

1 **Evaluation of Cell Type Annotation R Packages on Single Cell RNA-seq**
2 **Data**

3 Qianhui Huang¹, Yu Liu², Yuheng Du¹, Lana X. Garmire^{2*}

4 ¹ *Department of Biostatistics, University of Michigan, Ann Arbor, 48109, USA*

5 ² *Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor,*
6 *48105, USA.*

7 * To whom correspondence should be addressed. Email address: lgarmire@med.umich.edu

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24 **Abstract**

25 Annotating cell types is a critical step in single cell RNA-Seq (scRNA-Seq) data analysis. Some
26 supervised/semi-supervised classification methods have recently emerged to enable automated
27 cell type identification. However, comprehensive evaluations of these methods are lacking to
28 provide practical guidelines. Moreover, it is not clear whether some classification methods
29 originally designed for analyzing other bulk omics data are adaptable to scRNA-Seq analysis. In
30 this study, we evaluated ten cell-type annotation methods publicly available as R packages. Eight
31 of them are popular methods developed specifically for single cell research (Seurat, scmap,
32 SingleR, CHETAH, SingleCellNet, scID, Garnett, SCINA). The other two methods are
33 repurposed from deconvoluting DNA methylation data: Linear Constrained Projection (CP) and
34 Robust Partial Correlations (RPC). We conducted systematic comparisons on a wide variety of
35 public scRNA-seq datasets as well as simulation data. We assessed the accuracy through intra-
36 dataset and inter-dataset predictions, the robustness over practical challenges such as gene
37 filtering, high similarity among cell types, and increased classification labels, as well as the
38 capabilities on rare and unknown cell-type detection. Overall, methods such as Seurat, SingleR,
39 CP, RPC and SingleCellNet performed well, with Seurat being the best at annotating major cell
40 types. Also, Seurat, SingleR, CP and RPC are more robust against down-sampling. However,
41 Seurat does have a major drawback at predicting rare cell populations, and it is suboptimal at
42 differentiating cell types that are highly similar to each other, while SingleR and RPC are much
43 better in these aspects. All the codes and data are available at:
44 https://github.com/qianhuiSenn/scRNA_cell_deconv_benchmark.

45

46 **KEYWORDS:** scRNA-seq; cell type; annotation; classification; benchmark;

47

48

49

50

51

52 **Introduction**

53 Single cell RNA sequencing (scRNA-seq) has emerged as a powerful tool to enable the
54 characterization of cell types and states in complex tissues and organisms at the single-cell level
55 [1–5]. Annotating cell types amongst the cell clusters is a critical step before other downstream
56 analyses, such as differential gene expression and pseudo time analysis [6–9]. Conventionally, a
57 set of priorly known cell-type specific markers are used to label the cell types of the clusters
58 manually. This process is laborious and often is a rate-limiting step for scRNA-seq analysis. This
59 approach is also prone to bias and errors. The marker may not be specific enough to differentiate
60 the cell subpopulations in the same dataset, or it may not be generic enough to be applied from one
61 study to another. Automating the cell type labeling is critical to enhance reproducibility and
62 consistency among single cell studies.

63 Recently some annotation methods have emerged to systematically assign cell types in the new
64 scRNA-seq dataset, based on existing annotations from another dataset. Instead of using only top
65 differentiating markers, most methods project or correlate the new cells onto similar cells in the
66 well-annotated reference datasets, by leveraging the whole transcriptome profiles. These
67 annotation methods are developed rapidly, at the same time benchmark datasets that the
68 bioinformatics community agrees upon are lacking. These issues pose the urgent need to
69 comprehensively evaluate these annotation methods using datasets with different biological
70 variabilities, protocols and platforms. It is essential to provide practical guidelines for users.
71 Additionally, identification of limitations of each method through comparisons will also help boost
72 further algorithmic development, which in turn will benefit the scRNA-Seq community.

73 In this study, we evaluated ten cell annotation methods publicly available as R packages (**Table**
74 **1**). Eight of them are popular methods developed specifically for single cell research (Seurat [10],
75 scmap [11], SingleR [12], CHETAH [13], SingleCellNet [14], scID [15], Garnett [16], SCINA
76 [17]). Those methods can be further divided into two categories: Seurat, scmap, CHETAH,
77 SingleCellNet, and scID utilize the gene expression profile as a reference without prior knowledge
78 in signature sets, while Garnett and SCINA require additional pre-defined gene markers as the
79 input. Additionally, to potentially leverage existing deconvolution methods for other bulk omics
80 data, we also included two modified methods: Linear Constrained Projection (CP) and Robust
81 Partial Correlations (RPC) that are popular in DNA methylation analysis [18]. We conducted

82 systematic comparisons on six publicly available scRNA-seq datasets (**Table 2**) varying by
83 species, tissue and sequencing protocol, as well as six sets of simulation data with known truth
84 measure.

85 **Results**

86 **Intra-dataset accuracy evaluation**

87 We first tested the classification accuracy of 10 methods (**Table 1**) on six publicly available
88 scRNA-seq datasets (**Table 2**). These datasets include two human pancreatic islet datasets
89 (GSE85241 and GSE86469), two whole mouse datasets (Tabula Muris, or TM-Full), and two
90 peripheral blood mononuclear cells datasets (PBMC). Since Tabula Muris datasets are
91 heterogeneous in terms of tissue contents, to evaluate the tools' performance on homogeneous
92 data, we down sampled them separately into two mouse lung datasets (Tabula Muris-Lung, or TM-
93 Lung) by taking cells from lung tissue only. This results in eight real scRNA-seq datasets (**Table**
94 **1**). To avoid potential bias, we used the 5-fold cross validation scheme to measure the averaged
95 accuracy in the 1-fold holdout subset. We used three different performance measurement metrics:
96 overall accuracy, adjusted rand index (ARI), and V-measure [19–20] (see **Materials and**
97 **methods**). The evaluation workflow is depicted in Figure S1.

98 **Figure 1A-C** show the classification accuracy metrics on eight datasets. The top-five
99 performing annotation methods are Seurat, SingleR, CP, singleCellNet and RPC. Seurat has the
100 best overall classification performance in the 5-fold cross validation evaluation. On average, the
101 three accuracy metrics from Seurat are significantly higher (Wilcoxon paired rank test, $P < 0.05$)
102 than 9 other methods. SingleR has the second-best performance after Seurat, with all three metrics
103 higher than 8 other methods, among which 6 pair-wise method comparisons achieved statistical
104 significance (Wilcoxon paired rank test $P < 0.05$). Though slightly lower in average metric scores,
105 the classification performance of both singleCellNet and CP are comparable to SingleR.

106 In order to test the influence of cell-type number on tool's performance, we next evaluated the
107 TM-Full and TM-Lung results. As shown in **Table 2**, for two TM-Full datasets from 10X and
108 Smart-Seq2 platforms, which contain a large number of cell types (32 and 37 cell types,
109 respectively), we took a subset of cells from the lung tissue and created two TM-Lung datasets
110 that are relatively small in cell-type numbers, with 8 and 10 cell types for 10X and Smart-Seq2

111 platforms, respectively. Most methods perform well for both TM-Lung datasets with ARI > 0.9.
112 However, some of the methods had a drop of performance on the two TM-Full datasets. The
113 increased classification labels imposed a challenge. Garnett failed to predict on such large TM-
114 Full datasets. Additionally, SCINA, CHETAH and scmap have significantly lower classification
115 metrics on TM-Full datasets, compared to those on TM-Lung datasets. On the contrary, the
116 previously mentioned top-five methods are more robust despite the increase of complexity in TM-
117 Full datasets. Again, Seurat yields the best metric scores in both TM-Full datasets, demonstrating
118 its capability at analyzing complex datasets.

119 **Cross-dataset accuracy evaluation**

120 To evaluate the annotation tools in a more realistic setting, we conducted cross-dataset
121 performance evaluation on 10 datasets (5 pairs), where the referencing labels were obtained from
122 one dataset and the classification was done on another dataset of the same tissue type (**Table 2**).
123 Within a pair, we used FACS-sorted, purified dataset as the reference data, and the remaining one
124 as the query data (see **Materials and methods**). Among the 5 pairs of datasets, 4 are real
125 experimental data: PBMC cell pair with PBMC-sorted-ref and PBMC-3K-query; pancreas cell pair
126 with pancreas-celseq2-ref and pancreas-fluidigm-query; TM-Full pair with TM-Full-smartseq2-
127 ref and TM-Full-10X-query; TM-Lung pair with TM-Lung-smartseq2-ref and TM-Lung-10x-
128 query. The last pair is simulation datasets with the pre-defined truth, where the true assay without
129 dropouts (simulation_true-ref) was used as the reference and the raw assay with dropout mask
130 (simulation_raw-query) was used as the query.

131 **Figure 1D-F** shows the classification accuracy metrics on the above mentioned 5 pairs of query
132 and reference datasets. The top 3 performing annotation methods in the descending rank order are
133 Seurat, SingleR, and CP, the same as those in the same-dataset cross validation results (**Figure**
134 **1A-C**). In particular, they all perform very well on the simulation data with known truth measures,
135 as all three accuracy metrics are above 0.96. RPC is ranked 4th, slightly better than SCINA. Similar
136 to the 5-fold cross validation evaluation, methods such as scID, CHETAH, scmap and Garnett are
137 persistently ranked among the lowest-performing methods for accuracy. Interestingly,
138 singleCellNet, the method that performs relatively well (ranked 4th) in the same-dataset 5-fold
139 cross validation, is now consistently ranked on the 6th, behind RPC and SCINA, due to drop of
140 performance in TM-Full datasets. Besides annotation methods, the accuracy scores are also much

141 dependent on the datasets. For example, on complex PBMC datasets, even Seurat only reaches
142 0.76 for ARI. A further examination of the confusion matrices (Figure S2) for Seurat, SingleR,
143 CP, singleCellNet and RPC reveals that the challenge comes from distinguishing highly similar
144 cell types such as CD4+ T cells vs. CD8+ T cells or Dendritic cells vs. CD14+ Monocytes in the
145 PBMC datasets.

146 We also performed batch corrected cross-dataset accuracy evaluations on 4 pairs of
147 experimental datasets. For each pair of data, both reference and query datasets were aligned using
148 CCA [10,21]. The result is illustrated in Figure S3A-B. Most methods do not benefit from aligning
149 and integrating the datasets (Figure S3B). None of the other methods exceeds the performance of
150 Seurat in all three metrics after the batch correction (Figure S3A). The drop of performance in
151 those methods may be attributed to the fact that aligned datasets contain negative values after the
152 matrix correction and subtraction from the integration algorithm used in Seurat. In addition, some
153 algorithms require non-normalized data matrix as the input, while batch-corrected matrix from
154 Seurat is normalized, which may violate some models' assumptions.

155 Altogether, these results from both experimental and simulation data indicate that Seurat has
156 the best overall accuracy among the annotation methods in comparison, based on intra-dataset
157 prediction and cross-dataset prediction [10].

158 **The effect of cell type similarity**

159 Since it is challenging to distinguish highly similar cell populations using cross-data evaluation,
160 we next conducted simulations. We designed 20 simulation data sets composed of five cell groups
161 with varying levels of differential expression. Similar to others [22], we used Splatter [23] to pre-
162 define the same set of differential expression (DE) genes in simulation datasets, and only differed
163 the magnitudes of DE, from low, low-moderate, moderate to high (**Figure 2A**). When cell
164 populations are more separable, the classification task is easy for the majority of methods. As the
165 cell populations become less separable, all methods show a decrease in their performance (**Figure**
166 **2B-D**). The degrees of such decrease vary among the methods though. SingleR, RPC, Seurat,
167 singleCellNet and CP are in the first class that are relatively more robust than the other five
168 methods. SingleR and RPC are ranked the 1st and 2nd for their robustness against cell type
169 similarity, with all three metric scores above 0.9. Seurat is ranked the 4th after singleCellNet (the
170 3rd) when the samples are least separable (low DE), exposing its slight disadvantage. Garnett

171 failed to predict when cell-cell similarity is high (low DE). In this context, the pre-defined marker
172 genes may be ‘ambiguous’ to discriminate in multiple cell types, which may cause problems for
173 Garnett to train the classifier.

174 **The effect of increased classification labels on annotation performance**

175 The increased cell type classification labels imposed a challenge for some methods in inter and
176 intra-data predictions. We designed five simulation datasets each composed of an increased
177 number (N) of cell groups (N = 10, 20, 30, 40, 50) with a constant total cell numbers, gene
178 numbers, and level of differential expression among cell groups. Similar to the performance that
179 we observed on intra-data and inter-data classification experiments, the increased classification
180 grouping labels lead to dropping accuracy for most methods, except SingleR, which is extremely
181 robust without drop of performance (**Figure 2E-G**). RPC is consistently ranked 2nd regardless of
182 the cell group numbers. Seurat and CP are ranked the 3rd and 4rd for their robustness before N=30,
183 with small differences in accuracy metrics. However, after N=30, the accuracy of Seurat
184 deteriorates faster and is ranked 4rd instead. The performance issue in Seurat may be due to its
185 susceptibility towards cell-cell similarity. Since we keep a constant differential expression level
186 despite the increased cell grouping labels, more cell types have similar expression profiles and
187 they are more likely to be misclassified. On the other hand, Garnett failed to predict when
188 simulation data set has cell types $N > 20$. Therefore, the simulation study confirms the practical
189 challenge of increased cell labels in multi-label classification for most methods evaluated. SingleR
190 is the most robust method against increased complexity in both real dataset and simulation data
191 evaluations.

192 **The effect of gene filtering**

193 We also evaluated the stability of annotation methods in inter-dataset classification, by varying the
194 number of query input features. For this purpose, we used the human pancreas data pair (**Table 2**).
195 We randomly down sampled the features from Fluidigm data into 15,000, 10,000 and 5000 input
196 genes, based on the original log count distribution (**Figure 3A**). When the number of features
197 decreases, most methods show decreased metric scores as expected (**Figure 3B**). Seurat and
198 SingleR are the top 2 most robust methods over the decrease of feature numbers, and their ARI
199 scores remain high across all sampling sizes (ARI > 0.9). Again, methods such as Garnett, scID
200 and scmap are more susceptible to low feature numbers, since their performances decrease as the

201 feature number decreases. Therefore, using query data with fewer features than the reference data
202 may affect the prediction performance of those methods. Alternatively, we also downsized the
203 samples by reducing the number of raw reads before alignment and tag counting steps (**Figure**
204 **3C**). While most methods show fairly consistent accuracy scores with reduced raw reads as
205 expected, a couple of methods, such as singleCellNet and scID, are perturbed by this procedure
206 (**Figure 3D**).

207 **Rare population detection**

208 Identifying rare populations in single cells is a much biologically interesting aspect. We evaluated
209 the inter-dataset classification accuracy per cell population for the top 5 methods based on overall
210 accuracy and adjusted rand index (ARI) (Figure S4): Seurat, SingleR, CP, singleCellNet and RPC
211 (**Figure 1A-B**). We used a mixture of 9 cell populations with a wide variety of percentages (50%,
212 25%, 12.5%, 6.25%, 3.125%, 1.56%, 0.97%, 0.39%, 0.195%) in ten repeated simulation datasets
213 with different seeds (**Figure 4A**). When the size of the cell population is larger than 50 cells out
214 of 2000 cells, all five methods achieve high cell-type specific accuracy of over 0.8 (**Figure 4B**).
215 However, the classification performances drop drastically for Seurat and singleCellNet when the
216 cell population is 50 or less. On the other hand, most low-performing methods have fluctuated
217 performance and do not perform well in classifying the major cell populations (Figure S4B).
218 Interestingly, bulk-reference based methods such as SingleR, CP and RPC are extremely robust
219 against the size changes of a cell population. They employ averaged profiles as the references and
220 are not susceptible to low cell counts. One challenge for some other single cell methods is that
221 there are not enough cell counts from a low-proportion cell type. Some methods just remove or
222 ignore those cell types in the training phase (such as Garnett), or during alignment (such as Seurat)
223 by their threshold parameters of the algorithms.

224 **Unknown population(s) detection**

225 Among the scRNA-seq specific annotation tools, five methods (Garnett, SCINA, scmap,
226 CHETAH, scID) contain the rejection option that allows ‘unassigned’ labels. This is a rather
227 practical option, as the reference data may not contain all cell labels present in the query data. In
228 order to assess how accurate these methods are at labeling ‘unassigned’ cells, we used the scheme
229 of “hold-out one cell type evaluation” on the same simulation dataset pair used in cross-dataset
230 prediction. That is, we remove the signature of one cell type in the reference matrix while keeping

231 the query intact. The evaluation repeated five times for all five cell types. For each method, we
232 measured the average classification accuracies excluding the hold-out group (**Figure 4C**), and the
233 accuracy of assigning unlabeled class to the leave-out group in the query (**Figure 4D**). Among the
234 five methods compared, SCINA, scmap and scID all have metrics scores above the average level
235 of all tools tested for accuracy excluding the hold-out group (**Figure 4C**). However, SCINA has
236 better accuracy in rejecting cell groups existing in the query dataset but not in the reference (**Figure**
237 **4D**). Similar results were observed from “hold-out two cell type evaluation” (Figure S5). SCINA
238 has a relatively better balance between overall accuracy in existing cell types and precise rejection
239 in non-existing cell types.

240 The caveat here, however, is that none of the rejection-enabled methods are among the best
241 performing methods in terms of overall accuracy, stability and robustness to cell type similarities.
242 Since accuracy, stability and robustness are probably more important attributes to assess these
243 methods, the practical guide value based on the results of unknown population detection is limited.

244 **Time and memory comparison**

245 In order to compare the runtime and memory utilization of the annotation methods, we simulated
246 six data sets each composed of 20,000 genes, with 5 cell types of equal proportion (20%), in total
247 cell numbers of 5000, 10,000, 15,000, 20,000, 25,000, 50,000, respectively (see **Materials and**
248 **methods**). All methods show increases in computation time and memory usage when the number
249 of cells increases (**Figure 5**). Of the five top-performing methods in the intra-data and inter-data
250 annotation evaluations (**Figure 1**), singleCellNet and CP outperform others on speed (**Figure 5A**).
251 As the dataset size increases beyond 50,000 cells, methods such RPC require a runtime as large as
252 6 hours. For memory utilization, singleCellNet and CP consistently require less memory than other
253 top performing methods (**Figure 5B**). Notably, the best performing method Seurat (by accuracy)
254 requires memories as large as 100GB, when dataset size increases beyond 50,000 cells, which is
255 significantly larger than most other methods. In all, based on computation speed and memory
256 efficiency, singleCellNet and CP outperform others among the top-class accurate annotation
257 methods.

258 **Discussion**

259 In this study, we presented comprehensive evaluations of 10 computational annotation methods in
260 R packages, on single cell RNA-Seq data. Of the 10 methods, 8 of them are designed for single-
261 cell RNA-seq data, and 2 of them are our unique adaptation from methylation-based analysis. We
262 evaluated these methods on 6 publicly available scRNA-seq datasets as well as many additional
263 simulation datasets. We systematically assessed accuracy (through intra-dataset and inter-dataset
264 predictions), the robustness of each method with challenges from gene filtering, cell-types with
265 high similarity, increased cell type classification labels, and the capabilities on rare population
266 detection and unknown population detection, as well as time and memory utilization (**Figure 6**).
267 In summary, we found that methods such as Seurat, SingleR, CP, RPC and SingleCellNet
268 performed relatively well overall, with Seurat being the best-performing methods in annotating
269 major cell types. Methods such as Seurat, SingleR, RPC and CP are more robust against down-
270 sampling. However, Seurat does have a major drawback at predicting rare cell populations, as well
271 as minor issues at differentiating highly similar cell types and coping with the increased
272 classification labels, while SingleR and RPC are much better in these aspects.

273 During the preparation of the manuscript, another evaluation paper was published in a special
274 edition of Genome Biology [24]. We, therefore, address the differences between these two studies'
275 methodologies, before discussing our own findings in detail. First, rather than simply comparing
276 the methods claimed to be “single cell specific”, we uniquely repurpose two methods: Linear
277 Constrained Projection (CP) and Robust Partial Correlations (RPC). Although they were originally
278 developed for DNA methylation data deconvolution, their regression-based principle could be
279 adapted to scRNA-seq supervised/semi-supervised classification. We modified the final regression
280 coefficients as the probability of one specific cell type label, rather than the cell content as in DNA
281 methylation-based deconvolution. As the results indicated, CP and RPC has comparable prediction
282 with SingleR, the overall second best method. This shows the potential of repurposing existing
283 deconvolution methods from another bulk omics analysis. Secondly, for benchmark datasets, we
284 used fewer real experimental datasets. However, we uniquely included many simulated datasets
285 while the other study did not use any. We argue that it is important to have additional simulation
286 datasets, because evaluation based on manually annotated cell-type-specific markers in the
287 experimental data is prone to bias. On the contrary, one can introduce simulation datasets with
288 ‘ground truth’ and unbiasedly assess the tricky issues, such as identifying highly similar cell
289 populations or very rare cell populations. Thirdly, Seurat, the method with the best overall

290 accuracy in our study, is not included in the other study. The high annotation performance of Seurat
291 on intra-data and inter-data predictions in our study, is mostly due to the fact that it's a
292 classification method using an integrated reference. Its data transfer feature shares the same
293 anchors identification step as the data integration feature. However, unlike data integration, the
294 cell type classification method in Seurat does not correct the query expression data. On top of that,
295 its default setting projects the PCA structure of a reference onto the query, instead of learning a
296 joint structure with CCA [10,21]. This type of methods represents a new trend in single cell
297 supervised classification, evident by a series of scRNA-seq data integration methods (LIGER,
298 Harmony, scAlign *etc* [25–27]). Lastly, we only selected the packages in R with good
299 documentations, as R is still the most popular bioinformatics platform for open-source scRNA-
300 Seq analysis packages (e.g. the arguably most popular method Seurat, which the other study
301 omitted).

302 Although having slightly lower accuracy metrics scores than Seurat, SingleR and CP still have
303 very excellent performance in intra-data and inter-data prediction, with resilience towards gene
304 filtering and increased complexity in datasets. In addition, SingleR has better performance than
305 Seurat in predicting rare cell populations, dealing with increased cell type classification labels, as
306 well as differentiating highly similar cell types. This advantage of SingleR may benefit from its
307 method and the pseudo-bulk reference matrix. The averaged pseudo-bulk reference profile may
308 potentially remove the variation and noise from the original single cell reference profile, and it can
309 retain the expression profiles of all cell types and is not affected by the low count. SingleR uses
310 pseudo-bulk RNA-seq reference to correlate the expression profiles to each of the single cells in
311 the query data, and uses highly variable genes to find the best fit iteratively. For Seurat, the
312 annotation of the cell labels on query data is informed by the nearest anchor pairs. If two or more
313 cell types have similar profiles, their alignments may overlap which may cause misclassification.
314 Seurat also has some requirements on the minimum number of defined anchor pairs. In the case of
315 rare cell populations, the lack of the neighborhood information makes the prediction difficult.
316 Similar to other study [24], we also found that method that incorporates the prior-knowledge (e.g.
317 Garnett and SCINA) did not improve the classification performance over other methods that do
318 not have such requirements. This prior-knowledge is limited when cell-cell similarity is large. In
319 addition, as the number of cell types increases, the search for the marker genes will become
320 challenging, making these methods even less desirable.

321 Compared with intra-data prediction, inter-data prediction is more realistic but also more
322 challenging. Technical/platform and batch differences in inter-data prediction may impose major
323 challenges to the classification process, although the tissue and cell type contents are the same. In
324 our study, the CCA batch-correction preprocessing step did not improve the classification accuracy
325 for most methods. Among all experimental data used as the benchmark in this study, PBMC
326 datasets had the worst accuracy results (ARI=0.76 for the best method Seurat). Further inspection
327 of the confusion matrices revealed that the challenges come from distinguishing highly similar cell
328 types, which themselves may have some level of inaccuracy from the original experiments. If the
329 upstream unsupervised clustering methods are not sensitive enough to categorize similar cell
330 populations, this uncertainty may be carried through to the downstream cell annotation steps. This
331 again highlights the potential issue of evaluating the supervised/semi-supervised methods in single
332 cell data, where we are not certain about the ‘ground truth’ of the cell labels to begin with.
333 Recently, some studies used unsupervised classification methods through multi-omics integration,
334 and/or reconstruction of gene regulatory network [28,29], representing a new trend in this area.
335 As the multi-omics technology continues to advance [30], it will be of interest to evaluate these
336 methods, where both multi-omics and pre-defined marker information are available for the same
337 samples.

338 Overall, we recommend using Seurat for general annotation tasks for cell types that are
339 relatively separable and without rare population identification as the objective. However, for
340 datasets contain cell types with high similarities or rare cell populations, if a reference dataset with
341 clean annotations is available, SingleR, RPC and CP are preferable.

342

343 **Materials and methods**

344 **Real data sets**

345 Six real scRNA-seq data sets were downloaded and used for evaluations and validations (**Table**
346 **2**). The human pancreatic islet datasets were obtained from the following accession numbers:
347 GEO: GSE85241 (Celseq2) [10,31], GEO: GSE86469 (Fluidigm C1) [10,32]. The *Tabula Muris*
348 datasets Version 2 (10X Genomics and Smart-Seq2) were downloaded from FigShare:
349 <https://tabula-muris.ds.czbiohub.org/> [3]. The bead-purified PBMC dataset (10X Genomics) was

350 obtained from the Zheng dataset: <https://github.com/10XGenomics/single-cell-3prime-paper>, and
351 the PBMC-3K dataset (10X Genomics) was downloaded from
352 <https://support.10xgenomics.com/single-cell-gene-expression/datasets> [33]. These datasets differ
353 by species, tissue and sequencing protocol. For each of the datasets, we collected both raw counts
354 and cell-type annotations from the corresponding publications, except PBMC-3k, for which the
355 cell-type annotations were obtained through the standard single cell RNA-seq analysis and
356 classified using cell-type-specific marker genes. The extracted cell-type annotations for each
357 dataset were used as the ground truth for evaluations (Table S1).

358 *Data cleaning*

359 Datasets were paired in groups by tissue type (**Table 2**). Within a pair, we used the data generated
360 by Fluorescence-activated cell sorting (FACS) sorted method as reference data. Both reference
361 data and query data were further processed to make sure the cell types in reference data are larger
362 or equal to the cell types in the query data. When necessary, the query data were down sampled
363 following the original cell type count distribution. For the two Tabula Muris (TM-Full) datasets
364 from 10X and Smart-Seq2 platforms, which contain a large number of cell types (32 and 37 cell
365 types, respectively), we took a subset of cells from lung tissue and created two TM-Lung datasets
366 that have fewer cell types, 8 for 10X and 10 for Smart-Seq2 platform, respectively. As a result, we
367 have four pairs of experimental datasets: PBMC cell pair with PBMC-sorted-ref and PBMC-3K-
368 query; pancreas cell pair with pancreas-celseq2-ref and pancreas-fluidigm-query; TM-Full pair
369 with TM-Full-smartseq2-ref and TM-Full-10X-query; TM-Lung pair with TM-Lung-smartseq2-
370 ref and TM-Lung-10x-query.

371 *Data downsampling*

372 To explore the effects of different feature numbers and read depths on the performance of tools,
373 we randomly down sampled features (genes) from human pancreas-Fluidigm dataset into 5000,
374 10,000 and 15,000 input genes, following the original log count distribution. We repeated five
375 times for each downsampling scheme. Alternatively, we also down sampled the reads into 25%,
376 50%, 75% of the original read depths (with 2 repetitions) using *samtools* on BAM files, and then
377 realigned following the method provided by the original manuscript [32].

378 **Simulated Data Sets**

379 We simulated a dataset using Splatter, with 4000 genes and 2000 cells (Splatter parameters,
380 dropout.shape=-0.5, dropout.mid=1), and then split each dataset into 5 cell groups with proportions
381 10%, 30%, 30%, 10% and 20%. In addition, we also generated three additional simulation sets to
382 evaluate the robustness of tools. In the first set, we generated 10 simulation datasets each has
383 10,000 genes and 2,000 cells (use Splatter parameters dropout.shape=-0.5, dropout.mid=1, 10
384 different seeds), and then split each into 9 cell groups with proportions 50%, 25%, 12.5%, 6.25%,
385 3.125%, 1.56%, 0.97%, 0.39%, 0.195%, respectively. The second set contains 20 simulation
386 datasets, each composed of 10,000 genes and 2,000 cells splitting into 5 cell types with equal
387 proportions. These datasets have the same set of differentially expressed (DE) genes, but differ by
388 the magnitude of DE factors (*de.facScale* parameter in Splatter). We simulated each DE scale five
389 times with five different seeds. The DE scales and the parameterizations are: low: *de.facScale* =
390 *c*(0.1, 0.3, 0.1, 0.3, 0.2); low-moderate: *de.facScale* = *c*(0.3, 0.5, 0.3, 0.5, 0.4); moderate:
391 *de.facScale* = *c*(0.5, 0.7, 0.5, 0.7, 0.6); high: *de.facScale* = *c*(0.7, 0.9, 0.7, 0.9, 0.8). The third set
392 contains five simulation datasets each composed of an increased number (N) of cell groups (N =
393 10, 20, 30, 49, 50) with a constant total cell numbers (10,000), gene numbers (20,000), and level
394 of differential expression among cell groups. Each simulation dataset contains two paired assays.
395 The true assay without dropouts was used as the reference and the raw assay with dropout mask
396 was used as the query.

397 **Data Preprocessing**

398 *Cell and gene filtering*

399 We filtered out cells for which fewer than 200 genes were detected and any genes that were
400 expressed in fewer than 3 cells.

401 *Normalization*

402 For the annotation tools that require a normalized count matrix as input, we performed log-
403 normalization using a size factor of 10,000.

404 *Pseudo-bulk reference matrix*

405 For the annotation tools that use bulk rather than single-cell expression profiles as reference, we
406 took the average of the normalized count of each cell type group and made a pseudo-bulk RNA-
407 seq reference.

408 *Marker genes selection*

409 Some classification tools (SCINA and Garnett) require cell-type specific marker as the input.
410 When such marker information is neither provided by the corresponding tools nor retrievable by
411 public research, we extract them from the reference data by performing differential expression
412 analysis using Wilcoxon rank sum test (*FindAllMarkers* function from Seurat with parameters
413 only.pos = TRUE, min.pct = 0.25 and logfc.threshold = 0.25). Wilcoxon rank sum test is the most
414 common nonparametric test for a difference in mean expression between cell groups. The top 10
415 ranked marker genes for each cell type were used as the input for the corresponding tools.

416 **Supervised/Semi-supervised Annotation Methods**

417 We only considered pre-printed or published methods with detailed documentation on installation
418 and execution. We excluded any methods that required extensive running time, and where we were
419 unable to customize the reference dataset, or random and inconsistent predictions were produced.
420 In the end, ten cell annotation methods, publicly available as R packages, were evaluated in this
421 study. This includes eight methods (Seurat, scmap, SingleR, CHETAH, SingleCellNet, scID,
422 Garnett, and SCINA) commonly used to annotate scRNA-seq data. In addition, to investigate the
423 potential to repurpose deconvolution methods for other bulk omics analysis, we also included and
424 modified two methods originally designed for bulk DNA methylation that use a different type of
425 algorithms not yet reported in scRNA-seq specific tools: Linear Constrained Projection (CP) and
426 Robust Partial Correlations (RPC).

427 All parameters were set to default values following the author's recommendations or the
428 respective manuals (**Table 1**). For methods that allow "unknown" assignments (scmap, CHETAH,
429 scID, Garnett, and SCINA), we modified the parameter to force assignments where possible
430 (except for the evaluations where unknown assignments were allowed).

431 *Adaptation of CP and RPC methods for scRNA-Seq analysis*

432 In order to accommodate the methylation-based methods for scRNA-seq data, we made some
433 modifications. In original papers, both RPC and CP model the methylation profile of any given
434 sample as a linear combination of a given set of reference profiles representing underlying cell-
435 types present in the sample. Assume the number of underlying cell-types to be C , and each cell
436 type has a profile \mathbf{b}_c that constitutes the signature matrix \mathbf{H} [34–36]. Let \mathbf{y} be the profile of a given
437 sample and w_c be the weight estimation of cellular proportion of each cell type, and the underlying
438 model becomes:

$$439 \quad \mathbf{y} = \sum_{c=1}^C w_c \mathbf{b}_c + \epsilon$$

440 Both methods assume that reference profiles contain the major cell-types present in the sample
441 \mathbf{y} and sum of weights equal to 1. RPC estimates the weight coefficients using robust multivariate
442 linear regression or robust partial correlation, while CP uses a quadratic programming technique
443 known as linear constrained protection to estimate the weights [37].

444 In the modified version, we first converted the single cell RNA-seq reference data into pseudo-
445 bulk RNA-seq data matrix by taking the average of the normalized count of each cell type group.
446 Then we took the subset of pseudo-bulk RNA-seq data by keeping n features that exhibited high
447 cell-to-cell variations across C distinct cell types in the reference dataset, and had a small condition
448 number below 3 as the signature matrix \mathbf{H} [34]. We set the highly variable genes to 2000, using
449 *FindVariableFeatures* function from Seurat (Figure S6). We let \mathbf{y} be the profile of a given single
450 cell from the query data with the same 2000 genes from the signature matrix \mathbf{H} . While applying
451 both algorithms, we treated the estimated weight for each cell type as the probability and the cell
452 type with the highest weight was the identity of the corresponding single cell sample in the query
453 data. This conversion is based on the fact that \mathbf{y} no longer represents averages over many different
454 cell types, but only expression profile from only one cell type (since we have single cell data).

455 **Benchmarking**

456 *Five-fold cross validation and cross-dataset prediction*

457 For each dataset in four pairs of the real experimental datasets mentioned above, we used a 5-fold
458 cross validation where the four-fold data were used as the reference and the remaining one-fold as
459 the query. For the cross-dataset prediction, in addition to the four pairs of real datasets, we used

460 simulation datasets containing true assay (without dropouts) as the reference and raw assay (with
461 dropout mask) as the query.

462 In order to evaluate whether batch correction and data integration benefit the classification
463 performance, for each pair of real dataset, we aligned both reference and query dataset using CCA
464 [10,21] from the Seurat data integration function. Then we separated the aligned datasets and
465 performed the cross-dataset evaluation again.

466 *Performance evaluation on the effect of feature numbers and read depths*

467 To investigate the robustness of different methods with regards to feature numbers and read depths,
468 we used the down-sampled human pancreas Fluidigm data set as described in the data
469 downsampling section. In such evaluation, the human pancreas Celseq2 dataset was used as the
470 reference and the down-sampled human pancreas Fluidigm dataset was used as a query.

471 *Performance evaluation with effect of differential expression (DE) scale among cell groups*

472 In this assessment, we used 20 simulation data sets containing the same DE gene set but differing
473 only by DE factors as described earlier in the Simulated Data Sets section. Each simulation data
474 set contains two paired assays. The true assay (without dropouts) was used as the reference and
475 the raw assay (with dropout mask) was used as the query.

476 *Performance evaluation on the effect of increased classification labels*

477 In this evaluation, we designed five simulation data sets, each composed of an increased number
478 (N) of cell groups (N=10, 20, 30, 40, 50) with a constant total cell numbers, gene numbers, and
479 level of differential expression among cell groups. Each simulation data set contains two paired
480 assays. The true assay (without dropouts) was used as the reference and the raw assay (with
481 dropout mask) was used as the query.

482 *Rare and unknown population detection*

483 Each of the 10 simulation data sets in the rare population detection evaluation was composed of
484 10,000 genes and 2000 cells splitting into 9 cell types with proportions 50%, 25%, 12.5%, 6.25%,
485 3.125%, 1.56%, 0.97%, 0.39%, 0.195%. The simulation dataset in the unknown population
486 detection evaluation was composed of 4000 genes and 2000 cells splitting into 5 cell types. We
487 used the scheme of “hold-out one cell type evaluation” to evaluate prediction on the unknown

488 population, that is, removing the signature of one cell type in the reference matrix while predicting
489 the query. During each prediction, one cell group was removed from the reference matrix and the
490 query remained intact. We repeated the evaluation five times for all five cell types. We additionally
491 employed a “hold-out two cell type” experiment, in which we removed signatures of any
492 combination of two cell types in the reference matrix while keeping the query intact. The
493 evaluation was repeated ten times for all ten different combinations. Similarly, for each simulation
494 data set, the true assay (without dropouts) was used as the reference and the raw assay (with
495 dropout mask) was used as the query.

496 *Runtime and Memory Assessment*

497 In order to compare the computational runtime and memory utilization of annotation methods, we
498 simulate six datasets, with total cell numbers of 5000, 10,000, 15,000, 20,000, 25,000, and 50,000,
499 respectively, each composed of 20,000 genes, splitting into 5 cell types with the equal proportion.
500 The true assay (without dropouts) was used as the reference and the raw assay (with dropout mask)
501 was used as the query. Each execution was performed in a separate R session in our lab server (4
502 nodes (Dell PowerEdge C6420) of 2 X Intel(R) Xeon(R) Gold 6154 CPU @ 3.00GHz, 192GB
503 RAM, one node (Dell Poweredge R740) with 2 X Xeon(R) Gold 6148 CPU @ 2.40GHz, 192 GB
504 RAM, and two 16GB Nvidia V100 GPUs) with Slurm job scheduler. One processor and 100GB
505 memory were reserved for each job. From the job summary, we collected ‘Job Wall-clock time’
506 and ‘Memory Utilized’ for evaluation. We ran each method on each dataset five times to estimate
507 the average computation time.

508 **Evaluation Criteria**

509 The prediction results of the methods are evaluated using three different metrics: overall accuracy,
510 adjusted rand index, and V-measure. We used three different metrics to avoid possible bias in
511 evaluating the performance. The detailed explanations on these metrics were described earlier
512 [22,38,39]. Briefly, *Overall accuracy* is the percent agreement between the predicted label and the
513 true label. *Adjusted rand index* (ARI) is the ratio of all cell pairs that are either correctly classified
514 together or correctly not classified together, among all possible pairs, with adjustment for chance.
515 *V-measure* is computed as the harmonic mean of distinct homogeneity and completeness score. In
516 specific, homogeneity is used to assess whether each predicted cell type groups contains only

517 members of a single class, while completeness is used to assess whether all members of a given
518 class are assigned to the same predicted cell label.

519 **Code Availability**

520 All the codes and data are available at:
521 https://github.com/qianhuiSenn/scRNA_cell_deconv_benchmark.

522 **Authors' Contributions**

523 LG, and QH envisioned this project. QH implemented the project and conducted the analysis with
524 help from YL and YD. QH and LG wrote the manuscript. All authors have read and agreed on the
525 manuscript.

526 **Competing interests**

527 The authors declare no competing financial interests.

528 **Acknowledgements**

529 This research was supported by grants K01ES025434 awarded by NIEHS through funds provided
530 by the trans-NIH Big Data to Knowledge (BD2K) initiative (www.bd2k.nih.gov), R01 LM012373
531 and R01 LM012907 awarded by NLM, R01 HD084633 awarded by NICHD to L.X. Garmire.

532 **References**

- 533 [1] Plass M, Solana J, Wolf FA, Ayoub S, Misios A, Glažar P, et al. Cell type atlas and lineage
534 tree of a whole complex animal by single-cell transcriptomics. *Science* 2018;360.
535 <https://doi.org/10.1126/science.aag1723>.
- 536 [2] Cao J, Spielmann M, Qiu X, Huang X, Ibrahim DM, Hill AJ, et al. The single-cell
537 transcriptional landscape of mammalian organogenesis. *Nature* 2019;566:496–502.
- 538 [3] Tabula Muris Consortium, Overall coordination, Logistical coordination, Organ collection
539 and processing, Library preparation and sequencing, Computational data analysis, et al.
540 Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. *Nature*
541 2018;562:367–72.
- 542 [4] Yu P, Lin W. Single-cell Transcriptome Study as Big Data. *Genomics Proteomics
543 Bioinformatics* 2016;14:21–30.
- 544 [5] Mu Q, Chen Y, Wang J. Deciphering Brain Complexity Using Single-cell Sequencing.
545 *Genomics Proteomics Bioinformatics* 2019;17:344–66.

- 546 [6] Zappia L, Phipson B, Oshlack A. Exploring the single-cell RNA-seq analysis landscape
547 with the scRNA-tools database. *PLoS Comput Biol* 2018;14:e1006245.
- 548 [7] Zhu X, Yunits B, Wolfgruber T, Poirion O, Arisdakessian C, Garmire L. GranatumX: A
549 community engaging and flexible software environment for single-cell analysis. *bioRxiv*
550 2018:385591. <https://doi.org/10.1101/385591>.
- 551 [8] Bacher R, Kendzioriski C. Design and computational analysis of single-cell RNA-
552 sequencing experiments. *Genome Biol* 2016;17:63.
- 553 [9] Rostom R, Svensson V, Teichmann SA, Kar G. Computational approaches for interpreting
554 scRNA-seq data. *FEBS Lett* 2017;591:2213–25.
- 555 [10] Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM 3rd, et al.
556 Comprehensive Integration of Single-Cell Data. *Cell* 2019;177:1888–902.e21.
- 557 [11] Kiselev VY, Yiu A, Hemberg M. scmap: projection of single-cell RNA-seq data across data
558 sets. *Nat Methods* 2018;15:359–62.
- 559 [12] Aran D, Looney AP, Liu L, Wu E, Fong V, Hsu A, et al. Reference-based analysis of lung
560 single-cell sequencing reveals a transitional profibrotic macrophage. *Nat Immunol*
561 2019;20:163–72.
- 562 [13] de Kanter JK, Lijnzaad P, Candelli T, Margaritis T, Holstege FCP. CHETAH: a selective,
563 hierarchical cell type identification method for single-cell RNA sequencing. *Nucleic Acids*
564 *Res* 2019. <https://doi.org/10.1093/nar/gkz543>.
- 565 [14] Tan Y, Cahan P. SingleCellNet: A Computational Tool to Classify Single Cell RNA-Seq
566 Data Across Platforms and Across Species. *Cell Syst* 2019;9:207–13.e2.
- 567 [15] Boufeua K, Seth S, Batada NN. scID uses discriminant analysis to identify transcriptionally
568 equivalent cell types across single cell RNA-seq data with batch effect. *iScience*
569 2020:100914.
- 570 [16] Pliner HA, Shendure J, Trapnell C. Supervised classification enables rapid annotation of
571 cell atlases. *Nat Methods* 2019;16:983–6.
- 572 [17] Zhang Z, Luo D, Zhong X, Choi JH, Ma Y, Wang S, et al. SCINA: Semi-Supervised
573 Analysis of Single Cells in silico. *Genes* 2019;10:531.
- 574 [18] Teschendorff AE, Breeze CE, Zheng SC, Beck S. A comparison of reference-based
575 algorithms for correcting cell-type heterogeneity in Epigenome-Wide Association Studies.
576 *BMC Bioinformatics* 2017;18:105.
- 577 [19] Rand WM. Objective Criteria for the Evaluation of Clustering Methods. *J Am Stat Assoc*
578 1971;66:846–50.
- 579 [20] Rosenberg A, Hirschberg J. V-measure: A conditional entropy-based external cluster
580 evaluation measure. *Proceedings of the 2007 joint conference on empirical methods in*
581 *natural language processing and computational natural language learning (EMNLP-*
582 *CoNLL)*, 2007, p. 410–20.
- 583 [21] Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic
584 data across different conditions, technologies, and species. *Nat Biotechnol* 2018;36:411–20.
- 585 [22] Arisdakessian C, Poirion O, Yunits B, Zhu X, Garmire LX. DeepImpute: an accurate, fast,
586 and scalable deep neural network method to impute single-cell RNA-seq data. *Genome Biol*
587 2019;20:211.
- 588 [23] Zappia L, Phipson B, Oshlack A. Splatter: simulation of single-cell RNA sequencing data.
589 *Genome Biol* 2017;18:174.
- 590 [24] Abdelaal T, Michielsen L, Cats D, Hoogduin D, Mei H, Reinders MJT, et al. A comparison
591 of automatic cell identification methods for single-cell RNA sequencing data. *Genome*

- 592 Biology 2019;20. <https://doi.org/10.1186/s13059-019-1795-z>.
- 593 [25] Welch JD, Kozareva V, Ferreira A, Vanderburg C, Martin C, Macosko EZ. Single-Cell
594 Multi-omic Integration Compares and Contrasts Features of Brain Cell Identity. *Cell*
595 2019;177:1873–87.e17.
- 596 [26] Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, et al. Fast, sensitive and
597 accurate integration of single-cell data with Harmony. *Nat Methods* 2019;16:1289–96.
- 598 [27] Johansen N, Quon G. scAlign: a tool for alignment, integration, and rare cell identification
599 from scRNA-seq data. *Genome Biol* 2019;20:166.
- 600 [28] Duren Z, Chen X, Zamanighomi M, Zeng W, Satpathy AT, Chang HY, et al. Integrative
601 analysis of single-cell genomics data by coupled nonnegative matrix factorizations. *Proc*
602 *Natl Acad Sci U S A* 2018;115:7723–8.
- 603 [29] Zeng W, Chen X, Duren Z, Wang Y, Jiang R, Wong WH. DC3 is a method for
604 deconvolution and coupled clustering from bulk and single-cell genomics data. *Nat*
605 *Commun* 2019;10:4613.
- 606 [30] Ortega MA, Poirion O, Zhu X, Huang S, Wolfgruber TK, Sebra R, et al. Using single-cell
607 multiple omics approaches to resolve tumor heterogeneity. *Clin Transl Med* 2017;6:46.
- 608 [31] Muraro MJ, Dharmadhikari G, Grün D, Groen N, Dielen T, Jansen E, et al. A Single-Cell
609 Transcriptome Atlas of the Human Pancreas. *Cell Syst* 2016;3:385–94.e3.
- 610 [32] Lawlor N, George J, Bolisetty M, Kursawe R. Single-cell transcriptomes identify human
611 islet cell signatures and reveal cell-type-specific expression changes in type 2 diabetes.
612 *Genome* 2017.
- 613 [33] Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively
614 parallel digital transcriptional profiling of single cells. *Nat Commun* 2017;8:14049.
- 615 [34] Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of
616 cell subsets from tissue expression profiles. *Nat Methods* 2015;12:453–7.
- 617 [35] Venet D, Pecasse F, Maenhaut C, Bersini H. Separation of samples into their constituents
618 using gene expression data. *Bioinformatics* 2001;17 Suppl 1:S279–87.
- 619 [36] Abbas AR, Wolslegel K, Seshasayee D, Modrusan Z, Clark HF. Deconvolution of blood
620 microarray data identifies cellular activation patterns in systemic lupus erythematosus.
621 *PLoS One* 2009;4:e6098.
- 622 [37] Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, et
623 al. DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC*
624 *Bioinformatics* 2012;13:86.
- 625 [38] Zhu X, Ching T, Pan X, Weissman SM, Garmire L. Detecting heterogeneity in single-cell
626 RNA-Seq data by non-negative matrix factorization. *PeerJ* 2017;5:e2888.
- 627 [39] Poirion O, Zhu X, Ching T, Garmire LX. Using single nucleotide variations in single-cell
628 RNA-seq to identify subpopulations and genotype-phenotype linkage. *Nat Commun*
629 2018;9:4892.

630 **Figure Legends**

631 **Figure 1 Inter-data and cross-date accuracy comparison**

632 (A-C) Within data accuracy comparison, shown as heatmaps of three classification metrics, (A)
633 overall accuracy, (B) adjusted rand index (ARI), and (C) v-measure across eight real datasets. For

634 each dataset, a 5-fold cross validation is performed: using four folds as the reference and one-fold
635 as the query. **(D-F)** Between-data accuracy comparison, shown as heatmaps of three classification
636 metrics, **(D)** overall accuracy, **(E)** adjusted rand index (ARI), and **(F)** v-measure across four pairs
637 of experimental datasets and one pair of simulation datasets. PBMC cell pair: PBMC-sorted-ref
638 and PBMC-3K-query; pancreas pair: pancreas-celseq2-ref and pancreas-fluidigm-query; TM-Full
639 pair: TM-Full-smartseq2-ref and TM-Full-10X-query; TM-Lung pair: TM-Lung-smartseq2-ref
640 and TM-Lung-10x-query; simulation: true-ref and dropout-masked raw. TM-Lung datasets pair
641 was down sampled from TM-Full datasets pair by taking cells from lung tissue only. Within the
642 simulation datasets pair, the true assay without dropouts (true-ref) was used as the reference and
643 the raw assay with dropout mask (raw-query) was used as the query. The columns are datasets,
644 and the rows are annotation methods. The heatmap scale is shown on the figure, where the brighter
645 yellow color indicates a better classification accuracy score. On the right of each heatmap is a
646 boxplot to summarize the classification metrics among methods. Box colors represent different
647 methods as shown in the figure. The methods in the heatmap and the boxplot are arranged in
648 descending order by their average metrics score across all datasets. Some methods failed to
649 produce a prediction for certain data sets (indicated by grey squares).

650 ****: significantly higher ($P < 0.05$) than 9 other methods using pairwise Wilcoxon test.

651 ***: significantly higher ($P < 0.05$) than 8 other methods using pairwise Wilcoxon test.

652 **: significantly higher ($P < 0.05$) than 7 other methods using pairwise Wilcoxon test.

653 *: significantly higher ($P < 0.05$) than 6 other methods using pairwise Wilcoxon test.

654

655 **Figure 2 Effect of cell-cell similarity and increased classification labels on annotation tool**
656 **performance**

657 **(A)** PCA plots of simulation datasets generated by Splatter, each of which is composed of 10,000
658 genes and 2000 cells, splitting into 5 cell types with equal proportion, and contains the same
659 proportion of differentially expressed genes in each cell type. The datasets differ by changing the
660 magnitude of DE factors for those DE genes to simulate more or less differences between groups.

661 Based on the magnitude of DE factors in five cell groups, we generated 20 datasets with cell groups
662 similarity ranging from low, low-moderate, moderate to high DE (see **Materials and methods**).
663 Colors represent different cell types. True assay (without dropouts) is used as the query and raw
664 assay (with dropout) is used as the reference. **(B-D)** Plots showing three classification metrics to
665 evaluate each annotation method applied to the datasets in **(A)**. The x-axis is the DE scale for
666 differential expressed genes in each group, and the y-axis is the metric score. Results are shown as
667 mean+std over 5 repetitions. Line colors and point shapes correspond to different methods. The
668 metrics are: **(B)** overall accuracy, **(C)** adjusted rand index (ARI) and **(D)** v-measure. **(E-G)** Plots
669 illustrating three classification metrics to evaluate each annotation methods applied to five
670 simulation datasets, each of which is composed of an increased number (N) of cell groups (N =
671 10, 20, 30, 49, 50) with a constant total cell numbers (10,000), gene numbers (20,000), and level
672 of differential expression among cell groups. The x-axis is the number of cell types in each data
673 set, and the y-axis is the metric score. The metrics are: **(E)** overall accuracy, **(F)** adjusted rand
674 index (ARI) and **(G)** v-measure.

675 **Figure 3 Effects of feature (gene) numbers and read depths on annotation tool performance**

676 **(A)** The features (genes) in the human pancreas-fluidigm dataset are filtered by removing genes
677 that present in less than 3 cells, resulting in 19211 genes. The filtered features (genes) are randomly
678 down sampled into 5000, 10,000 and 15,000 input genes, following the original log count
679 distribution. Such down-sampling was repeated 5 times. **(C)** The BAM file reads in the human
680 pancreas-fluidigm dataset are randomly down sampled into 25%, 50%, 75% of the original read
681 depths. **(B)(D)** Plots depicting the three classification metrics (overall accuracy, adjusted rand
682 index and v-measure) of each method applied to the down sampling approaches in **(A)** and **(C)**
683 respectively. The x-axis is the down sampling size for feature numbers or reads, and the y-axis is
684 the metrics score. Results are shown as mean+std over 5 repetitions. Line colors and point shapes
685 correspond to different methods. SCINA failed when the number of input features reached 5000,
686 thus no point is shown.

687

688 **Figure 4 Performance comparison on rare cell type and unknown cell types detection**

689 All datasets are generated by Splatter. **(A)** Cell population distribution of simulation data (10
690 repeats), composed of 10,000 genes and 2000 cells, split into 9 cell types with proportions of 50%,
691 25%, 12.5%, 6.25%, 3.125%, 1.56%, 0.97%, 0.39%, 0.195%, respectively. **(B)** Plot illustrating
692 cell-type specific accuracy across 9 cell groups in **(A)**, for the five annotation methods that exceed
693 0.8 in overall accuracy and adjusted rand index (ARI). The x-axis is the cell groups in the
694 descending order for their cell proportions, and the y-axis is the cell-type specific classification
695 score. Results are shown as mean+std over 5 repetitions. **(C)** Performance metrics (overall
696 accuracy, adjusted rand index and v-measure) of another simulation data set, composed of 4000
697 genes and 2000 cells splitting into 5 cell types. True assay (without dropouts) is used as the
698 reference and the raw assay (with dropout) is used as the query. During each prediction, one cell
699 group is removed from the reference matrix and the query remains intact. The x-axis lists methods
700 with rejection options (e.g. allowing ‘unlabeled’ samples), and the y-axis is the classification
701 metrics score excluding the hold-out group. **(D)** Boxplots showing the accuracy of methods in **(C)**,
702 when assigning ‘unlabeled’ class to the leave-out group in the query.

703 **Figure 5 Speed and memory usage comparison**

704 Speed and memory comparison on six pairs of simulation data with increasing numbers of cells
705 (5000, 10,000, 15,000, 20,000, 25,000, 50,000). True assay (without dropouts) is used as the
706 reference and the raw assay (with dropout) is used as the query. Both reference and query contain
707 the same number of cells. Color depicts different annotation methods. **(A)** Natural log of running
708 time (y-axis) vs. cell size (x-axis) over five repetitions in each data point. **(B)** Natural log of peak
709 memory usage (y-axis) vs. cell size (x-axis) over five repetitions in each data point.

710 **Figure 6 Benchmark Summary**

711 Summary of the classification performance in each evaluation. Each row is a method and each
712 column is evaluations from intra-data and inter-data prediction (Intra-Inter), cell-cell similarity
713 (DE-Scale), increased classification labels, downsampling of genes, downsampling of reads, rare
714 group detection, unknown population detection (rejection), time and memory utilization. The
715 heatmap shows individual method’s rank based on averaged metric scores over overall accuracy,
716 adjusted rand index, and v-measure for each evaluation indicated in the bottom column. Time and
717 memory are ranked by utilization. Grey box indicated that the method does not participate in the

718 evaluation. The methods in the heatmap are arranged in ascending order by their average rank over
719 inter-data and intra-data prediction.

720 *Note:* pop_overall: averaged metric scores for all simulations in rare population detection.
721 Low_count: averaged metrics scores for classifying cell types < 1.56% in population.
722 rej_exe_overall: averaged classification metrics score excluding the hold-out group. rej_overall:
723 accuracy of assigning ‘unlabeled’ class to the leave-out group in the query.

724

725 **Tables**

726 **Table 1 List of Single-Cell RNA-sequencing/Methylation Cell Annotation tools**
727 **benchmarked in this study**

728 **Table 2 Datasets used in this study**

729 **Supplementary material**

730 **Supplementary Figures**

731 **Supplementary Figure 1 Benchmark Workflow**

732 Illustration of the workflow for this study consists of 1) preprocessing 2) prediction 3) evaluations.

733 **Supplementary Figure 2 Cell-type specific accuracy for top 5 performing methods on PBMC**
734 **cross-dataset prediction**

735 Confusion matrix of cell-type specific accuracy for PBMC inter-dataset predictions among top
736 performing annotation methods (Seurat, SingleR, CP, singleCellNet, RPC). The x-axis is the
737 predicted label from each algorithm, and the y-axis is the true label in the query data.

738 **Supplementary Figure 3 Inter-data prediction using aligned reference and query matrix**

739 For each of the four pairs of experimental data used in cross-data evaluation, we aligned both
740 reference and query dataset using CCA from the Seurat data integration function. Then we
741 separated the aligned datasets and performed the cross-dataset evaluation again. (A) Inter-data
742 accuracy comparison, shown as heatmaps of three classification metrics (overall accuracy,

743 adjusted rand index (ARI), and v-measure). **(B)** Boxplots illustrating the averaged metrics scores
744 before and after alignment for each method. The x-axis is the methods, and the y-axis is the
745 classification metrics score.

746 **Supplementary Figure 4 Rare population detection evaluation for remaining 5 methods**

747 **(A)** Boxplots illustrating the averaged overall accuracy and adjusted rand index over all the rare
748 population detection simulation data. The x-axis is the methods evaluated, and the y-axis is the
749 metric score. **(B)** Rare population detection results for the five methods with lower overall
750 accuracy and ARI. The x-axis is the cell groups in the descending order for their cell proportions,
751 and the y-axis is the cell-type specific classification score.

752 **Supplementary Figure 5 Hold-out two cell type rejection evaluation**

753 “Hold-out two cell type” experiment was performed on the same simulation dataset pair used in
754 cross-dataset prediction. In this experiment, signatures of any combination of two cell types were
755 removed in the reference matrix while keeping the query intact. The evaluation was repeated ten
756 times for all ten different combinations. **(A)** The x-axis lists methods with rejection options (e.g.
757 allowing ‘unlabeled’ samples), and the y-axis is the classification metrics score excluding the hold-
758 out groups. **(B)** Boxplots showing the accuracy of methods in **(A)**, when assigning ‘unlabeled’
759 class to leave-out groups in the query.

760 **Supplementary Figure 6 The optimal number of highly variable genes (HVG) to be used in** 761 **CP and RPC algorithms**

762 The highly variable genes are identified from reference dataset and ranked by standardized
763 variance from mean-variance feature selection methods with variance-stabilizing transformation.
764 **(A)** The boxplot depicts the overall accuracy averaged over five pairs of inter-dataset predictions
765 (pbmc, pancreas, tabula-Full, tabula-Lung, and simulation) with the top 100, 200, 500, 1000, 2000,
766 and 5000 highly variable genes as input features for CP and RPC methods. The x-axis is the
767 number of highly variable features, and the y-axis is the overall accuracy. Methods are reflected
768 by different box colors. **(B)** The boxplot represents the condition number of the pseudo-bulk
769 reference matrix averaged over four combinations of cross-dataset predictions with the top 100,

770 200, 500, 1000, 2000, and 5000 highly variable genes as input features. The x-axis is the number
771 of highly variable features, and the y-axis is the condition number.

772

773 **Supplementary Tables**

774 **Supplementary Table 1 Composition of cell-types in each real dataset**

Table 1 List of Single-Cell RNA-sequencing/Methylation Cell Deconvolution tools benchmarked in this study.

| Software | Method/Algorithm | Bulk/Single Reference | Require Pre-defined Marker Genes | Allow Unknown | Version Under R 3.6.0 | Reference |
|-----------------|---|------------------------------|---|----------------------|------------------------------|------------------|
| SingleR | Correlation-based with Iterative Tuning | Bulk | No | No | SingleR_1.0.0 | [12] |
| CP | Reference-based method using Constrained Projection | Bulk | No | No | EpiDISH_2.0.2 | [18] |
| RPC | Reference-based Robust Partial Correlations | Bulk | No | No | EpiDISH_2.0.2 | [18] |
| Garnett | Elastic net Multinomial Regression | Single | Yes | Yes | garnett_0.1.4 | [16] |
| SCINA | Bimodal Distribution assumption for marker genes | Single | Yes | Yes | SCINA_1.1.0 | [17] |
| Seurat | Define anchor with CCA, L2-norm and MNN | Single | No | No | Seurat_3.0.1 | [10] |
| singleCellNet | Multi-Class Random Forest | Single | No | No | singleCellNet_0.1.0 | [14] |
| CHETAH | Correlation-based with Hierarchical Classification | Single | No | Yes | CHETAH_1.1.2 | [13] |

| | | | | | | |
|-------|---|--------|----|-----|-----------------|------|
| scmap | K-nearest-neighbor classification with cosine similarity | Single | No | Yes | scmap_1.6.0 | [11] |
| scID | Fisher's Linear Discriminant Analysis-like Framework | Single | No | Yes | scID_0.0.0.9000 | [15] |

Table 2 Datasets used in this study

| Dataset Name | Protocol | No. of Cells | No. of Genes | No. of Cell Types | Species/Tissue/Description | Reference |
|--------------------------|-----------------|---------------------|---------------------|--------------------------|---|------------------|
| PBMC-Sorted | 10X | 91,649 | 18,986 | 7 | Human Peripheral Blood Mononuclear Cells | [33] |
| PBMC-3K | 10X | 2467 | 13,714 | 6 | Human Peripheral Blood Mononuclear Cells | 10X Genomics |
| Pancreas-Sorted | CEL-Seq2 | 2285 | 34,363 | 13 | Human Pancreas | [10, 31] |
| Pancreas | Fluidigm C1 | 638 | 34,363 | 13 | Human Pancreas | [10, 32] |
| Tabula Muris-Sorted | Smart-Seq2 | 24,622 | 22,252 | 37 | Mouse | [3] |
| Tabula Muris | 10X | 20,000 | 17,866 | 32 | Mouse | [3] |
| Tabula Muris Lung-Sorted | Smart-Seq2 | 1563 | 22,253 | 10 | Mouse Lung | [3] |
| Tabula Muris Lung | 10X | 1303 | 17,866 | 8 | Mouse Lung | [3] |
| Simulation1_true | Splatter | 2000 | 4000 | 5 | Simulation data for inter-data prediction | |
| Simulation1_raw | Splatter | 2000 | 4000 | 5 | - | |

| | | | | | |
|------------------------------|----------|------|--------|---|--|
| Simulation2_true | Splatter | 2000 | 10,000 | 9 | Simulation data with descending cell proportion for each cell group, repeat with 10 random seeds. |
| Simulation2_raw | Splatter | 2000 | 10,000 | 9 | - |
| Simulation_Low_true | Splatter | 2000 | 10,000 | 5 | Simulation data with low differential expression scale for each cell group, repeat with 5 random seeds. |
| Simulation_Low_raw | Splatter | 2000 | 10,000 | 5 | - |
| Simulation_Low_Moderate_true | Splatter | 2000 | 10,000 | 5 | Simulation data with low-moderate differential expression scale for each cell group, repeat with 5 random seeds. |
| Simulation_Low_Moderate_raw | Splatter | 2000 | 10,000 | 5 | - |
| Simulation_Moderate_true | Splatter | 2000 | 10,000 | 5 | Simulation data with moderate differential expression scale for each cell group, repeat with 5 random seeds. |
| Simulation_Moderate_raw | Splatter | 2000 | 10,000 | 5 | - |

| | | | | | |
|----------------------|----------|---|--------|--------------------|--|
| Simulation_High_true | Splatter | 2000 | 10,000 | 5 | Simulation data with high differential expression scale for each cell group, repeat with 5 random seeds. |
| Simulation_High_raw | Splatter | 2000 | 10,000 | 5 | - |
| Simulation3_true | Splatter | 10,000 | 20,000 | 10;20;30; 40;50 | Simulation data with increased cell type labels from 10 to 40 cell types. |
| Simulation3_raw | Splatter | 10,000 | 20,000 | 10;20;30; 40;50 | - |
| Simulation4_true | Splatter | 5000/10,000/ 15,000/20,000/ 25,000/50,000 | 20,000 | 5 | Simulation data with increased cell number from 5000 to 50,000. |
| Simulation4_raw | Splatter | 5000/10,000/ 15,000/20,000/ 25,000/50,000 | 20,000 | 5 | - |

Note: Raw data is true simulation data with the addition of dropouts. Sorted data were generated from Fluorescence-activated cell sorting (FACS) sorted method

Supplementary Table 1 Composition of cell-types in each real dataset

| PBMC-Sorted | PBMC-3K | Pancreas-celseq2 | Pancreas-Fluidigm | Tabula Muris-Lung-smartseq2 | Tabula Muris-Lung-10X | Tabula Muris-Full-smartseq2 | Tabula Muris-Full-10X |
|-------------------------------|-------------------------------|-------------------------------|-------------------------------|------------------------------------|---------------------------------|------------------------------------|--------------------------------------|
| B cells (10,084) | B cells (342) | Acinar (274) | Acinar (21) | B cells (57) | B cells (140) | B cells (2029) | B cells (5615) |
| CD14+ Monocytes (2,465) | CD14+ Monocytes (2,465) | Activated stellata (90) | Activated stellata (16) | Ciliated columnar cell (25) | - | Basal cell (1340) | Basal cell (27) |
| CD34+ (6,312) | - | Alpha (843) | Alpha (239) | Classical monocyte (90) | Classical monocyte (4) | Basal cell of epidermis (1648) | Basal cell of epidermis (2020) |
| CD4+ T cells (42,166) | CD4+ T cells (479) | Beta (445) | Beta (258) | Leukocyte (35) | Leukocyte (9) | Basophil (25) | - |
| CD8+ T cells (22,138) | CD8+ T cells (308) | Delta (203) | Delta (25) | Lung endothelia cell (693) | Lung endothelia cell (24) | Bladder cell (695) | Bladder cell (192) |

| | | | | | | | |
|-------------------------|----------------------------|---------------------|-----------------------|--------------------------------|---------------------------------|----------------------------------|-------------------------------------|
| Dendritic cells (99) | Dendritic cells (33) | Ductal (258) | Ductal (36) | Monocyte (65) | - | Bladder urothelial cell (683) | Bladder urothelial cell (141) |
| NK cells (8,385) | NK cells (155) | Endothelial (21) | Endothelial (14) | Myeloid cell (85) | Myeloid cell (2) | Blood cell (206) | Blood cell (153) |
| - | - | Epsilon (4) | Epsilon (1) | Natural killer cell (37) | Natural killer cell (113) | Cardiac muscle cell (133) | - |
| - | - | Gamma (110) | Gamma (18) | Stromal cell (423) | Stromal cell (888) | Ciliated columnar cell (25) | - |
| - | - | Macrophage (15) | Macrophag e (1) | T cell (53) | T cell (123) | Classical monocyte (90) | Classical monocyte (4) |
| - | - | Mast (6) | Mast (3) | - | - | DN1 thymic pro-T cell (32) | - |
| - | - | Quiescent stella | Quiescent stella | - | - | Endocardial cell (165) | Endocardial cell (1) |

| | | | | | | | |
|---|---|---------|---------|---|---|---|---------------------------------|
| - | - | (12) | (1) | - | - | Endothelial cell | Endothelial cell |
| | | Schwann | Schwann | | | (3319) | (971) |
| - | - | (4) | (5) | - | - | Endothelial cell of hepatic sinusoid | - |
| | | | | | | (182) | |
| - | - | - | - | - | - | Epithelial cell | Epithelial cell |
| | | | | | | (201) | (99) |
| - | - | - | - | - | - | Fibroblast | Fibroblast |
| | | | | | | (2189) | (4) |
| - | - | - | - | - | - | Granulocyte | Granulocyte |
| | | | | | | (761) | (73) |
| - | - | - | - | - | - | Granulocytopoietic cell | Granulocytopoietic cell |
| | | | | | | (221) | (14) |
| - | - | - | - | - | - | Hematopoietic precursor cell | Hematopoietic precursor cell |
| | | | | | | (265) | (24) |

| | | | | | | | |
|---|---|---|---|---|---|--|---|
| - | - | - | - | - | - | Hepatocyte (391) | Hepatocyte (374) |
| - | - | - | - | - | - | Immature B cell (344) | Immature B cell (3) |
| - | - | - | - | - | - | Immature T cell (1337) | Immature T cell (222) |
| - | - | - | - | - | - | Keratinocyte (330) | Keratinocyte (1035) |
| - | - | - | - | - | - | Kidney collecting duct epithelial cell (121) | Kidney collecting duct epithelial cell (24) |
| - | - | - | - | - | - | Late pro-B cell (306) | Late pro-B cell (14) |
| - | - | - | - | - | - | Leukocyte (683) | Leukocyte (19) |
| - | - | - | - | - | - | Luminal epithelial cell of mammary gland | Luminal epithelial cell of mammary gland |

(578) (35)

Lung endothelial cell Lung endothelial cell

(693) (24)

Macrophage Macrophage

(395) (208)

Mesenchymal cell Mesenchymal cell

(830) (5200)

Mesenchymal stem cell Mesenchymal stem
cell

(499) (169)

Monocyte Monocyte

(331) (42)

Myeloid cell Myeloid cell

(1208) (2)

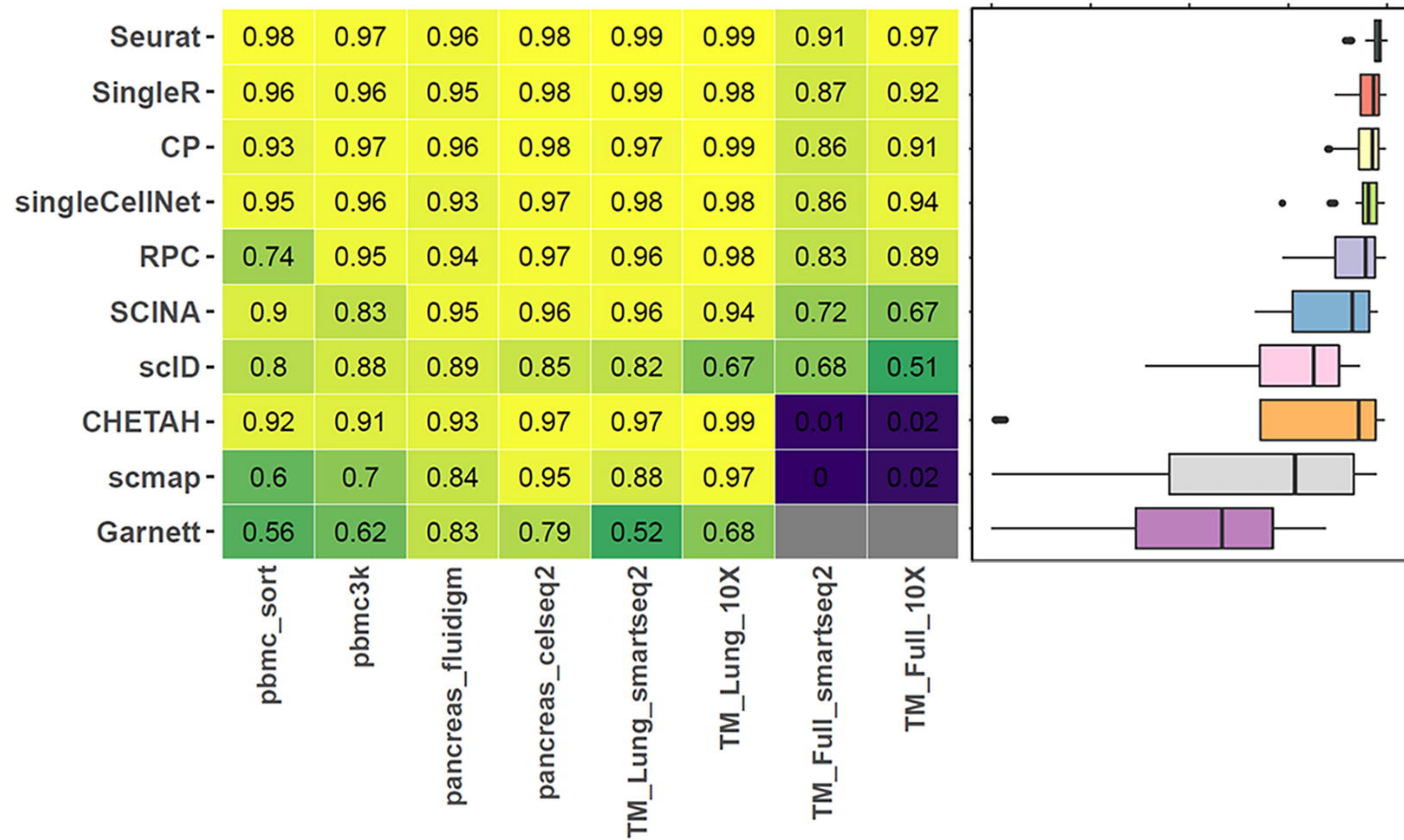
Natural killer cell Natural killer cell

(171) (142)

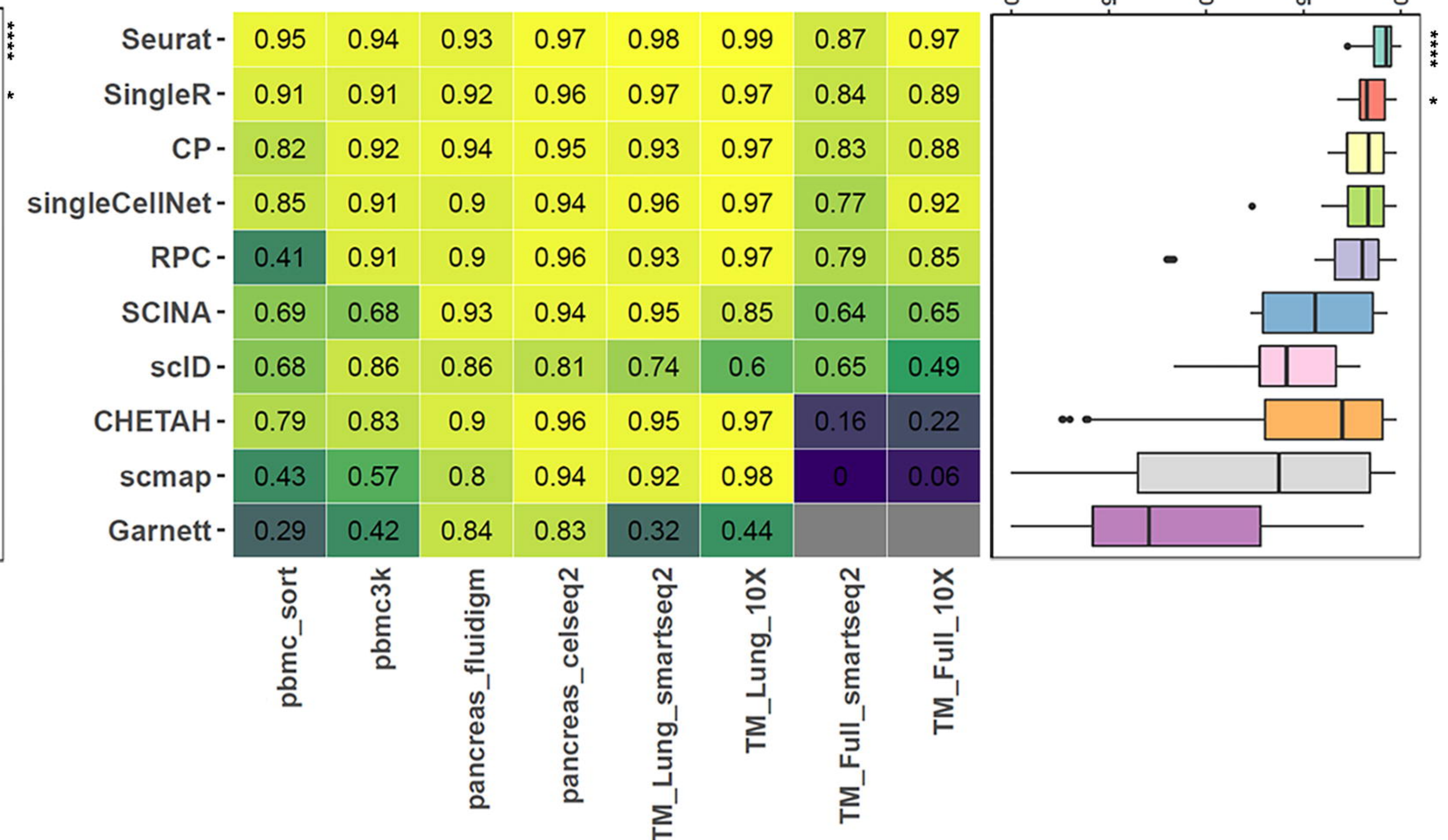
| | | | | | | | |
|---|---|---|---|---|---|---|--|
| - | - | - | - | - | - | Skeletal muscle satellite cell (540) | Skeletal muscle satellite cell (11) |
| - | - | - | - | - | - | Stromal cell (863) | Stromal cell (1153) |
| - | - | - | - | - | - | T cell (793) | T cell (1985) |

Note: Values are indicated as Cell Type (Cell count)

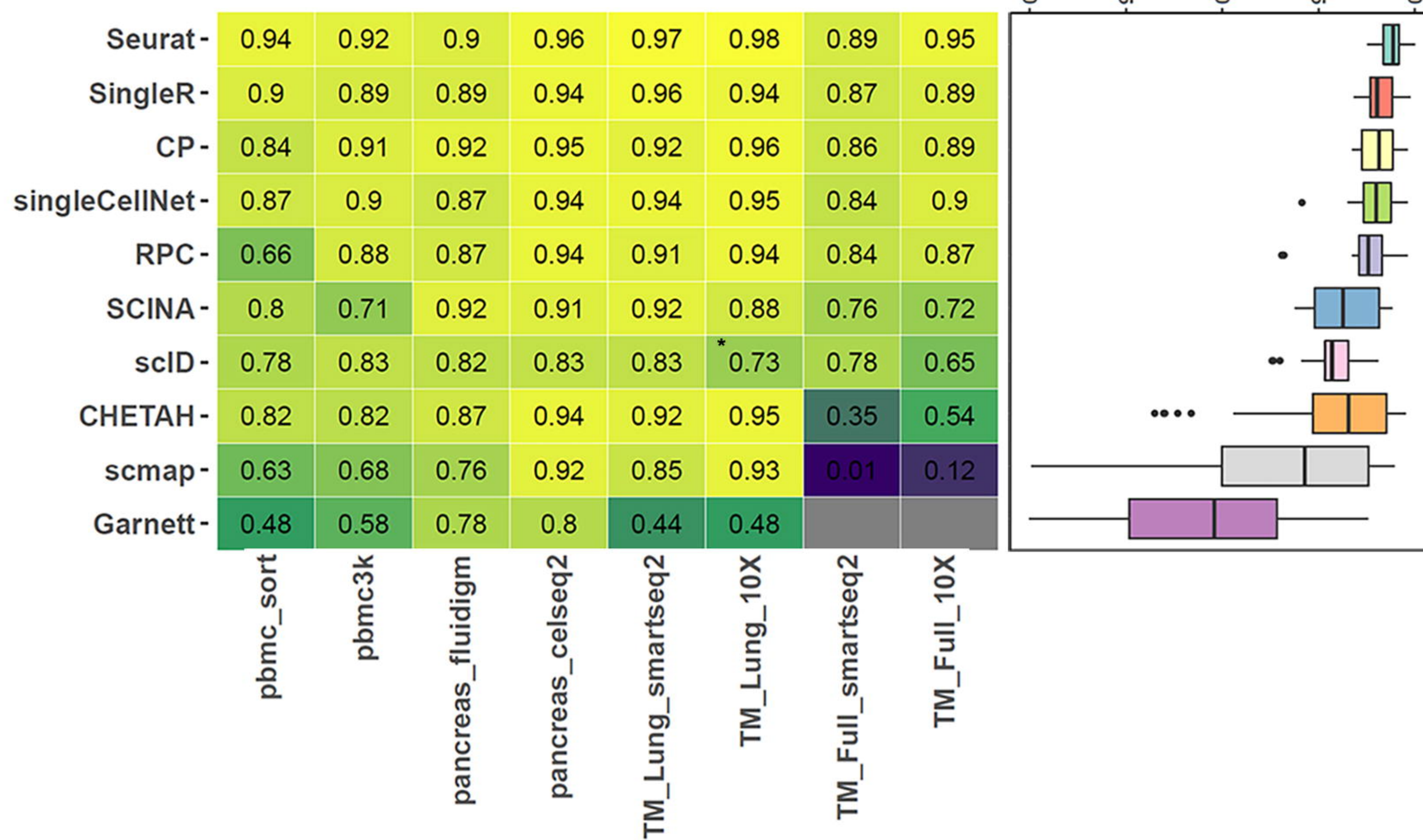
A Five-Fold Overall Accuracy



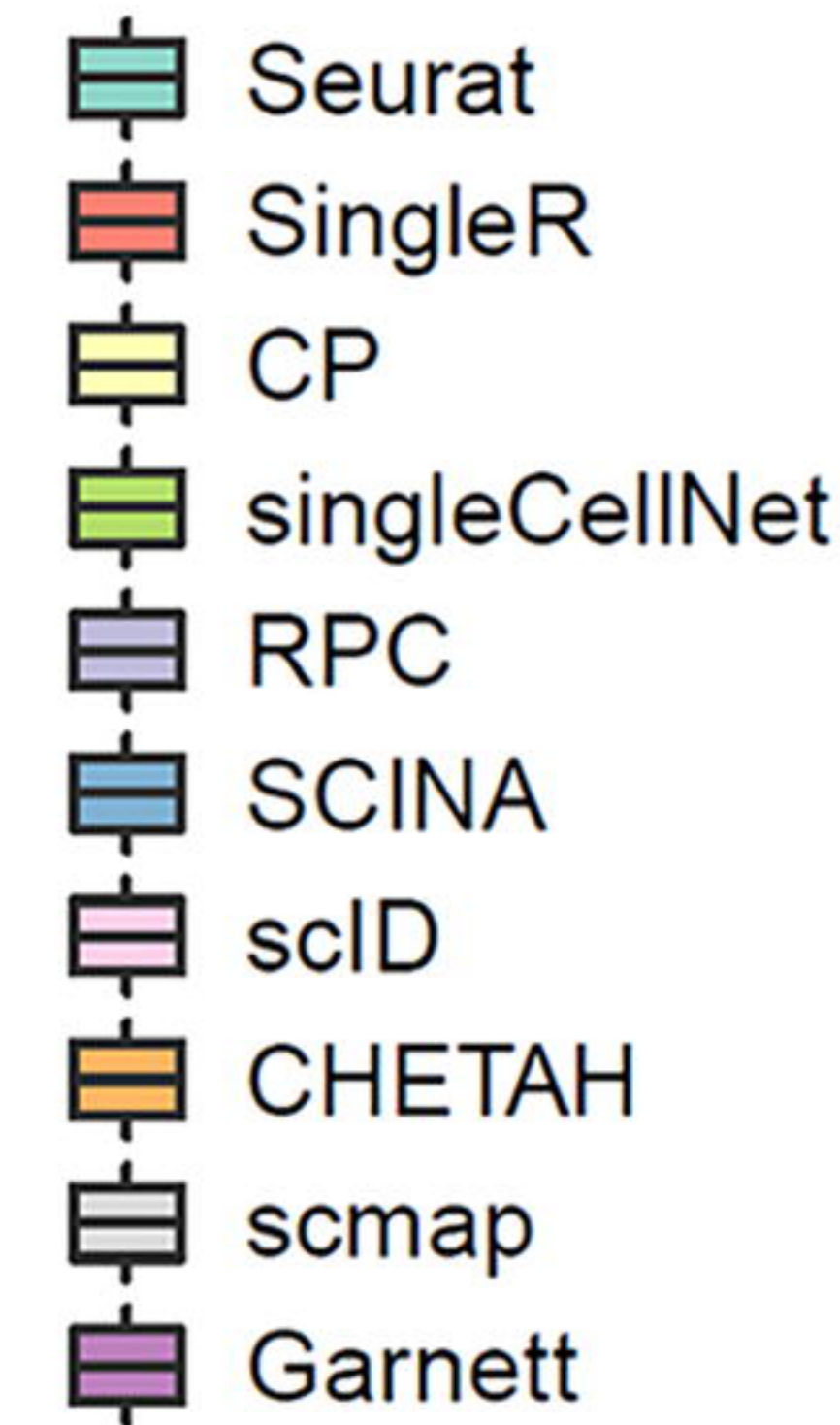
B Five-Fold ARI



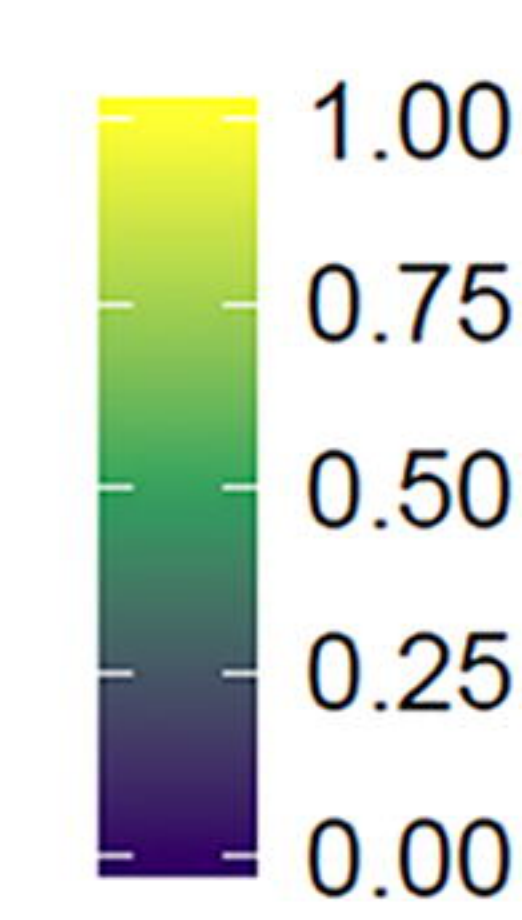
C Five-Fold V-measure



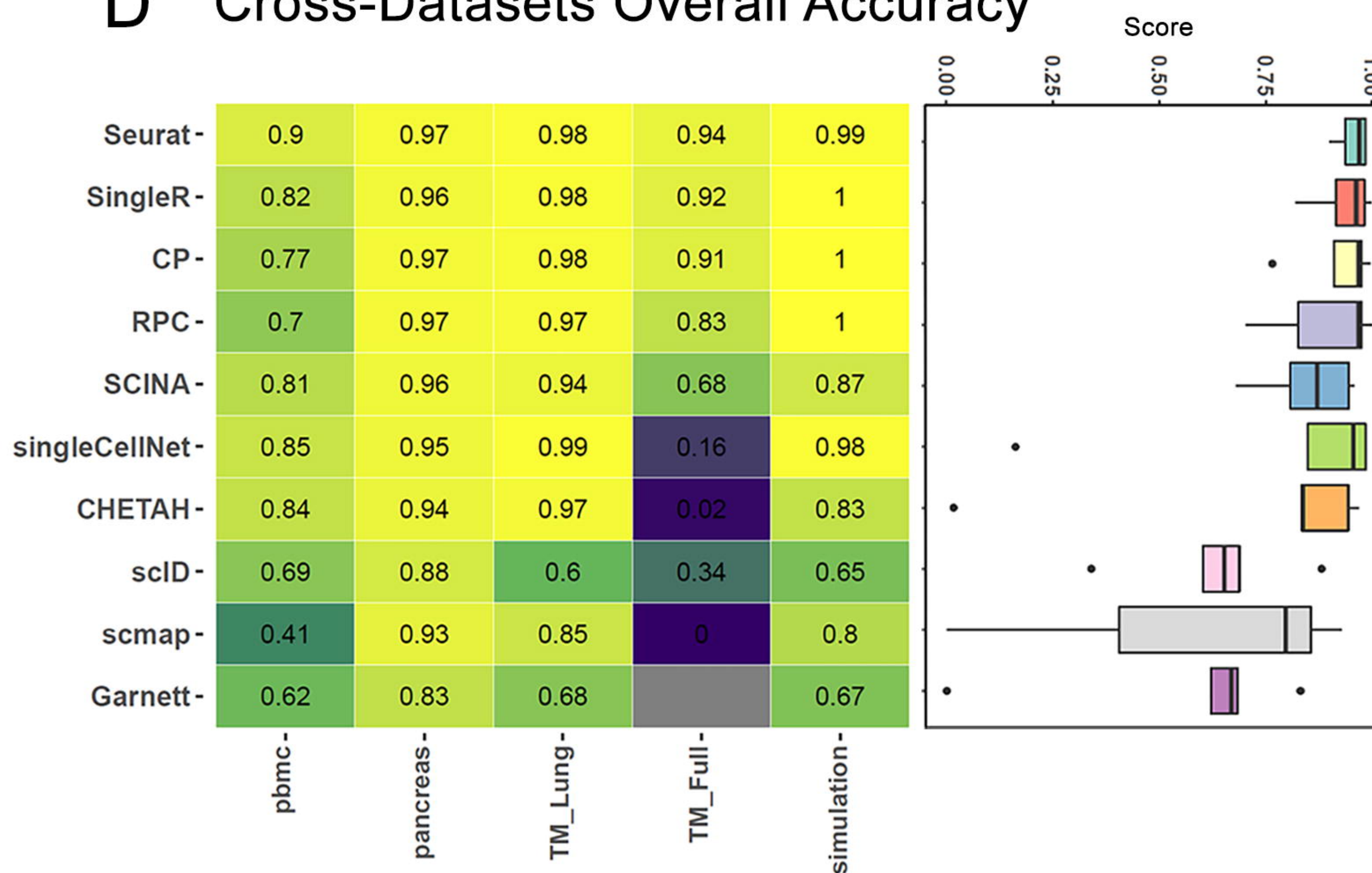
Methods



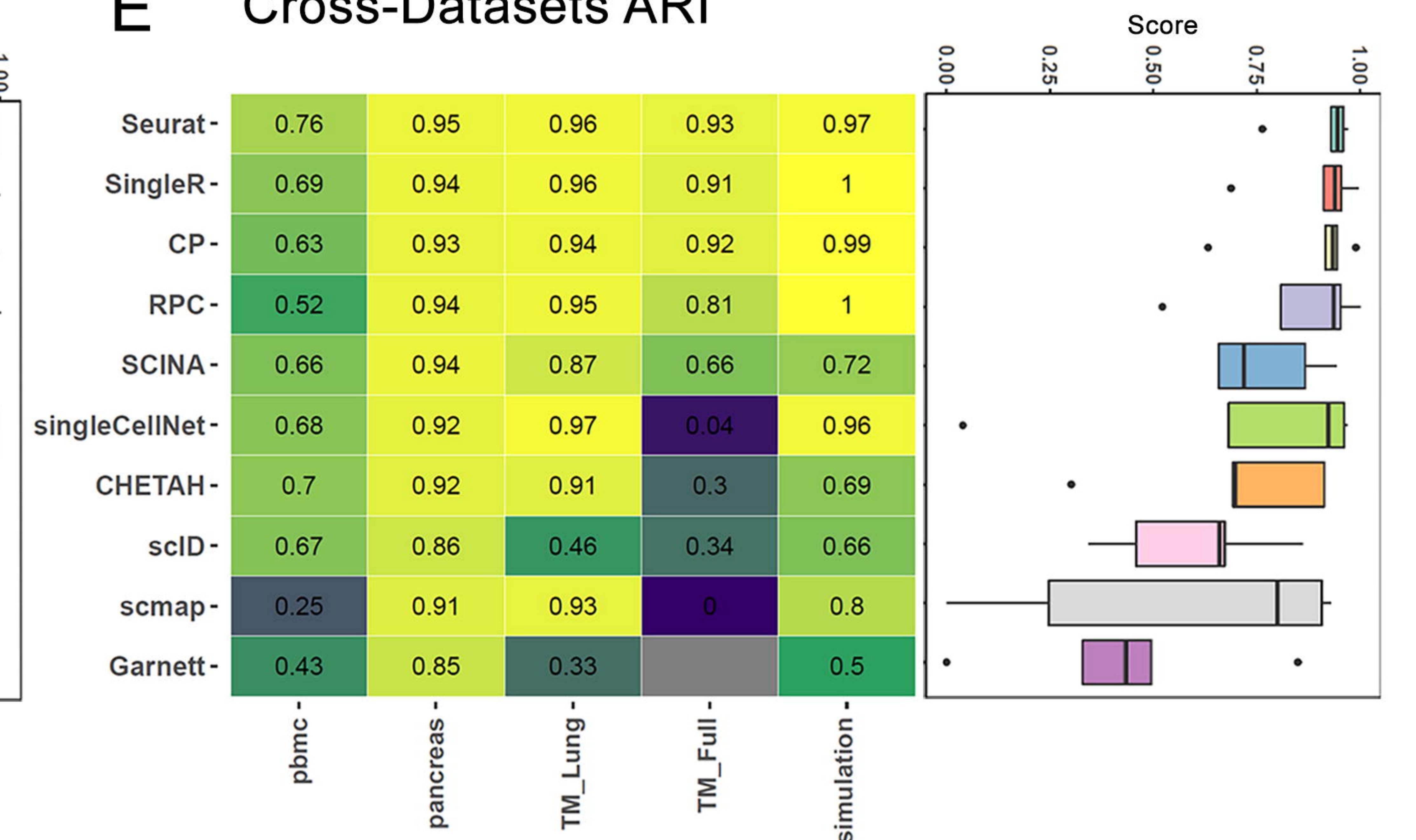
Score



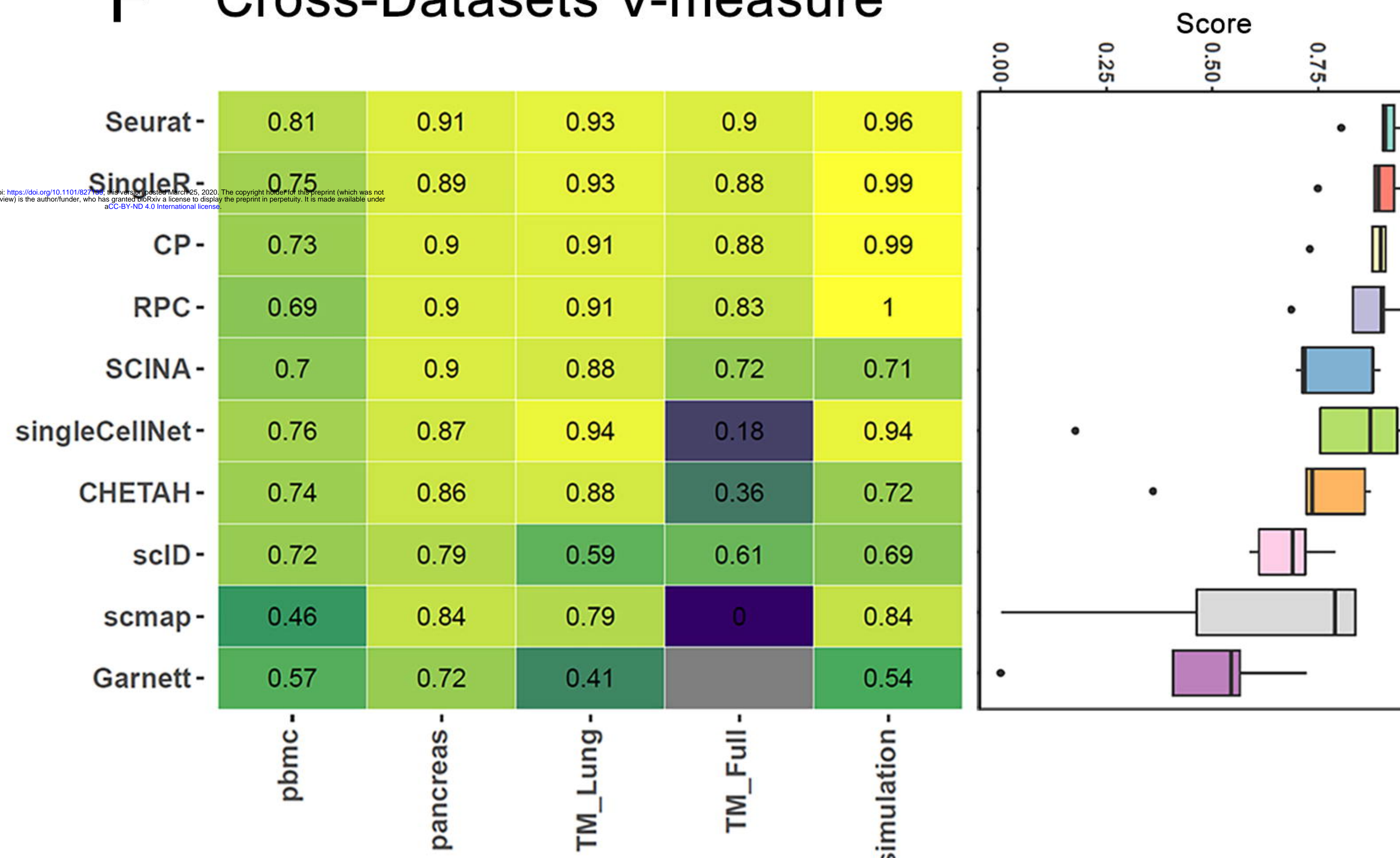
D Cross-Datasets Overall Accuracy



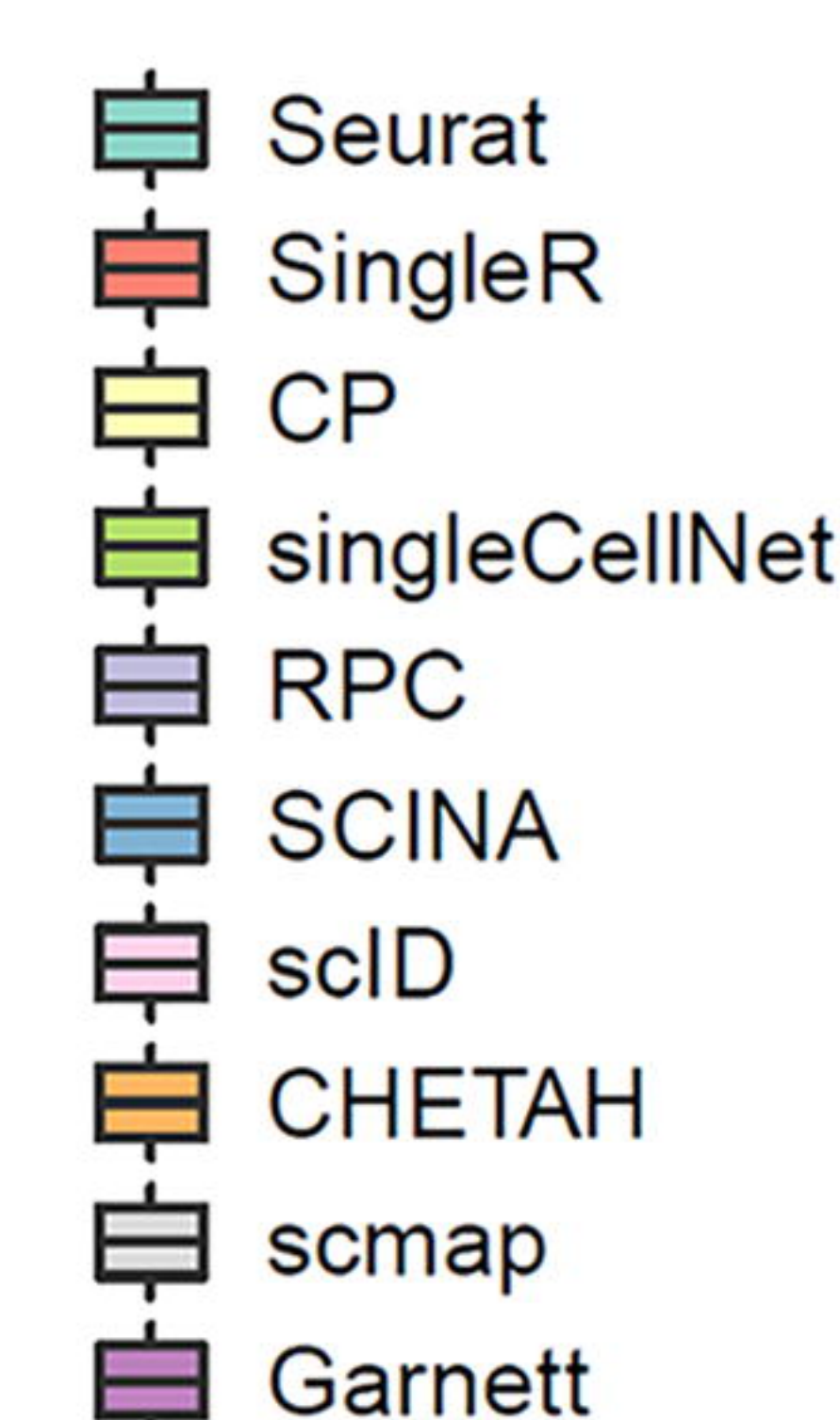
E Cross-Datasets ARI



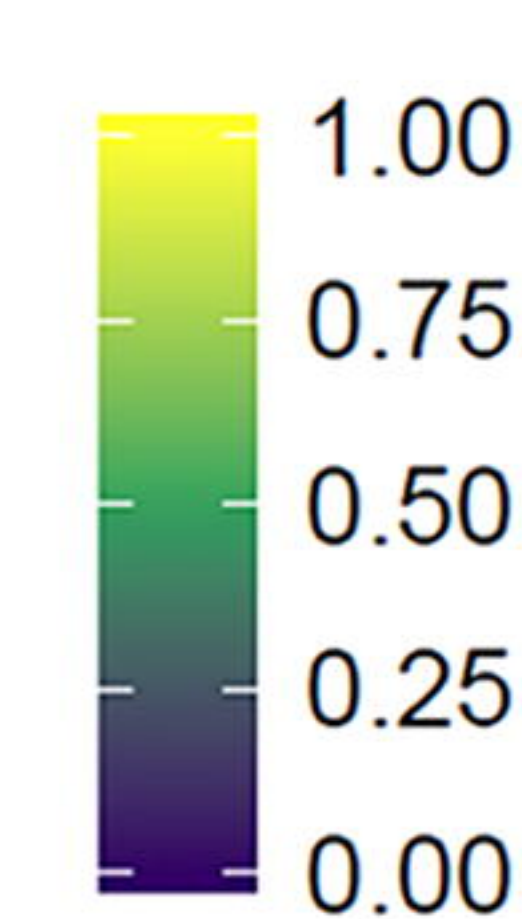
F Cross-Datasets V-measure

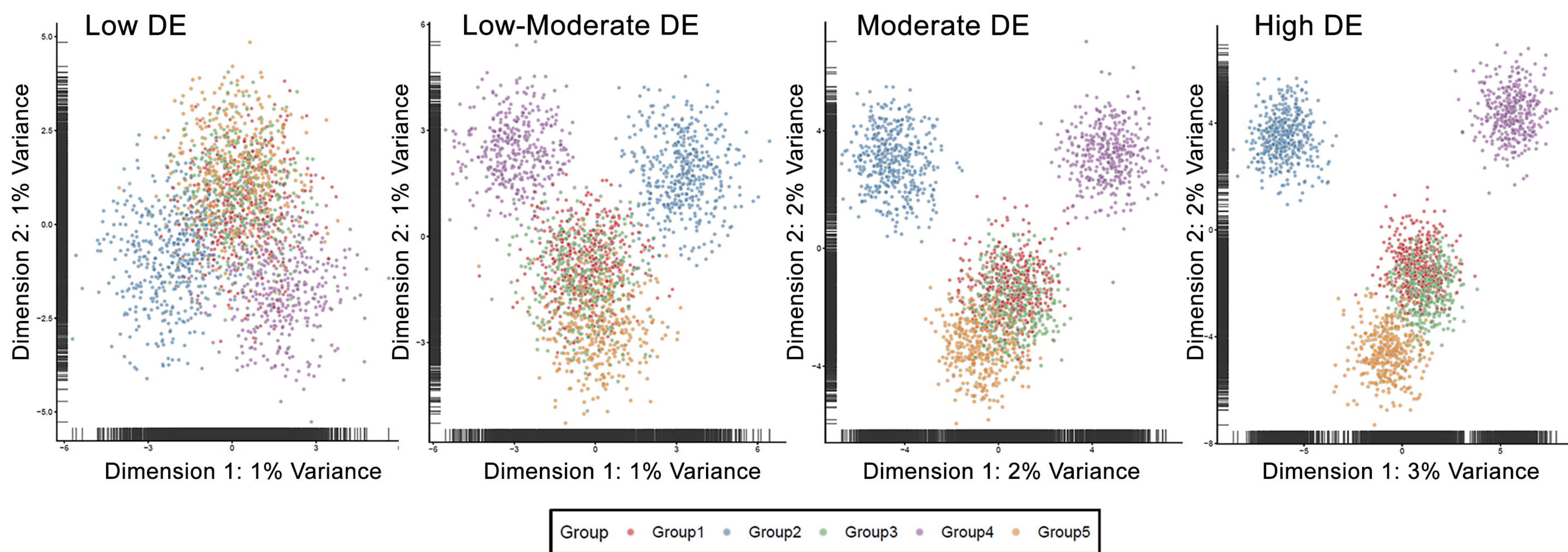
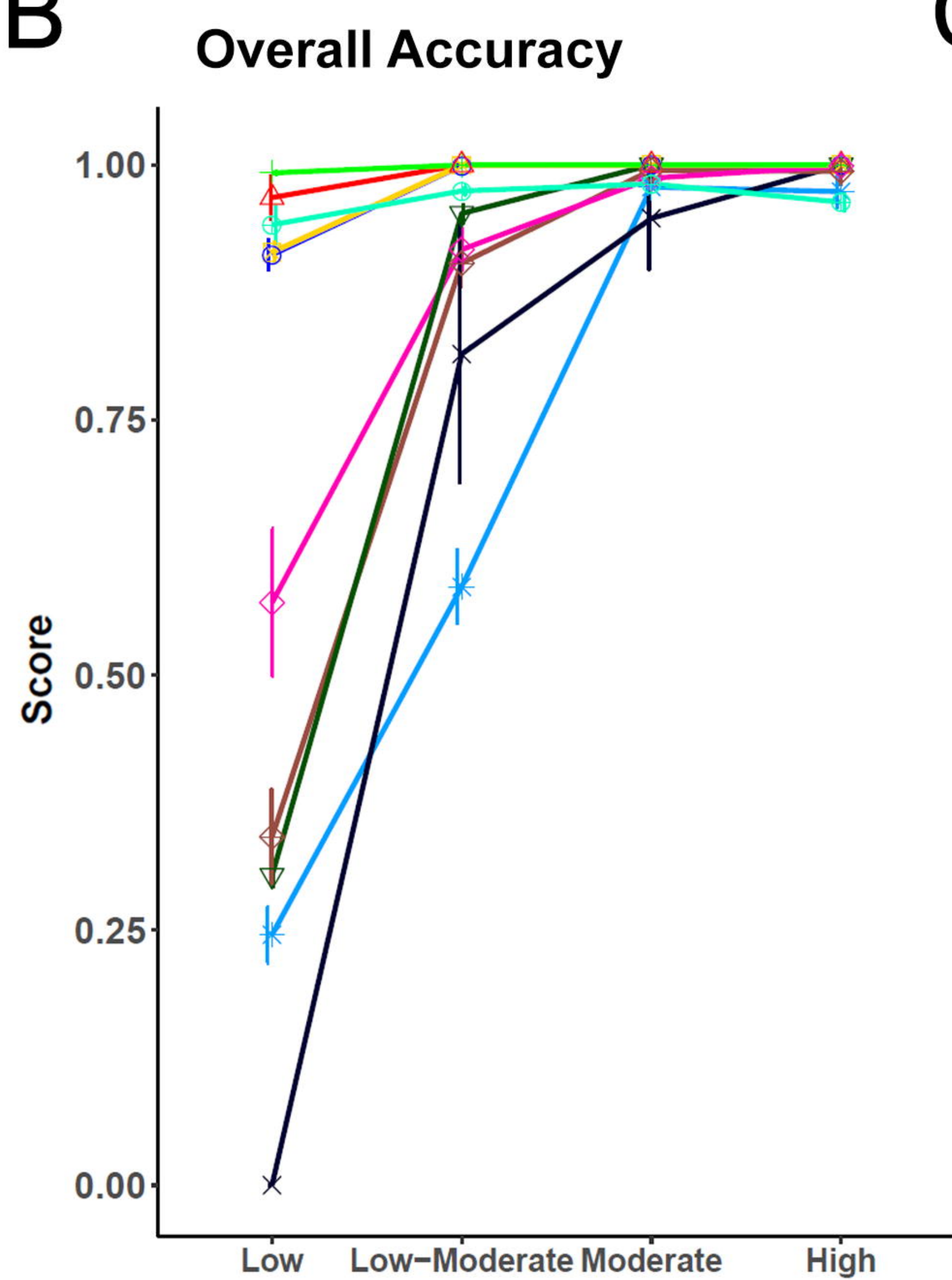
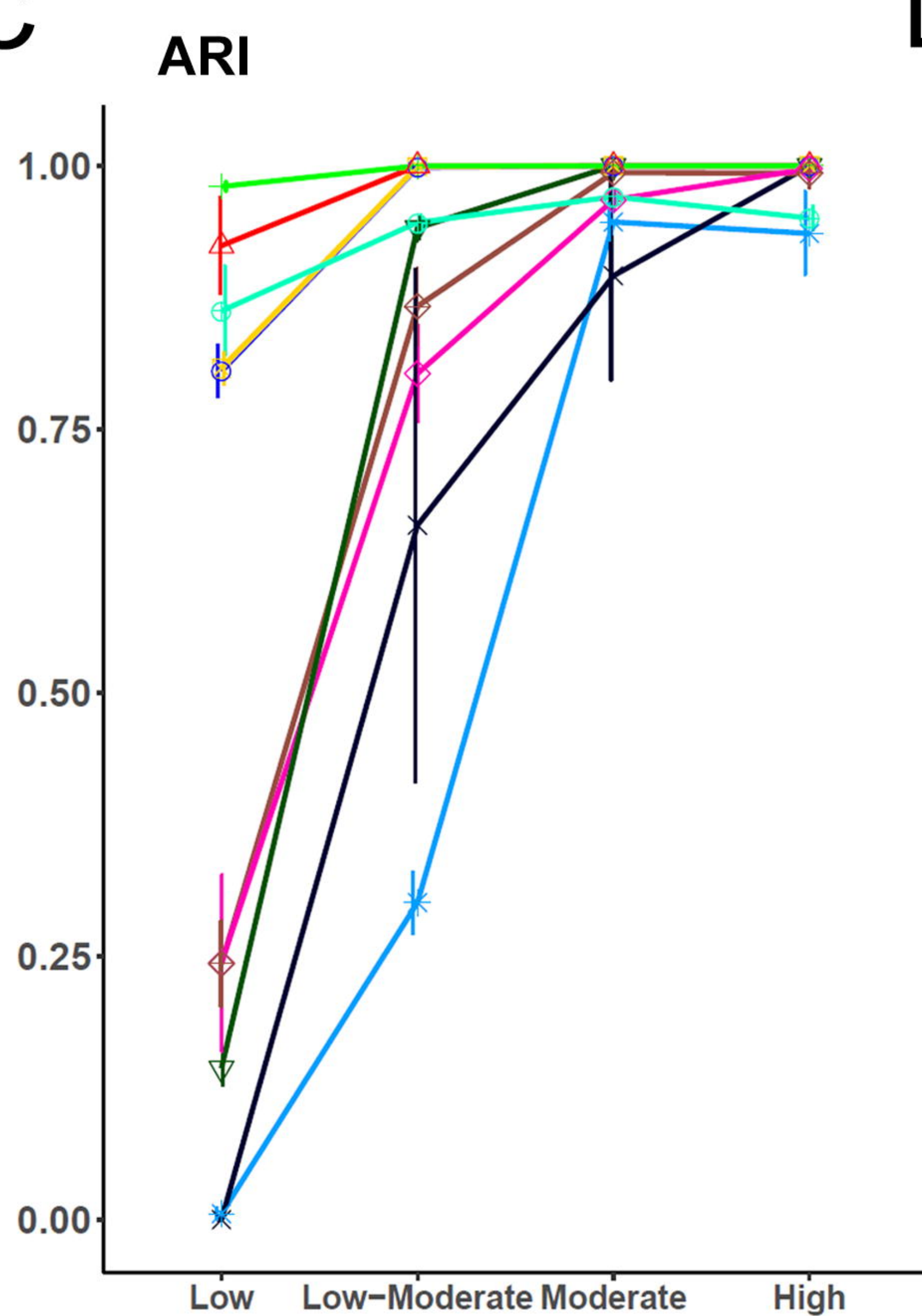
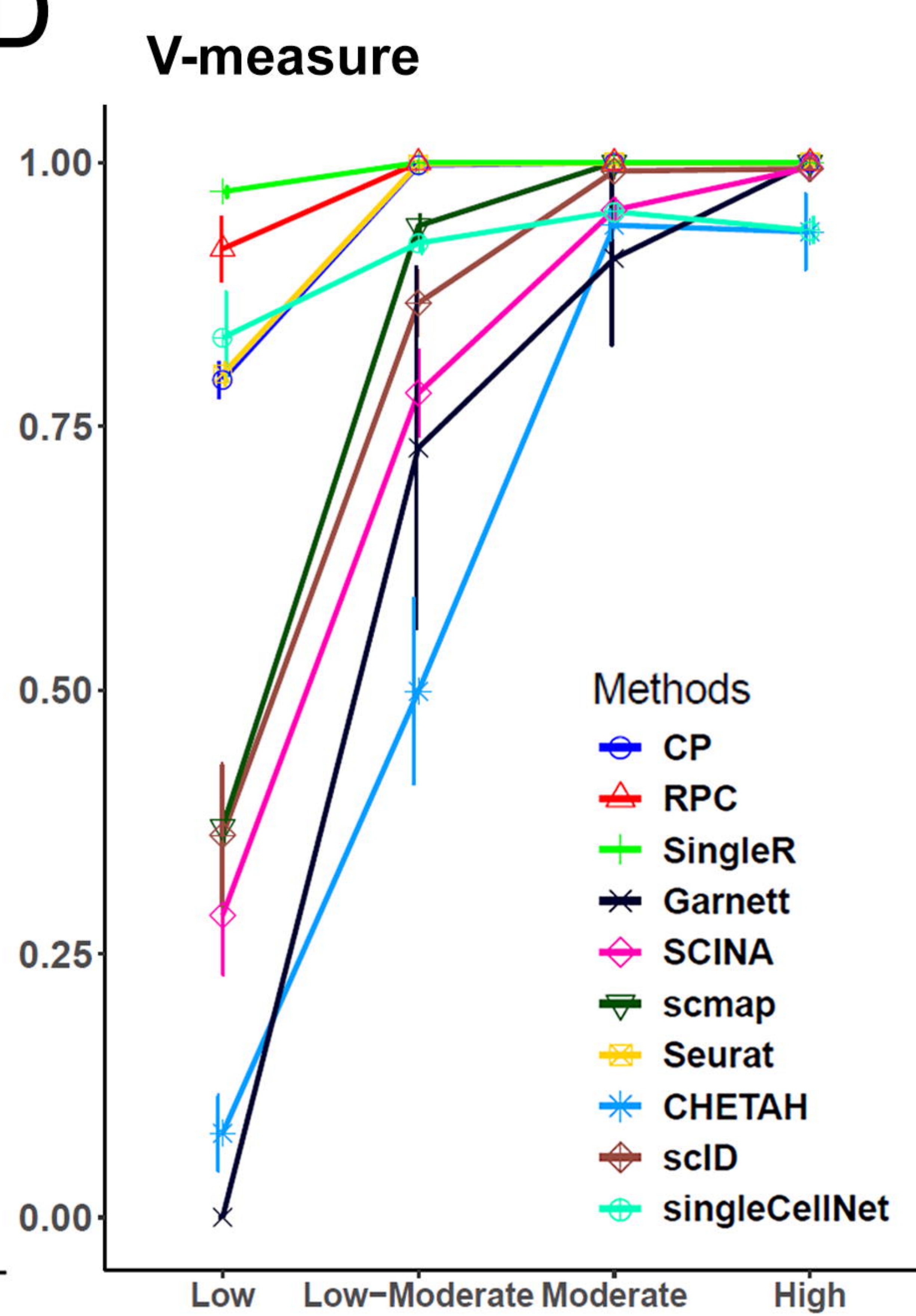
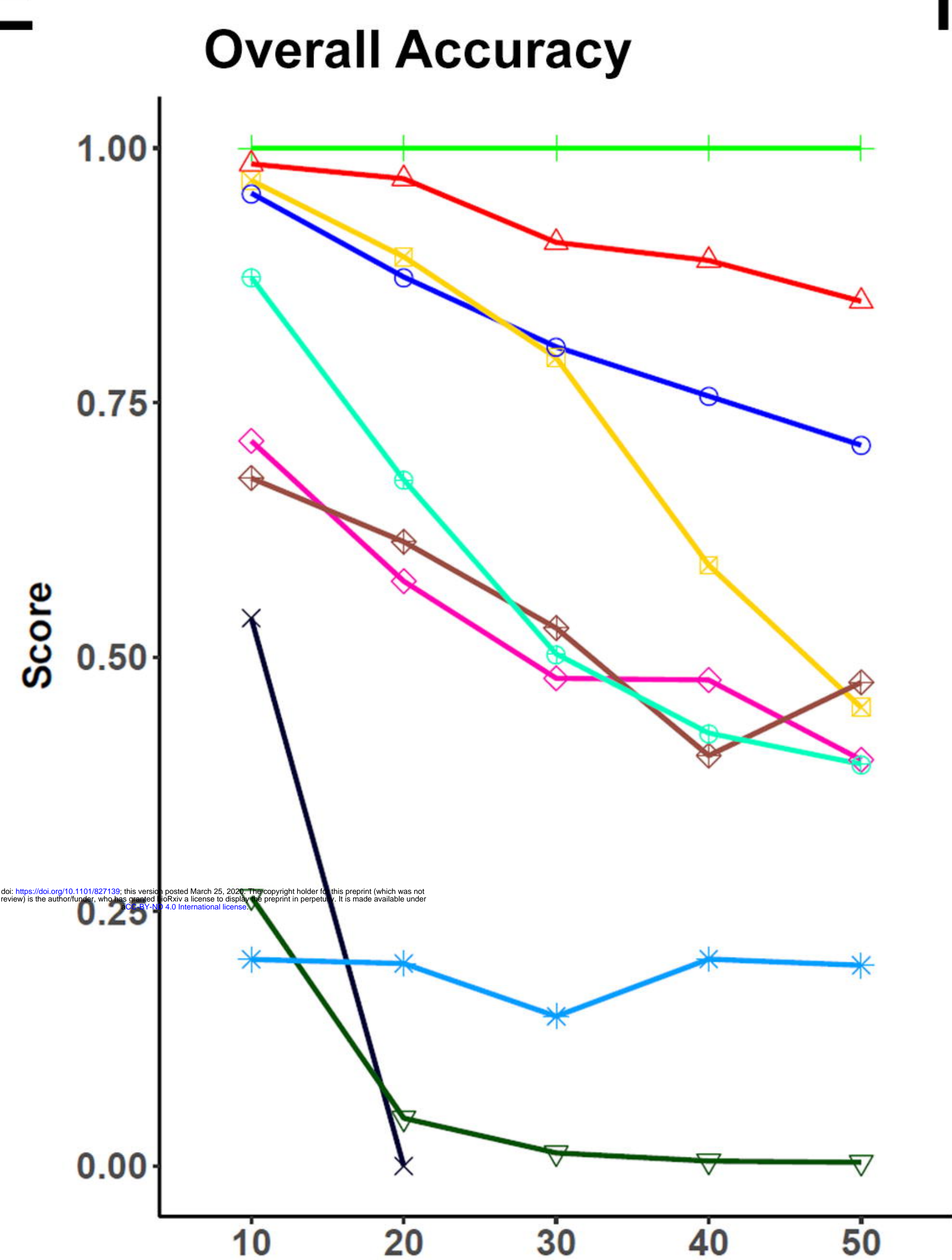
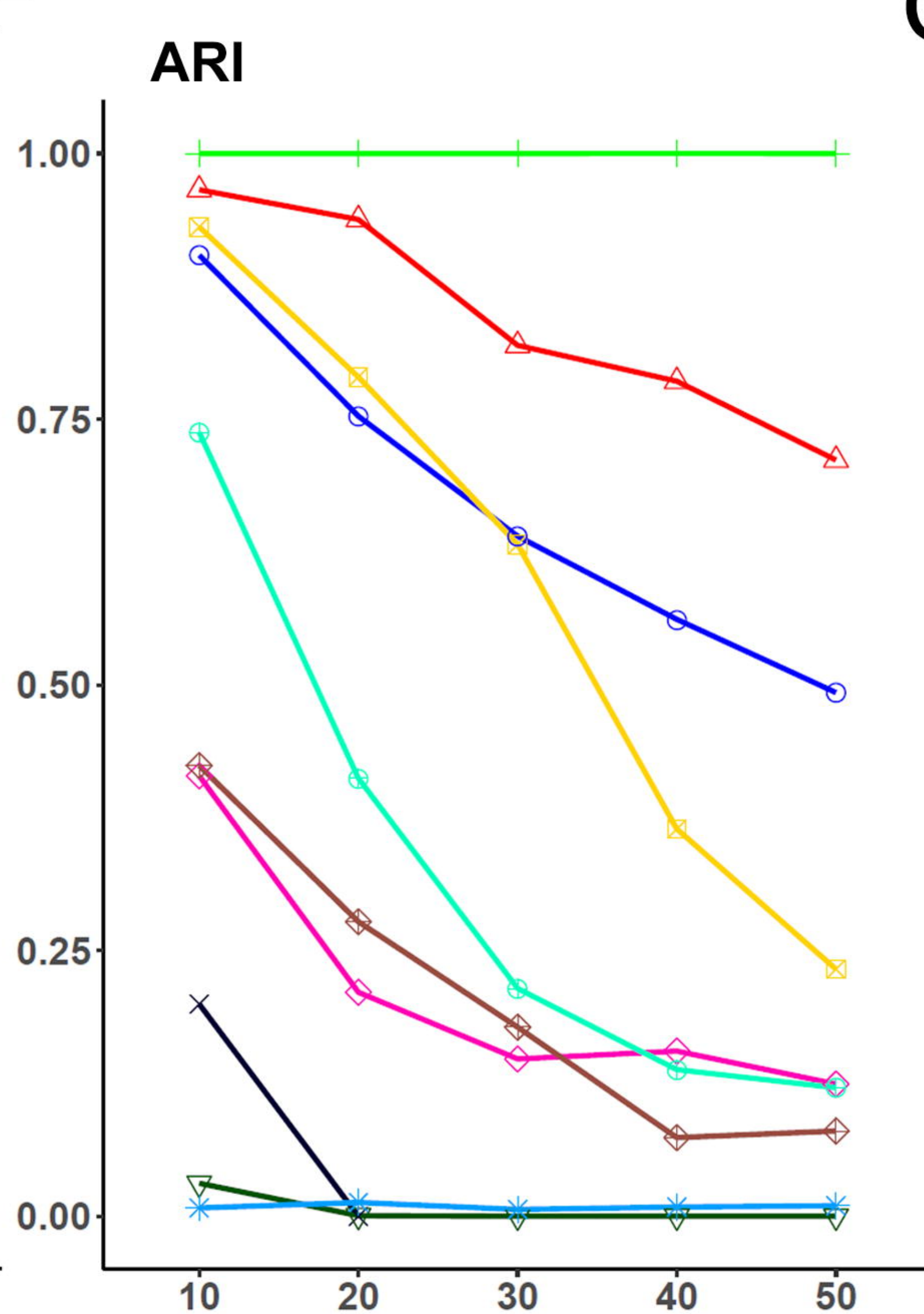
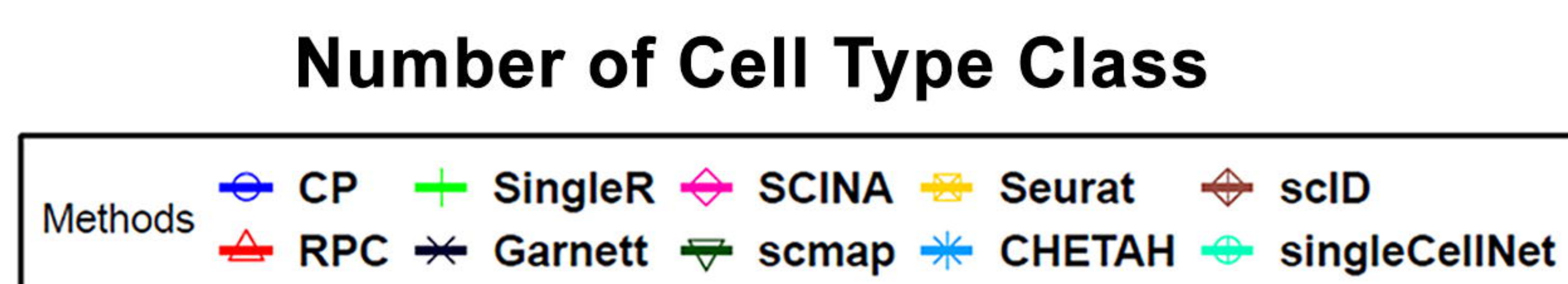
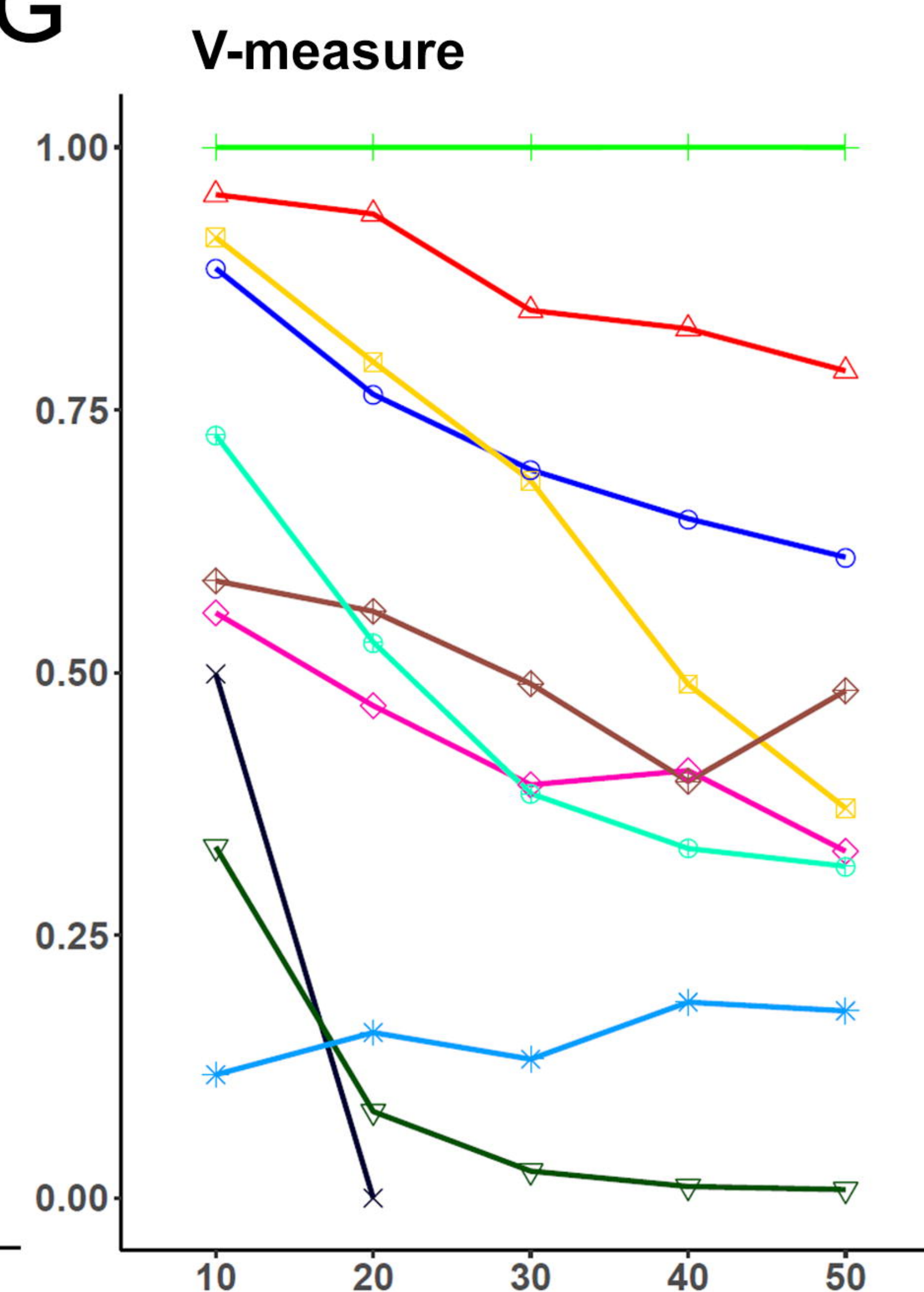


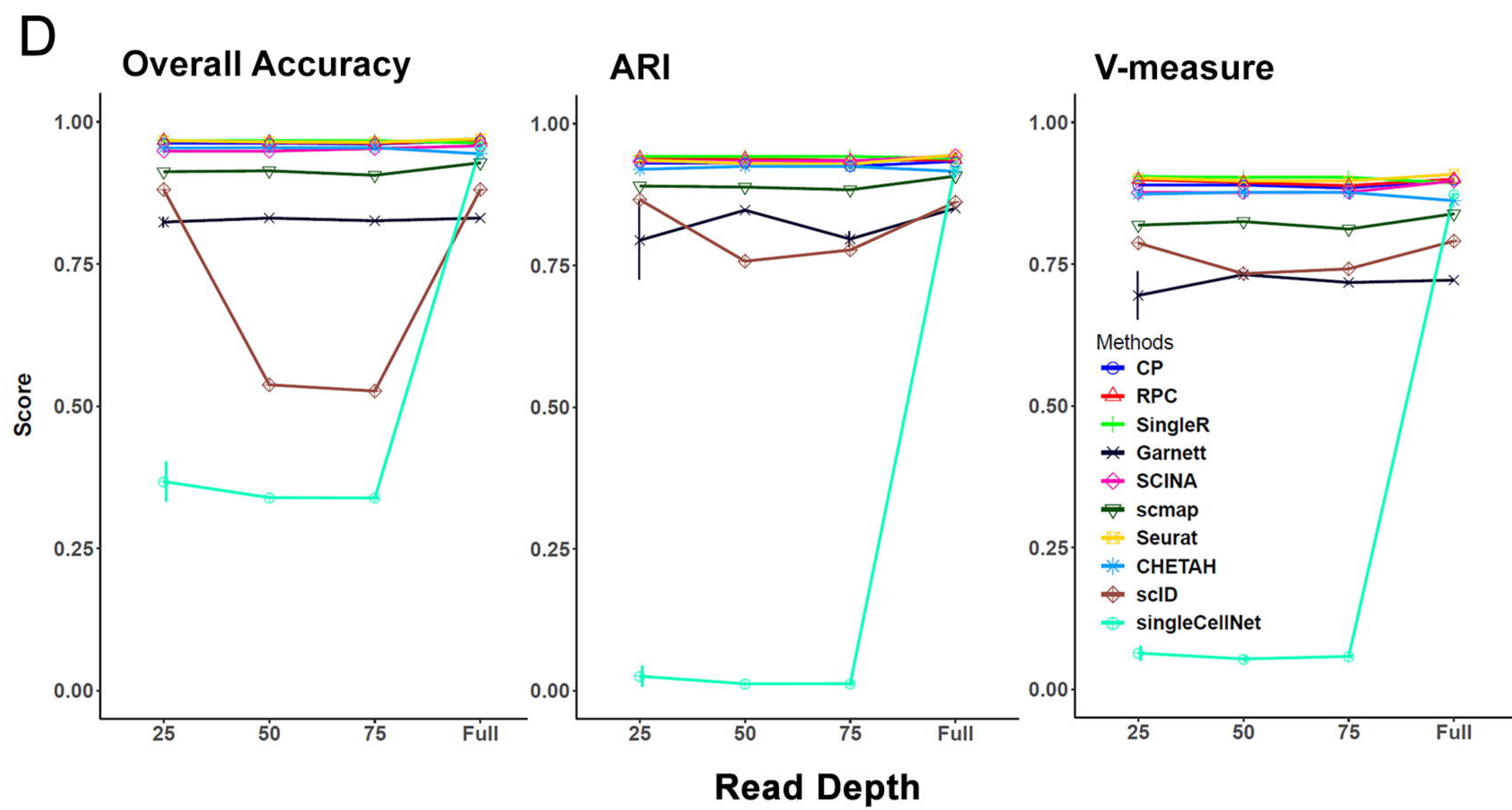
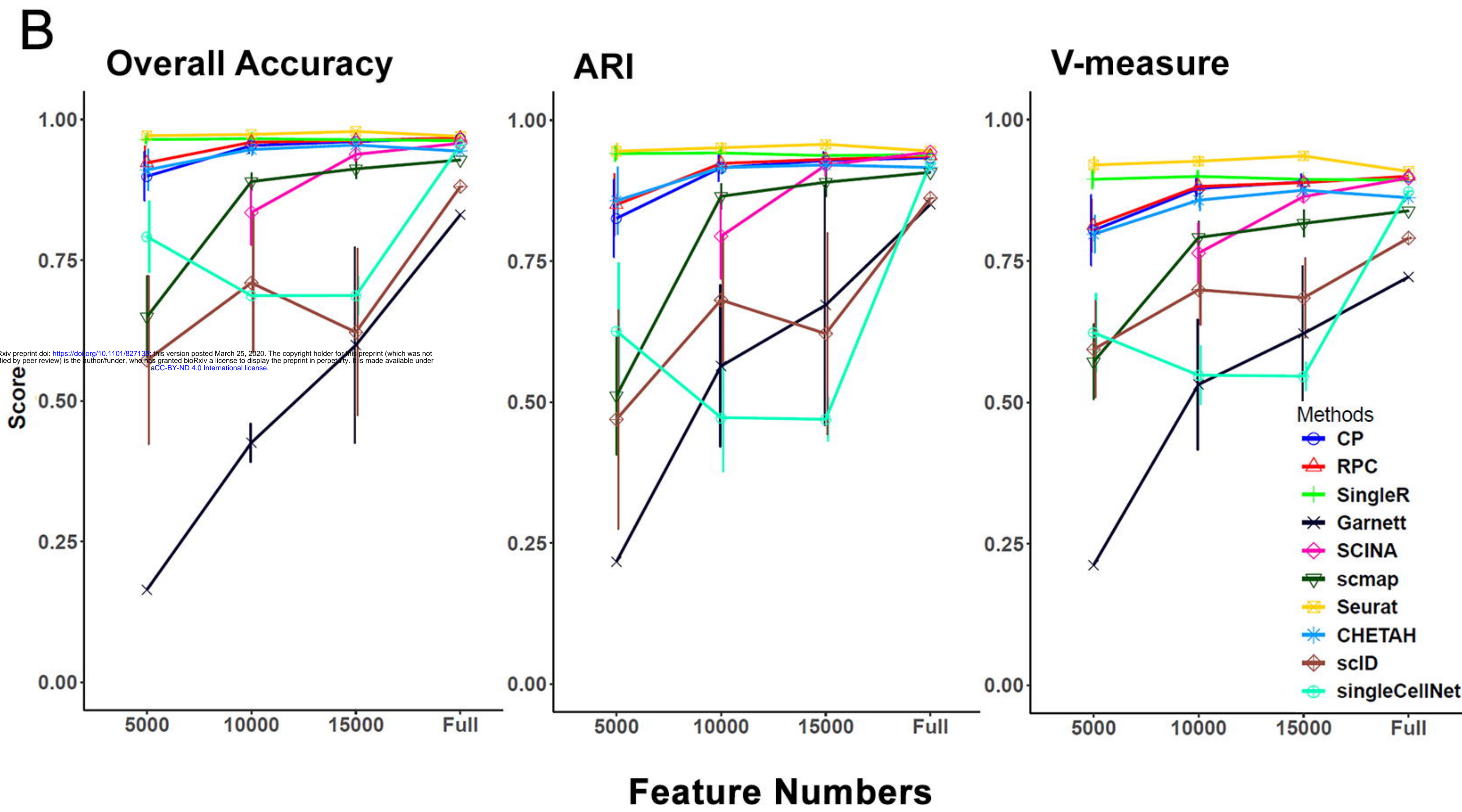
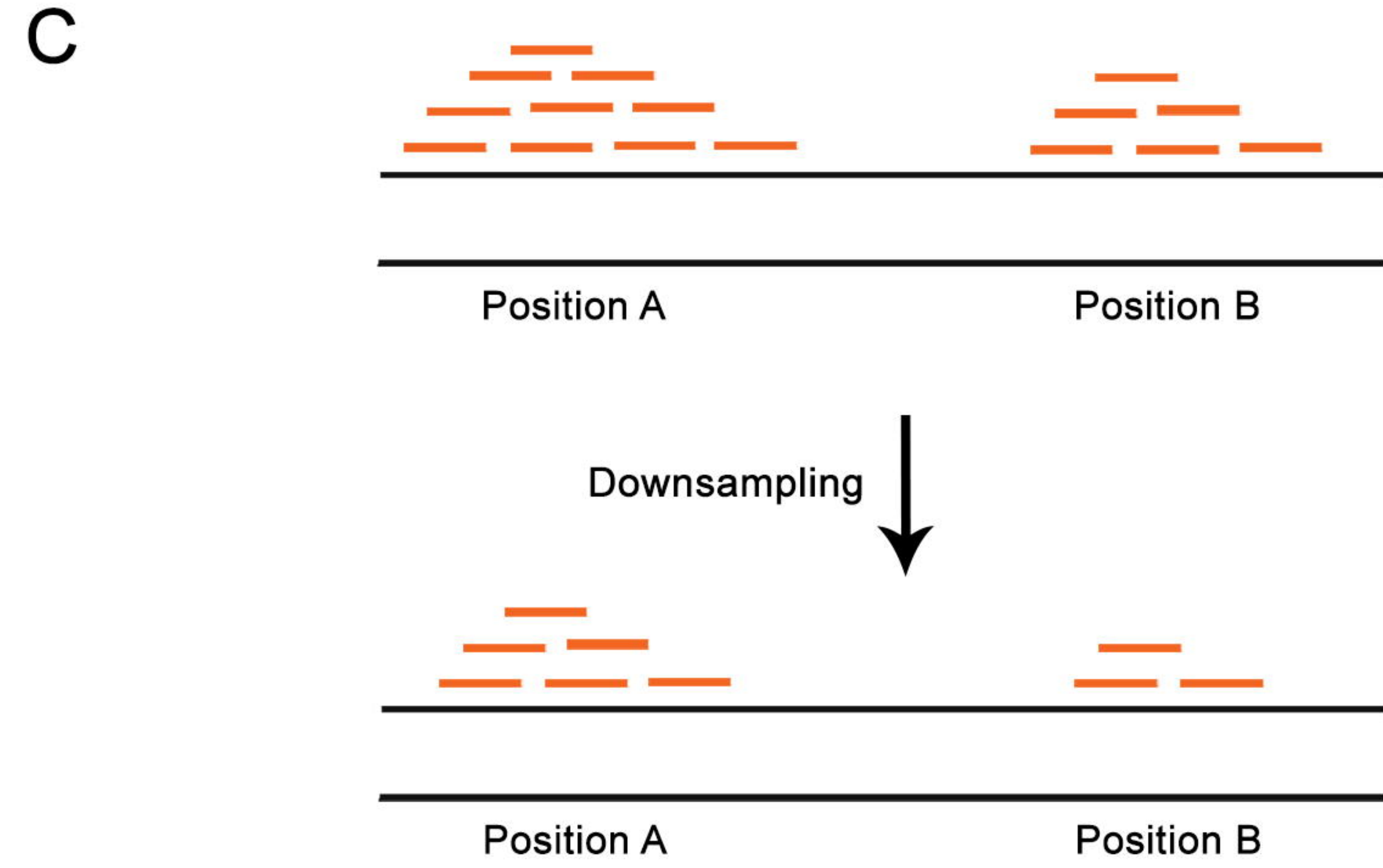
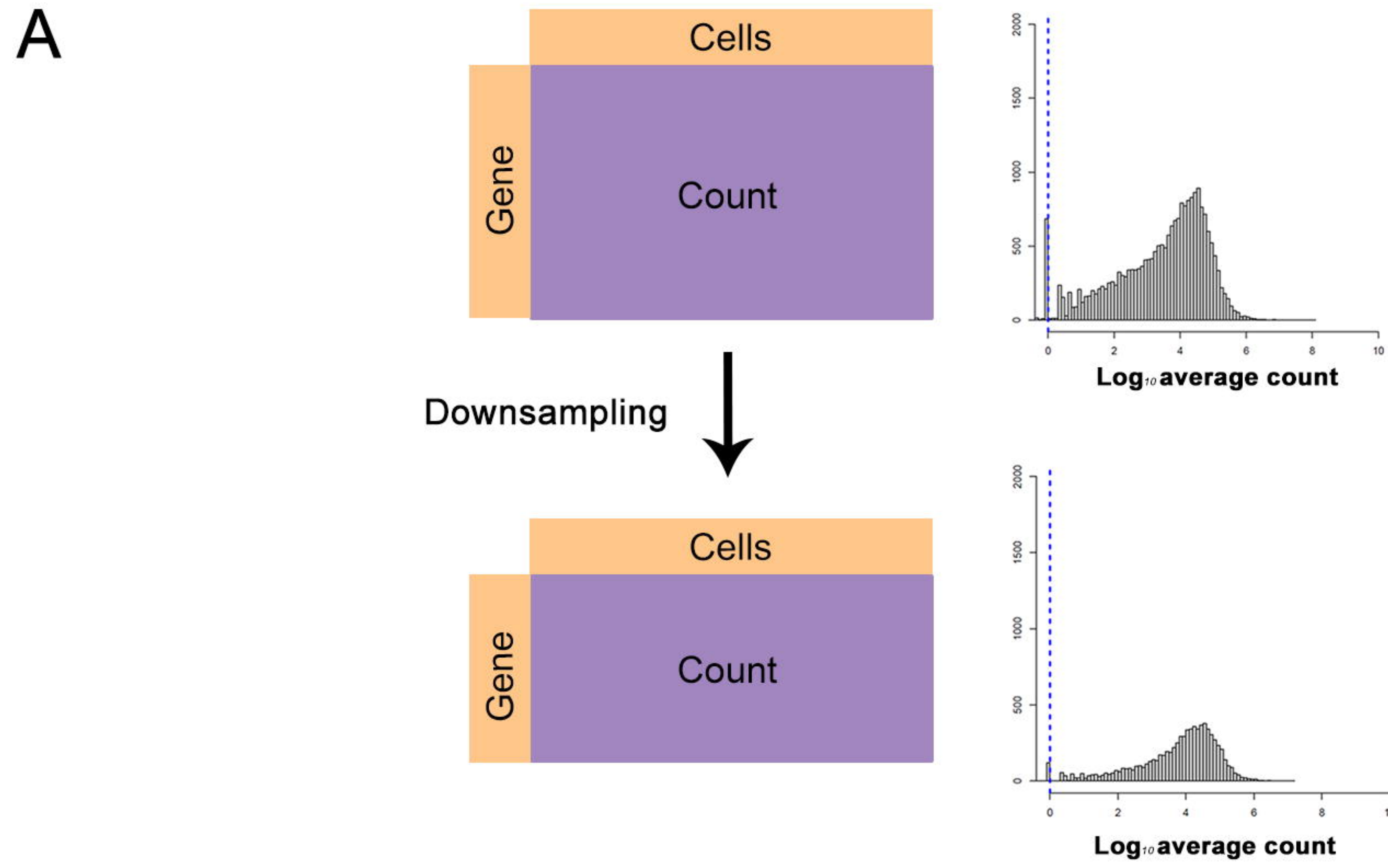
Methods

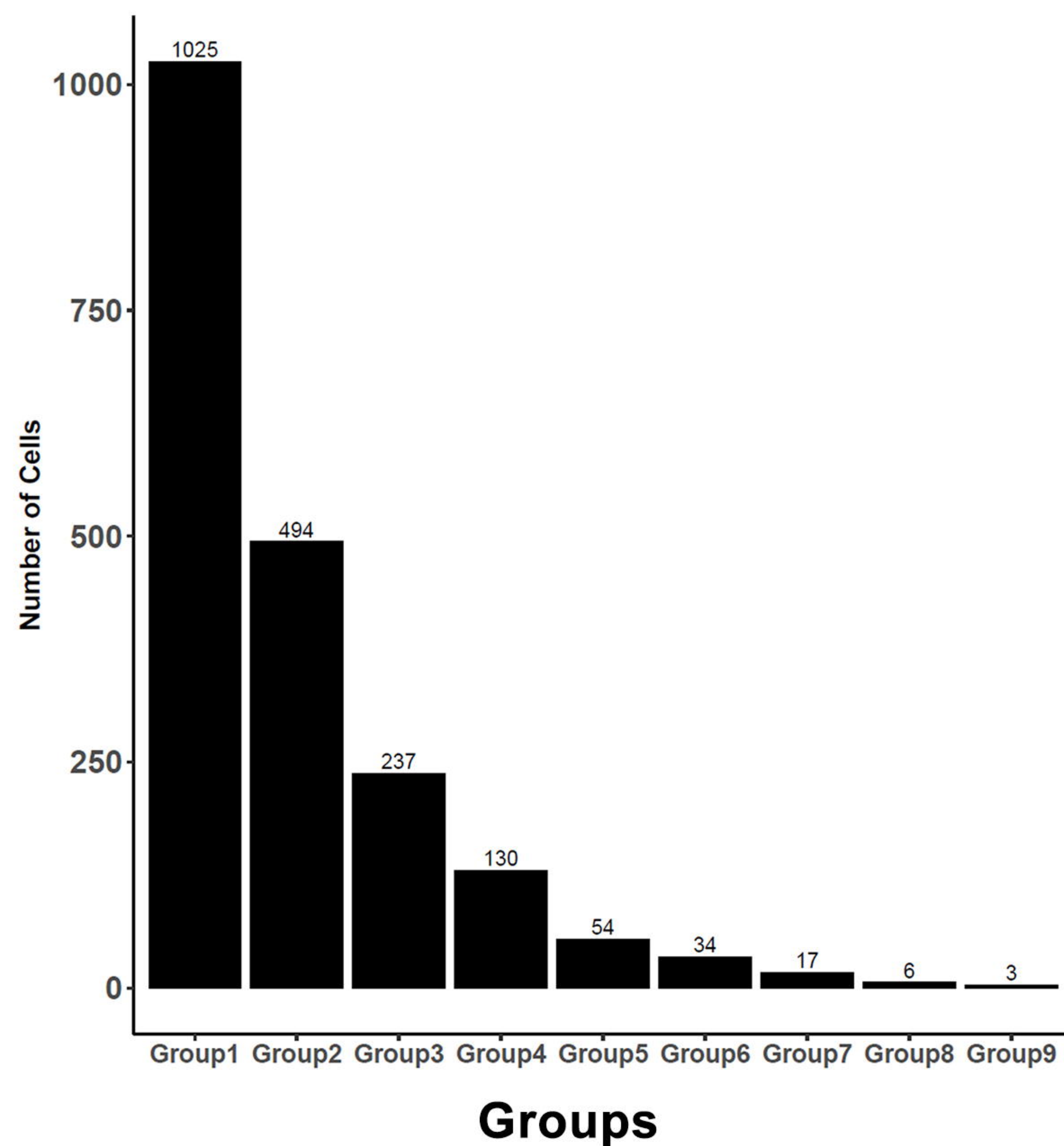
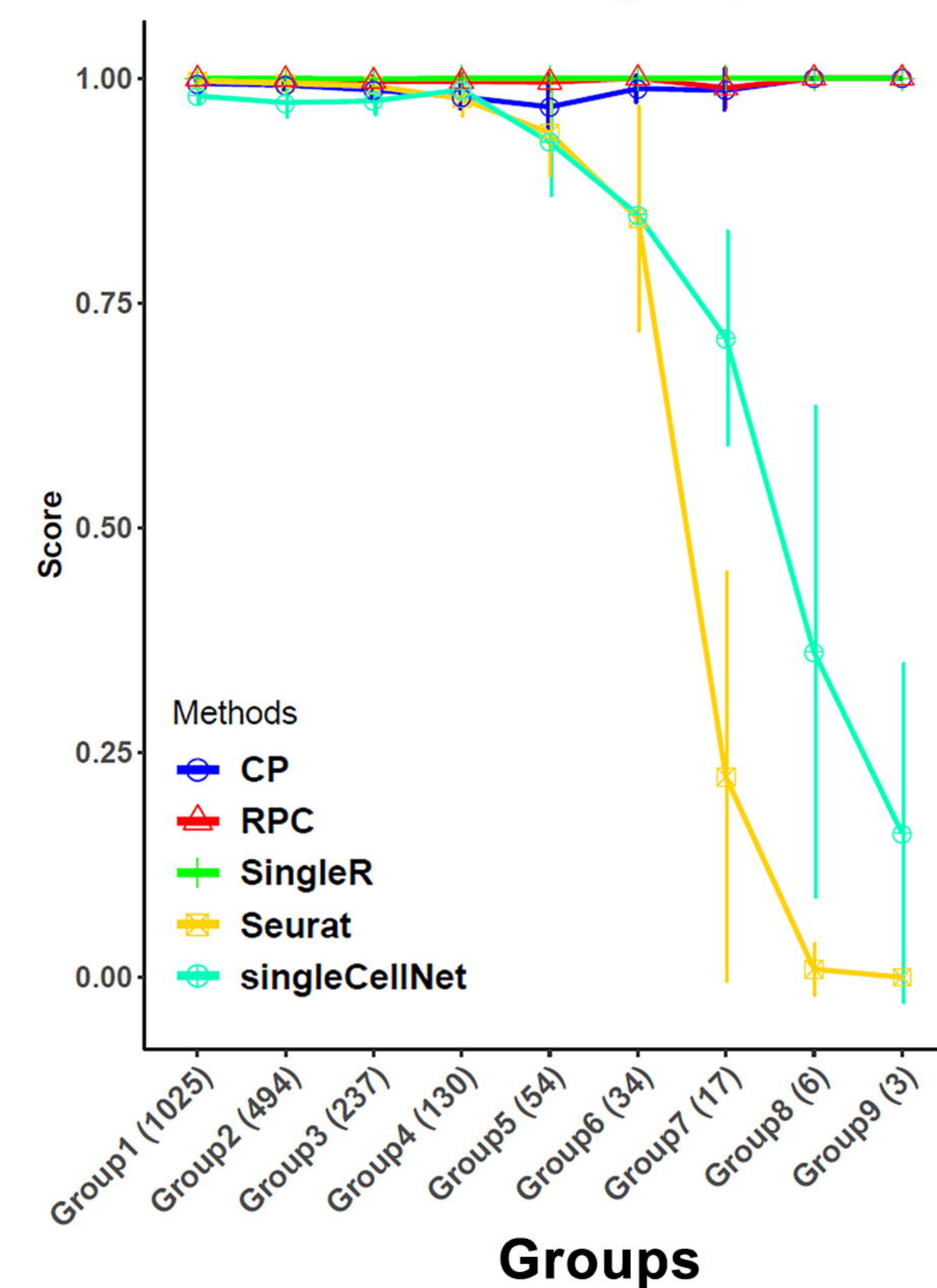
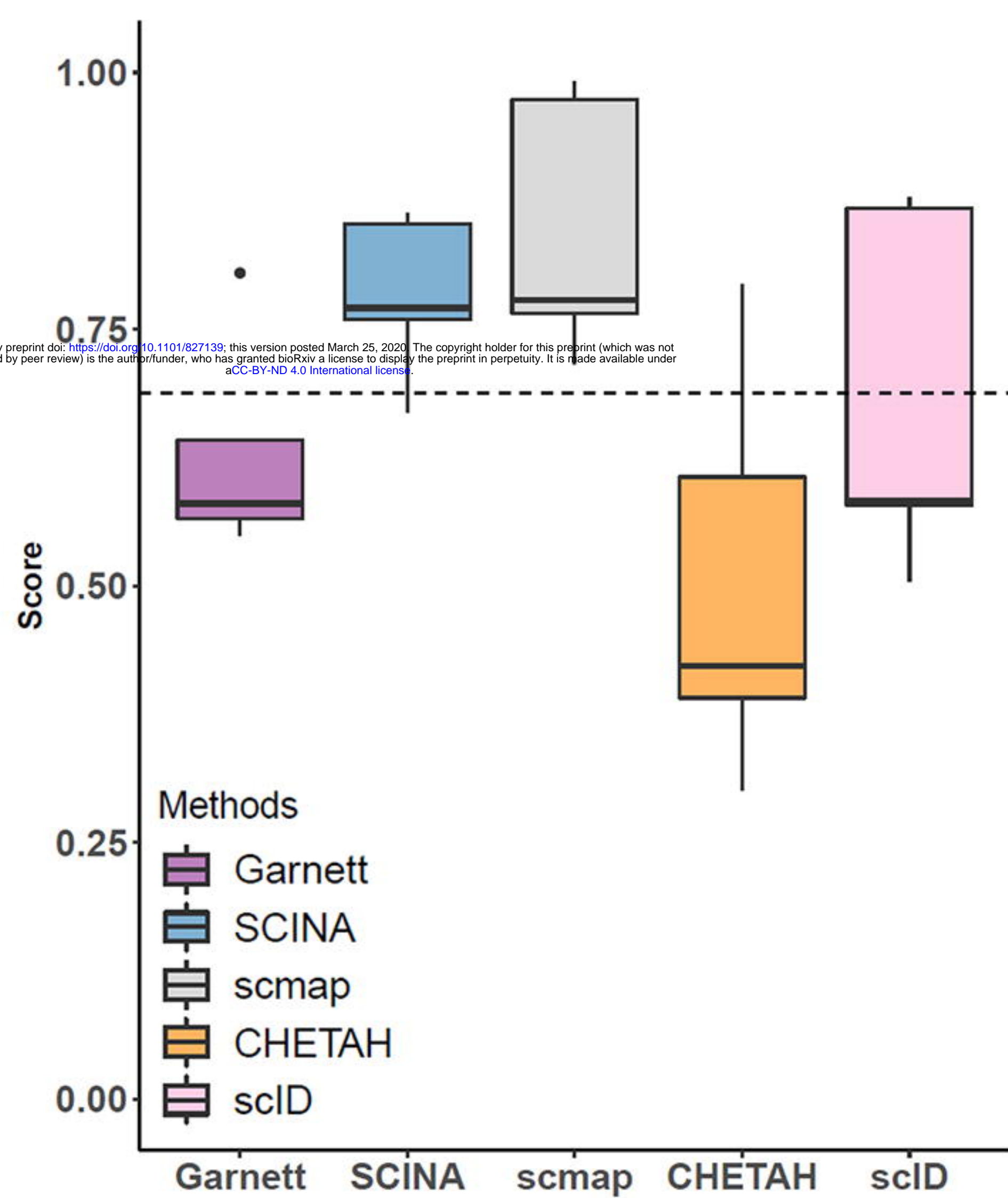
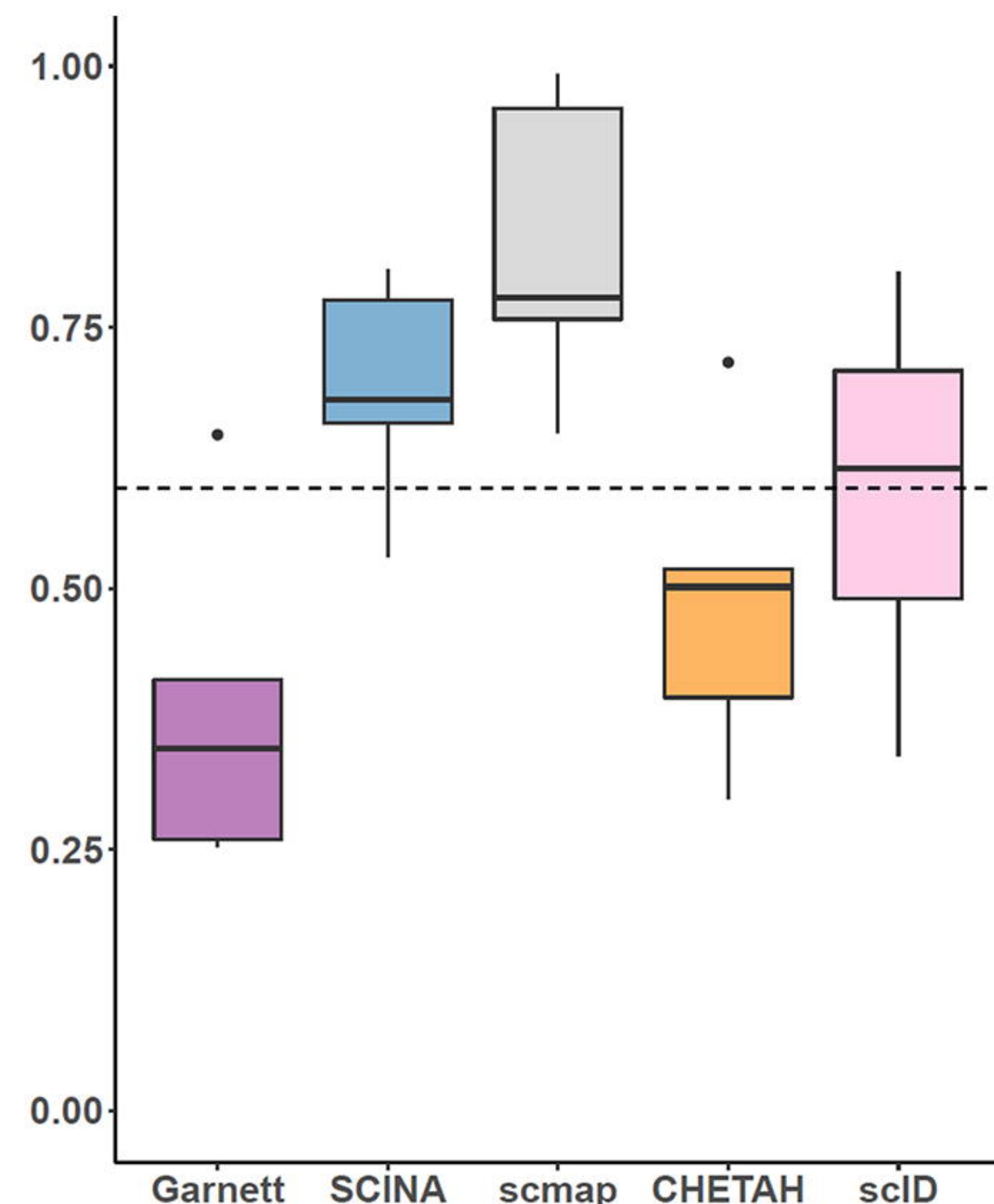
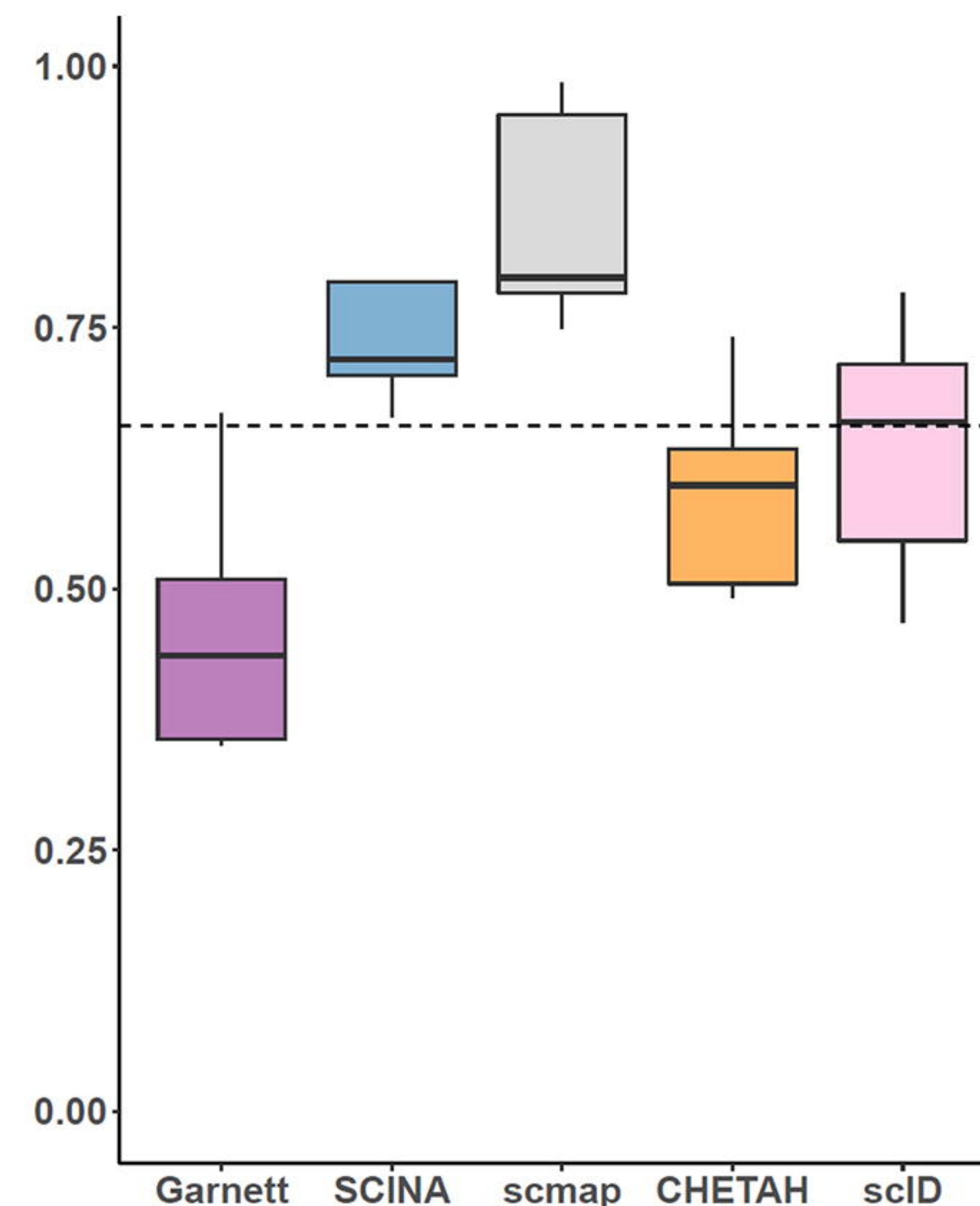
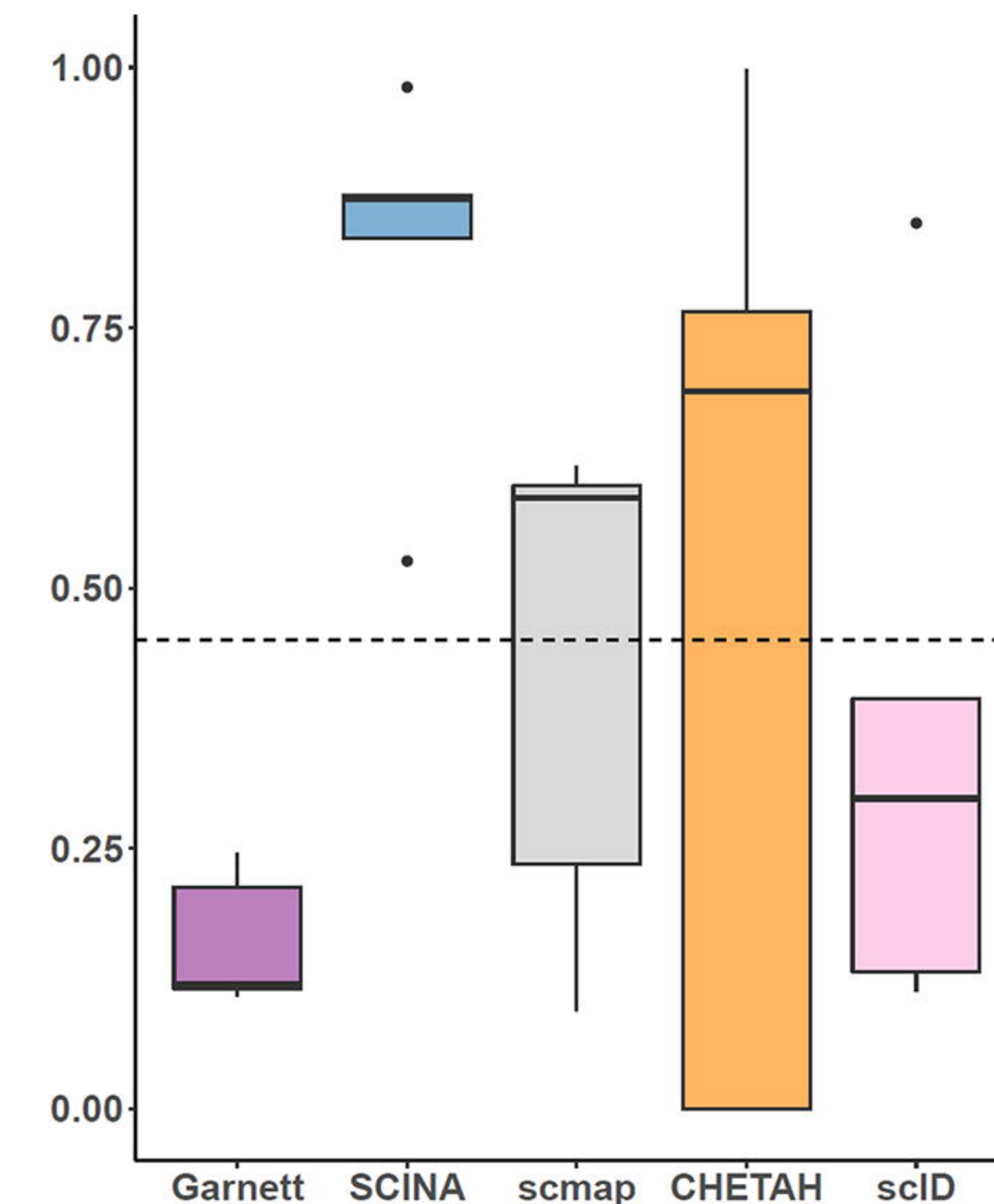


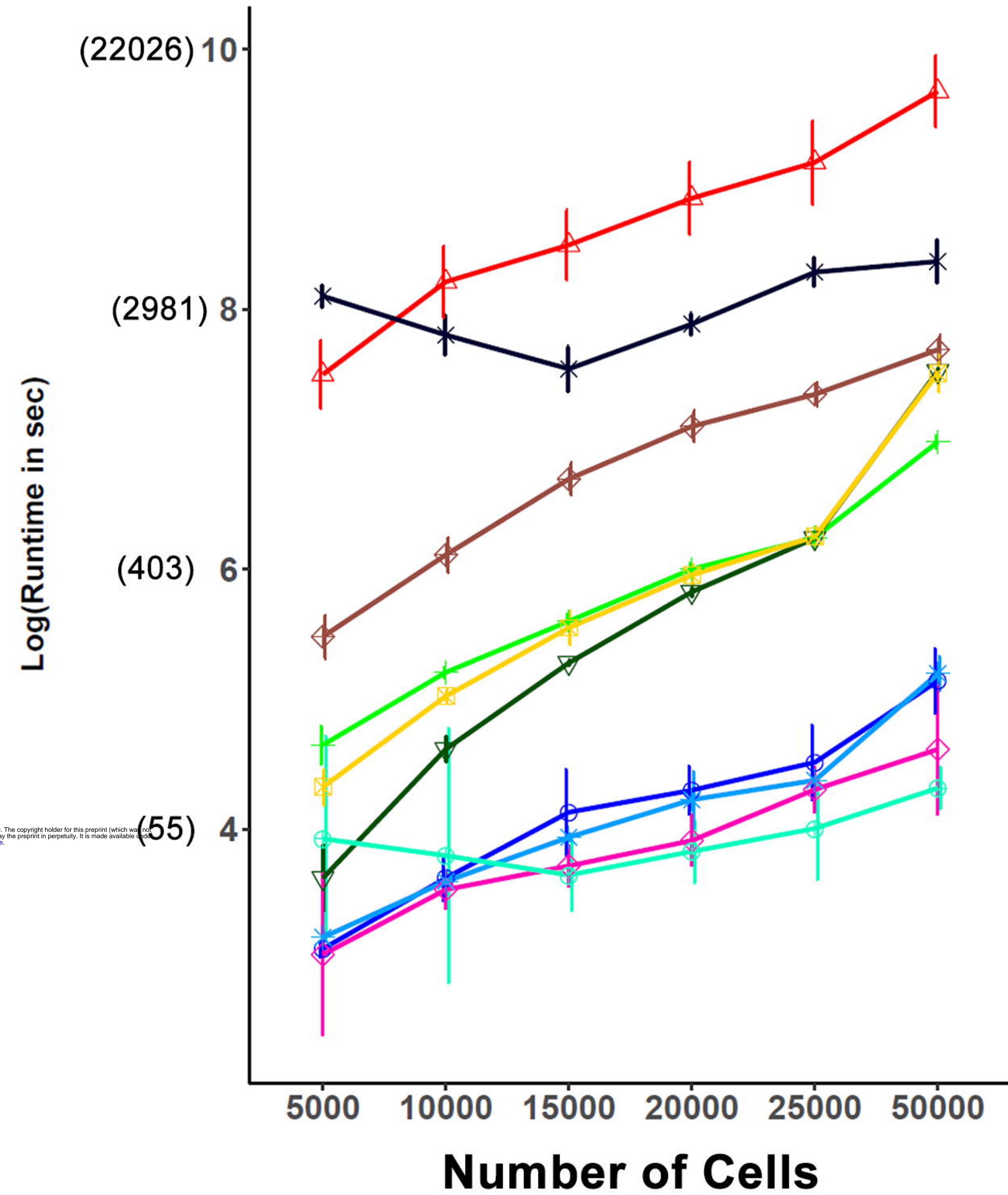
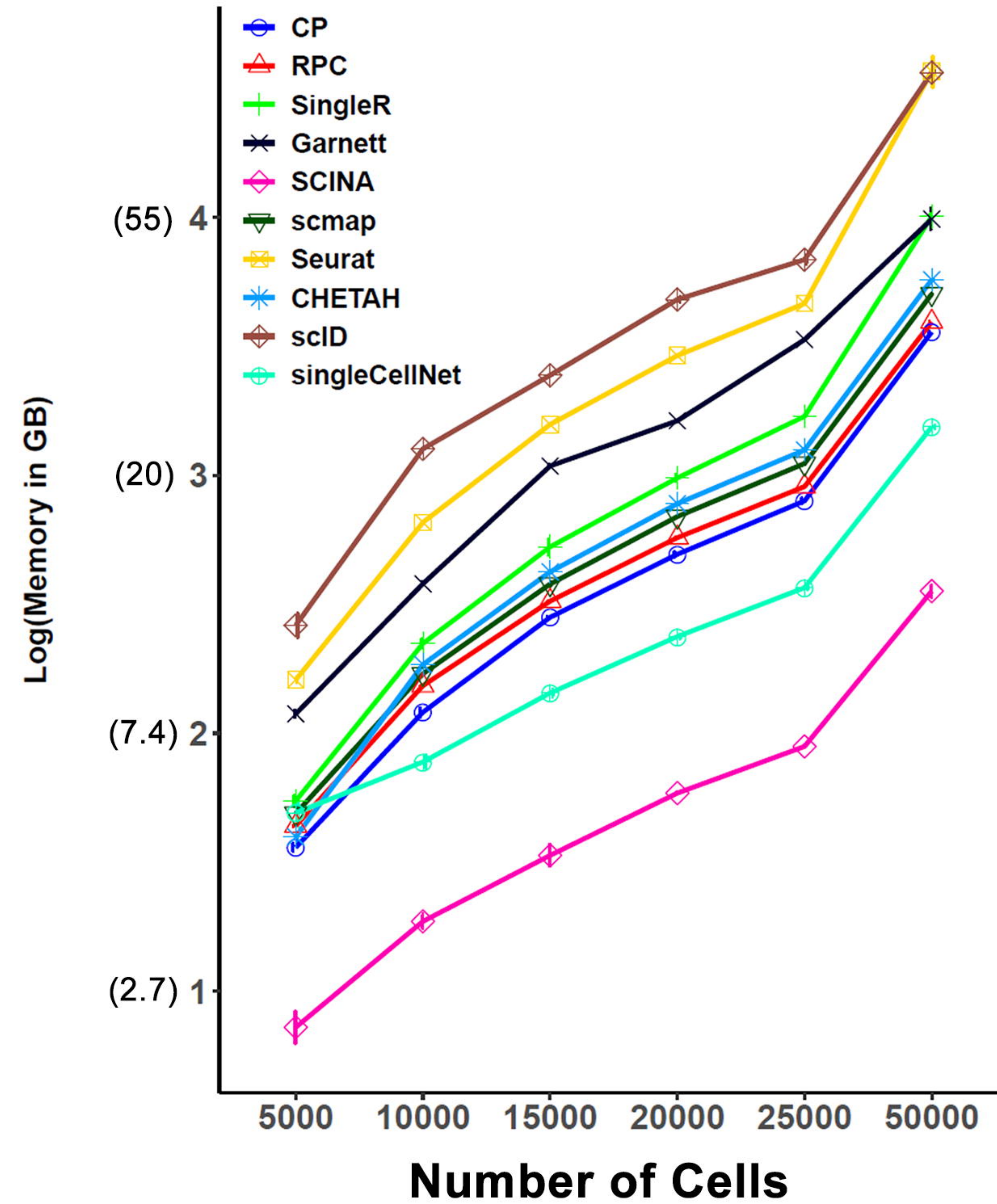
Score

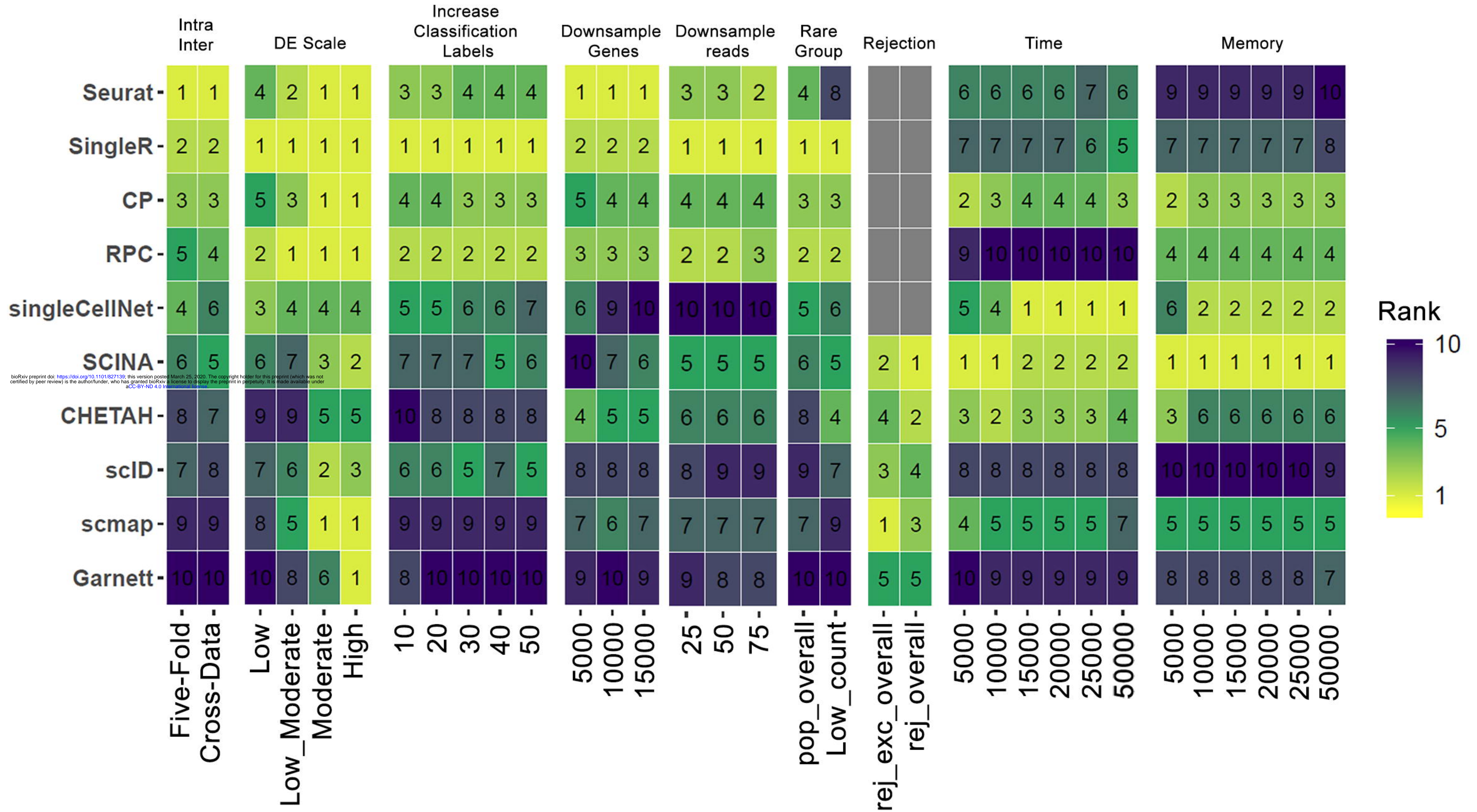


A**B****C****D****E****F****G**

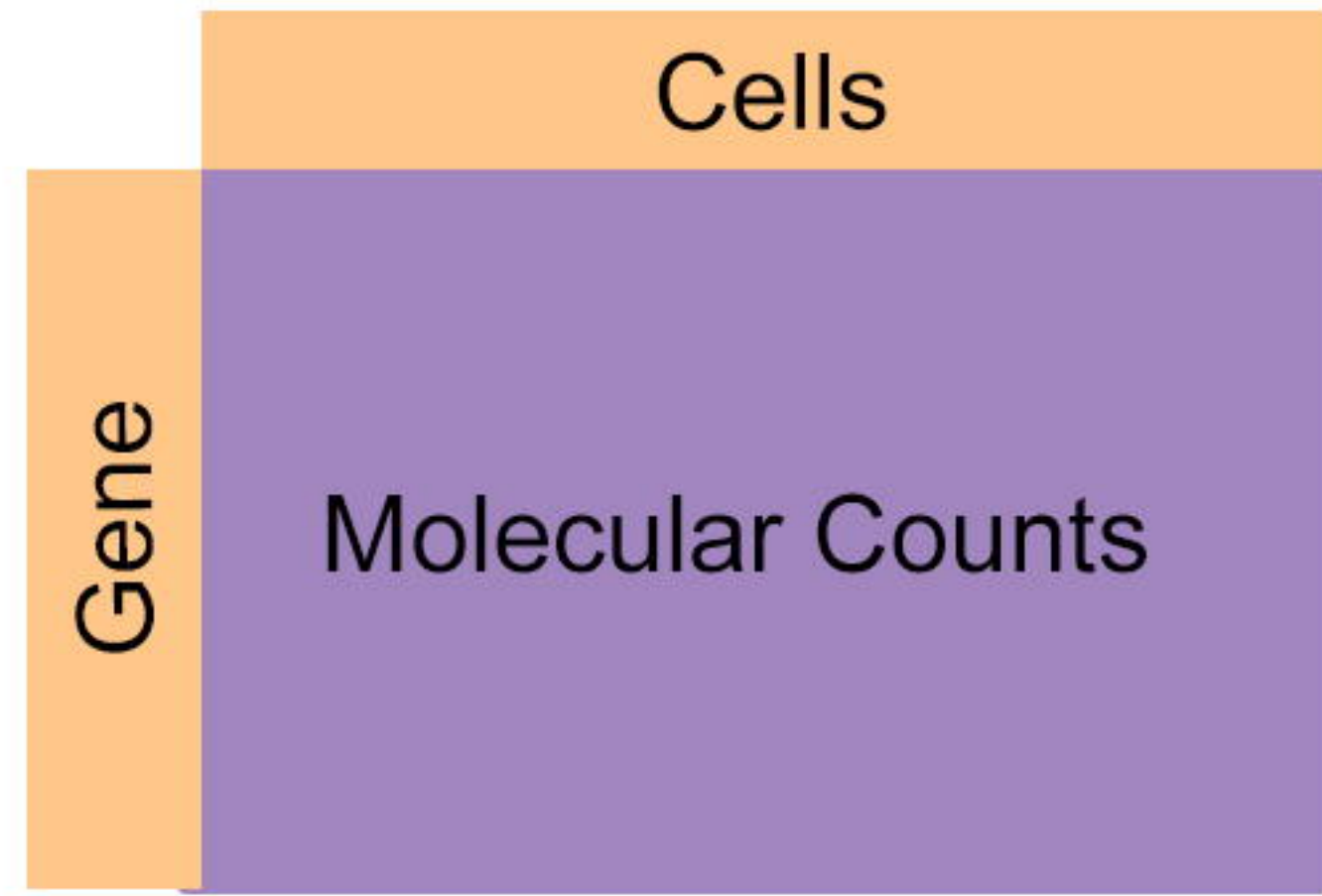


A**Group population scale in the simulation data****B****Cell-type specific accuracy across cell groups****C****Overall Accuracy****ARI****V-measure****D****Accuracy of the leave-out group to be unassigned****Methods with Rejection Option****Methods with Rejection Option**

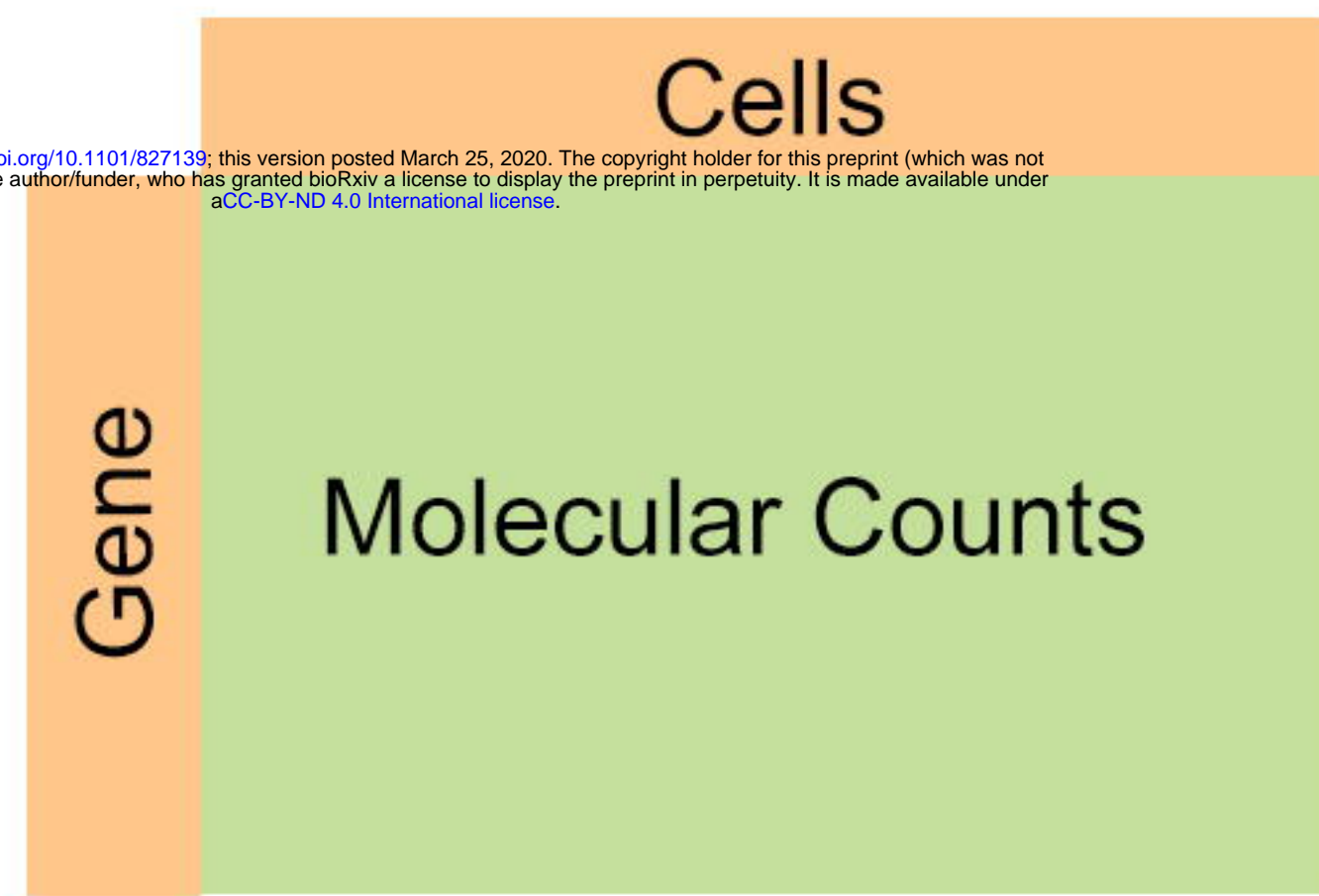
A**Runtime Comparison****B****Memory Utilization Comparison**



Data Preprocessing

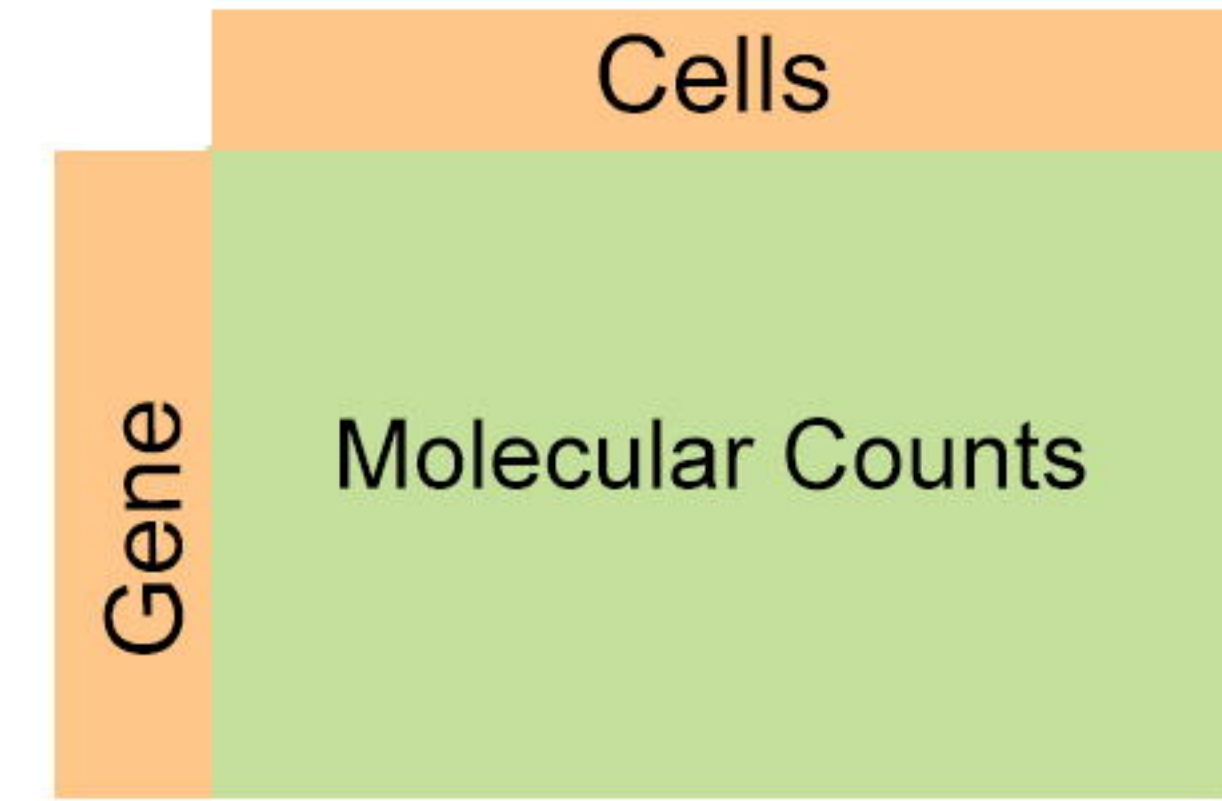


Cell/Gene Filter
Optional Normalization

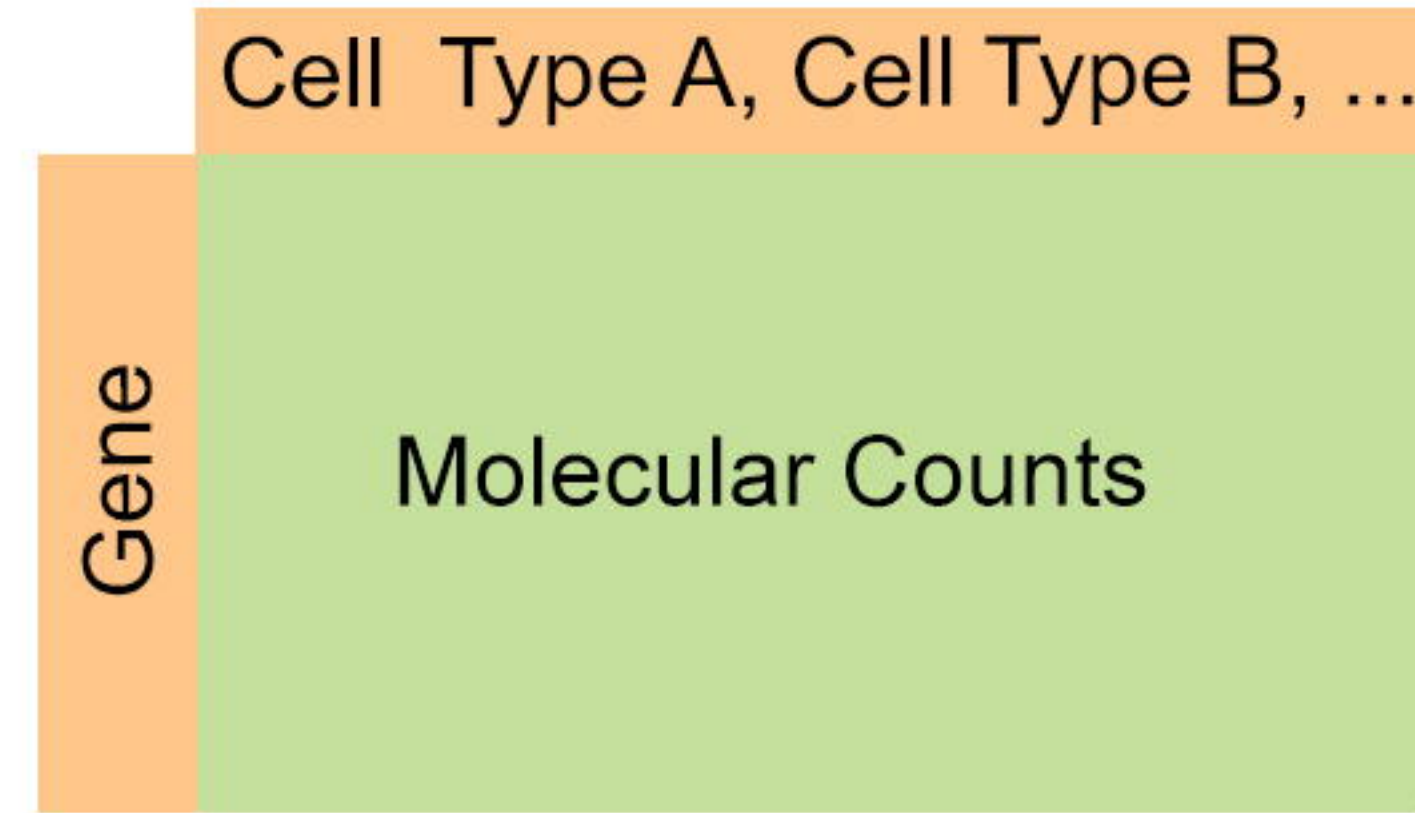


bioRxiv preprint doi: <https://doi.org/10.1101/827138>; this version posted March 25, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

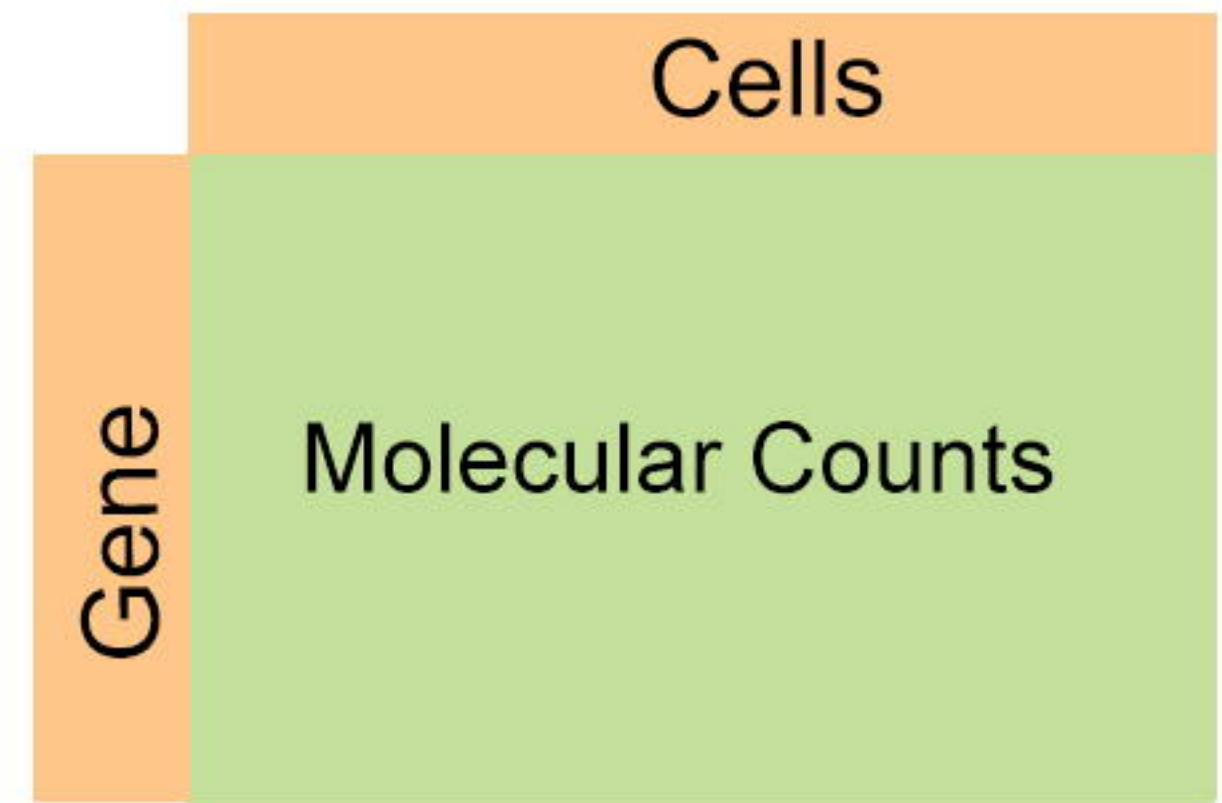
Reference



OR



Query



Different Algorithm

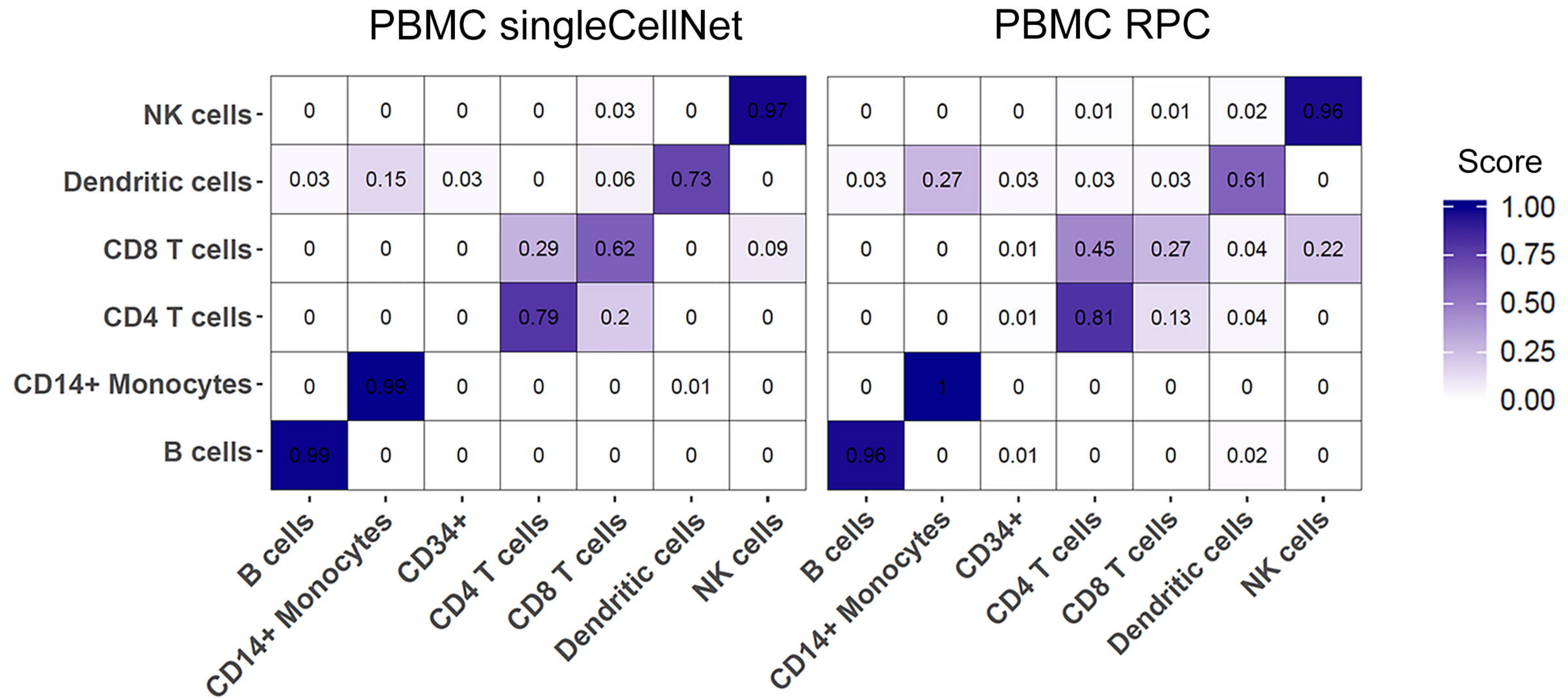
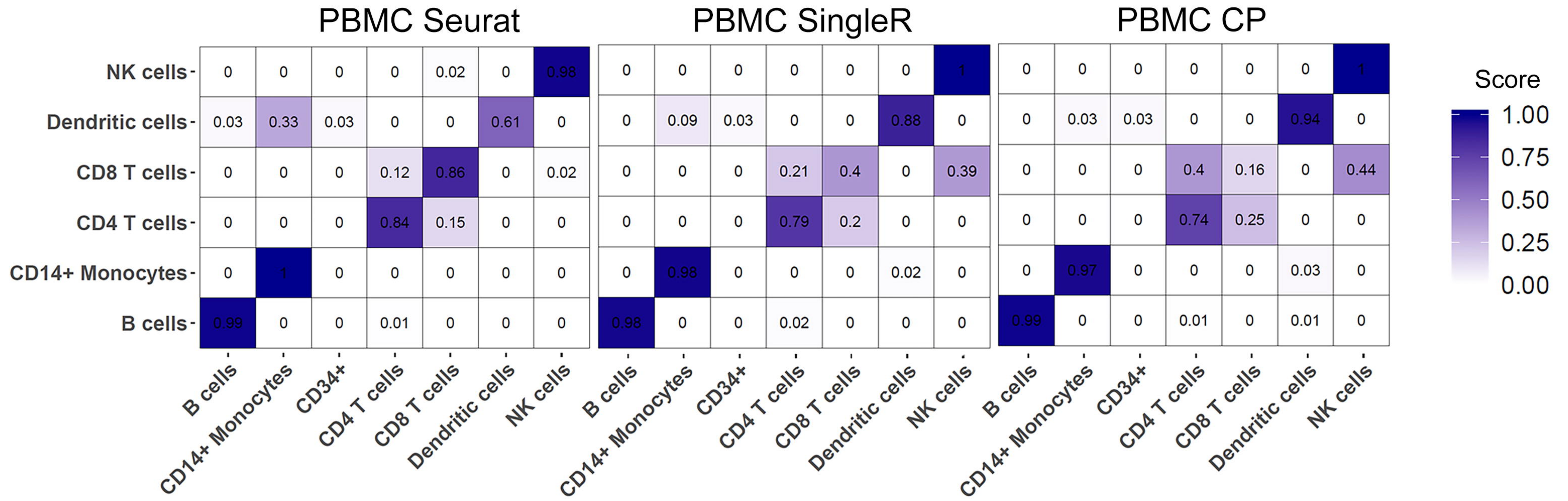


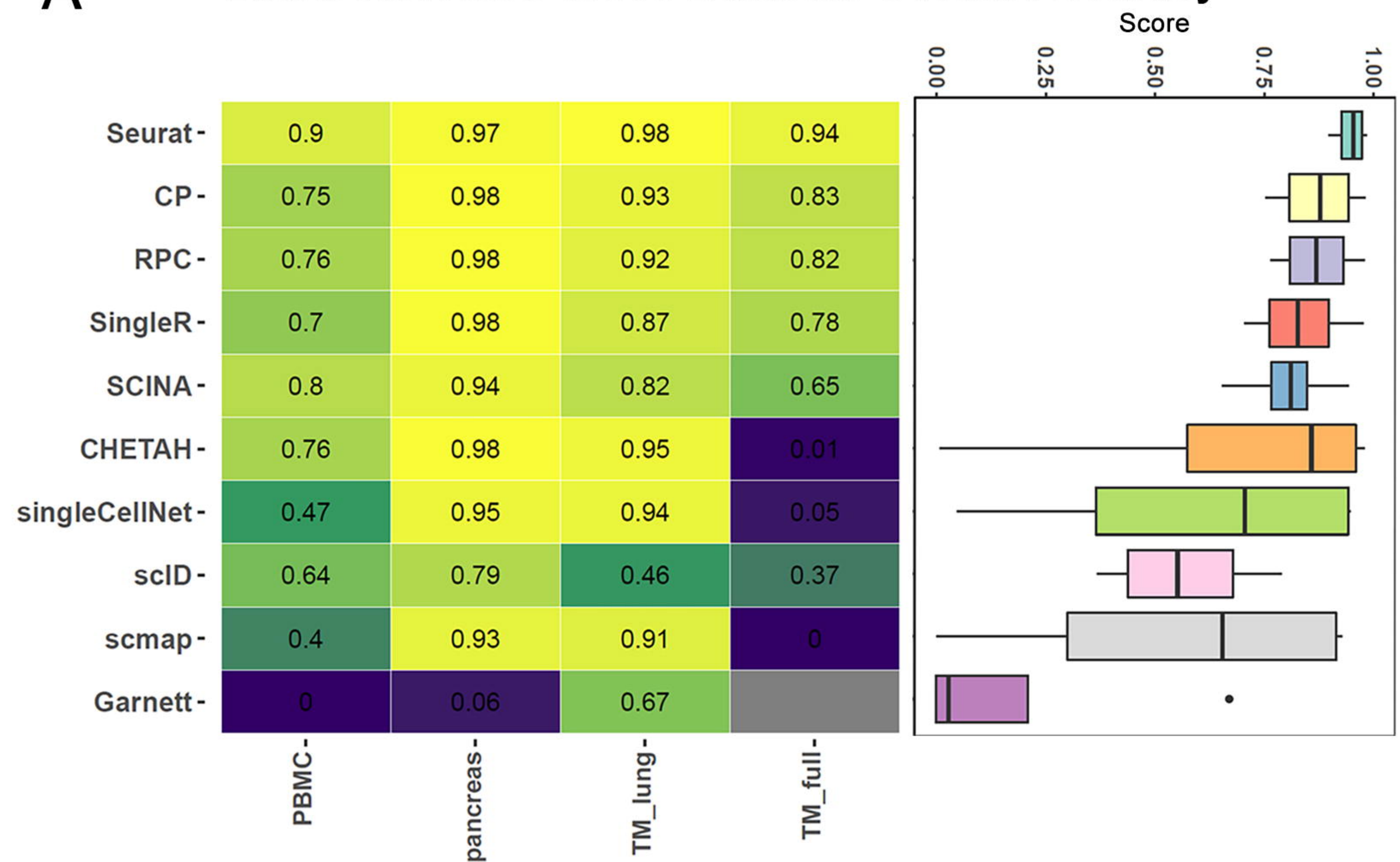
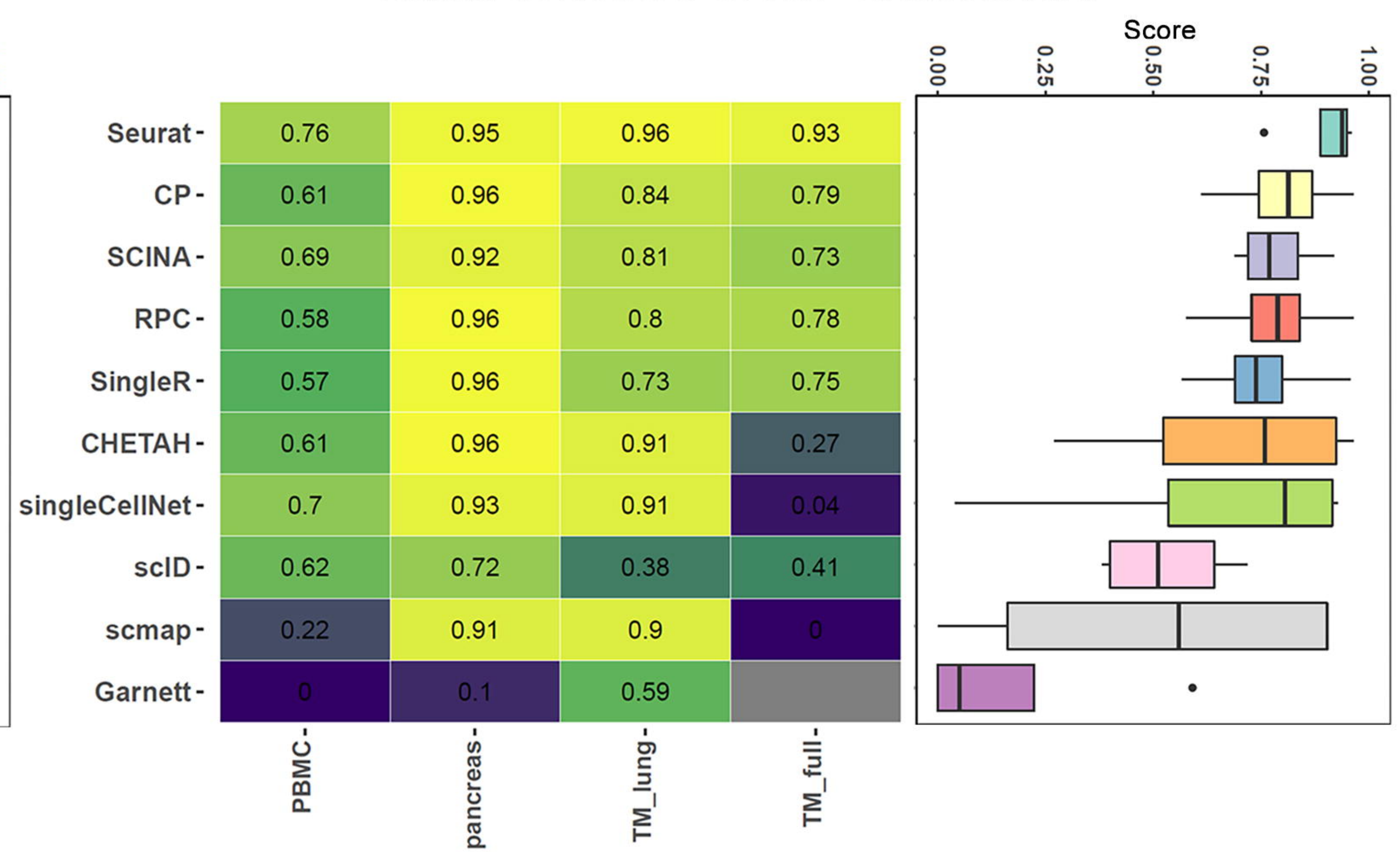
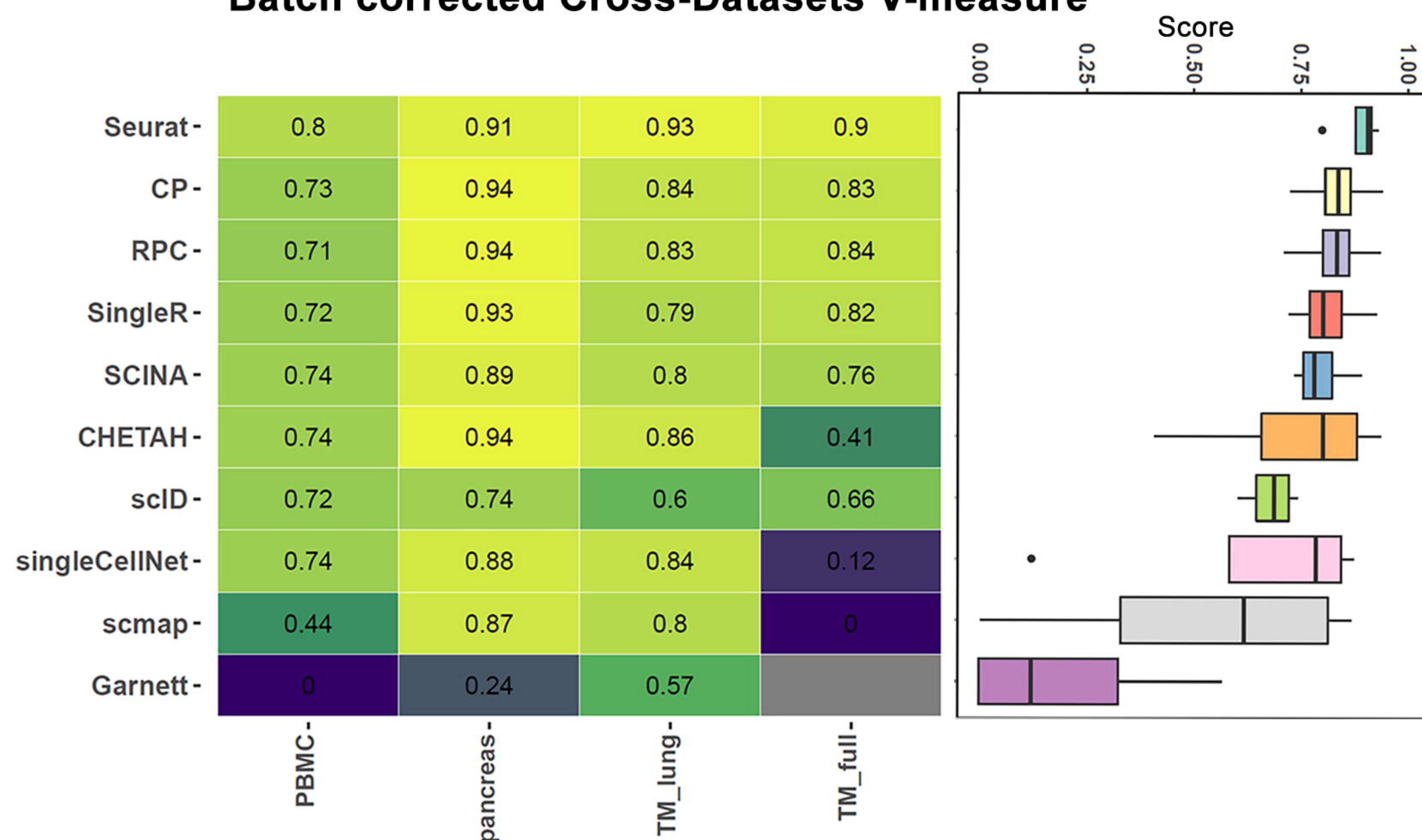
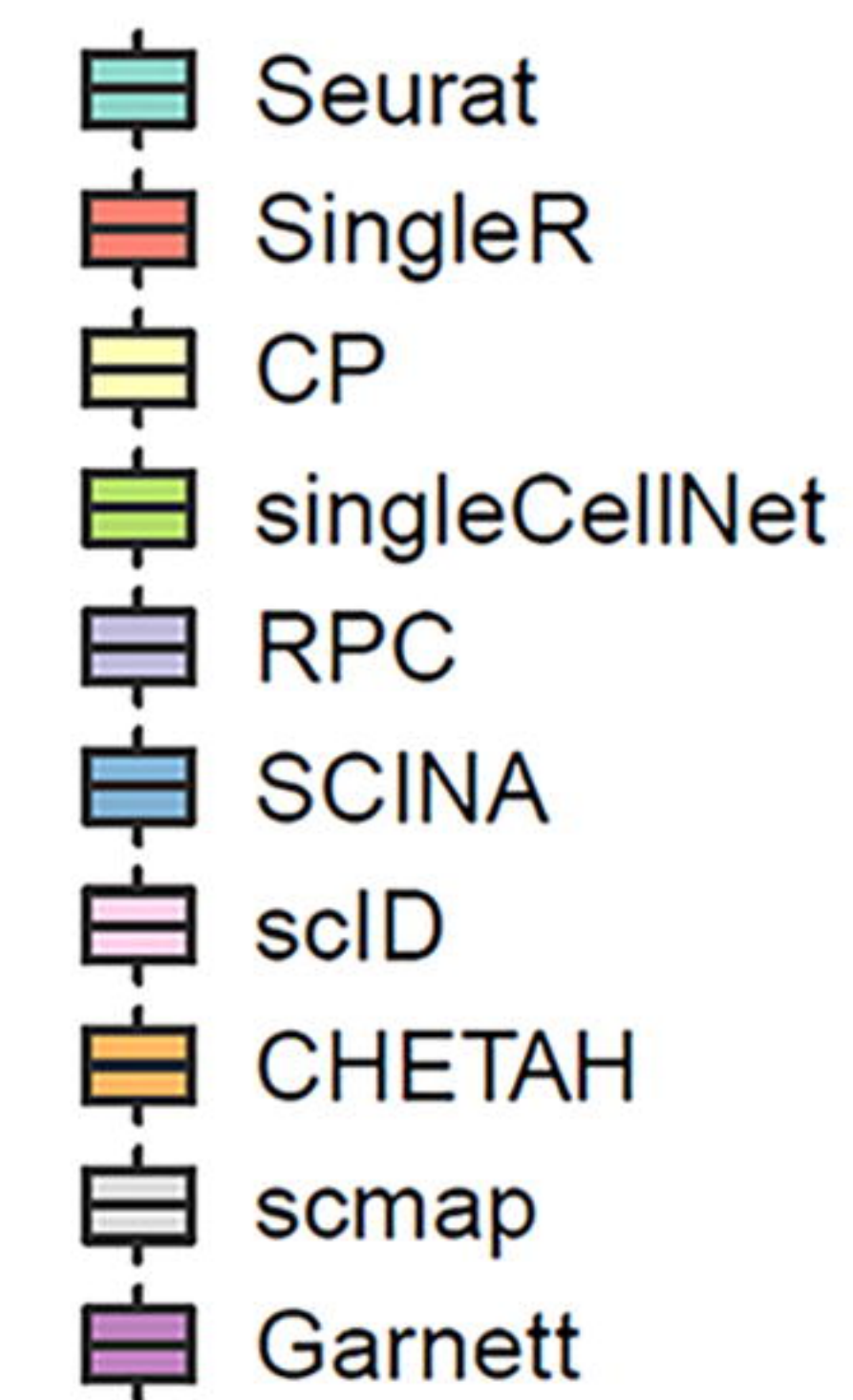
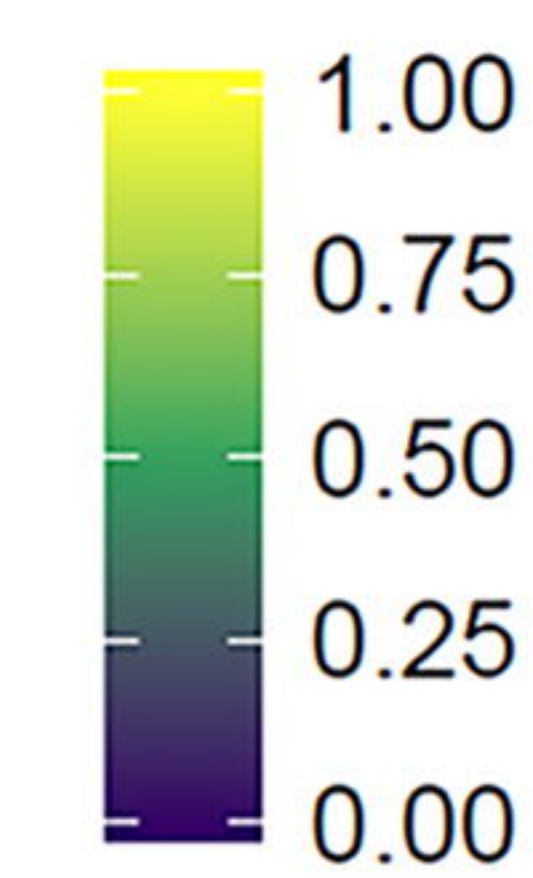
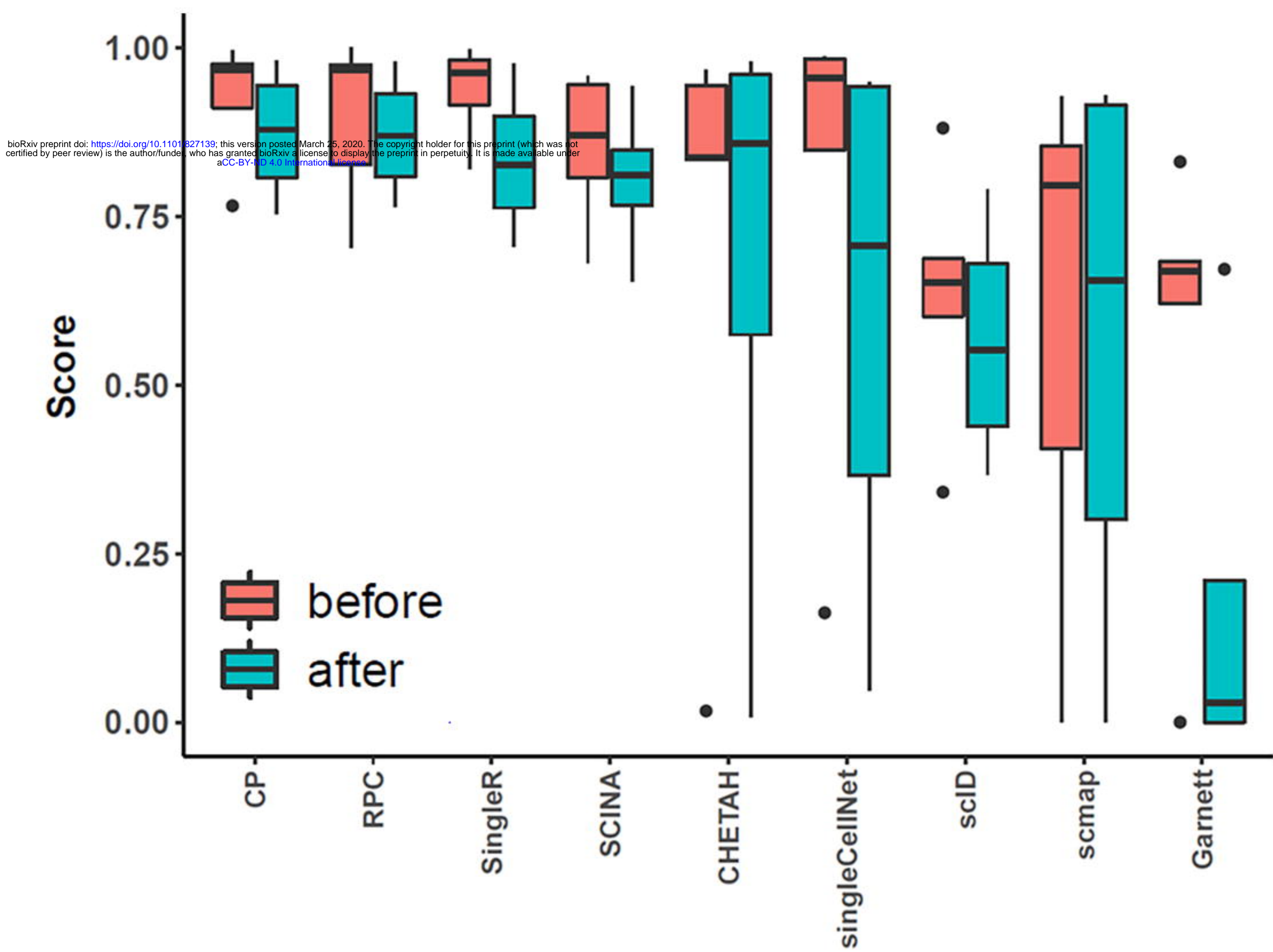
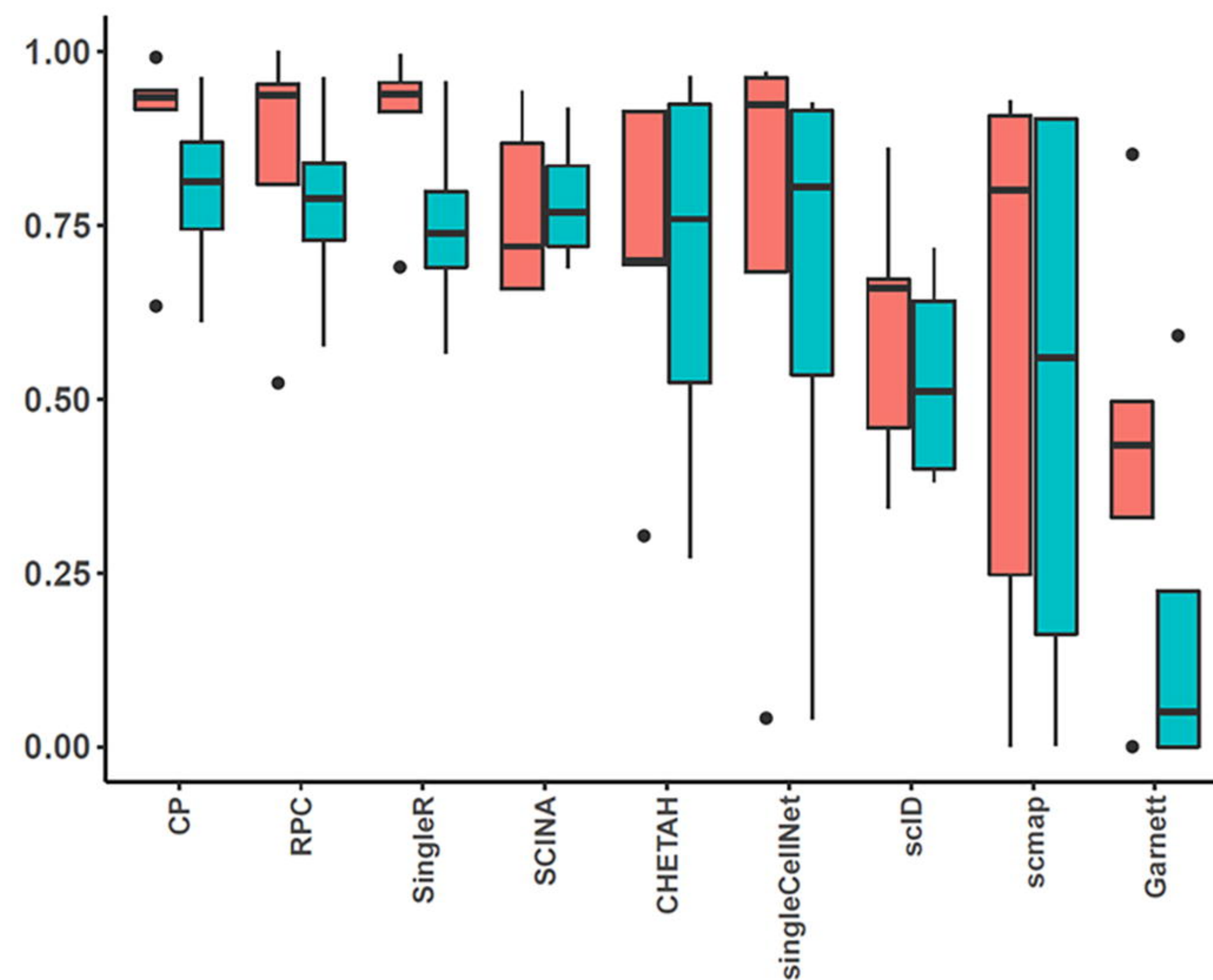
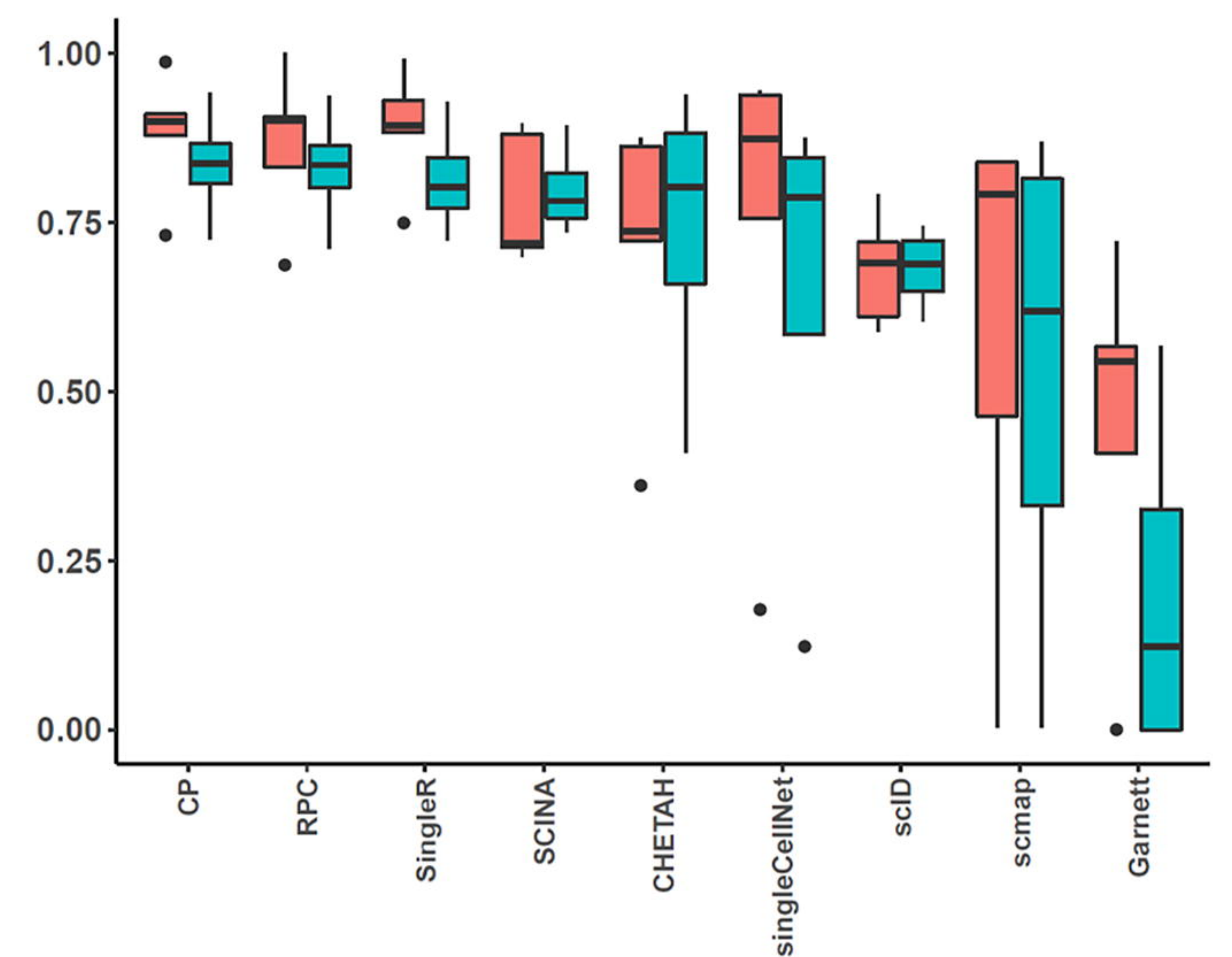
| Sample ID | Prediction |
|-----------|------------|
| Cell1 | Group5 |
| Cell10 | Group1 |
| Cell100 | Group2 |
| Cell1000 | Group2 |
| Cell1001 | Group2 |
| Cell1002 | Group2 |
| Cell1003 | Group2 |
| Cell1004 | Group5 |
| Cell1005 | Group1 |
| Cell1006 | Group3 |
| Cell1007 | Group2 |
| Cell1008 | Group3 |
| Cell1009 | Group3 |

Assessment

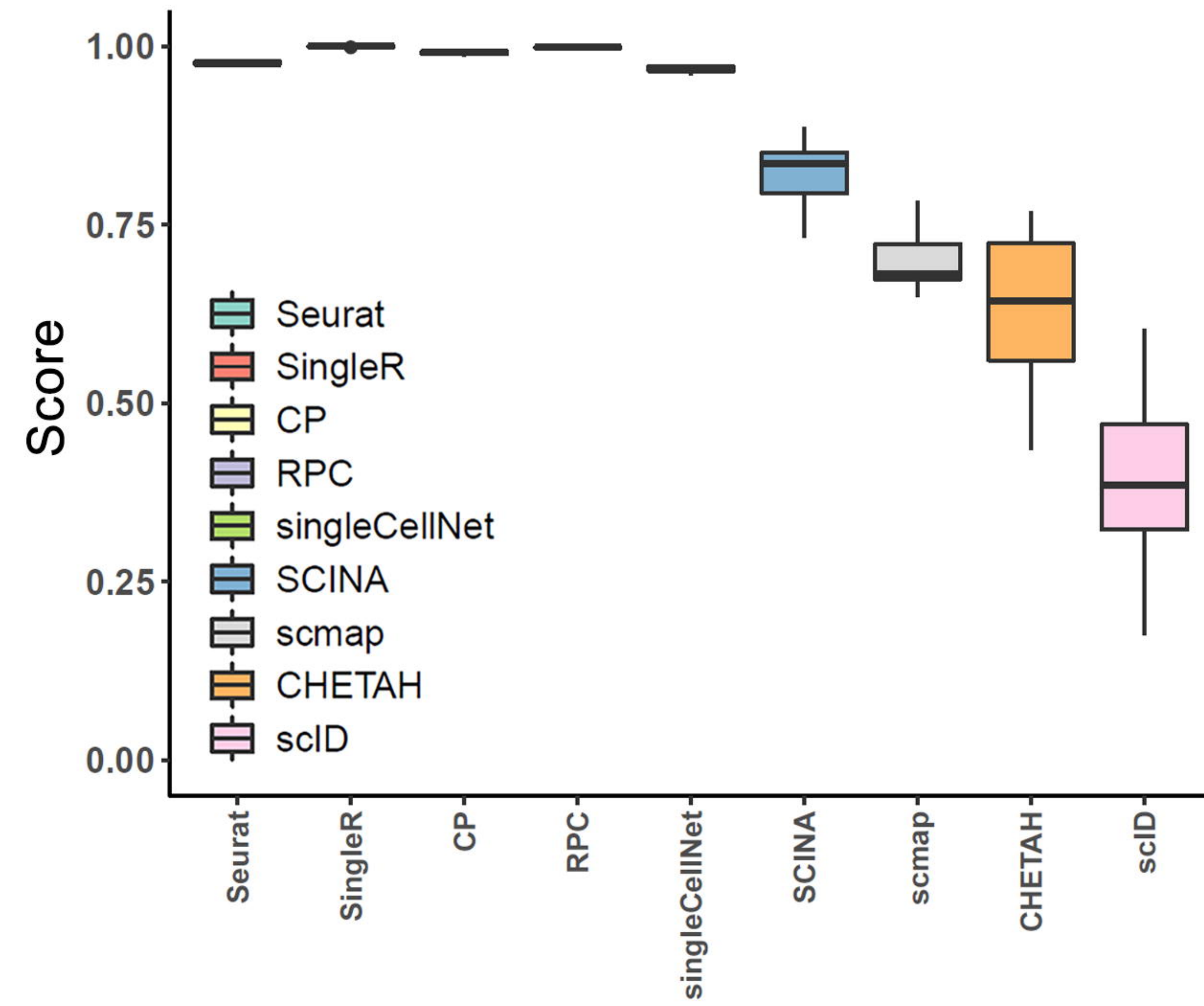


ARI
V-measure
Multi-class
Confusion Matrics

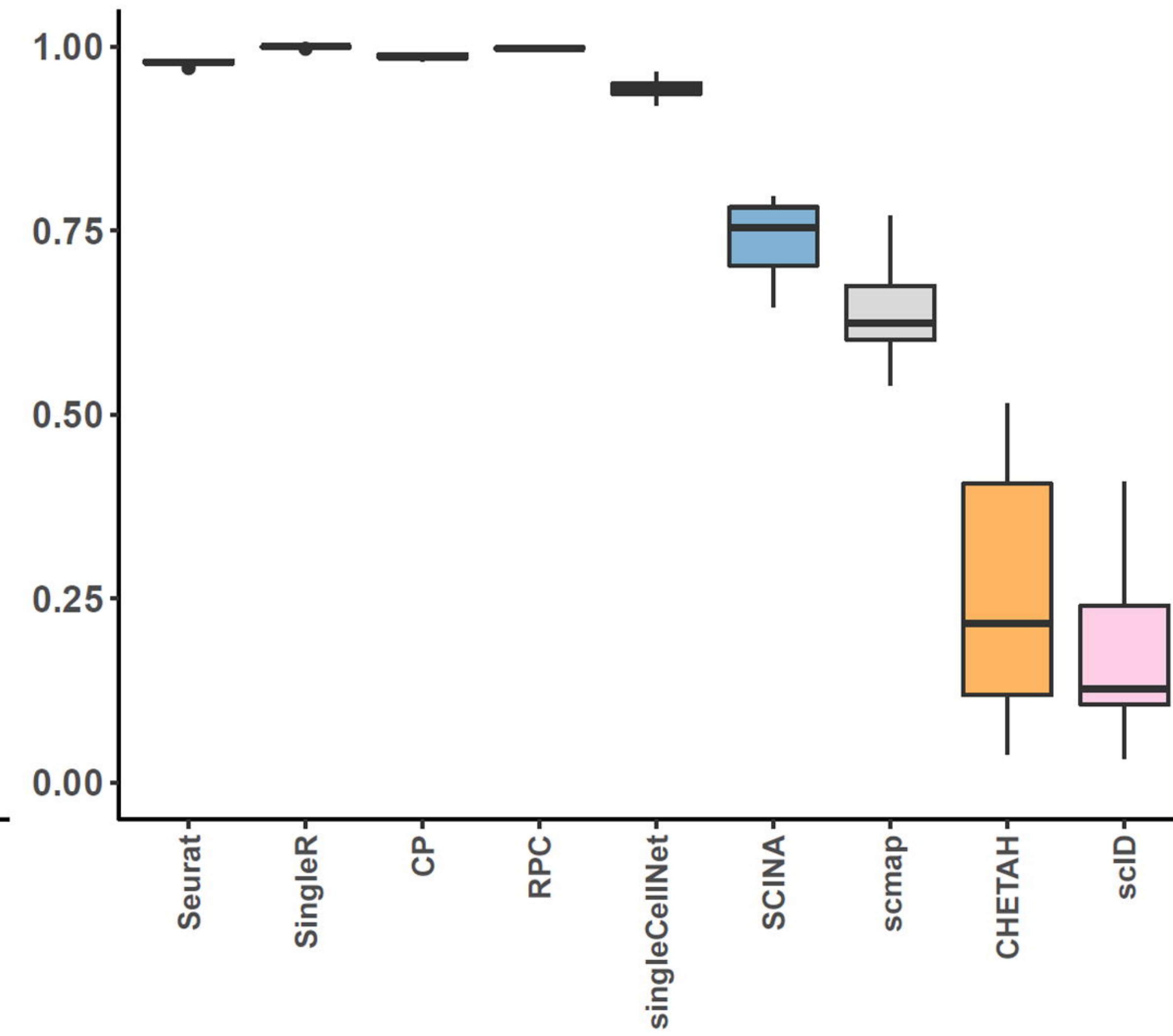


A **Batch corrected Cross-Datasets Overall Accuracy****Batch corrected Cross-Datasets ARI****Batch corrected Cross-Datasets V-measure****Methods****Score****B** **Overall Accuracy Before and After Batch Correction****ARI Before and After Batch Correction****V measure Before and After Batch Correction**

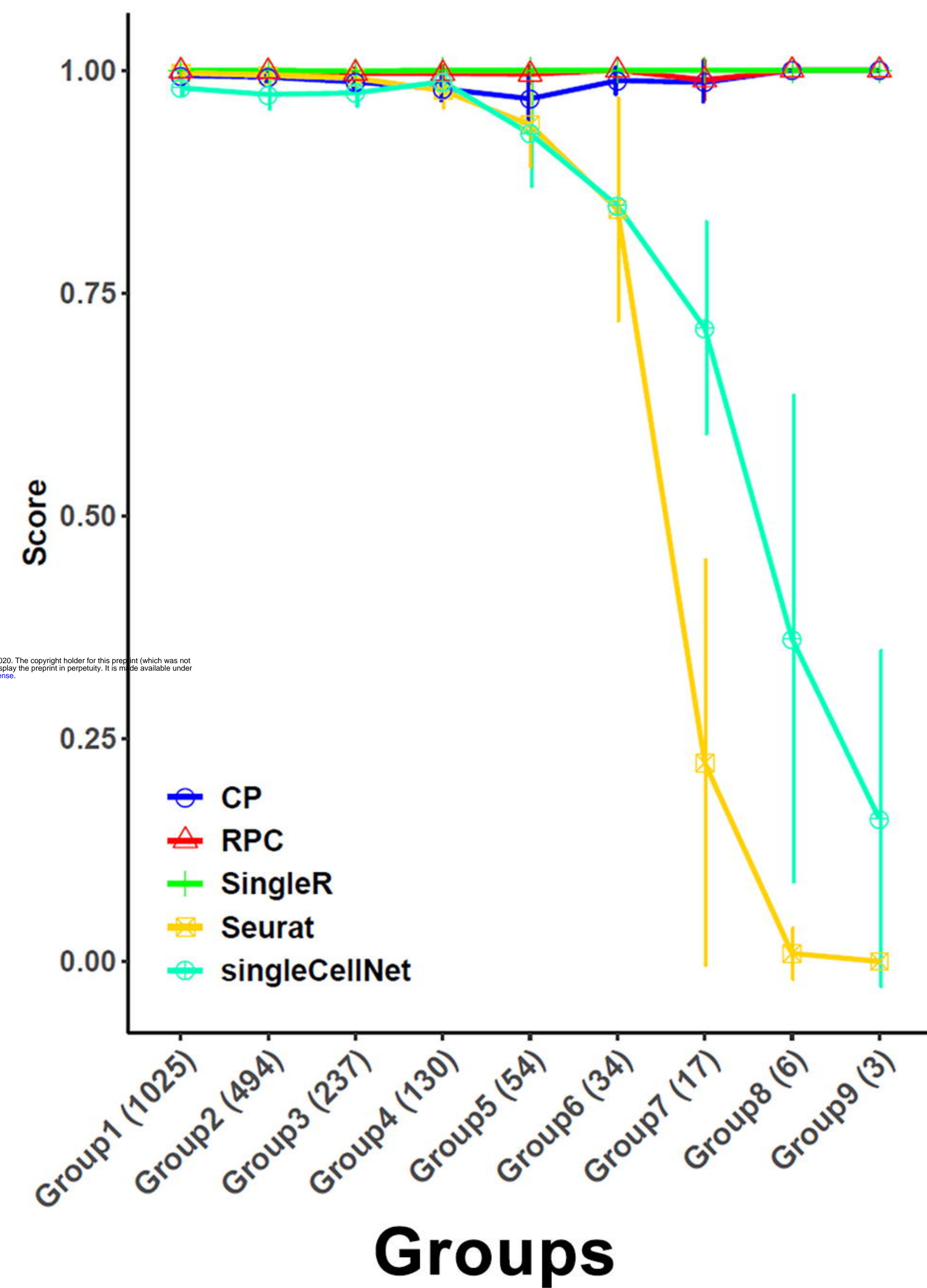
A Overall Accuracy



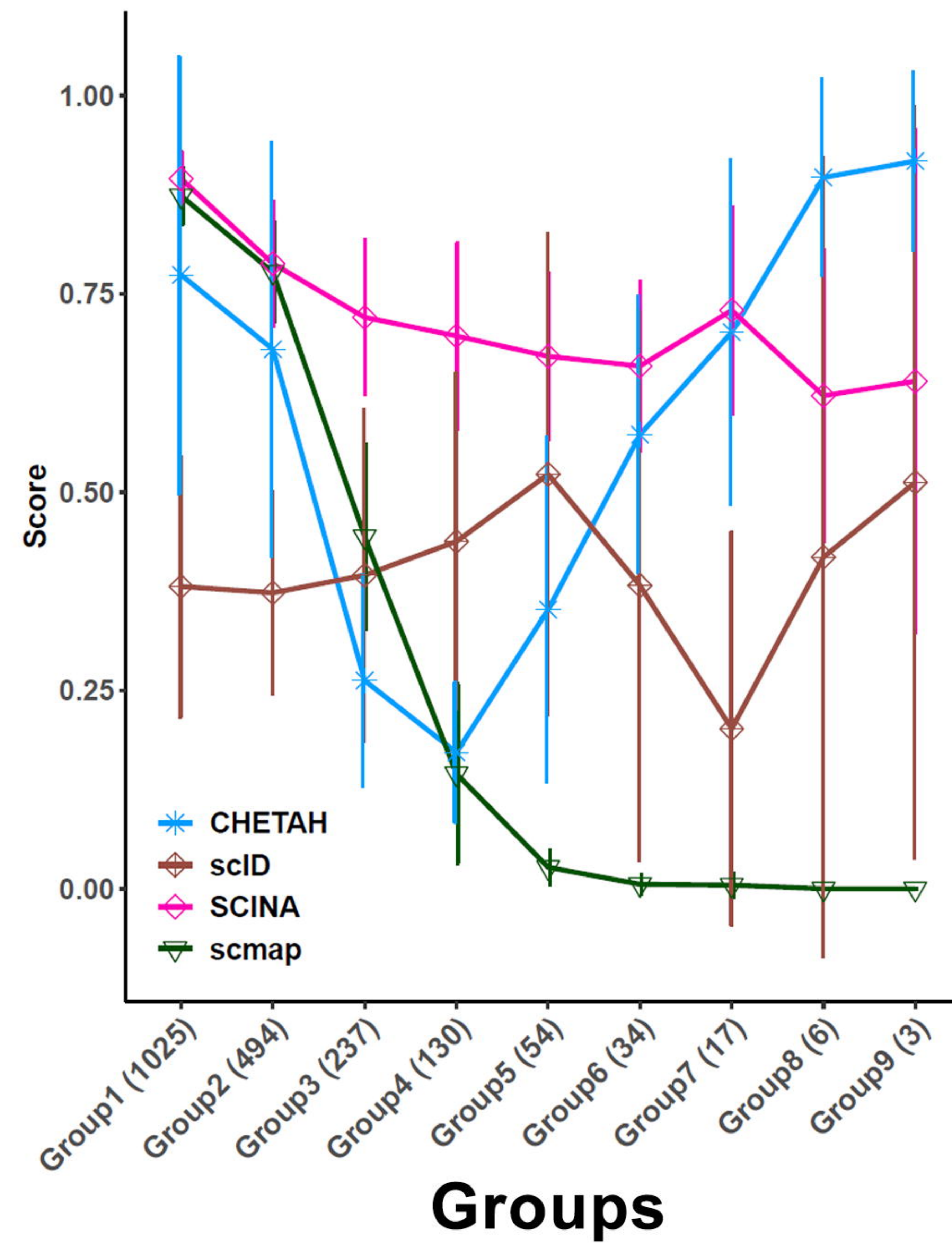
ARI

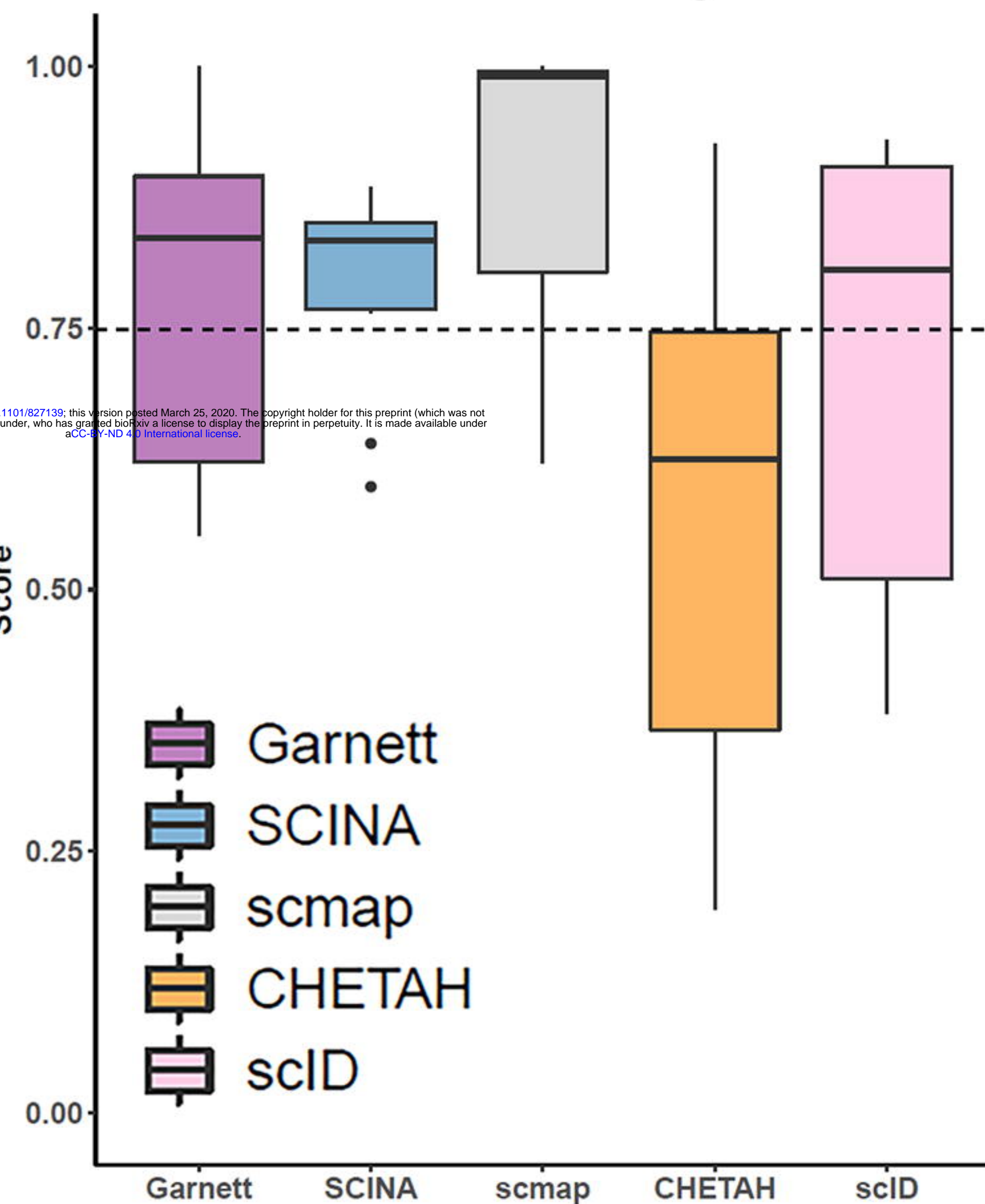
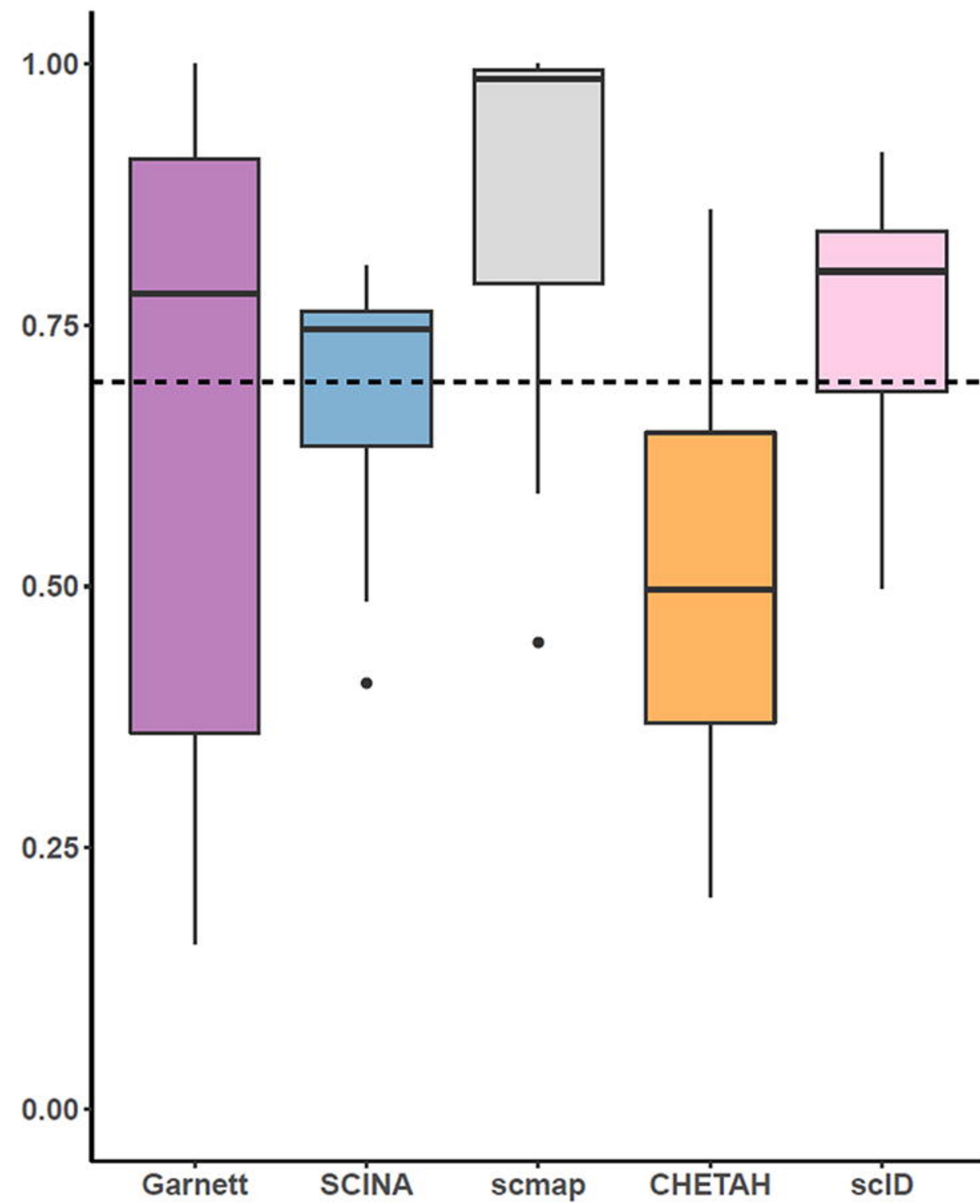
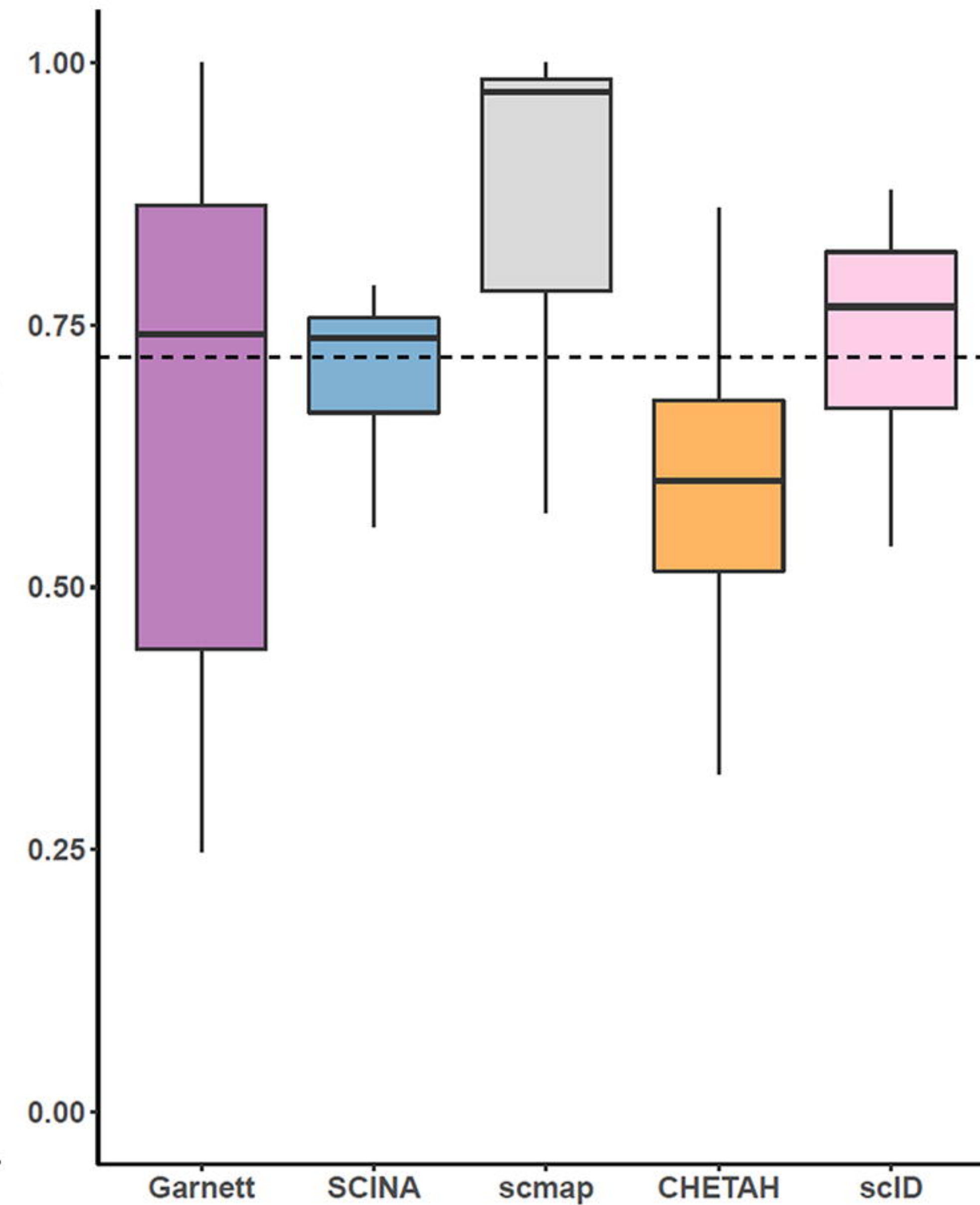
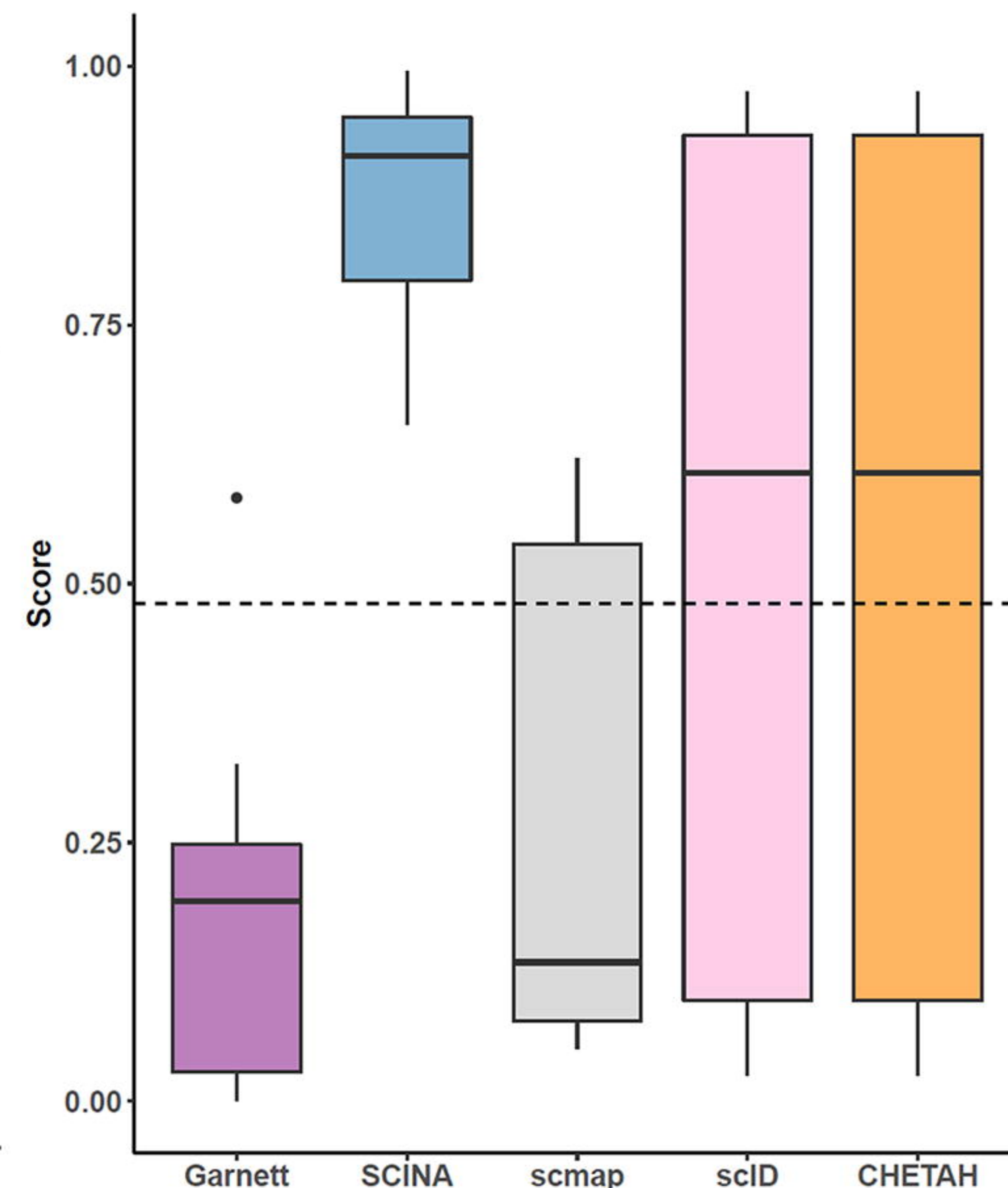


B Top 5 performing methods cell-type specific accuracy

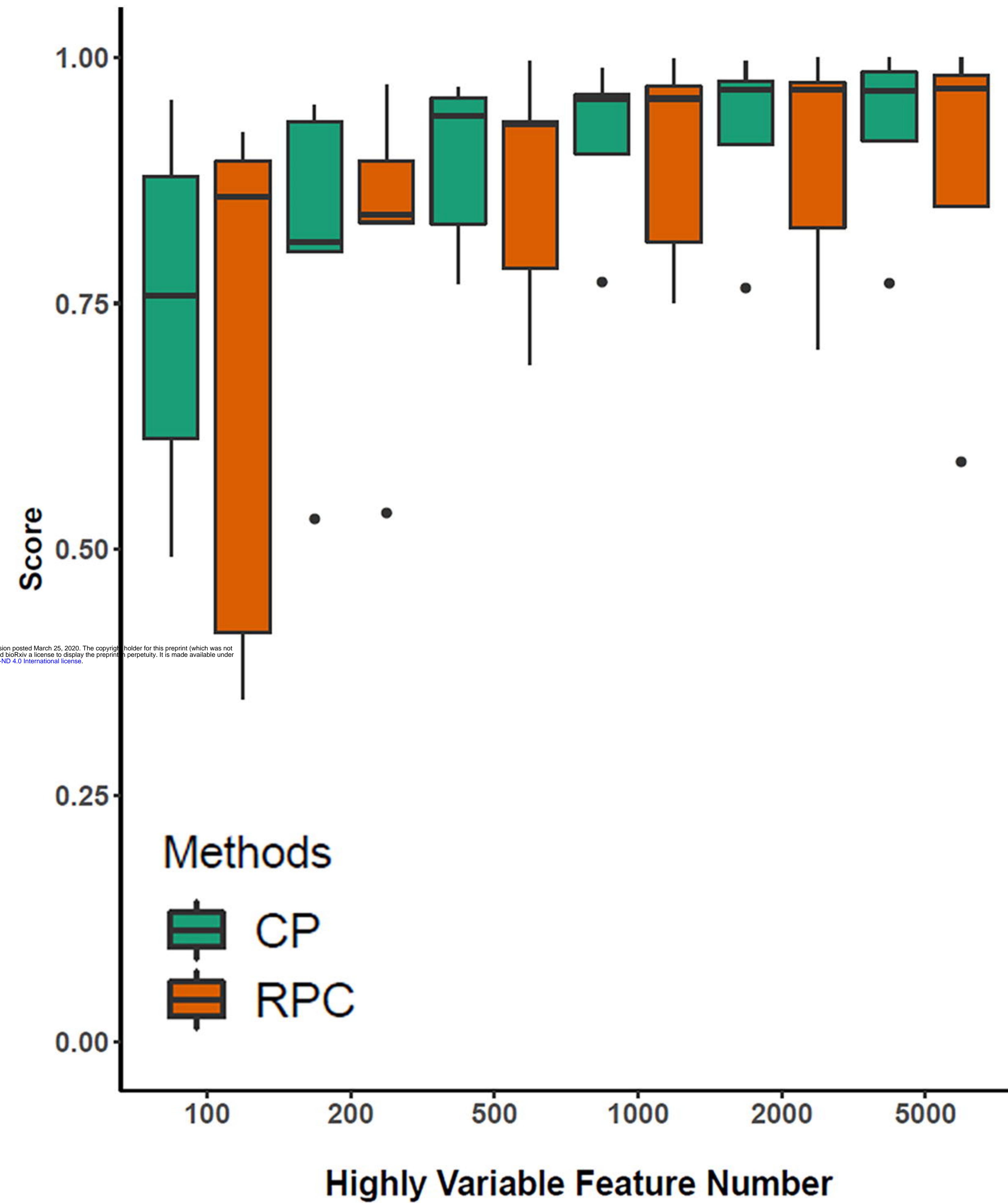


Remaining methods cell-type specific accuracy



A**Overall Accuracy****ARI****V-measure****B****Accuracy of leave-out groups to be unassigned****Methods with Rejection Option**

A Overall Accuracy using different number of HVG



B Matrix Condition Number using different number of HVG

