Clusters with correlation < 0.50 are assigned as Neg.Cell by clustifyr. E) Comparison of

280   classification power using different feature selection methods (M3Drop, Seurat variable gene

281   selection, selection of high variance genes from reference dataset, or no variable gene

282   selection).

283

284   **FIGURE 3. clustifyr implements multiple workflows for cell type classification.** A)

285   Comparison of ranked correlation vs gene list metrics for alpha, beta, and delta cells in

286     pancreatic dataset. B) Consensus cell type calls from using a reference scRNA-seq dataset and

287     gene list methods on alpha, beta, and delta cells in the pancreatic data.

288

# REFERENCES

290     1. Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel
291     digital transcriptional profiling of single cells. Nat Commun. 2017;8:14049.

292     2. Chen G, Ning B, Shi T. Single-Cell RNA-Seq Technologies and Related Computational Data
293     Analysis. Front Genet. 2019;10:317.

294     3. Luecken MD, Theis FJ. Current best practices in single-cell RNA-seq analysis: a tutorial. Mol
295     Syst Biol. 2019;15:e8746.

296     4. Farrell JA, Wang Y, Riesenfeld SJ, Shekhar K, Regev A, Schier AF. Single-cell reconstruction
297     of developmental trajectories during zebrafish embryogenesis. Science [Internet]. 2018;360.
298     Available from: https://www.ncbi.nlm.nih.gov/pubmed/29700225

299     5. The Tabula Muris Consortium, Overall coordination, Logistical coordination, Organ collection
300     and processing, Library preparation and sequencing, Computational data analysis, Cell type
301     annotation, Writing group, Supplemental text writing group & Principal investigators. Single-cell
302     transcriptomics of 20 mouse organs creates a Tabula Muris. Nature. 2018;562:367–72.

303     6. Kiselev VY, Andrews TS, Hemberg M. Challenges in unsupervised clustering of single-cell
304     RNA-seq data. Nat Rev Genet. 2019;20:273–82.

305     7. Vallejos CA, Risso D, Scialdone A, Dudoit S, Marioni JC. Normalizing single-cell RNA
306     sequencing data: challenges and opportunities. Nat Methods. 2017;14:565–71.

307     8. Kiselev VY, Yiu A, Hemberg M. scmap: projection of single-cell RNA-seq data across data
308     sets. Nat Methods. 2018;15:359–62.

309     9. Aran D, Looney AP, Liu L, Wu E, Fong V, Hsu A, et al. Reference-based analysis of lung
310     single-cell sequencing reveals a transitional profibrotic macrophage. Nat Immunol.
311     2019;20:163–72.

312     10. Pliner HA, Shendure J, Trapnell C. Supervised classification enables rapid annotation of cell
313     atlases. Nat Methods. 2019;16:983–6.

314     11. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic
315     data across different conditions, technologies, and species. Nat Biotechnol. 2018;36:411–20.

316     12. Lun AT, McCarthy DJ, Marioni JC. A step-by-step workflow for low-level analysis of single-
317     cell RNA-seq data with Bioconductor. F1000Res. 2016;5:2122.

318     13. Sergushichev AA. An algorithm for fast preranked gene set enrichment analysis using
319     cumulative statistic calculation [Internet]. bioRxiv. 2016 [cited 2019 Nov 14]. p. 060012.
320     Available from: https://www.biorxiv.org/content/10.1101/060012v1

321   14. Abdelaal T, Michielsen L, Cats D, Hoogduin D, Mei H, Reinders MJT, et al. A comparison of
322   automatic cell identification methods for single-cell RNA sequencing data. Genome Biol.
323   2019;20:194.

324   15. Andrews TS, Hemberg M. M3Drop: dropout-based feature selection for scRNASeq.
325   Bioinformatics. 2019;35:2865–7.

326   16. Heng TS, Painter MW, Immunological Genome Project, Consortium. The Immunological
327   Genome Project: networks of gene expression in immune cells. Nat Immunol. 2008;9:1091–4.

328   17. Novershtern N, Subramanian A, Lawton LN, Mak RH, Haining WN, McConkey ME, et al.
329   Densely interconnected transcriptional circuits control cell states in human hematopoiesis. Cell.
330   2011;144:296–309.

331   18. Ding J, Adiconis X, Simmons SK, Kowalczyk MS, Hession CC, Marjanovic ND, et al.
332   Systematic comparative analysis of single cell RNA-sequencing methods. bioRxiv.
333   2019;632216.

334   19. Han X, Wang R, Zhou Y, Fei L, Sun H, Lai S, et al. Mapping the Mouse Cell Atlas by
335   Microwell-Seq. Cell. 2018;172:1091–107 e17.

336   20. Li C, Liu B, Kang B, Liu Z, Liu Y, Ren X, et al. SciBet: a fast classifier for cell type
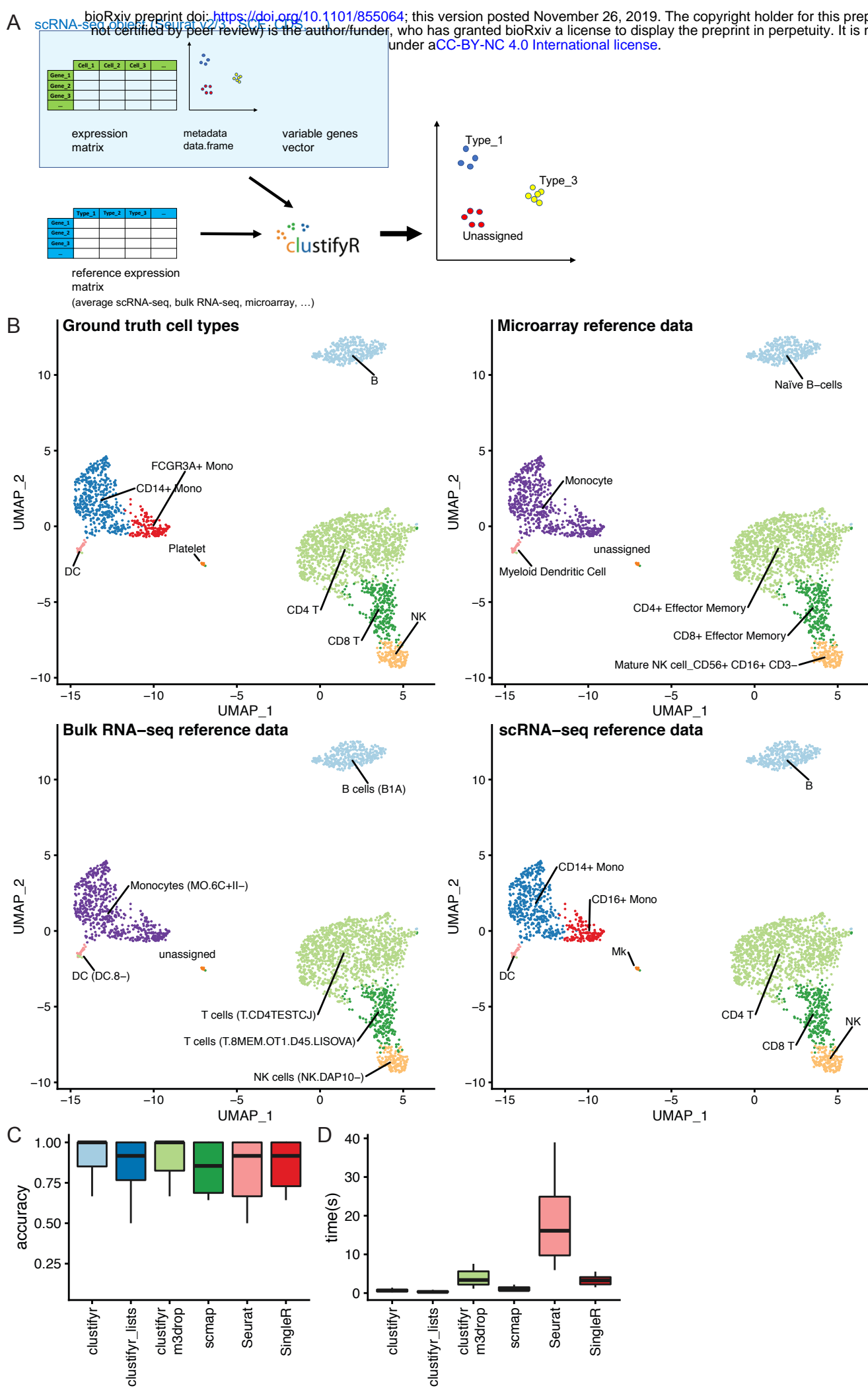337   identification using single cell RNA sequencing data. bioRxiv. 2019;645358.

338   21. Duò A, Robinson MD, Soneson C. A systematic performance evaluation of clustering
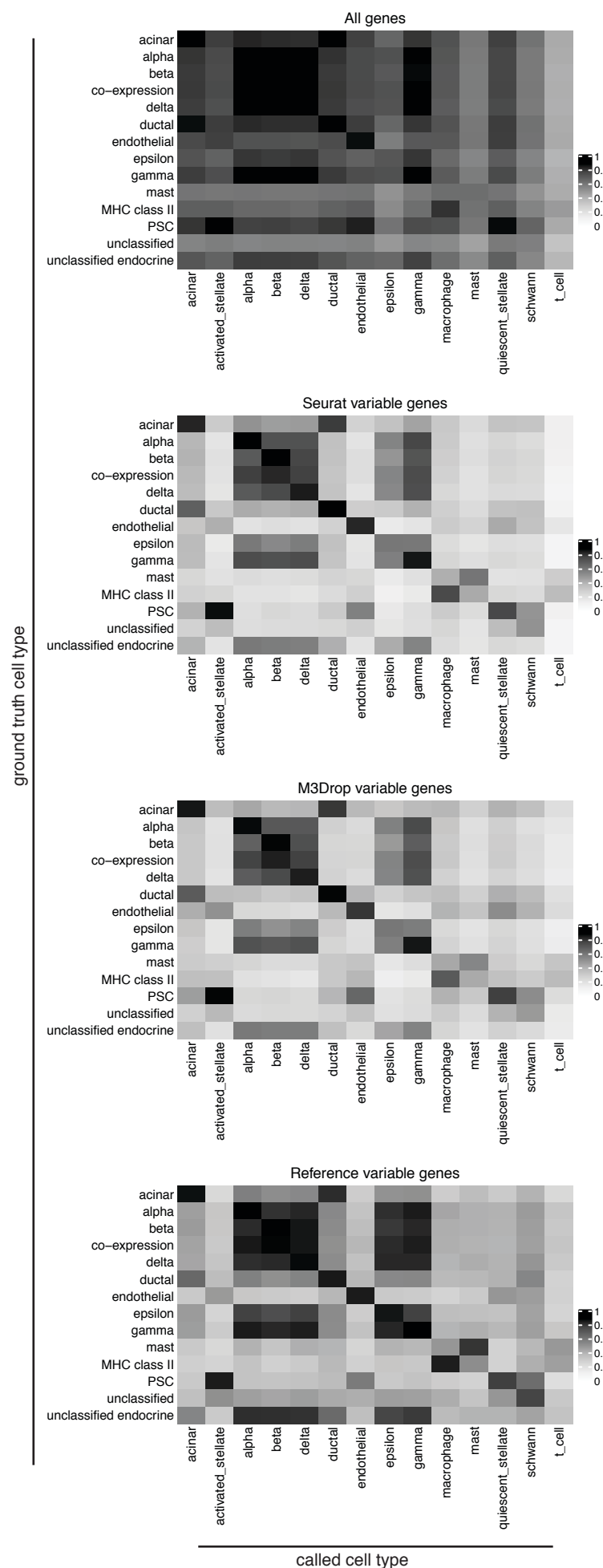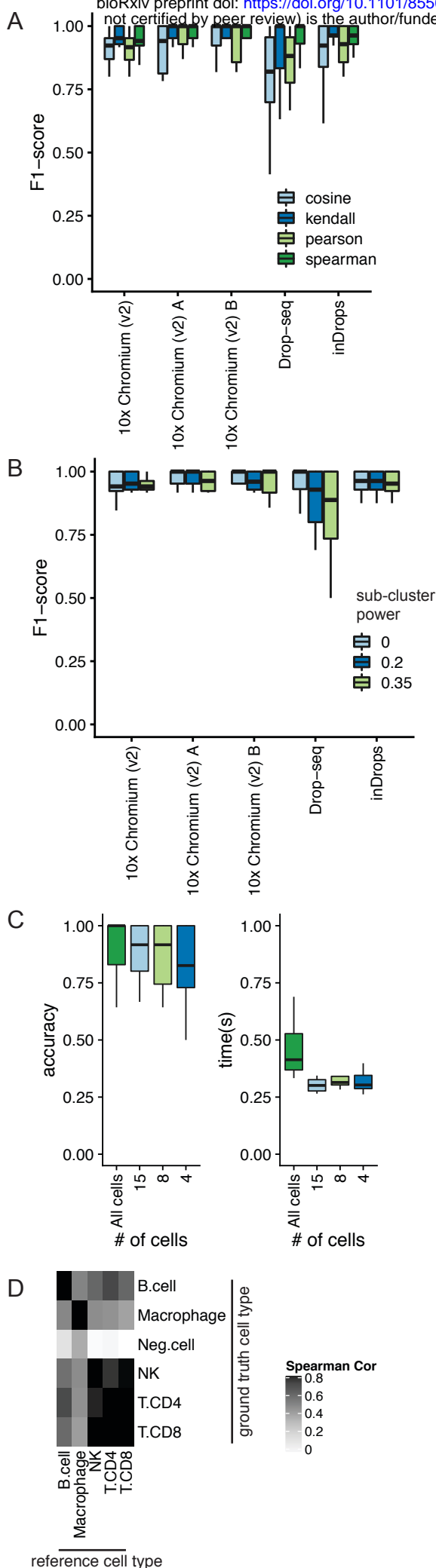339   methods for single-cell RNA-seq data. F1000Res. 2018;7:1141.

340   22. Soneson C, Robinson MD. Bias, robustness and scalability in single-cell differential
341   expression analysis. Nat Methods. 2018;15:255–61.

342   23. Cao J, Spielmann M, Qiu X, Huang X, Ibrahim DM, Hill AJ, et al. The single-cell
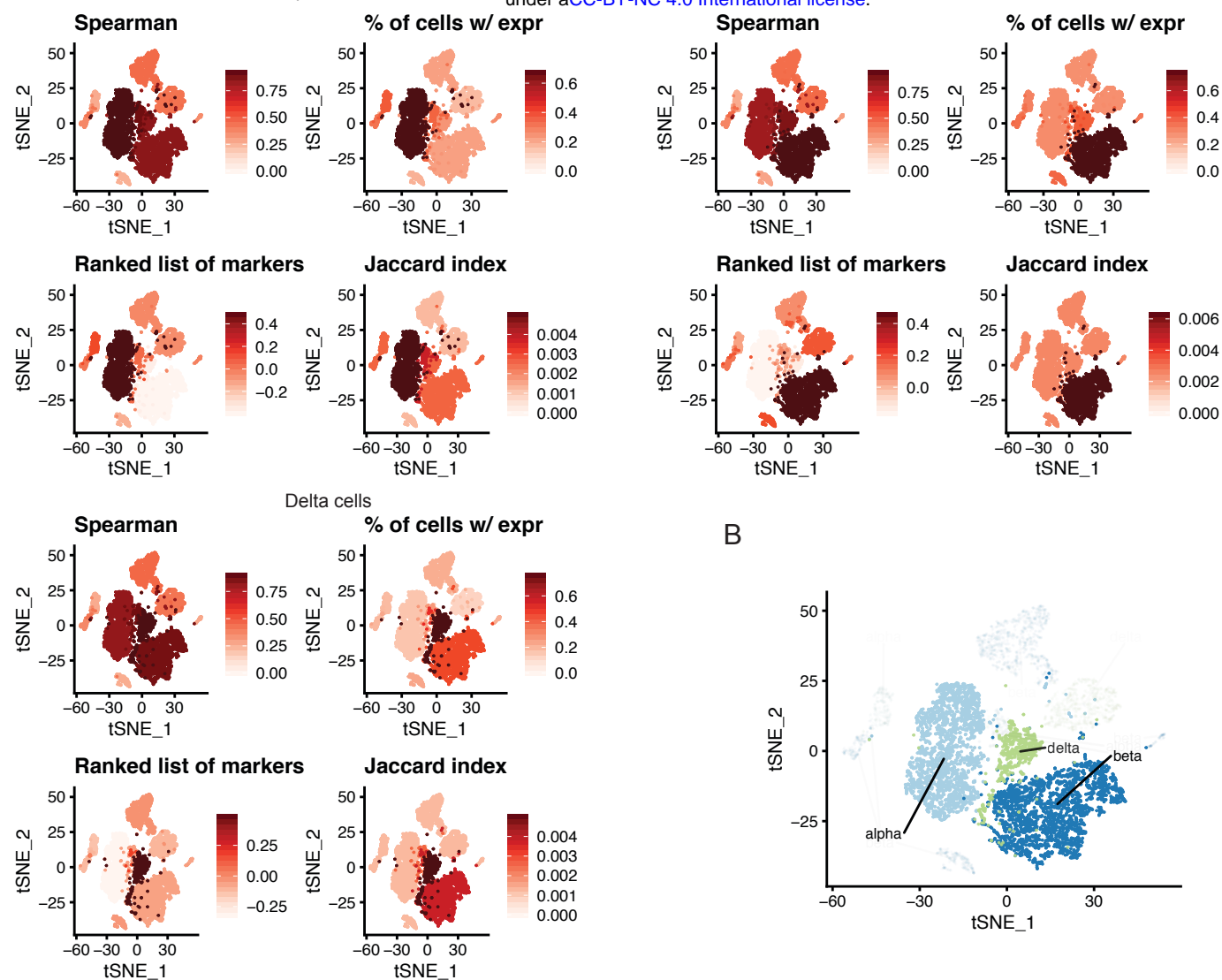343   transcriptional landscape of mammalian organogenesis. Nature. 2019;566:496–502.

344   24. Richer AL, Riemondy KA, Hardie L, Hesselberth JR. Simultaneous measurement of
345   biochemical phenotypes and gene expression in single cells. bioRxiv. 2019;820233.

Fig 1.

Fig 2.

Fig 3.

**A**

Alpha cells



Beta cells



Delta cells



**B**