1    **Anti-correlated Feature Selection Prevents False Discovery of Subpopulations in**

2    **scRNAseq**

3    Scott R Tyler,[1,2,*] Ernesto Guccione[2,3,4], and Eric E Schadt[1,*]

4    [1]Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New

5    York, NY, USA.

6    [2] Department of Oncological Sciences, Tisch Cancer Institute, Icahn School of Medicine at

7    Mount Sinai, New York, NY, USA

8    [3] Center for Therapeutics Discovery, Department of Oncological Sciences and Pharmacological

9    Sciences, Tisch Cancer Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA

10    [4] Bioinformatics for Next Generation Sequencing (BiNGS) Shared Resource Facility, Icahn

11    School of Medicine at Mount Sinai, New York, NY, USA

12    *Corresponding authors

13

14

**Abstract**

16    **While sub-clustering cell-populations has become popular in single cell-omics,**

17    **negative controls for this process are lacking. Popular feature-selection/clustering**

18    **algorithms fail the null-dataset problem, allowing erroneous subdivisions of homogenous**

19    **clusters until nearly each cell is called its own cluster. Using 45,348 scRNAseq analyses**

20    **of real and synthetic datasets, we found that anti-correlated gene selection reduces or**

21    **eliminates erroneous subdivisions, increases marker-gene selection efficacy, and**

22    **efficiently scales to 245k cells without the need for high-performance computing.**

23

**Results**

25    A frequent first task in performing cell-type identification from scRNAseq is feature

26    selection to identify genes that are cell-type specific markers based on various statistical

27    properties. Current approaches include measures of the relationship between a gene's mean and

28    variance (i.e., overdispersion)[1-3] and a gene's mean and dropout rate[4]. An open problem however

29    is how algorithms handle the "null-dataset;" that is, when there is only a single cell-type present.

30    Given the popularity of sub-clustering (i.e., iteratively subdividing the initially identified

31    clusters)[5-8], it is important to know that these groups are not being erroneously subdivided, thus

32    producing false subtypes[9]. While novel sub-populations of interest should always be validated via

33    bench-biology methods, an algorithmic assurance that one is not being misled can save money

34    and years of effort attempting to validate erroneously discovered "novel sub-populations." Given

35    the imperfections in clustering algorithms[10], sub-clustering itself can be a valid practice,  because

36    a single round of clustering may be insufficient to fully divide a dataset into its constituent groups.

37    However, we must have confidence that such algorithms will correctly identify single populations,

38    preventing the false discovery of nonexistent sub-populations. In the case of a single cell

39    population, either 1) a feature selection algorithm would accurately report that there are no genes

40    that define sub-populations, or 2) the clustering algorithm would determine that only a single

41    cluster is present.

42    We sought to devise an algorithm to identify cell-type marker-genes that would not only

43    identify subpopulations of cell-types with high accuracy, but also solve the null-dataset problem.

44    We thus began from first principles, asking the question: "what is a cell-type?". Traditional

45    molecular biology has defined cell-types based on distinct cellular functions that are concordant

46    with expression of distinct sets of genes: "marker-genes" (**Fig. 1a**), that often include hierarchical

47    mutually exclusive gene expression. For example, in the pancreas the gene *NEUROD1* is a pan-

48    endocrine marker, expressed in many different cell-types but should be mutually exclusively

49    expressed from exocrine marker-genes[11]. If we accept this definition of cell-type and -lineage

50    specific genes, we can algorithmically discover marker-genes from scRNAseq, as these genes

51    will show a statistical excess of negative correlations with other genes (**Fig. 1b**). Given this

52    premise, if only a single cell-identity is present in a dataset, we would expect an absence of an

53    anti-correlation pattern since the cells of other cell-identities would not be present (**Fig. 1c**).

54    Indeed, looking at known marker-genes from different cell types in the pancreas (i.e. *AMY2A*

55    expressed in acinar cells and *SST* expressed in delta cells), we see the expected anti-correlation

56    pattern between *AMY2A* and *SST* (**Fig. 1d**), which disappears when examining subsets

57    comprised of only a single cell type (**Fig. 1e**). Notably, the anti-correlation pattern holds for

58    lineage-markers as well as cell-type markers (**Fig. 1f**).

59    Using these observations, we constructed an algorithm that identifies genes with an

60    excess of negative correlations relative to what would be expected if the gene were un-patterned,

61    as empirically measured with a bootstrap shuffled null background (**Fig. 1g,h**). We then select

62    genes that have an excess of negative correlations, controlling for false positives by setting an

63    appropriate false discovery rate (FDR) (**Fig. 1i**). Overall, this procedure selects the genes that

64    have significantly more negative correlations with other genes than would be expected by chance

65    (See **Methods** for details). While others have performed small-scale experiments using positive

66    correlations for feature selection, it was deemed infeasible due to computational run-time[12]; here

67    we create an open-source, efficient implementation in python to overcome this barrier, but focus

68    attention on negative correlation patterns as opposed to positive.

69          Given our reasoning that the anti-correlation pattern should go away when examining data

70    representing only a single cell-type (**Fig. 1c**), with preliminary support for our rationale in a single

71    dataset (**Fig. 1e**), we hypothesized that anti-correlation-based feature selection would be

72    sufficient to solve the null-dataset problem, while status quo algorithms may not adequately solve

73    for this problem. With the null-dataset, no "cell-type or cell-state specific genes" should be

74    identified as this is a single population of cells. We tested this hypothesis by performing feature

75    selection and affinity propagation (AP)-based clustering on two datasets composed of scRNAseq

76    from homeostatic cell line culture from NIH3T3 (**Fig. 1j**) or HEK293T cells (**Fig. 1k**), which we

77    anticipate would capture the biologically relevant variation in only a single clustering round, and

78    any attempt to *further* subdivide beyond that should be algorithmically blocked. Indeed, the anti-

79    correlation algorithm allowed for only a single round of clustering, while the other algorithms tested

80    allowed for further subdivisions (**Fig. 1j,k**).

81          While this preliminary evidence suggests that anti-correlation-based feature selection

82    solves the issue of false positives from sub-clustering homogenous populations, real-world

83    datasets do not harbor a "ground-truth." We therefore simulated a single cluster using Splatter

84    which produces negative binomially distributed gene expression matrices[13]. We performed

85    feature selection using the noted algorithms[1-4] and passed these features to four different

86    clustering algorithms including Affinity Propagation, K-means+Elbow-rule, K-means+Silhouette,

87    and locally weighted Louvain modularity (See **methods** for algorithm details). In all cases, the

88    anti-correlation-based method for feature selection detected no valid features within a single

89    population of cells, thus addressing the null-dataset problem, while all other feature selection and

90    clustering algorithm combinations failed the null-dataset problem, selecting noisy features that

91    resulted in at least several clusters (**Fig. 1l**). Note that most feature selection algorithms frequently

92    require the user to manually set the number of "discoveries" or selected features, which is likely

93    a key contributor to this failure of the null-dataset problem when using standard feature selection

94    approaches.

95         Without an algorithmic check to prevent erroneous sub-clustering, one could recursively

96    divide a dataset until it is fully subdivided (each individual cell representing its own cluster), here

97    dubbed "recursion-to-completion" (**Fig. 2a**). In practice, this would indicate that someone

98    analyzing a scRNAseq dataset could always decide to sub-cluster a "cluster of interest" and report

99    a "novel subpopulation" of cells, resulting in false discoveries. To test the robustness of each

100   feature selection algorithm to the recursion-to-completion problem, we selected four publicly

101   available datasets from differing species and platforms including droplet-based UMI approaches

102   (**Fig. 2b**) and full-length transcript single-cell and -nucleus RNAseq (sNucSeq) (**Fig. 2c**)[14]. Again,

103   we found that standard overdispersion- and dropout-based feature selection methods enabled

104   recursion-to-completion, often finding hundreds of clusters, while anti-correlation-based feature

105   selection were robust to this problem. Anti-correlation showed fewer rounds of recursion ($P \leq 0.05$

106   for TukeyHSD post-hocs), and fewer overall clusters ($P \leq 1e-3$ for TukeyHSD post-hocs) relative

107   to other methods (**Figure 2d-e**). This demonstrates that anti-correlation-based feature selection

108   is robust to differing technologies, species, and sequencing type, retaining the ability to minimize

109   false sub-divisions.

110        To verify these results with known ground-truth, we simulated 4 clusters, and allowed each

111   algorithm to iteratively sub-cluster until either no features were returned, or only a single cluster

112   was identified. Consistent with our findings from real-world datasets, anti-correlation-based

113   feature selection protected against erroneous sub-clustering, while other approaches allowed for

114   several rounds of recursive sub-clustering, yielding hundreds to thousands of final 'clusters' (fewer

115   average rounds of sub-clustering: $P=1.08e-6, F=52.9$, main-effects 1-way ANOVA; $P \leq 6.2e-6$ for

116   TukeyHSD post-hocs; fewer total clusters: $P=7.2e-10, F=238.2$, main-effects 1-way ANOVA;

117   $P \leq 1.3e-9$ for TukeyHSD post-hocs); **Extended Data Fig. 1a**). These simulated data demonstrate

118    that anti-correlated feature selection guards against erroneously splitting a single population of

119    cells, while the algorithms tested here enable false discoveries of what appear to be "novel sub-

120    types."

121         We next sought to determine the overall accuracy of these feature selection algorithms,

122    where ground-truth differentially expressed genes (DEGs) should be selected by feature selection

123    algorithms, and non-DEGs should not be selected. To this end, we used Splatter to simulate

124    datasets comprised of 4, 6, 8, and 10 clusters. Our anti-correlation algorithm had the best

125    accuracy, F1-score, Mathew's Correlation Coefficient (MCC), precision, true negative rate, FPR,

126    and false discovery rate (FDR) compared to other feature selection algorithms (**Extended Data**

127    **Fig. 1b**). However, anti-correlation-based feature selection had average recall (also called

128    sensitivity or false negative rate); this is explained however, by Splatter's wide-spread co-

129    expression of *all* genes in *all* clusters (**Extended Data Fig. 2a**). In other words, using Splatter, *all*

130    *clusters* express the "marker-genes" of *all other* clusters, therefore blunting the anti-correlations

131    of marker-genes seen in practice (**Fig. 1**), thus reducing the apparent sensitivity. SERGIO

132    however is a gene regulatory network (GRN) based scRNAseq simulation approach that more

133    accurately represents empirical scRNAseq datasets[15] and does not induce co-expression of all

134    marker genes in all clusters (**Extended Data Fig. 2b**). Using this simulation paradigm anti-

135    correlation-based feature selection outperformed other approaches by *every* metric including

136    recall/sensitivity (**Extended Data Fig. 1c**). Furthermore, using seven pancreatic datasets,[2, 16-20]

137    the anti-correlated genes were either tied for, or had significantly higher p-value significance rank,

138    precision, and recall for pancreatic specific genes based on gProfiler/Human Protein Atlas tissue

139    enrichment compared to other algorithms (**Extended Data Fig. 1d**)[21, 22].

140         To assess the practical scalability of anti-correlation-based feature selection, we re-

141    processed and ran a larger dataset (245,389 cells) from a *Tabula Muris* data-release[23]. The full

142    feature selection process took 60.95 minutes, while calculating the cell-cell correlations, distance,

143    and clustering were far more computationally intense taking several days (see **Methods** for

6

144  clustering details) (**Fig 2f**). These findings show that anti-correlation-based feature selection

145  should not be a major limiting factor for large datasets.

146      We also sought to demonstrate our feature-selection approach's utility in safe sub-

147  clustering in practice; to this end, we focused on a cluster whose marker genes included

148  insulin/amylin (*INS1/2*, *IAPP*) and glucagon (*GCG*), the markers for pancreatic beta and alpha

149  cells, respectively, indicating that this cluster was insufficiently divided in the first clustering round.

150  We performed sub-clustering with anti-correlation, identifying leukocyte, alpha-, beta-, and delta-

151  cell populations. We further sub-clustered the insulin high population, and unexpectedly found the

152  rare[24] population of pancreatic-polypeptide (*Ppy/Pyy*) expressing PP-cells (**Fig. 2g**), a cluster

153  comprising only 0.01% of the original dataset. Attempting to further sub-divide PP-cells yielded

154  no usable features, thus showing that anti-correlation-based feature selection can facilitate

155  extremely sensitive sub-clustering to identify rare biologically meaningful populations from large

156  datasets, while also preventing errant subdivisions.

157      As seen in the final sub-cluster round, however, while anti-correlation-based feature-

158  selection is biologically accurate and answers the question: "Should this cluster be sub-

159  clustered?", it does not ensure that downstream algorithms will select the correct number of

160  clusters; this remains an outstanding problem as previously reported[9]. However, passing the first

161  step of successfully identifying a homogeneous population, through anti-correlation-based feature

162  selection, provides confidence that meaningful structure existed in the parent population.

163      Overall, these results demonstrate that anti-correlation-based feature selection solves

164  the null-dataset and recursion-to-completion problems, outperforms others in overall feature

165  selection accuracy, and works with both UMI and full-length sequencing methods. These

166  properties can prevent wasted time and money for bench-practitioners attempting to validate

167  novel sub-populations by providing an algorithmic check to false discoveries in scRNAseq.

168  Lastly, our open source python package (titled anticor_features) is open-source, pip installable,

169  and compatible with SCANPY/AnnData[25] to enable broad adoptability.

170    **Code and Data Availability**

171    All code used for implementing the anti-correlation-based feature selection approach is available

172    as a stand-alone package:

173    https://bitbucket.org/scottyler892/anticor_features

174    and is also pip installable:

175             python3 -m pip install anticor_features

176    All code for running simulations and comparisons used in this study are available at:

177    https://bitbucket.org/scottyler892/anti_correlation_vs_overdispersion/

178

179     **Methods**

180     *Example of anti-correlation principle on pancreatic dataset*

181         A previously published scRNAseq dataset and annotations were used for scatter plots of

182     *AMY2A* for acinar cells, *SST* for delta cells, and *NEUROD1* for endocrine cells (**Fig. 1d-f**)[2].

183

184     *Normalization of scRNAseq datasets to be used for benchmarking*

185         Due to large variation (often orders of magnitude differences) in total UMI counts across

186     cells and it's downstream effects on cell-to-cell distance metrics, we normalized each cell within

187     UMI based datasets through bootstrapped UMI downsampling as described here:

188     https://bitbucket.org/scottyler892/pyminer_norm. In brief, a cutoff is selected for both the number

189     of observed genes in a cell as well as the number of total UMI observed in a cell. Cells not meeting

190     these criteria are removed, and all other cells are normalized through UMI downsampling. UMI

191     downsampling is done through simulating the transcriptome of a given cell, and randomly

192     selecting N transcripts, where N is the desired number of total UMI for each cell to have, in this

193     case 95% of the cutoff used for total UMI count. Thus, each cell is randomly sampled to the same

194     UMI depth.

195         To normalize full-length sequencing datasets with TPM or similar units, we created a

196     variant of quantile normalization we call truncated quantile normalization. First a cutoff ($g$) is

197     selected for the number of genes to be expressed in each cell in the final normalized dataset.

198     Next, cells with fewer than $g$+1 genes expressed are removed, then for each cell, the

199     transcriptome is subtracted by the expression value of gene $g$+1 for that cell, thus setting the $g$+1

200     gene's expression to zero, leaving the remaining top $g$ expressed genes with >0 expression in all

201     cells. All negative values are then set to 0. For ties at the expression-level of $g$ that would result

202     in differing number of observed genes, genes are randomly selected to be preserved or set to

203     zero stochastically. This yields a vector for each cell for whom the top expressed $g$ genes are

204     kept, but shifted downwards in a manner that does not introduce an artificially large gap between

9

205    the lowest expressed gene (*g*) and zero. These top *g* genes for each cell are then quantile

206    normalized. This process is implemented in the pyminer_norm pip package, and can be called

207    from the command-line on tsv files:

208        python3 -m pyminer_norm.quantile_normalize -i in_file.tsv -o out_file_qNorm.tsv -n 2000

209    to perform truncated quantile normalization on the top 2000 genes for each cell.

210

211    *NIH3T3 and HEK293T cell line datasets*

212        This    dataset    was    downloaded    from    10x    Genomics'    website    at

213    (https://support.10xgenomics.com/single-cell-gene-expression/datasets/3.0.2/1k_hgmm_v3).

214    The cells of mouse or human origin were separated into distinct datasets for our purposes here

215    based on the sum of reads that mapped to each species' transcriptome, while doublets were

216    excluded. In the case of both human and mouse references, cells were kept that had >3162

217    counts mapping to hg19 or mm10 for HEK293T and NIH3T3 respectively, cells were also only

218    kept if they had >1000 genes observed. The remaining cells were then downsampled to 3003

219    counts for each dataset to normalize for variable count depth that otherwise spanned two orders

220    of magnitude.

221

222    *Affinity Propagation*

223        Our    implementation    of    affinity    propagation    was    based    on    the    sklearn

224    sklearn.cluster.AffinityPropagation function, in which the preference vector is initialized to the row-

225    wise minimum of the input matrix; in this case, the negative squared Euclidean distance of the

226    Spearman correlations across all cells. We observed that as datasets scale, the original affinity

227    propagation algorithm fragments single populations into many small populations that were similar

228    to each other. We therefore follow the original affinity propagation results with an analysis that

229    calculates the distance (in affinity space) between cluster centers (also called exemplars). The

230    standard deviation of within-cluster affinities is then calculated. For each cluster-cluster pair from

231  the original affinity propagation cluster results, we then determine the number of combined

232  standard deviations required to traverse half the Euclidean distance in affinity space between two

233  cluster centers. This measure is the number of standard deviations needed to reach the waypoint

234  between two cluster centers. Because these are standard deviation measures, we can convert

235  these to transition probabilities, as with a Z-score, using the scipy.stats.norm.sf function. This

236  creates a cluster x cluster matrix of transition probabilities; this probability matrix is then subjected

237  to dense weighted Louvain modularity. Final clusters are determined by the results of this

238  procedure, where AP clusters that were determined by Louvain modularity to belong to the same

239  community are merged. All code and cluster for the affinity propagation with merged procedure

240  can be accessed through running PyMINEr with the appended arguments: " -ap_clust -ap_merge"

241  at the command line or interactively via the pyminer.pyminer.pyminer_analysis function using the

242  arguments: ap_clust=True, ap_merge=True.

243

244  *Clustering – K-means with Elbow and K-means with silhouette*

245  First each dataset (already log transformed) was subset for the genes selected by the

246  given feature selection algorithm, then genes were min-max linear normalized between 0 and 1.

247  K-means clustering was performed using the sklearn.cluster KMeans function. For the elbow rule,

248  the sum of squared Euclidean distances of samples to their cluster center was used in conjunction

249  with the given k value. We took the elbow to be the value of k which yielded the minimum distance

250  to the origin.

251  For the silhouette method, we calculated the average silhouette score with the

252  sklearn.metrics silhouette_score function, and sample level silhouettes calculated with the

253  silhouette_samples function. The number of clusters was selected by moving from k=1 to k_max,

254  testing for whether there existed a cluster whose maximum sample level silhouette was less than

255  the average silhouette score for the whole dataset (as determined by the silhouette_score

256  function).

257

258    *Clustering – Locally weighted Louvain modularity*

259    We created a kNN graph embedding and subjected it to Louvain modularity as follows:

260    1. Calculate Spearman correlation of all cells against all other cells (matrix: **S**).

261    2. Calculate the inverse squared Euclidean distance matrix from the Spearman matrix

262    (matrix: **D**), divided by the square-root of the number of cells. In this matrix, cells that are

263    more similar to each have higher values, and cells that are dissimilar have lower values,

264    inversely proportional to the squared Euclidean distance.

265    3. For each cell, *i*, (i.e.: row in matrix **D**) subtract the upper 95$^{th}$ percentile (or top 200$^{th}$ closest

266    cell, whichever yields fewer connections) of distance vector (**D**$_i$), then mask all negative

267    values to zero, thus creating a weighted local distance matrix (matrix: **L**).

268    4. To ensure that all cells are on an equivalent scale, each row in **L** is divided by it's maximum

269    ( **L**$_i$ = **L**$_i$ / max(**L**$_i$) ).

270    5. The normalized local distance matrix **L** serves as the weighted adjacency matrix for

271    building the network for weighted Louvain modularity.

272

273    The locally weighted adjacency matrix was subjected to Louvain modularity as implemented in

274    the python pip package: python-louvain.

275

276    *Implementation of other feature selection algorithms*

277    Because each feature selection algorithm expects slightly different processing methods

278    relative to each other (either normalized and log-transformed, or count data), we followed author

279    guidance in implementation.

280

281    *PyMINEr's overdispersion pipeline*: is contained within the originally published full PyMINEr

282    pipeline, but is also callable within python as follows:

283    feature_table = do_over_dispers_feat_select(ids=cell_ids,

284    ID_list=gene_ids,

285    in_mat=exprs)

286

287    *Seurat's overdispersion*: Per author guidelines, we log-normalized the input expression matrix

288    and selected features as follows:

289    obj<- NormalizeData(CreateSeuratObject(exprs))

290    obj <- FindVariableFeatures(obj)

291    var_feat <- VariableFeatures(obj)

292

293    *Original Brennecke algorithm*: We used the implementation of the original overdispersion-based

294    feature selection algorithm as implemented in the M3Drop package as follows:

295    Brennecke_HVG <- BrenneckeGetVariableGenes(exprs, fdr = 0.05, minBiolDisp = 0.5)

296

297    *M3Drop*: Unlike other most other feature selection algorithms, M3Drop allows for either a pre-

298    specified FDR, or a pre-specified percentage of the transcriptome to select. In our testing using

299    the FDR approach (which could theoretically solve that the null-dataset problem), we found that

300    each dataset required fine tuning of this cutoff to provide reasonable results, and in the case of

301    full-length transcript based approaches did not select any genes even in the full datasets, which

302    are known to be biologically complex. We therefore sought a more realistic implementation that

303    did not require manual tuning for each dataset, and therefore implemented the "percentage"

304    approach within M3Drop so that a standard call yielded meaningful results regardless of dataset,

305    without necessitating a manual inspection for hyperparameter selection for all datasets, which

306    could also be seen as tuning hyperparameters to fit our expectations of the data. The

307    implementation was as follows:

308    results <- M3DropGetExtremes(exprs, percent=0.05, suppress.plot=TRUE)

13

309   Using the genes within the results$right section as the genes with an excess of zeros for the final

310   selected genes.

311

312   *Details of anti-correlation feature selection algorithm*

313       We aimed to develop an algorithm that identifies genes that have "too many" negative

314   correlations below a dynamically selected cutoff that make the selected genes more negatively

315   correlated with other genes than one would expect from random chance. To this end we began

316   with a False Positive Rate (FPR) of 0.001, for identifying a cutoff at which correlations should be

317   counted as a "discovery" (D, where more significant), or "non-discovery" (ND, where less

318   significant). Using a bootstrap shuffled null background, in which all discoveries (D) are false,

319   because true positives (TP) are known to be equal to zero:

320   $$FP + TP = N(D)$$

321   Where D is all discoveries, more significant that the cutoff. Therefore because this is measured

322   from a bootstrap shuffled null background (i.e.: TP = 0):

323   $$FP = N(D)$$

324   Using this knowledge, we created the null background of gene-gene Spearman correlations is

325   generated through randomly sampling 5,000 genes, shuffling within-genes, such that a gene-

326   gene correlation plot would have its x-y pairing shuffled, calculating pairwise Spearman

327   correlations.

328

329   **<u>Definitions</u>**

330   **$E_o$**: the original expression matrix

331   *rand*: an integer vector of the length 5000 for the random samples within the space of 1..n, where

332   n is the number of genes

333   **$E_r$**: The random subset matrix that is permuted as defined below:

334   For i..$N$(*rand*):

335     $\mathbf{E}_{r,i} = permute(\mathbf{E}_{o,\ rand[\ i\ ]})$

336     Where $\mathbf{E}_r$ provides a N(cell) x N(rand) matrix, which is a within-gene bootstrap shuffled

337     version of a subset of the transcriptome, therefore unpairing the gene-gene pairs for measuring

338     the null background of Spearman correlations.

339     In our testing, using a greater number of randomly selected genes, *N(rand)*, for the

340     permutation based null-background did alter the null-distributions, as these distributions were

341     stable at this sampling depth, and did not notably change the selected cutoffs. Note that the

342     method of rank transformation for Spearman correlation effects the outcome; here we perform

343     dense-rank transformation. Non-dense rank transformations frequently result in large gaps within

344     the distributions because of ties. This is particularly important with count-based datasets where

345     ties are frequent.

346     The null Spearman background matrix (**B**) was the symmetric 5000 x 5000 comparison of

347     this sample (5000 choose 2 combinations).

348     For i=1..*N(rand)* and j=1..*N(rand)*:

349     $$\mathbf{B}_{i,j} = Spearman(\mathbf{E}_{r,i}, \mathbf{E}_{r,j})$$

350

351     Next, this **B** background matrix, of null Spearman rho values, is filtered for only values $\mathbf{B}_{i,j} <0$, thus

352     creating a negative correlation null-background; this is needed because the null background for

353     values $\mathbf{B}_{i,j} >0$ and values $\mathbf{B}_{i,j} <0$ follow different distributions (**Extended Data Fig. 2c**), indicating

354     the necessity to measure them independently. Self-comparisons and duplicate comparisons were

355     also removed.

356     For i=1..*N(rand)*, and j=i+1..*N(rand)*:

357     $$\boldsymbol{b} = (\mathbf{B}_{i,j} \in \mathbf{B} \mid \mathbf{B}_{i,j} < 0 \mid i > j)$$

358     Conceptually, this filtering is also important because the estimated number of false

359     positives (*FP*) for a given gene *i* is dependent on the number of genes that are actually randomly

15

360    distributed, or truly correlated. For example, gene X is co-regulated within a module of 2000

361    genes, while gene Y is not genuinely correlated with any other genes. Given that the number of

362    genes is static and zero sum, this true positive co-regulation removes those genes from possible

363    false positive negatively correlated genes.

364

365        This null background *vector* (*b*) is used to calculate an the cutoff ($C_{neg}$) that most closely

366    matches the desired FPR (default=1 in 1000 false positives), with a discovery considered as a

367    Spearman rho value < $C_{neg}$ in the gene-gene correlation matrix (**S**) calculated from the unshuffled

368    original expression matrix ($\mathbf{E}_o$), This cutoff is used for the estimated false discovery rate (*FDR*) for

369    the original intact unshuffled dataset.

370    Given that:

371
$$FPR = \frac{FP}{FP + TN}$$

372    and

373
$$N(b) = FP + TN$$

374    Because TP = FN = 0, given that *b* was generated from a bootstrap shuffled null. We therefore

375    find that:

376
$$FPR = \frac{FP}{N(b)}$$

377
$$N(b) * FPR = FP$$

378    Therefore, to identify the appropriate cutoff ($C_{neg}$), that yields the FPR(=1e-3 by default), we

379    simply take the Spearman rho value of *b* that is located within the sorted background vector that

380    gives the ratio of false positives to true negatives.

381
$$b_{sort} = sort(b)$$

382    Such that for i=1..$N(b_{sort})$-1, $b_{sort,i} < b_{sort,i+1}$

16

383    We then calculate the $C_{neg}$ cutoff, but taking the value at the index that gives the expected ratio

384    of false positives to true negatives as determined by the FPR hyperparameter (default=1e-3)

385                                    $$C_{neg} = b_{sort,\lfloor FP \rfloor}$$

386    Next, we use this empirically determined cutoff ($C_{neg}$), applying it to classify "discoveries" of

387    negative correlations in the correlation matrix **S** as calculated from the original, non-shuffled

388    dataset ($\mathbf{E}_o$). Where a discovery is defined as a Spearman rho value $\mathbf{S}_{i,j}$ less than the $C_{neg}$ cutoff.

389            Again it is important to note two things: 1) the null distribution of Spearman correlations,

390    are in fact two separate distributions concatenated around zero, for the null distribution of rho

391    values <0, and the null distribution of rho values >0 (**ED. Fig. 2c**); and 2) that variable abundance

392    of True Positives within the positive correlation domain will decrease the total number of

393    comparisons that fall within the negative correlation domain of these distributions; these two

394    distributions are therefore in competition with one another, meaning that they must be quantified

395    independently. For these reasons, when applying the empirically measured cutoff ($C_{neg}$) from the

396    shuffled transcriptome, we must apply it only to the correlations falling below zero.          To

397    apply this cutoff ($C_{neg}$) to the original expression matrix ($\mathbf{E}_o$), we first calculate the symmetric

398    gene-gene Spearman rho matrix (**S**).

399    Next, the number of total ($T$) Spearman rhos values <0 within **S** is tabulated for the application of

400    our cutoff ($C_{neg}$):

401                            $$T_i = N\big(\mathbf{S}_{i,j} \in \mathbf{S} \mid \mathbf{S}_{i,j} < 0\big)$$

402    For i=1..n, where n is the number of genes.

403    Note also, that $T_i$ sums to the total number of discoveries (D) and non-discoveries (ND).

404                        $$T_i = N(D_i) + N(ND_i) = TP + TN + FP + FN$$

405    Where:

406                    $$N(D_i) = N\big(\mathbf{S}_{i,j} \in \mathbf{S} \mid \mathbf{S}_{i,j} < C_{neg}\big) = FP + TP$$

407                $$N(ND_i) = N\big(\mathbf{S}_{i,j} \in \mathbf{S} \mid \mathbf{S}_{i,j} < 0 \mid \mathbf{S}_{i,j} > C_{neg}\big) = FN + TN$$

17

408

409    Further, the discoveries are comprised of both false positives (FP) and true positives (TP),

410    however, *which* individual values within the discovery class is a FP or TP is unknown. Using the

411    FPR however, we can estimate the number of *expected* FPs given the total number of

412    comparisons <0 for the given gene ($T_i$). In other words, if this gene were random in its negative

413    correlations, then only a specific number of false positives would be expected ($\widehat{FP_i}$), using $\boldsymbol{C}_{neg}$

414    as a cutoff.

415    $$\widehat{FP_i} = T_i * FPR$$

416    Therefore, with *FDR* defined as:

417    $$FDR = \frac{FP}{(FP + TP)}$$

418    We can estimate the *FDR* for each gene, determining if it has an over abundance of negative

419    correlations compared to what is expected from the null distribution:

420    $$\widehat{FDR_i} = \frac{\widehat{FP_i}}{N(D_i)}$$

421    We then select genes that have a >15x excess in discoveries relative the expected number of

422    false positives under the null distribution assumption. This corresponds to an estimated $\widehat{FDR} =$

423    0.066 (1/15). This yields the set of all excessively negatively correlated genes (A):

424    $$A = \left\{ gene_i \in genes \,\middle|\, \widehat{FDR_i} < 0.066 \right\}$$

425    Lastly, given that spurious positivity is still possible and even expected, we add one last layer of

426    protection against false discoveries. The positive/negative status of a single gene likely does not

427    define a truly "novel subtype" – particularly in a technique such as single-cell -omics where

428    stochastic dropout from random sampling is expected. We therefore apply an additional filter from

429    the premise that the genes whose expression patterns separate meaningful populations should

430    also be positively correlated with other genes that are following similar regulatory patterns. To

431     select this population of genes, we find genes that have greater than 10 positive correlations

432     above the positive correlation cutoff ($C_{pos}$), as calculated similarly to ($C_{neg}$) as described above.

433     $$M = \left\{ gene_i \in genes \mid N\left( \mathbf{S_{i,j}} \mid \mathbf{S_{i,j}} > C_{pos} \mid i \neq j \right) > 10 \right\}$$

434     The final included features are the intersect of A and M:

435     $$F = A \cap M$$

436     Overall, this means that genes must contain both an excess of negative correlations, and be a

437     member of a "module" of at least 10 genes that move in concert.

438

439     *Recursion benchmarks*

440         An initial run of locally weighted Louvain modularity was performed, then the given dataset

441     was subset to contain only the cells of a given cluster in the prior round of clustering. Next, feature

442     selection and locally weighted Louvain modularity was applied again, recursively until either each

443     cell was called its own "cell-type"/cluster or produced "cell-types"/clusters with ≤5 cells.

444         Circular recursion graphs were displayed using networkx[26], with layout determined by the

445     graphviz_layout(prog='twopi') layout[27].

446

447     *In silico recursive clustering benchmark*

448         Four clusters were simulated using Splatter[13], and all algorithms were allowed to

449     recursively select features, which were then subjected to locally weighted Louvain modularity until

450     one of the following conditions were met: no features were selected, the clustering algorithm only

451     found a single cluster, or the results of clustering formed groups of 5 or fewer cells.

452

453     *Real-world recursive clustering benchmark*

454         The above described recursion procedure was applied to the previously released mouse

455     heart scRNAseq dataset,[28] and human PBMC dataset[29] for UMI based technologies, and mouse

456   hippocampus single nucleus RNAseq[14] and human dendritic cell/monocyte[30] datasets were used

457   for full length transcript sequencing based approaches. Each dataset was normalized as

458   described above and is available in the repository site containing this benchmark:

459   https://bitbucket.org/scottyler892/anti_correlation_vs_overdispersion in the data folder. The same

460   recursive clustering procedure was followed as described for the *in silico* recursion benchmark

461   above.

462

463   *Feature selection accuracy based on Splatter and Sergio simulations*

464      For both simulation paradigms, we simulated 4, 6, 8, and 10 clusters. 2500 cells were

465   simulated with 10000 genes, of which 2000 were intended to be differentially regulated across

466   clusters. Once simulations were completed, the datasets were downsampled down to 95% of the

467   cell with the lowest total counts in the given dataset, using the pyminer_norm python package[31].

468      Splatter simulations were generated using the bin/simulate_data.R with the above

469   described clusters, cells, and gene parameters. SERGIO simulations were generated from the

470   bin/generate_sergio_sim.py script, which was called from the bin/simulate_data.R file. For each

471   cluster, a single "master-regulator" gene was used to induce high expression of its child nodes in

472   the GRN. The non-differentially regulated genes were random negative binomial distributions

473   added to the network with the np.random.negative_binomial function.

474      Similar to performing pathway analyses, a proper background list of genes is necessary

475   for quantifying enrichment. For example there may be a simulated low-expression gene that was

476   "differentially expressed" in ground-truth, however, was only expressed in two cells after

477   simulation of the low expressed gene. In this situation, this gene it would not be realistically

478   possible to "detect" this gene as differentially expressed even if ground truth clusters were known.

479   Therefore to generate a background of detectably differentially expressed genes, were performed

480   differential expression analysis by 1-way ANOVA (aov function) using the known ground truth

481   cluster labels. This gives us a list of detectably differentially expressed genes to use as the ground

482    truth desired genes for feature selection, while non-detectably differentially expressed were all

483    treated as not desired for selection. This parallels pathway analysis in that, if a gene is not

484    detectably expressed, it should not be included in the custom background.

485

486    *Pancreatic datasets for feature selection*

487    The seven pancreatic datasets[2, 16-20] used for feature selection efficacy benchmarking

488    were processed as previously described[2]; the available post-processing datasets were used as-

489    is. These datasets are also now re-packaged in the data zip contained within the benchmark

490    repository. To assess efficacy, three primary metrics were used via gProfiler analysis using the

491    human protein atlas "HPA" pathways which indicates genes are enriched for certain tissues and

492    sub-tissue niches[21, 22]. For each dataset, a custom background was used, comprised of the genes

493    expressed in the given dataset. For each analysis, the HPA results were filtered to include only

494    the pancreatic tissues and niches, the pancreatic HPA pathway that was the most significant was

495    counted as a method's best pancreatic match. The -log10(p-values), precision, and recall for this

496    best match was used for comparisons. To adjust for the wide range and skewed distributions in

497    significance across datasets and methods, we rank transformed the -log10(p-values); precision

498    and recall however are all on a scale between 0 and 1, and were therefore analyzed directly.

499    Significance was determined with the aov and TukeyHSD functions to measure the main effects

500    and post-hocs respectively. The aov function was called with the formula: metric ~ method +

501    dataset.

502

503    *Tabula Muris dataset*

504    The senescent Tabula Muris dataset[23] was used to demonstrate the scalability of our

505    analytic pipeline. This dataset was previously filtered to contain only cells with ≥2500 UMI counts.

506    We therefore downsampled the dataset such that all cells contained 2500 UMI, and log2

507    transformed it for analysis. The downsampling process was performed using the bio-pyminer-

508    norm package that is pip installable:

509        python3 -m pip install bio-pyminer-norm

510    The process of downsampling is reported in detail at the repository website:

511    https://bitbucket.org/scottyler892/pyminer_norm

512

513    Subclustering rounds were first feature selected with the anti-correlation package that we

514    released here, using default parameters:

515        from anticor_features.anticor_features import

516        anti_cor_table = get_anti_cor_genes(exprs, feature_ids, species = "mmusculus")

517    Locally weighted Louvain modularity was used for clustering as described above. Note that while

518    the default functionality of our feature selection package automatically removes ribosomal,

519    mitochondrial, and hemoglobin related genes, for fair comparison with other methods, these

520    genes were left in for possible selection when comparing to other algorithms. This can be

521    customized using the pre_remove_pathways argument. The default removal list are genes

522    contained in the following pathways (all related to ribosomal, mitochondrial, and hemoglobin):

523    "GO:0044429","GO:0006390","GO:0005739","GO:0005743","GO:0070125","GO:0070126","GO:

524    0005759","GO:0032543","GO:0044455","GO:0005761","GO:0005840","GO:0003735","GO:0022

525    626","GO:0044391","GO:0006614","GO:0006613","GO:0045047","GO:0000184","GO:0043043"

526    ,"GO:0006413","GO:0022613","GO:0043604","GO:0015934","GO:0006415","GO:0015935",

527    "GO:0072599","GO:0071826","GO:0042254","GO:0042273","GO:0042274","GO:0006364","GO:

528    0022618","GO:0005730","GO:0005791","GO:0098554","GO:0019843","GO:0030492"

529    Alternatively, if the user whishes to exclude specific features, these can be included in the

530    pre_remove_features list argument; however, this was left empty for all of the work presented

531    here.

## References

1.  Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e1821 (2019).
2.  Tyler, S.R. et al. PyMINEr Finds Gene and Autocrine-Paracrine Networks from Human Islet scRNA-Seq. *Cell Reports* **26**, 1951-1964.e1958 (2019).
3.  Brennecke, P. et al. Accounting for technical noise in single-cell RNA-seq experiments. *Nature methods* **10**, 1093 (2013).
4.  Andrews, T.S. & Hemberg, M. M3Drop: dropout-based feature selection for scRNASeq. *Bioinformatics* **35**, 2865-2867 (2018).
5.  Madissoon, E. et al. scRNA-seq assessment of the human lung, spleen, and esophagus tissue stability after cold preservation. *Genome Biology* **21**, 1 (2019).
6.  Cui, Y. et al. Single-Cell Transcriptome Analysis Maps the Developmental Track of the Human Heart. *Cell Reports* **26**, 1934-1950.e1935 (2019).
7.  Kaplan, N. et al. Single-Cell RNA Transcriptome Helps Define the Limbal/Corneal Epithelial Stem/Early Transit Amplifying Cells and How Autophagy Affects This Population. *Investigative Ophthalmology & Visual Science* **60**, 3570-3583 (2019).
8.  Ayyaz, A. et al. Single-cell transcriptomes of the regenerating intestine reveal a revival stem cell. *Nature* **569**, 121-125 (2019).
9.  Kiselev, V.Y., Andrews, T.S. & Hemberg, M. Challenges in unsupervised clustering of single-cell RNA-seq data. *Nature Reviews Genetics* **20**, 273-282 (2019).
10.  Kleinberg, J. An impossibility theorem for clustering. *Advances in neural information processing systems*, 463-470 (2003).
11.  Liu, H. et al. Systematically labeling developmental stage-specific genes for the study of pancreatic β-cell differentiation from human embryonic stem cells. *Cell Research* **24**, 1181-1200 (2014).
12.  Andrews, T.S. & Hemberg, M. Dropout-based feature selection for scRNASeq. *bioRxiv*, 065094 (2018).
13.  Zappia, L., Phipson, B. & Oshlack, A. Splatter: simulation of single-cell RNA sequencing data. *Genome Biology* **18**, 174 (2017).
14.  Habib, N. et al. Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn neurons. *Science* **353**, 925-928 (2016).
15.  Dibaeinia, P. & Sinha, S. SERGIO: A Single-Cell Expression Simulator Guided by Gene Regulatory Networks. *Cell Systems* **11**, 252-271.e211 (2020).
16.  Li, J. et al. Single-cell transcriptomes reveal characteristic features of human pancreatic islet cell types. *EMBO reports* **17**, 178-187 (2016).
17.  Muraro, M.J. et al. A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell Syst* **3**, 385-394.e383 (2016).
18.  Segerstolpe, Å. et al. Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and Type 2 Diabetes. *Cell metabolism* **24**, 593-607 (2016).
19.  Wang, Y.J. et al. Single-Cell Transcriptomics of the Human Endocrine Pancreas. *Diabetes* **65**, 3028-3038 (2016).
20.  Xin, Y. et al. RNA Sequencing of Single Human Islet Cells Reveals Type 2 Diabetes Genes. *Cell metabolism* **24**, 608-615 (2016).
21.  Raudvere, U. et al. g:Profiler: a web server for functional enrichment analysis and conversions of gene lists (2019 update). *Nucleic Acids Research* **47**, W191-W198 (2019).
22.  Uhlén, M. et al. Tissue-based map of the human proteome. *Science* **347**, 1260419 (2015).
23.  Almanzar, N. et al. A single-cell transcriptomic atlas characterizes ageing tissues in the mouse. *Nature* **583**, 590-595 (2020).
24.  Brereton, M.F., Vergari, E., Zhang, Q. & Clark, A. Alpha-, Delta- and PP-cells: Are They the Architectural Cornerstones of Islet Structure and Co-ordination? *The journal of*

*histochemistry and cytochemistry : official journal of the Histochemistry Society* **63**, 575-591 (2015).

25. Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biology* **19**, 15 (2018).

26. Hagberg, A., Swart, P. & S Chult, D. (Los Alamos National Lab.(LANL), Los Alamos, NM (United States), 2008).

27. Ellson, J., Gansner, E.R., Koutsofios, E., North, S.C. & Woodhull, G. in Graph drawing software 127-148 (Springer, 2004).

28. Genomics, x. (2018).

29. Genomics, x. (2018).

30. Villani, A.-C. et al. Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors. *Science (New York, N.Y.)* **356**, eaah4573 (2017).

31. Tyler, S.R., Bunyavanich, S. & Schadt, E.E. PMD Uncovers Widespread Cell-State Erasure by scRNAseq Batch Correction Methods. *bioRxiv*, 2021.2011.2015.468733 (2021).

599

**Figure 1: Anti-correlation algorithm premise and passage of the null-dataset problem. a**, The logic behind anti-correlation-based feature selection. Marker-genes will be expressed at

602    higher levels in their lineage/cell-type compared to cells outside of that lineage or cell-type. **b**,
603    As a scatter plot where expression of marker A is plotted against marker B, cells of type A and B
604    will form an L-shaped anti-correlation pattern, while cell-type C would express low levels of both
605    marker A and B. **c**, This anti-correlation pattern would disappear when examining a single
606    population of cells. **d**, The anti-correlation pattern of marker-genes appears in an example
607    dataset,[2] where high expression of *AMY2A* in acinar cells forms an anti-correlation pattern with
608    *SST* in delta cells of the pancreas. **e**, The anti-correlation pattern between *AMY2A* and *SST*
609    disappears when only subset for delta cells. **f**, The anti-correlation pattern is also present in
610    lineage-marking-genes as shown by the pattern of *AMYA2* and *NEUROD1*, which labels all
611    endocrine cells of the pancreas. **g**, The anti-correlation-based feature selection algorithm first
612    calculates a null background of Spearman correlations based on bootstrap shuffled gene-gene
613    pairs to calculate a background. **h**, Next the cutoff value closest matching the desired false
614    positive rate (FPR) is determined. Displayed is a histogram of the bootstrap shuffled null-
615    background of Spearman correlations less than zero. **i**, Lastly genes which show more
616    significant negative correlations (x-axis) than expected by chance (black line), given the gene's
617    number of total negative correlations (y-axis), are selected: i.e. those to the right of the cutoff
618    line. These are then used to calculate the False Discovery Rate (FDR) for each gene (See
619    **Methods** for details). **j-k**, Heatmaps of selected features, and the total number of subclusters for
620    each method of feature selection paired with AP clustering, when algorithms were allowed to
621    sub-divide iteratively for homeostatic cell line scRNAseq: (**j**) NIH3T3, (**k**) HEK293T. **l**, Boxplots
622    indicating the total number of clusters identified by each method of feature selection (box colors)
623    and clustering (noted in panels) showed that anti-correlation-based features selection (arrows)
624    identified no features, indicating a single population in all cases, while other methods produced
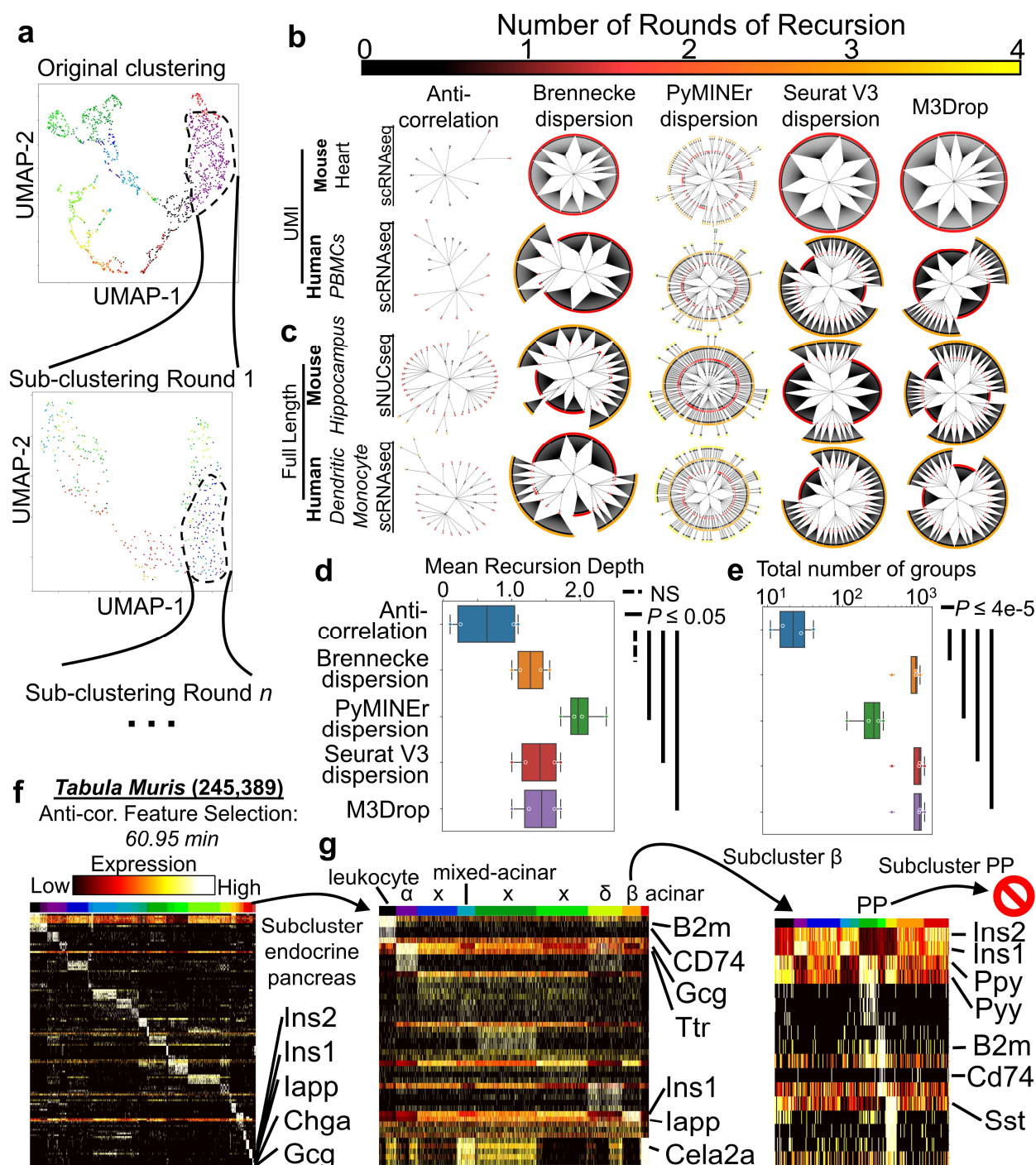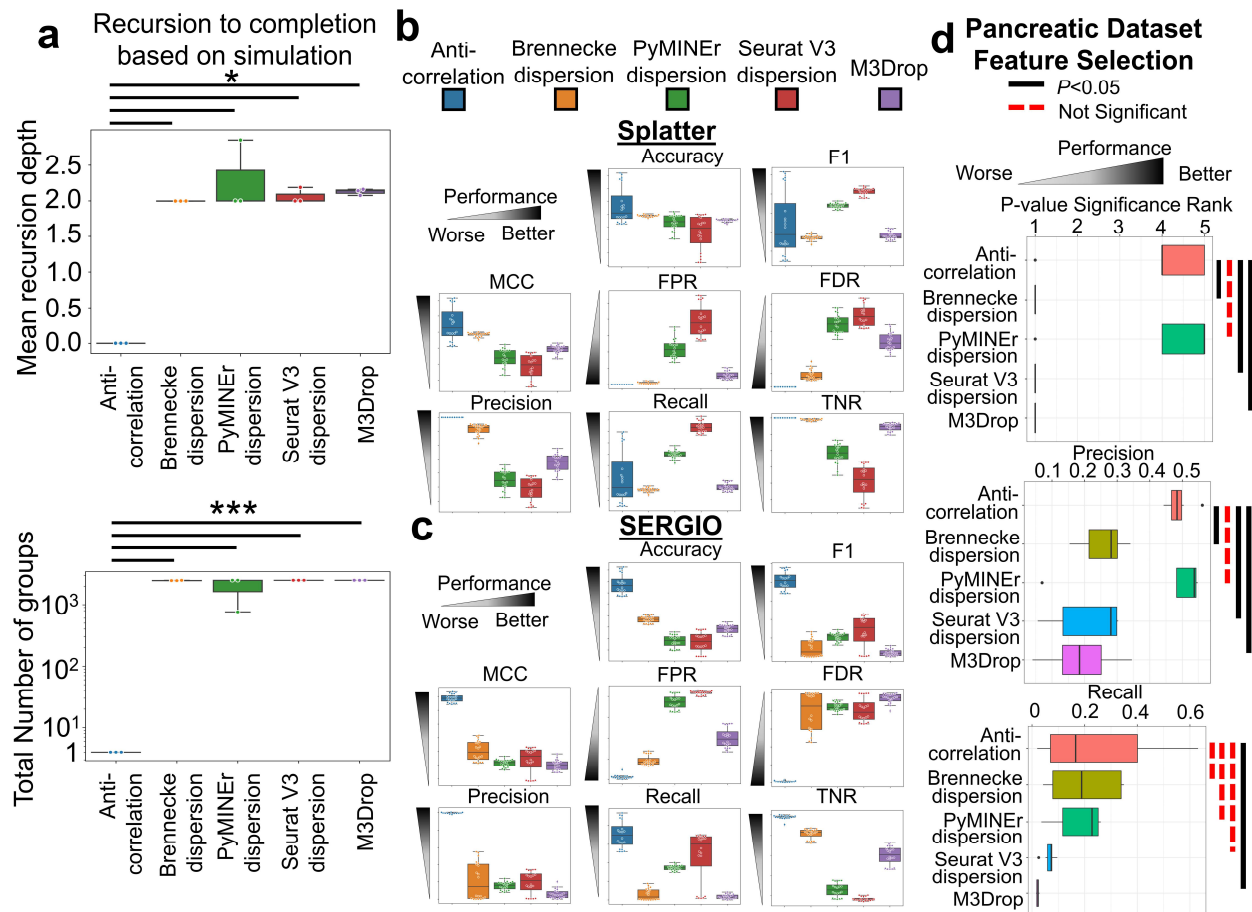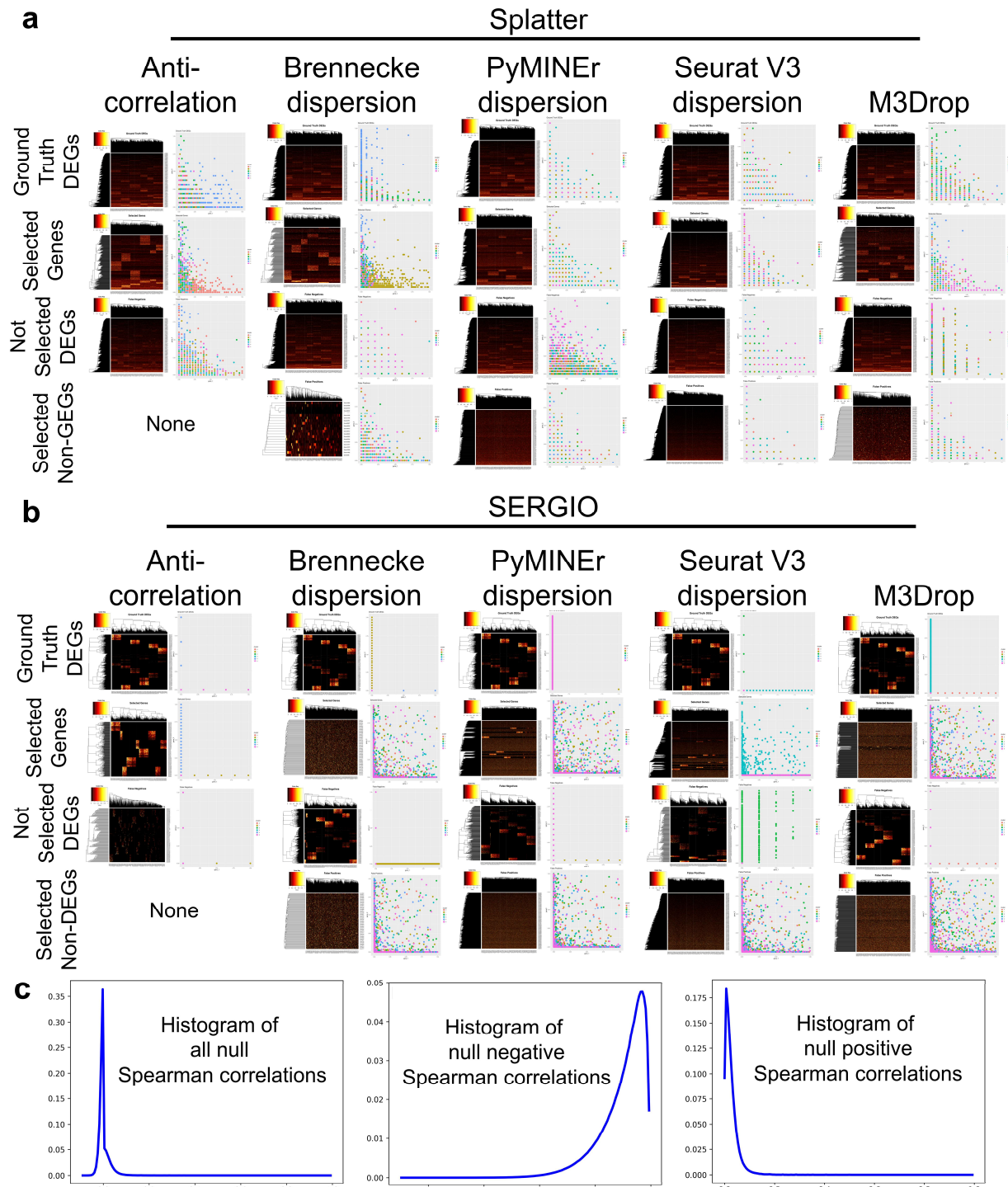625    more clusters, thus failing the null-dataset problem.

**Figure 2: Recursion-to-completion in real datasets and anti-correlation algorithm scaling.**
**a**, A schematic of sub-clustering is shown in the form of UMAP projections of the original dataset (left panel), and a sub-clustering iteration of a population found in the first round of feature selection and clustering (right panel). **b-c**, In real datasets of varying technologies, status quo algorithms fail the recursion-to-completion problem while the anti-correlation-based approach prevented recursion-to-completion. Recursive clustering plots where each point indicates a cluster at a given recursive clustering recursion-depth as denoted in successive

634     rings and color. **d**, Boxplots of the mean recursion depth for each of the final sub-clusters for
635     each noted method. **e**, Boxplots of the total number of groups obtained through iterative sub-
636     clustering. **f**, A heatmap of the top 5 marker genes per cluster are shown for the 26 primary
637     lineages from the full senescent *Tabula Muris* dataset[23], with the last cluster representing a
638     mixture of endocrine pancreas. **g**, When subclustered with anti-correlated feature selection, cell-
639     type droplets (x) as well as classically described leukocyte, α, δ, β, and acinar populations were
640     discovered. Subclustering β cells discovered mixed-lineage droplets with δ and leukocyte cells
641     as well as the rare PP-cell population, but additional subclustering of PP-cells was prevented by
642     anti-correlation-based feature selection.

**Extended Data Figure 1: Anti-correlation-based feature selection outperforms other methods in recursion-to-completion and feature selection efficacy. a**, Using Splatter simulation of four clusters, all algorithms were allowed to select features and perform locally weighted Louvain modularity-based clustering recursively. Shown are boxplots indicating the mean recursive depth and total number of clusters on a log scale. The anti-correlation algorithm did not allow for any recursive clustering, resulting in fewer clusters identified (*:$P \leq 6.2e\text{-}6$; ***:$P \leq 1.7e\text{-}9$; all ANOVA/TukeyHSD post-hoc comparisons against anti-correlation). **b-c**, Taken as a classification problem in which a feature selection algorithm's task is to select detectably differentially expressed genes across clusters, we quantified each algorithm's accuracy, F1 score, Mathew's Correlation Coefficient (MCC), false positive rate (FPR), false discovery rate (FDR), precision, true negative rate (TNR), and recall. **b**, Boxplots of classification metrics (panels) by feature selection approach (colored boxes) using Splatter simulations[13]. **c**, Boxplots indicating the performance of each features selection method (colored boxes) for each metric (panels), using SERGIO, gene regulatory network based simulations[15]. **d**, Using 7 pancreatic datasets[2, 16-20], each algorithm's selected features was analyzed for significance with pancreatic tissue enrichment via gProfiler and the human protein atlas[21, 22]; displayed are boxplots of the "best" pancreatic pathway by p-value comparing this pathway's rank p-value, precision, and recall.

**Extended Data Figure 2: Examples of Splatter and SERGIO simulations, and feature selection**. **a,b**, For both simulation paradigms (**a**) Splatter and (**b**) SERGIO, heatmaps are shown for the ground truth differentially expressed genes (DEGs), the selected-genes, non-selected DEGs, and selected genes that are not differentially expressed. Next to the heatmaps are gene-gene scatter plots of randomly selected genes from the indicated class (row) for the feature selection algorithms (columns). Points indicate an individual cell's expression of random

31

668    gene-x and gene-y for the designated gene class and algorithm, colorized by the simulated

669    cluster. (**a**) Splatter DEGs show widespread co-expression of DEGs within all clusters, while (**b**)

670    SERGIO allows for cluster specific expression of DEGs. (**c**) An example histogram of null

671    distribution patterns of Spearman rhos on shuffled datasets shows that, even on shuffled data

672    with no true positives, negative rhos follow a different distribution than positive rho values.