1  # CiteFuse enables multi-modal analysis of CITE-seq data

2  Hani Jieun Kim[1,2,3,+], Yingxin Lin[1,2,+], Thomas A. Geddes[2,4], Jean Yang[1,2] and Pengyi Yang[1,2,3,*]

3  [1] School of Mathematics and Statistics, The University of Sydney, NSW, Australia

4  [2] Charles Perkins Centre, The University of Sydney, NSW, Australia

5  [3] Computational Systems Biology Group, Children's Medical Research Institute, Faculty of Medicine
6  and Health, The University of Sydney, Westmead, NSW, Australia

7  [4] School of Life and Environmental Sciences, The University of Sydney, NSW, Australia

8  [+] These authors contributed equally

9  * Corresponding author (pengyi.yang@sydney.edu.au)

10  ## Abstract

11  Multi-modal profiling of single cells represents one of the latest technological advancements in

12  molecular biology. Among various single-cell multi-modal strategies, cellular indexing of transcriptomes

13  and epitopes by sequencing (CITE-seq) allows simultaneous quantification of two distinct species: RNA

14  and surface marker proteins (ADT). Here, we introduce CiteFuse, a streamlined package consisting of

15  a suite of tools for pre-processing, modality integration, clustering, differential RNA and ADT expression

16  analysis, ADT evaluation, ligand-receptor interaction analysis, and interactive web-based visualization

17  of CITE-seq data. We show the capacity of CiteFuse to integrate the two data modalities and its relative

18  advantage against data generated from single modality profiling. Furthermore, we illustrate the pre-

19  processing steps in CiteFuse and in particular a novel doublet detection method based on a combined

20  index of cell hashing and transcriptome data. Collectively, we demonstrate the utility and effectiveness

21  of CiteFuse for the integrative analysis of transcriptome and epitope profiles from CITE-seq data.

22  **Keywords:** CITE-seq, ADT, single-cell, integration, multi-modality, multi-omic, doublet detection,

23  ligand-receptor interaction

## Introduction

The latest advancement in multi-modal profiling of single cells promises to revolutionise our understanding in cellular biology that was previously inconceivable through bulk profiling technologies (Datlinger et al., 2017; Macaulay et al., 2015; Mohammed et al., 2017). Among various single-cell multi-modal strategies, cellular indexing of transcriptomes and epitopes by sequencing (CITE-seq) (Stoeckius et al., 2017) and its variants such as RNA expression and protein sequencing (REAP-seq) (Peterson et al., 2017) represent a class of approaches that allows simultaneous quantification of global gene expression and cellular proteins using single-cell RNA-sequencing (scRNA-seq) and antibody-derived tags (ADTs), respectively, on single cells. Further extensions such as multiplexed detection of proteins, transcriptomes, clonotypes and CRISPR perturbations enable additional modalities to be profiled on single cells (Mimitou et al., 2019).

While the surface proteins of individual cells measured by ADTs are also transcriptomically profiled by scRNA-seq, the measurements of these two different molecule species produced from the same genes do not necessarily correlate with each other, presumably because of post-transcriptional and post-translational gene regulation (See, Lum, Chen, & Ginhoux, 2018). Therefore, computational integration of single cell multi-modal profiling data may allow a more accurate characterisation of cells (e.g., cell type identification) (Buettner et al., 2015) and provide new biological insights that may be observable from neither a single data source (Lin et al., 2019) nor modality (Stuart et al., 2019).

Here we present CiteFuse, a computational framework that implements a suite of methods and tools for CITE-seq data from pre-processing through to integrative analytics. This includes doublet detection, network-based modality integration, cell type clustering, differential RNA and ADT expression analysis, ADT evaluation, ligand-receptor interaction analysis, and interactive web-based visualisation of the analyses (**Figure 1A**). Using both simulations and an experimental CITE-seq dataset generated from PBMCs (Mimitou et al., 2019), we demonstrate the integrative capacity of CiteFuse in various scenarios and its advantage over analysing each individual source and modality of data. CiteFuse represents the first method specifically designed to systematically integrate RNA and ADT modalities of single cells in CITE-seq data. We anticipate its increasing utility given the rapidly accumulating volume of multi-omic and multi-modality single cell data generated using CITE-seq from various biological studies (Mimitou et al., 2019; Stoeckius et al., 2017). Finally, CiteFuse is implemented as an R package

53  (http://SydneyBioX.github.io/CiteFuse/) as well as a user-friendly web application

54  (http://shiny.maths.usyd.edu.au/CiteFuse/), allowing users to upload and analyse their CITE-seq

55  datasets.

## Results

**CiteFuse gains information from multi-modal integration of CITE-seq data**

58  To take advantage of the complementary information present in multi-modal CITE-seq data, CiteFuse

59  integrates mRNA and ADT expression by constructing networks across single cells for each data

60  modality and fusing these networks using a similarity network fusion algorithm (Wang et al., 2014)

61  (**Figure 1A, blue tile**). It subsequently uses a spectral clustering algorithm to cluster the cells based on

62  the fused matrix. To test whether there is any advantage in using the fused multi-modal expression

63  matrix over the single-modal matrices, we performed a comparison between the different modalities

64  and across existing clustering algorithms with simulated CITE-seq data (Zhang et al., 2019) (**Figure

65  S1A**). We demonstrate that in both "easy" and "hard" scenarios (see Methods), CiteFuse clusters cells

66  more accurately than directly applying spectral clustering on the two single-modal data types (**Figure

67  S1B**). Moreover, we demonstrate that CiteFuse performs better compared to several established

68  clustering procedures, including SIMLR (Wang et al., 2017), PCA + *k*-means, and Seurat (Satija et al.,

69  2015) with either RNA or ADT expression matrix (**Figure S1B**).

70  To test if the information gain from multi-modal analysis using CiteFuse observed from the simulation

71  study translates into real-world data analysis, we next applied CiteFuse to a recent human PBMC CITE-

72  seq dataset (Mimitou et al., 2019) (**Figure 1B**). We show that clustering using CiteFuse on multi-modal

73  data and directly applying spectral clustering on single-modal (ADT or RNA) data lead to different

74  clustering outcomes (**Figure S2A**). We found that CiteFuse can generate four CD4+ T-cell clusters

75  (**Figures 1C and S2B**), of which three are CD4+ memory T-cells (clusters 2, 9, and 16) expressing high

76  level of S100A4 (a marker of memory T-cells) and one is CD4+ naive T-cells (cluster 14) expressing

77  high level of SELL (a marker of naive T-cells) (Elyahu et al., 2019; Haining et al., 2008) (**Figure S2C**).

78  In contrast, clustering using ADT alone leads to over-partitioning of CD4+ T-cells into five clusters and

79  clustering using RNA alone leads to under-partitioning of these cells into three clusters (**Figures S2B**

80  **and S2C**).

81    Moreover, we observed that clustering using RNA-alone fails to partition CD27+ and CD27- populations

82    of memory T-cells, whilst clustering using CiteFuse or ADT-alone can discriminate these two

83    populations, albeit to different resolutions (**Figures 1D and S2B**). A closer examination of the CD27-

84    CD4+ memory T-cell subpopulations (**Figure 1C**; light and dark blue clusters in CiteFuse; light blue

85    cluster in ADT alone) reveals that only CiteFuse can discriminate between CD27- DR+ (light blue) and

86    CD27- DR- (dark blue) memory T-cell subpopulations (Fonseka et al., 2018) (**Figures 1D and S2D**),

87    revealing that only CiteFuse has the capacity to finely map T-cell subpopulations and further

88    demonstrates the gain in information CiteFuse benefits from multi-modal analysis.

89    **CiteFuse detects both cross- and within-sample doublets**

90    Identification and removal of doublets from scRNA-seq data derived from microfluidic technology is

91    essential for downstream analysis. Cell hashing is a multiplexing technique commonly used in CITE-

92    seq for pooling multiple samples (Stoeckius et al., 2018). Because a key principle in cell hashing is the

93    selection of ubiquitously and highly expressed surface markers, against which distinct hashtag

94    oligonucleotide (HTO)-conjugated antibodies are raised, the high number of the ubiquitous epitopes

95    raises the possibility of utilising HTO-derived expression to detect within-sample doublets marked by

96    anomalous HTO expression. To this end, CiteFuse takes advantage of the matched matrices for RNA,

97    ADT, and HTO expression generated from a CITE-seq experiment (**Figure 2A**) and implements a

98    stepwise approach to detect and filter both cross- and within-sample doublets (**Figure 2B**). In the first

99    step, a Gaussian mixture model is used to identify cross-sample doublets that have more than one

100   hashtag (i.e. stained by orthogonal HTOs) (**Figure S3A**). Next, by leveraging the ubiquitous nature of

101   HTO expression, CiteFuse detects within-sample doublets from DBSCAN clustering of single cells

102   based on two features—total number of captured unique molecular identifiers (UMIs) and total HTO

103   expression (**Figure 2B**). Data are filtered in step one based on the mixture modelling step for cross-

104   sample doublets and then based on a baseline HTO threshold calculated through the Gaussian mixture

105   model for within-sample doublets (see Methods).

106   We benchmarked our doublet filtering approach with alternative methods, HTODemux (Stoeckius et al.,

107   2018) and Scrublet (Wolock et al., 2019), on the PBMC dataset (Mimitou et al., 2019) by demonstrating

108   that doublets/multiplets detected through CiteFuse show comparably high number of unique genes and

109   UMIs (**Figure S3B**). Notably, we show that the within-sample doublets identified through CiteFuse

110    represent outlier cells that have both high total UMIs and high HTO expression (**Figure S3C**). We show

111    that our approach captures most doublets detected through HTODemux and Scrublet but also identifies

112    additional ones that may have been missed by HTODemux and Scrublet (**Figure S3D**). When we

113    quantified the total UMIs and number of unique genes in cells exclusively identified by each method

114    (**Figure S3E**), we found that doublets exclusively detected by HTODemux and Scrublet show

115    characteristics that resemble singlets whereas those only detected by CiteFuse resemble doublets

116    (**Figures S3B and S3E**).

117    Strikingly, we observed the most improved separation of clusters on the first two principal components

118    of HTO expression before and after filtering of doublets detected by CiteFuse (**Figure 2C**), suggesting

119    our CiteFuse pipeline enables more accurate filtering of both within- and cross-sample doublets when

120    HTO libraries are available.

121    **CiteFuse doublet filtering preserves the separation between T-cell subpopulations**

122    To evaluate the impact of filtering method on the downstream analysis, we applied CiteFuse clustering

123    on data either unfiltered (4292 cells) or filtered using the different doublet detection methods—

124    HTODemux (3753 cells), Scrublet (3968 cells), and CiteFuse (3612 cells). Visualisation of the clusters

125    on UMAP revealed very different clustering outcomes by each filtering method, revealing that filtering

126    method can have a large impact on downstream analysis (**Figure S4A**).

127    We demonstrate the impact of filtering method on downstream analysis by evaluating the capacity of

128    the unfiltered and filtered datasets to define CD4+ and CD8+ T-cell types, two major groups of T

129    lymphocytes. We found that the CiteFuse-filtered dataset leads to the best separation of CD4+ (clusters

130    2, 9, 14, and 16) and CD8+ (clusters 3, 7, 10, and 15) T-cell populations on the basis of purity scores

131    (**Figure S4B**). Moreover, our results showed that the CiteFuse-filtered dataset can further discriminate

132    CD27+ and CD27- subpopulations within CD4+ and CD8+ T-cells (**Figure S4C-E**). Surprisingly, we

133    observed that HTODemux- and Scrublet-filtered datasets have low capacity to discriminate between

134    CD4+ and CD8+ T-cells, let alone CD27+ and CD27- subpopulations within each of the major T-cell

135    populations (**Figure S4C-E**).

136    **CiteFuse enables the evaluation of ADTs and visualisation of ADT-RNA networks**

137    The selection of a set of ADTs for CITE-seq may be an expensive process, requiring in many cases

138    optimisation through flow cytometry for antibody concentration and selection. To maximise the selection

139    of ADTs for subsequent CITE-seq experiments, CiteFuse implements a set of evaluation tools that

140    enables CITE-seq end-users to assess ADTs for relative importance and potential redundancy (**Figure**

141    **2D**). This includes correlating and visualising ADTs based on their expressions (**Figure S5A**) as well

142    as calculating the relative importance of individual ADTs based on CiteFuse clustering outcome using

143    a random forest model (see Methods) (**Figure S5B**). For example, in the PBMC CITE-seq dataset, we

144    found that CD223 and IgG1 are the two ADTs receiving the lowest importance scores and therefore

145    may not provide much additional information for cell type clustering. Indeed, we observed minimum

146    changes in the clustering outcome (ARI=0.99) even without the two ADTs (**Figure S5C**). We find that

147    more ADTs can be excluded (Subsets 2-3) with minimal effect on clustering results. In addition to ADT

148    evaluation, CiteFuse can also perform cluster-specific differential gene expression analysis to detect

149    and compare differentially expressed RNA and ADT (**Figure 2E**) and generate visualisation of ADT-

150    RNA correlation networks unique to each cluster, allowing users to evaluate relationships between ADT

151    and RNA in an intra-cluster manner (**Figure 2F**).

152    **CiteFuse facilitates accurate identification of ligand-receptor interactions**

153    Most studies on ligand-receptor interaction in single-cell biology rely solely on mRNA expression

154    (Vento-Tormo et al., 2018), thereby making an implicit assumption that the level of mRNA expression

155    is a proxy for the cell-surface protein expression. Yet studies have shown that the levels of mRNA and

156    proteins of the same gene can vary widely (Gry et al., 2009; Liu, Beyer, & Aebersold, 2016). In case of

157    cell-surface proteins, this is further complicated by the amount of proteins translocated to plasma

158    membrane. CITE-seq opens the possibility to use protein expression at the cell-surface to predict

159    ligand-receptor interactions. To this end, we predicted ligand-receptor interactions based on mRNA

160    expression of the ligand and ADT expression of the receptor, after normalisation and scaling of the

161    mRNA and ADT expression data (see Methods) (**Figures 2G and S6A**). We compared the ligand-

162    receptor interactions identified by CiteFuse with those identified from the conventional approach where

163    the expression of RNA alone is used as a readout for both ligand and receptor expression (**Figure S6A**).

164    We found that the overlap in interactions between the conventional approach and CiteFuse was variable

165    across clusters, but generally a large portion of the ligand-receptor interactions identified through the

166   conventional approach (referred to as RNA-specific) were not identified as interactions through

167   CiteFuse (**Figure S6B**). We also observed in each cluster a fraction of interactions that were identified

168   only by CiteFuse (referred to as CiteFuse-specific) (**Figure S6B**).

169   We then hypothesised that the large proportion of interactions in the conventional approach that are not

170   detected by CiteFuse may be because of false positive predictions. To investigate this, we calculated

171   the normalised log expression of the ADT and mRNA of all receptors that were identified in a ligand-

172   receptor interaction for each category (CiteFuse-specific, RNA-specific, and Common). We found that

173   although the mRNA expression of the receptors was comparable between the categories the ADT

174   expression of these receptors was much lower in the RNA-specific group than the other two groups

175   (**Figure S6C**). Notably, we found that a strong positive correlation of ADT and mRNA expression

176   (ranked relative to each cluster; see Methods) for receptors identified in a ligand-receptor interaction in

177   the Common and CiteFuse-specific categories but no correlation for those in the RNA-specific category

178   (**Figure S6D**). Similarly, we show that the mRNA expression of ligands detected in the RNA-specific

179   category have higher rankings than those detected in the other two categories (**Figure S6E**). These

180   data show that interactions identified through the conventional approach, which relies on RNA

181   expression alone, may introduce false interactions. These false interactions may potentially be driven

182   by high RNA expression that is not reciprocated in the cell-surface protein expression and thus

183   demonstrates the need to utilise both mRNA and ADT expression in ligand-receptor interaction

184   predictions (**Figure S6F**).

## Methods

### Integration of CITE-seq data through similarity network fusion and spectral clustering

187   To integrate multi-modal CITE-seq data, CiteFuse first normalises the ADT expression through centred

188   log-ratio (CLR) transformation. It next calculates cell-to-cell similarity matrices from ADT expression

189   using *perb* similarity metric from the *propr* package (Quinn, Richardson, Lovell, & Crowley, 2017) and

190   RNA expression using Pearson's correlation on highly variable genes identified with the *scran* package

191   (Lun, McCarthy, & Marioni, 2016). The two similarity matrices are scaled using an exponential similarity

192   kernel and then fused by a similarity network fusion algorithm (Wang et al., 2014).

193    CiteFuse performs spectral clustering (Ng, Jordan, & Weiss, 2002) to identify clusters from the fused

194    similarity matrix. Spectral clustering on single-modal matrices from CITE-seq data were performed for

195    comparison. As well as spectral clustering, CiteFuse also provides the additional option of Louvain

196    clustering (Blondel, Guillaume, Lambiotte, & Lefebvre, 2008), which is an algorithm based on the shared

197    nearest neighbours, which CiteFuse identifies from the fused similarity matrix. Finally, UMAP or tSNE

198    can be applied to the fused similarity matrix to visualise the multi-modal data.

199    **CITE-seq data simulation and evaluation of CiteFuse**

200    To evaluate the integrative capacity of CiteFuse, we simulated CITE-seq data with SymSim (Zhang, Xu,

201    & Yosef, 2019) and assessed the difference in clustering outcome between the modality of data and

202    also by different clustering methods.

203    For each simulation, we generated a dataset of 500 single cells among which were six cell types where

204    total numbers of RNA and ADT were 10,000 and 100, respectively. The following parameter settings

205    for sigma ($\sigma$), which controls within-population variability, and minimum population size (min_pop) were

206    used to simulate CITE-seq data of different levels of difficulty.

207    • Simulation 1 (easy): $\sigma$ (RNA) = 0.8; $\sigma$ (ADT) = 0.2; and min_pop = 50

208    • Simulation 2 (hard): $\sigma$ (RNA) = 0.9; $\sigma$ (ADT) = 0.4; and min_pop = 20

209    We generated 10 datasets for each simulation setting and benchmarked CiteFuse against spectral

210    clustering on single-modal matrices and three different clustering methods: *k*-means clustering on PCA

211    reduced dimension (PCA + *k*-means), SIMLR (Wang, Zhu, Pierson, Ramazzotti, & Batzoglou, 2017)

212    and Seurat (Satija, Farrell, Gennert, Schier, & Regev, 2015). For PCA + *k*-means, *k*-means clustering

213    was performed on the first 10 principal components. For *k*-means clustering and SIMLR, the number of

214    clusters was set as six so to be consistent with the simulation set-up. While for Seurat, we set the

215    resolution parameters between 1.5 and 2 such that the number of communities detected by Louvain

216    clustering is consistent with the number of cell types in the simulations. The concordance in clustering

217    outcome was evaluated as the adjusted rand index (ARI), where a higher index indicates better

218    clustering performance.

219    **CITE-seq data from healthy human PBMCs**

220      To demonstrate our method, we used the recently published CITE-seq data (Mimitou et al., 2019).

221      Specifically, we used the ECCITE-seq dataset from PBMC samples isolated from the blood of healthy

222      human controls. The samples from the human healthy PBMC datasets were pooled from 10x libraries

223      with four distinct barcodes, representing the four hashtag oligonucleotides (HTO) used in the cell

224      hashing.

225      **Calculation of signature scores for T-cell subpopulations**

226      To calculate the signature scores for the various immune populations, we averaged the expression of

227      the following sets of genes that were previously defined as marker genes for the respective cell types

228      of interest:

229          (1) S100A4, CRIP1, and AHNAK were used to define memory CD4+ T-cells (Elyahu et al., 2019;

230              Haining et al., 2008);

231          (2) TCF, ID3, CCR7, and SELL were used to define naive CD4+ T-cells (Elyahu et al., 2019;

232              Haining et al., 2008);

233          (3) GNLY, GZMB, PRF1, GZMA, NKG7, HLA-DRB1, and HLA-DPA1 were used to define CD4+

234              CD27- DR+ T-cells (Fonseka et al., 2018);

235      **CiteFuse doublet detection approach**

236      CiteFuse implements a stepwise procedure to identify both the cross-sample doublets and within-

237      sample doublets from CITE-seq data when cell hashing data is available.

238          (1) Cross-sample doublet identification

239              First, we fit a two-component Gaussian mixture model to each log-transformed HTO expression.

240              The intersection point defined from the mixture model is used to categorise each cell in terms

241              of whether the HTO is either highly or lowly expressed. The cells found to have a single highly

242              expressed HTO are considered as singlets whilst those that have two or more highly expressed

243              HTOs are considered as doublets or multiplets. Cells without any highly expressed HTOs are

244              considered as empty droplets.

245          (2) Within-sample doublet identification

246              Data filtered by cross-sample doublets are next subject to within-sample doublet identification

247              using a density-based spatial clustering and noise detection algorithm (DBSCAN) on an HTO-

248    specific matrix comprising of two features—total number of UMIs and log-transformed HTO

249    expression. The two parameters used in the DBSCAN for this study are eps = 190 and minPts

250    = 50. This procedure is repeated for each HTO and the smallest cluster from DBSCAN

251    clustering is assigned as within-sample doublets.

252    We benchmarked our doublet detection method against two existing methods: HTODemux (Stoeckius

253    et al., 2018) from the Seurat package and Scrublet (Wolock, Lopez, & Klein, 2019). We used the default

254    parameter settings for sim_doublet_ratio and n_neighbors to construct the KNN classifier to simulate

255    doublets with Scrublet by following their online tutorial (https://github.com/AllonKleinLab/scrublet/) and

256    set an expected doublet rate of 0.04. We compared the total number of UMI and the number of unique

257    expressed genes for each cell by each method (HTODemux, Scrublet, and CiteFuse). To compare the

258    effect of filtering method on the downstream analysis, we performed spectral clustering on the output

259    of the similarity network fusion and calculated the purity score of CD8+ cells against CD4+ cells in

260    individual clusters for each filtering method.

261    **Calculation of purity score**

262    To calculate the purity of CD4+ and CD8+ T-cell populations, we first identified CD4+ and CD8+ T-cells

263    by creating a Gaussian mixture model on expression of CD4, CD8, and CD11c. For CD4+ T-cells, we

264    created a Gaussian mixture model of CD4 and CD11c expression to define the CD4+ CD11c-

265    population. For CD8+ T-cells, the same approach was employed but with CD8 and CD11c expression.

266    Using the threshold calculated from the mixture model, cells were assigned as either CD4+ negative or

267    positive cells and CD8+ negative or positive cells. Next, using the CD4+ and CD8+ T-cell labels, we

268    calculated the purity of each cluster for either CD4 or CD8 T-cells. A purity of 1 denotes a cluster

269    composed purely of either CD4 or CD8 T-cells, and a purity score of 0 denotes a cluster devoid of either

270    cell type.

271    **Analysis and visualisation of differentially expressed RNA and ADT**

272    To identify the differentially expressed mRNA and ADTs for each cluster, we used the Wilcoxon rank

273    sum test to compare the log-transformed expression of mRNA and ADT for each cluster against all

274    other clusters. The p-values were adjusted using the Benjamini and Hochberg method (Benjamini &

275    Hochberg, 1995).

276     For the selection of RNA and ADT markers for a given cluster, we considered the following three criteria:

277     (1) An adjusted *P*-value of lower than 0.05;

278     (2) The mean expression of RNA and ADT in the cells of the cluster is greater than the mean

279         expression of RNA and ADT in cells of all other clusters; and

280     (3) The proportion of cells in the cluster expressing the RNA and ADT is greater than the proportion

281         of cells expressing the RNA and ADT across all other clusters by at least 10%.

282     CiteFuse enables two exploration methods to visualise the results of differential expression analysis for

283     both RNA and ADT in a single plot:

284     (1) DEcomparisonPlot

285         The DEcomparisonPlot visualises the positive log10 transformed adjusted P-values as a dot of

286         the RNA and the negative log10 transformed adjusted p-values of its corresponding ADT signal

287         on the same y-axis.

288     (2) DEbubblePlot

289         We used the circlepack plot to visualise the RNA and ADT markers, where each marker is

290         represented by a circle and the size of the circle represents the magnitude of the negative log10

291         *P*-value. The circles representative of RNA and ADT markers from the same clusters are then

292         grouped into a larger circle, representing individual clusters. The circlepack plots are generated

293         using the R package *ggraph* (Pedersen, 2017).

294     **ADT-RNA correlation network construction**

295     To construct the ADT-RNA co-expression network, we calculated the Pearson's correlation between

296     mRNA and ADT expression. Other correlation calculation methods, such as the Spearman and Kendall

297     correlation, are also available as options in our CiteFuse package. ADT-RNA pairs with high absolute

298     correlation (above a default setting of 0.6) are used to construct the ADT-RNA correlation network. The

299     networks are visualised using R packages, igraph (Csardi & Nepusz, 2006) and visNetwork (Almende

300     & Thieurmel, 2016).

301     **Evaluation of ADT importance**

302    To evaluate the importance for each ADT towards the clustering outcome, we trained a random forest

303    model on a subset of randomly sampled cells (80% of total), using the clustering labels from the

304    similarity network fusion of the PBMC CITE-seq data. After 50 repeated fitting of the random forest

305    model, we quantified the feature importance in terms of the mean decrease in Gini index as a surrogate

306    of the importance of each ADT towards clustering outcome. We defined ADT importance score as the

307    median of the feature importance of all runs. A higher score indicates greater importance of the ADT.

308    Next, to identify potentially redundant ADTs that do not contribute significantly towards clustering

309    outcome, we sorted the ADTs by importance and drew cut-offs in accordance to the local maximums

310    of the difference in importance scores. We then retained the subset of ADTs the with importance scores

311    greater than the cut-offs and performed similarity network fusion analysis. We calculated the adjusted

312    rand index (ARI) to measure the concordance in clustering outcome for each subset of ADTs against

313    that of the full dataset.

314    **Ligand-receptor interaction prediction**

315    One of the key challenges in analysing ligand-receptor relationships between two modalities is the

316    difference in scaling and distribution. To address this, we first scaled each feature into a range of 0 to

317    1 through min-max normalisation. Specifically, for every value of a feature $x$ across all single cells, the

318    normalised expression $z$ is calculated by

319    $z = \frac{x - min(x)}{max(x) - min(x)}.$

320    Another challenge we encountered was the difference in distribution between the two modalities: we

321    observed that the distribution of mRNA expression tends to be more zero-inflated than ADT expression.

322    Because comparing unequal distributions has the potential to introduce bias, especially during ligand-

323    receptor predictions when the mean expression is compared, we thus performed another step of

324    transformation on the ADT expression to force the low-expression values to zero. For the normalised

325    expression $z$, with $z \in [0,1]$, the transformed expression is calculated by

326    $z = \begin{cases} 0, & z < t \\ z, & z \geq t \end{cases},$

327    where $t$ is set as 0.5 by default.

328 Lastly, we performed a similar procedure to the method from Vento-Tormo et al. to predict ligand-

329 receptor interactions (Efremova, Vento-Tormo, Teichmann, & Vento-Tormo, 2019). For each ligand-

330 receptor interaction pair originating from a cluster expressing the ligand and another cluster expressing

331 the receptor, we performed a permutation test on the mean of the average RNA expression from the

332 ligand cluster and the mean of the ADT expression from the receptor cluster. Only ligand-receptor pairs

333 with a *P*-value of lower than 0.05 were defined as significant pairs.

334 **Calculation of average and relative ranking of RNA and ADT expression**

335 For the analysis of the ligand-receptor interactions identified through CiteFuse and the conventional

336 approach using only mRNA expression, we calculated the concordance of mRNA and ADT expression

337 of receptors. Because the same gene may be predicted to be involved as a receptor in a ligand-receptor

338 interaction in multiple clusters, we performed a cluster-specific analysis as the expression and

339 correlation of the mRNA and ADT of the receptor is likely to be different between clusters. Therefore,

340 we evaluated concordance between mRNA and ADT in a cluster-specific and relative manner by

341 calculating the ranking of mRNA and ADT expression in the cluster of interest in relation to all other

342 clusters. We then plotted the relative ranking of mRNA and ADT expression against one another. For

343 ligands, we also calculated a cluster-specific ranking based on their mRNA expression.

344 **Data and code availability**

345 All data used in this study are available under accession numbers GSE126310. Sources for code

346 used in this study are available from http://SydneyBioX.github.io/CiteFuse/.

## Author contributions

348 H.J.K. and Y.L. conceived the study with input from J.Y.H.Y. and P.Y.; H.J.K. and Y.L. developed the

349 computational methods, tools, and the R package and led the data analysis and interpretation with input

350 from J.Y.H.Y. and P.Y.; T.A.G. contributed to the development of computational methods; Y.L.

351 implemented the Shiny app with input from H.J.K., J.H.Y.L. and P.Y.; H.J.K., Y.L., and P.Y. wrote the

352 manuscript with input from J.Y.H.Y.; All authors revised, edited, and approved the final version of the

353 manuscript.

## Acknowledgments

## Declaration of interests

364     The authors declare that they have no competing interests.

## References

366     Almende, B. V., & Thieurmel, B. (2016). visNetwork: Network Visualization using "vis.js" Library.

367         *CRAN*.

368     Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful

369         Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B (Methodological)*.

370         https://doi.org/10.1111/j.2517-6161.1995.tb02031.x

371     Blondel, V. D., Guillaume, J. L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities in

372         large networks. *Journal of Statistical Mechanics: Theory and Experiment*.

373         https://doi.org/10.1088/1742-5468/2008/10/P10008

374     Buettner, F., Natarajan, K. N., Casale, F. P., Proserpio, V., Scialdone, A., Theis, F. J., … Stegle, O.

375         (2015). Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data

376         reveals hidden subpopulations of cells. *Nature Biotechnology*, *33*(2), 155–160.

377         https://doi.org/10.1038/nbt.3102

378     Csardi, G., & Nepusz, T. (2006). The igraph software package for complex network research.

379         *InterJournal Complex Systems*, 1695.

380     Datlinger, P., Rendeiro, A. F., Schmidl, C., Krausgruber, T., Traxler, P., Klughammer, J., … Bock, C.

381    (2017). Pooled CRISPR screening with single-cell transcriptome readout. *Nature Methods*,

382    *14*(3), 297–301. https://doi.org/10.1038/nmeth.4177

383  Efremova, M., Vento-Tormo, M., Teichmann, S. A., & Vento-Tormo, R. (2019). CellPhoneDB v2.0:

384    Inferring cell-cell communication from combined expression of multi-subunit receptor-ligand

385    complexes. *BioRxiv*. https://doi.org/10.1101/680926

386  Elyahu, Y., Hekselman, I., Eizenberg-Magar, I., Berner, O., Strominger, I., Schiller, M., … Monsonego,

387    A. (2019). Aging promotes reorganization of the CD4 T cell landscape toward extreme

388    regulatory and effector phenotypes. *Science Advances*. https://doi.org/10.1126/sciadv.aaw8330

389  Fonseka, C. Y., Rao, D. A., Teslovich, N. C., Korsunsky, I., Hannes, S. K., Slowikowski, K., …

390    Raychaudhuri, S. (2018). Mixed-effects association of single cells identifies an expanded

391    effector CD4+ T cell subset in rheumatoid arthritis. *Science Translational Medicine*.

392    https://doi.org/10.1126/scitranslmed.aaq0305

393  Gry, M., Rimini, R., Strömberg, S., Asplund, A., Pontén, F., Uhlén, M., & Nilsson, P. (2009).

394    Correlations between RNA and protein expression profiles in 23 human cell lines. *BMC

395    Genomics*. https://doi.org/10.1186/1471-2164-10-365

396  Haining, W. N., Angelosanto, J., Brosnahan, K., Ross, K., Hahn, C., Russell, K., … Stegmaier, K.

397    (2008). High-throughput gene expression profiling of memory differentiation in primary human T

398    cells. *BMC Immunology*. https://doi.org/10.1186/1471-2172-9-44

399  Lin, Y., Ghazanfar, S., Wang, K. Y. X., Gagnon-Bartsch, J. A., Lo, K. K., Su, X., … Yang, J. Y. H.

400    (2019). ScMerge leverages factor analysis, stable expression, and pseudoreplication to merge

401    multiple single-cell RNA-seq datasets. *Proceedings of the National Academy of Sciences of the

402    United States of America*. https://doi.org/10.1073/pnas.1820006116

403  Liu, Y., Beyer, A., & Aebersold, R. (2016). On the Dependency of Cellular Protein Levels on mRNA

404    Abundance. *Cell*. https://doi.org/10.1016/j.cell.2016.03.014

405  Lun, A. T. L., McCarthy, D. J., & Marioni, J. C. (2016). A step-by-step workflow for low-level analysis

406    of single-cell RNA-seq data with Bioconductor. *F1000Research*.

407    https://doi.org/10.12688/f1000research.9501.2

408 Macaulay, I. C., Haerty, W., Kumar, P., Li, Y. I., Hu, T. X., Teng, M. J., … Voet, T. (2015). G&amp;T-

409  seq: parallel sequencing of single-cell genomes and transcriptomes. *Nature Methods*, *12*, 519.

410 Mimitou, E. P., Cheng, A., Montalbano, A., Hao, S., Stoeckius, M., Legut, M., … Smibert, P. (2019).

411  Multiplexed detection of proteins, transcriptomes, clonotypes and CRISPR perturbations in

412  single cells. *Nature Methods*. https://doi.org/10.1038/s41592-019-0392-0

413 Mohammed, H., Hernando-Herraez, I., Savino, A., Scialdone, A., Macaulay, I., Mulas, C., … Reik, W.

414  (2017). Single-Cell Landscape of Transcriptional Heterogeneity and Cell Fate Decisions during

415  Mouse Early Gastrulation. *Cell Reports*. https://doi.org/10.1016/j.celrep.2017.07.009

416 Ng, A. Y., Jordan, M. I., & Weiss, Y. (2002). On spectral clustering: Analysis and an algorithm. In

417  *Advances in Neural Information Processing Systems*.

418 Pedersen, T. L. (2017). R: Package 'ggraph.' *Cran*.

419 Peterson, V. M., Zhang, K. X., Kumar, N., Wong, J., Li, L., Wilson, D. C., … Klappenbach, J. A.

420  (2017). Multiplexed quantification of proteins and transcripts in single cells. *Nature*

421  *Biotechnology*, *35*, 936.

422 Quinn, T. P., Richardson, M. F., Lovell, D., & Crowley, T. M. (2017). Propr: An R-package for

423  Identifying Proportionally Abundant Features Using Compositional Data Analysis. *Scientific*

424  *Reports*. https://doi.org/10.1038/s41598-017-16520-0

425 Satija, R., Farrell, J. A., Gennert, D., Schier, A. F., & Regev, A. (2015). Spatial reconstruction of

426  single-cell gene expression data. *Nature Biotechnology*. https://doi.org/10.1038/nbt.3192

427 See, P., Lum, J., Chen, J., & Ginhoux, F. (2018). A single-cell sequencing guide for immunologists.

428  *Frontiers in Immunology*. https://doi.org/10.3389/fimmu.2018.02425

429 Stoeckius, M., Hafemeister, C., Stephenson, W., Houck-Loomis, B., Chattopadhyay, P. K., Swerdlow,

430  H., … Smibert, P. (2017). Simultaneous epitope and transcriptome measurement in single cells.

431  *Nature Methods*, *14*(9), 865–868. https://doi.org/10.1038/nmeth.4380

432 Stoeckius, M., Zheng, S., Houck-Loomis, B., Hao, S., Yeung, B. Z., Mauck, W. M., … Satija, R.

433  (2018). Cell Hashing with barcoded antibodies enables multiplexing and doublet detection for

434        single cell genomics. *Genome Biology*, *19*(1), 224. https://doi.org/10.1186/s13059-018-1603-1

435   Stuart, T., Butler, A., Hoffman, P., Hafemeister, C., Papalexi, E., Mauck, W. M., … Satija, R. (2019).

436        Comprehensive Integration of Single-Cell Data. *Cell*, *177*(7), 1888-1902.e21.

437        https://doi.org/10.1016/j.cell.2019.05.031

438   Vento-Tormo, R., Efremova, M., Botting, R. A., Turco, M. Y., Vento-Tormo, M., Meyer, K. B., …

439        Teichmann, S. A. (2018). Single-cell reconstruction of the early maternal–fetal interface in

440        humans. *Nature*. https://doi.org/10.1038/s41586-018-0698-6

441   Wang, B., Mezlini, A. M., Demir, F., Fiume, M., Tu, Z., Brudno, M., … Goldenberg, A. (2014).

442        Similarity network fusion for aggregating data types on a genomic scale. *Nature Methods*, *11*(3),

443        333–337. https://doi.org/10.1038/nmeth.2810

444   Wang, B., Zhu, J., Pierson, E., Ramazzotti, D., & Batzoglou, S. (2017). Visualization and analysis of

445        single-cell rna-seq data by kernel-based similarity learning. *Nature Methods*.

446        https://doi.org/10.1038/nMeth.4207

447   Wolock, S. L., Lopez, R., & Klein, A. M. (2019). Scrublet: Computational Identification of Cell Doublets

448        in Single-Cell Transcriptomic Data. *Cell Systems*, *8*(4), 281-291.e9.

449        https://doi.org/10.1016/j.cels.2018.11.005

450   Zhang, X., Xu, C., & Yosef, N. (2019). Simulating multiple faceted variability in single cell RNA

451        sequencing. *Nature Communications*, *10*(1). https://doi.org/10.1038/s41467-019-10500-w

452

453 **Figures and legends**



454

455 **Figure 1. An overview of CiteFuse and application to clustering of PBMC CITE-seq data.** (A) A
456 summary of the key components and functions implemented in CiteFuse. (B) UMAP visualisation of
457 human PBMC CITE-seq data (Mimitou et al., 2019). (C) Clustering outputs (represented by colours of
458 points) of CD4+ T-cells using multi-modality (CiteFuse), or single-modality (antibody-derived tag [ADT]
459 or RNA alone). (D) Expression of key markers of sub-cell types in CD4+ T-cells.

460

**Figure 2. Doublet detection and downstream analysis using CiteFuse.** (A) A schematic representation of the CITE-seq experiment and the cell hashing data generated by using hashtag oligonucleotide (HTO). (B) The doublet detection approach implemented in CiteFuse. This includes cross-sample doublet identification using Gaussian mixture modelling and a novel within-sample doublet identification method using a combined index of cell hashing and transcriptome data. (C) PCA visualisation of HTO expression before and after filtering of doublets using HTODemux (Stoeckius et al., 2018), Scrublet (Wolock et al., 2019), or CiteFuse. (D-G) Key downstream analytical tools implemented in CiteFuse.

469 **Supplementary figures and legends**



470

471 **Figure S1. Evaluation of CiteFuse and other alternative methods using simulations (related to**
472 **Figure 1)**. (A) A schematic summary of different methods and data modalities used for clustering cells.
473 (B) Ten simulations were conducted for an easy and a hard scenario, respectively. Y-axis shows the
474 adjusted rand index (ARI) calculated for clustering outputs from using various methods and data
475 modalities on each of the two scenarios were presented as boxes.

476

**Figure S2. Clustering of CITE-seq data using single- or multi-modality (related to Figure 1).** (A) UMAP of the fused expression matrix, ADT-alone and RNA-alone expression matrix of the human PBMC CITE-seq data (Mimitou et al., 2019). Clustering outcomes are highlighted by coloured points for both multi-modality (CiteFuse) and single-modality (ADT or RNA) approaches. (B) Centred log-ratio (CLR, y-axis) transformed ADT expression of CD4 and CD27 epitopes in clusters defined from CiteFuse, ADT-alone, and RNA-alone and (C) log RNA expression of S100A4, a marker of CD4+ memory T-cells, and SELL, a marker of naive CD4+ T-cells, in clusters defined from each approach. Clusters correspond to memory CD27+, CD27- DR+, CD27- DR-, and naive cells are highlighted by red arrows. (D) CLR-transformed expression of ADT (CD4 and CD27; first two panels) and log RNA expression of a set of signature genes for memory, naive, or CD27- HLA-DR+ CD4+ memory cells (third, fourth, and fifth panels) highlighted on UMAP of fused similarity matrix. A brighter colour denotes higher expression.
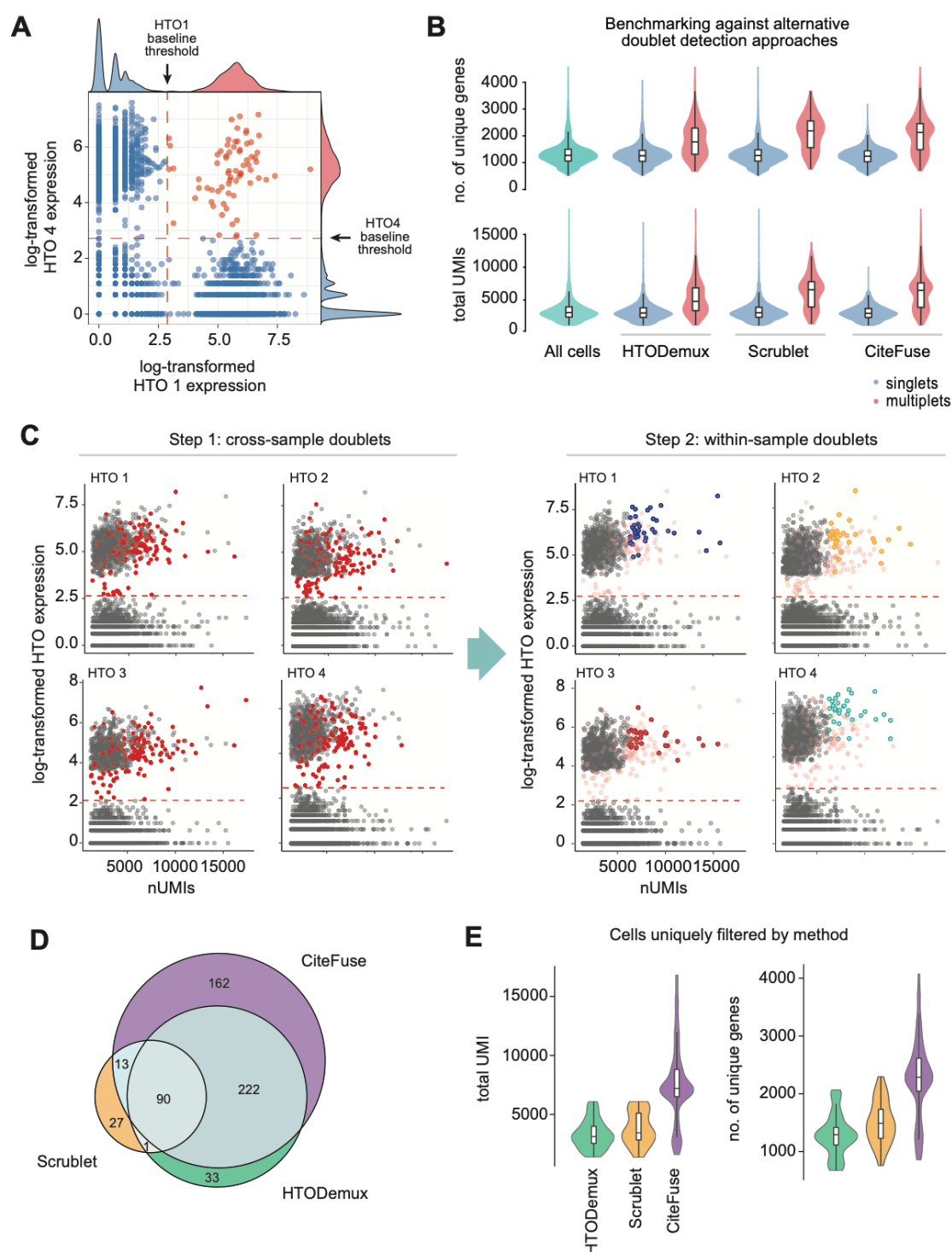
488

**Figure S3. Cross-sample and within-sample doublet detection of CiteFuse (related to Figure 2).**
(A) Gaussian mixture modelling of log-transformed hashtag oligonucleotide (HTO) expression to
identify cross-sample doublets (red points). (B) Total number of unique molecular identifiers (nUMI)
and total number of genes expressed in all cells (both filtered and unfiltered) and HTODemux-,
Scrublet-, and CiteFuse-identified singlets and doublets/multiplets. (C) A scatter plot of nUMI and log-
transformed HTO expression for each HTO (1-4) highlighted by cross-sample doublets (red; left
panel) and within-sample doublets (color-coded by HTO sample; right panel). (D) A Venn diagram of
doublets depicting the overlap in identified doublets between the three filtering methods. (E) nUMI and
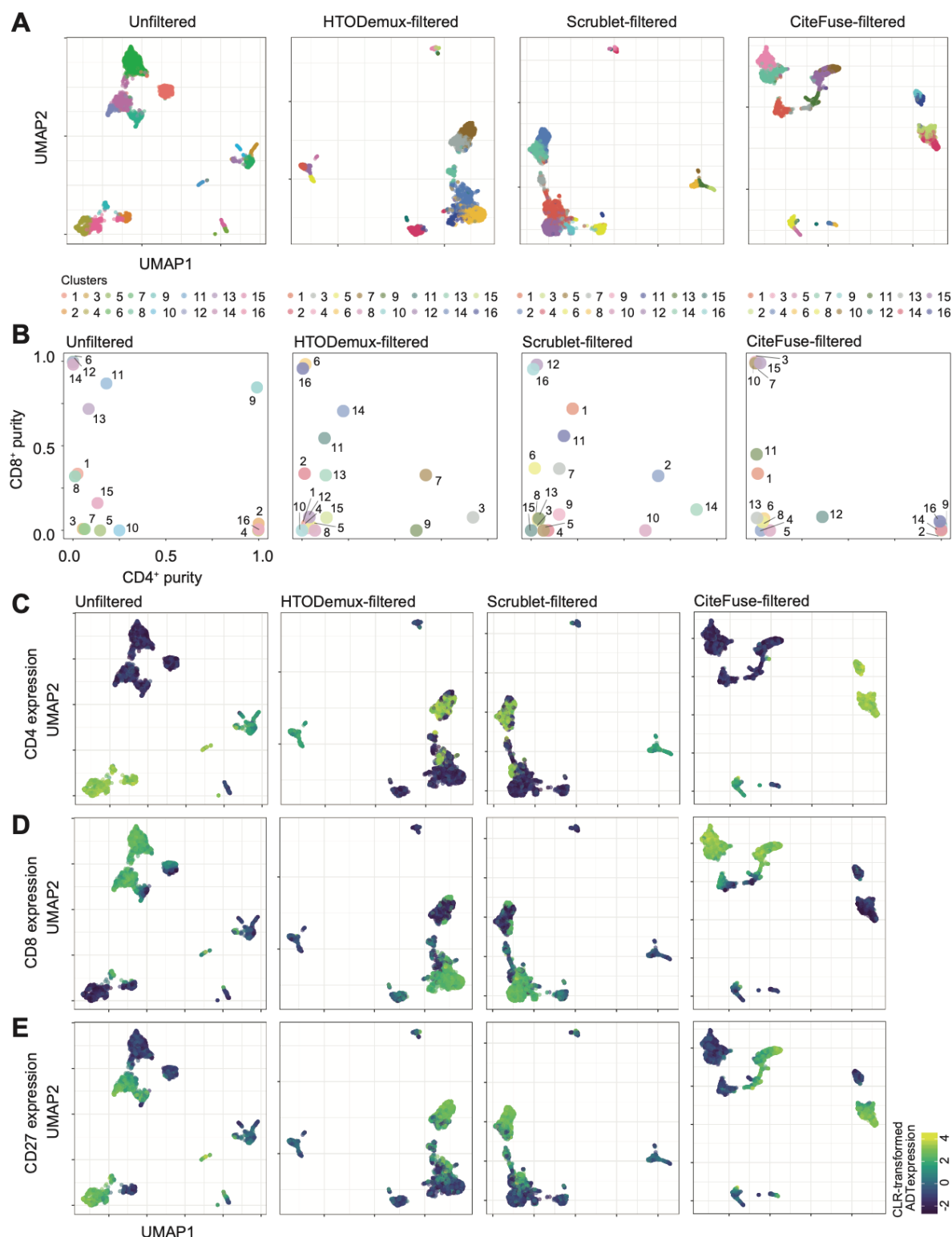total number of genes expressed in doublets uniquely identified by each filtering method.

**Figure S4. Clustering results from unfiltered and doublet filtered data (related to Figure 2).** (A) UMAPs of the unfiltered, HTODemux-filtered, Scrublet-filtered, and CiteFuse-filtered matrix. Clusters generated by fused matrix of both unfiltered and filtered data are highlighted in different colours. (B) Purity scores of CD8+ cells (y-axis) against CD4+ (and CD11c-) (x-axis) cells in individual clusters by unfiltered data or data filtered by each of the three methods. CLR-transformed ADT expression of (C) CD4 (D) CD8 and (E) CD27 highlighted on UMAPs from (A).
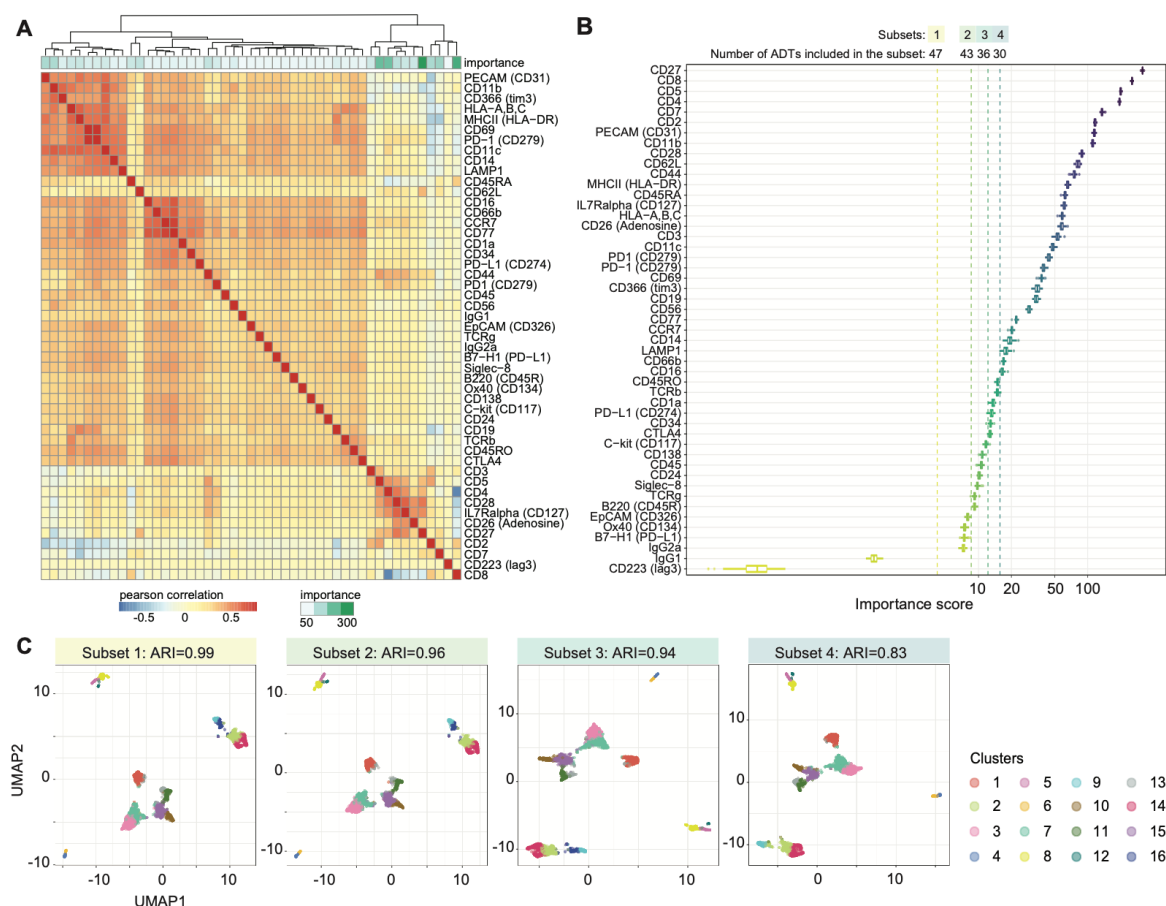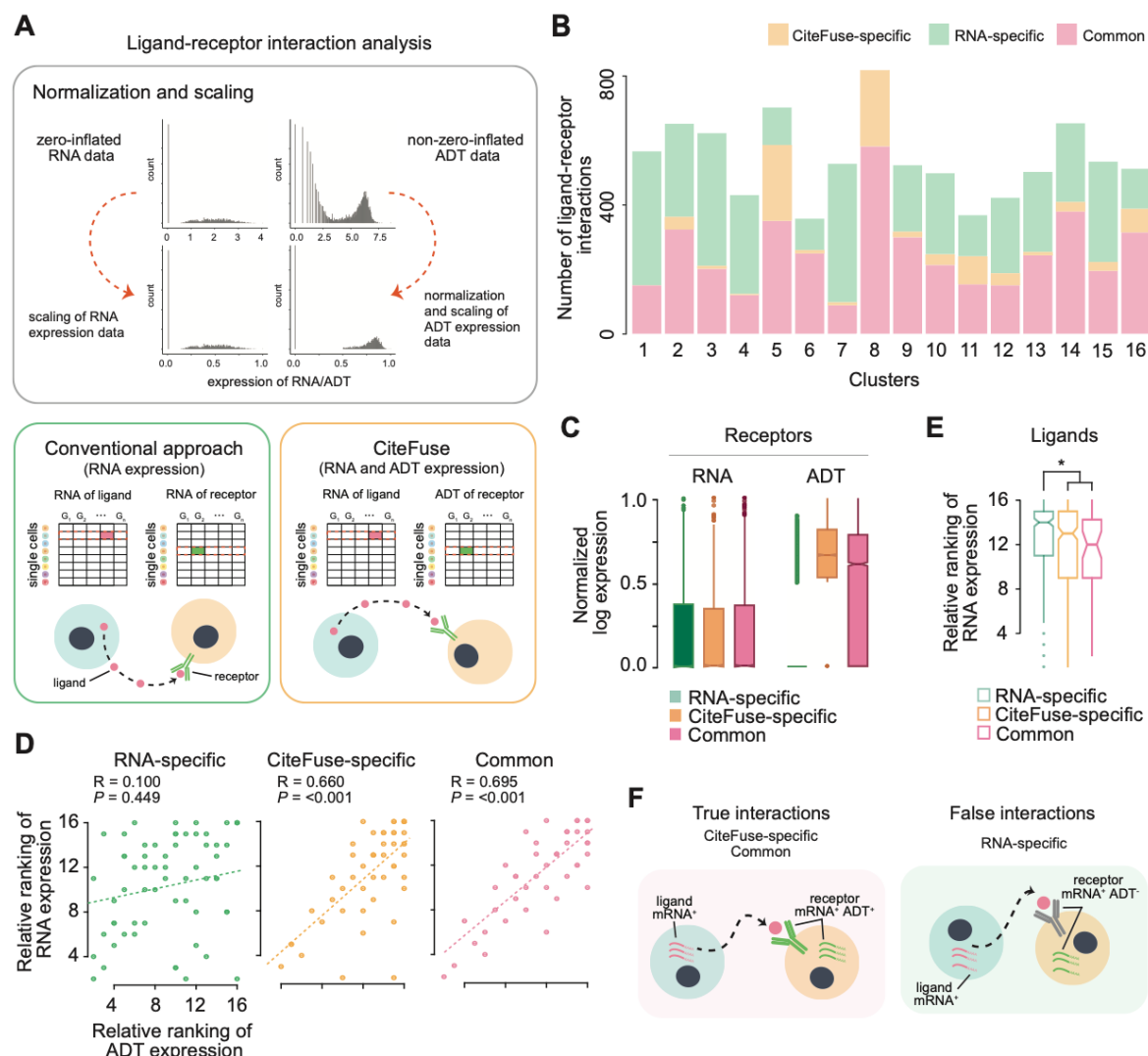
505

**Figure S5. Evaluation of ADTs on CiteFuse clustering outputs (related to Figure 2).** (A) Heatmap of pairwise correlation of ADT expression. Importance score of each ADT was generated by fitting a random forest on CiteFuse clustering outputs of fused matrix (see Methods). (B) Importance scores (x-axis) of ADT towards CiteFuse clustering outputs calculated as the average Gini index after 10 repeated fitting of random forest model. (C) UMAP of CiteFuse with various subsets of ADTs (in decreasing order from left to right panels) and adjusted rand index (ARI) of clustering outcomes against the full ADT set.

512

**Figure S6. Ligand-receptor interaction prediction with CiteFuse (related to Figure 2).** (A) A schematic illustrating the pre-processing step in CiteFuse (scaling of the RNA and ADT expression data and normalization of ADT expression data) and the two types of ligand-receptor interaction prediction methods: 1) conventional approach based on only RNA expression data and 2) CiteFuse approach based on both RNA and ADT expression to predict ligand-receptor interactions. (B) Number of ligand-receptor interactions predicted for each cluster by both conventional approach and CiteFuse (Common), or only by conventional approach (RNA-specific) or CiteFuse (CiteFuse-specific). (D) Scatter plot of relative ranking of RNA and ADT expression across clusters for all receptors identified in a ligand-receptor interaction in each of the three categories (i.e. RNA-specific, CiteFuse-specific, and Common). (E) Relative ranking of RNA expression of ligands in the three categories predicted by conventional approach and/or CiteFuse. (F) Schematic illustration of true and false ligand-receptor interactions and their mRNA and ADT expression where "+" and "-" denote high and low expression, respectively.