



## A machine learning method for the discovery of minimum marker gene combinations for cell-type identification from single-cell RNA sequencing

Brian D Aeversmann, Yun Zhang, Mark Novotny, et al.

*Genome Res.* published online June 4, 2021

Access the most recent version at doi:[10.1101/gr.275569.121](https://doi.org/10.1101/gr.275569.121)

---

<b>P&lt;P</b>	Published online June 4, 2021 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---

CRISPR and RNAi Genetic Screening.  
Your new superpower.

LEARN MORE

CELLECTA

---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

1 **Title:** NS-Forest: A machine learning method for the discovery of minimum marker gene  
2 combinations for cell type identification from single-cell RNA sequencing

3 **Authors:** Brian Aevermann<sup>1</sup>, Yun Zhang<sup>1</sup>, Mark Novotny<sup>1</sup>, Mohamed Keshk<sup>1</sup>, Trygve  
4 Bakken<sup>2</sup>, Jeremy Miller<sup>2</sup>, Rebecca Hodge<sup>2</sup>, Boudewijn Lelieveldt<sup>3,4</sup>, Ed Lein<sup>2</sup>, Richard H.  
5 Scheuermann<sup>1,5,6\*</sup>

6 **Affiliations:** <sup>1</sup>J. Craig Venter Institute, La Jolla, CA, USA; <sup>2</sup>Allen Institute for Brain  
7 Science, Seattle, WA, USA; <sup>3</sup>Department of Radiology, Leiden University Medical  
8 Center, Leiden, The Netherlands; <sup>4</sup>Department of Intelligent Systems, Delft University of  
9 Technology, Delft, The Netherlands; <sup>5</sup>University of California San Diego, La Jolla, CA,  
10 USA; <sup>6</sup>La Jolla Institute for Immunology, La Jolla, CA, USA

11 \*Contact info:  
12 Richard H. Scheuermann, Ph.D.  
13 Director, La Jolla Campus  
14 J. Craig Venter Institute  
15 4120 Capricorn Ln.  
16 La Jolla, CA 92037  
17 [rscheuermann@jcvl.org](mailto:rscheuermann@jcvl.org)  
18 858-200-1876

19  
20

21 **Abstract**

22 Single-cell genomics is rapidly advancing our knowledge of the diversity of cell  
23 phenotypes, including both cell types and cell states. Driven by single-cell/nucleus RNA  
24 sequencing (scRNA-seq), comprehensive cell atlas projects characterizing a wide range  
25 of organisms and tissues are currently underway. As a result, it is critical that the  
26 transcriptional phenotypes discovered are defined and disseminated in a consistent and  
27 concise manner. Molecular biomarkers have historically played an important role in  
28 biological research, from defining immune cell-types by surface protein expression to  
29 defining diseases by their molecular drivers. Here we describe a machine learning-based  
30 marker gene selection algorithm, NS-Forest version 2.0, which leverages the non-linear  
31 attributes of random forest feature selection and a binary expression scoring approach to  
32 discover the minimal marker gene expression combinations that optimally captures the  
33 cell type identity represented in complete scRNA-seq transcriptional profiles. The marker  
34 genes selected provide an expression barcode that serves as both a useful tool for  
35 downstream biological investigation and the necessary and sufficient characteristics for  
36 semantic cell type definition. The use of NS-Forest to identify marker genes for human  
37 brain middle temporal gyrus cell types reveals the importance of cell signaling and non-  
38 coding RNAs in neuronal cell type identity.

39

## 40 Introduction

41 Cells are the fundamental functional units of life. In multicellular organisms, different cell  
42 types play different physiological roles in the body. The identity and function of a cell - the  
43 cell phenotype - is dictated by the subset of genes/proteins expressed in that cell at any  
44 given point in time. Abnormalities in this expressed genome are disorders that form the  
45 physical basis of disease (Scheuermann et al. 2009). Thus, understanding normal and  
46 abnormal cellular phenotypes is key for diagnosing disease and identifying therapeutic  
47 targets.

48 Single cell transcriptomic technologies that measure cell transcriptional phenotypes using  
49 single cell/single nucleus RNA sequencing (scRNA-seq) are revolutionizing cell biology.  
50 The expression profiles produced by these technologies can be used to define cell types  
51 and their states based on the genes they express. For simplicity, throughout the text we  
52 will use the term “cell type” to refer to these distinct cell phenotypes that include discrete  
53 canonical cell types and distinct cell states. Numerous atlas projects designed to provide  
54 a comprehensive enumeration of normal cell types and states are currently underway,  
55 including the Human Cell Atlas (Regev et al. 2017), California Institute for Regenerative  
56 Medicine (CIRM) (Darmanis et al 2015; Enge et al. 2017; Nowakowski et al. 2017),  
57 LungMAP (Schiller et al. 2019), Pancreas atlas (Muraro et al. 2016), Heart atlas (Asp et  
58 al. 2019), and NIH Brain initiative (Mott et al. 2018). By leveraging these atlases of normal  
59 cell types defined using specimens from healthy patients as references, the role of  
60 expression deviations in disease are now being investigated (Levitin et al. 2018; Al-  
61 Dalahmah et al. 2020; Chaudhry et al. 2019).

62 Despite the incredible promise of single cell transcriptomic analysis, representations of  
63 these cell type clusters and their transcriptional phenotypes have not been adequately  
64 formalized in a standardized way to ensure effective dissemination in accordance with  
65 FAIR principles (Wilkinson et al. 2016). One approach for formalizing this type of  
66 knowledge representation and dissemination is to use the semantic framework provided  
67 by biomedical ontologies. For cell types defined by single cell transcriptomics, the Cell  
68 Ontology (CL) is an established biomedical ontology that could be used to address FAIR-  
69 compliant cell phenotype dissemination (Bard et al. 2005; Diehl et al. 2011; Meehan et

70 al. 2011; Bakken et al. 2017). With the rapid expansion in both datasets and cell types  
71 being defined using scRNA-seq, the challenge will be to make the generation of these  
72 semantic knowledge representations scalable.

73 Toward a scalable dissemination solution, we previously proposed to define cell types  
74 based on the minimum combination of necessary and sufficient features that capture cell  
75 type identity and uniquely characterize a discrete cell phenotype (Bakken et al. 2017). In  
76 the case of cell types identified by scRNA-seq experiments, these features would  
77 correspond to the combination of marker genes unique to a given gene expression cluster  
78 that provides high sensitivity and high specificity for cell type classification.

79 In this regard, determining marker gene combinations for cell type clusters is different  
80 from differential expression analysis (DE). Commonly used scRNA-seq analysis tools -  
81 Seurat (Stuart et al. 2019) and Scanpy (Wolf et al. 2018) – are often used for differential  
82 gene analysis. After cluster analysis, genes are evaluated by comparing expression in  
83 cells in a target cluster versus expression in all other cells using, for example, the  
84 Wilcoxon Rank Sum test. However, the resulting ranked set of genes cannot be used to  
85 determine the best individual marker or the best marker combinations from either the p-  
86 value rank or fold difference in expression. In contrast, marker gene determination should  
87 explicitly test for classification power and ability to discriminate a gene expression cluster  
88 of interest.

89 The ideal marker gene would show a “binary expression” pattern. These are markers that  
90 are expressed at high levels in all individual cells of a given cell type and not expressed  
91 in the cells of any other cell type. These binary expression markers are particularly useful  
92 in many downstream assays such as RT-PCR (Aevermann et al. 2021), or spatial  
93 transcriptomics where low level expression in non-target cells could be problematic.  
94 However, candidate marker genes identified by traditional differential expression analysis  
95 do not necessarily enrich for binary expression. Candidate marker genes produced by  
96 these approaches are often expressed at high levels in the target cluster and lower but  
97 measurable levels in off-target clusters. We refer to these markers as quantitative  
98 markers as their discriminatory power is derived from specific expression level thresholds,  
99 and so their utility would be dependent on the analytical sensitivity of the assay performed.

100 In other cases, a single binary marker may not be available for the cell type cluster in  
101 question, requiring the identification of marker combinations for optimal classification.

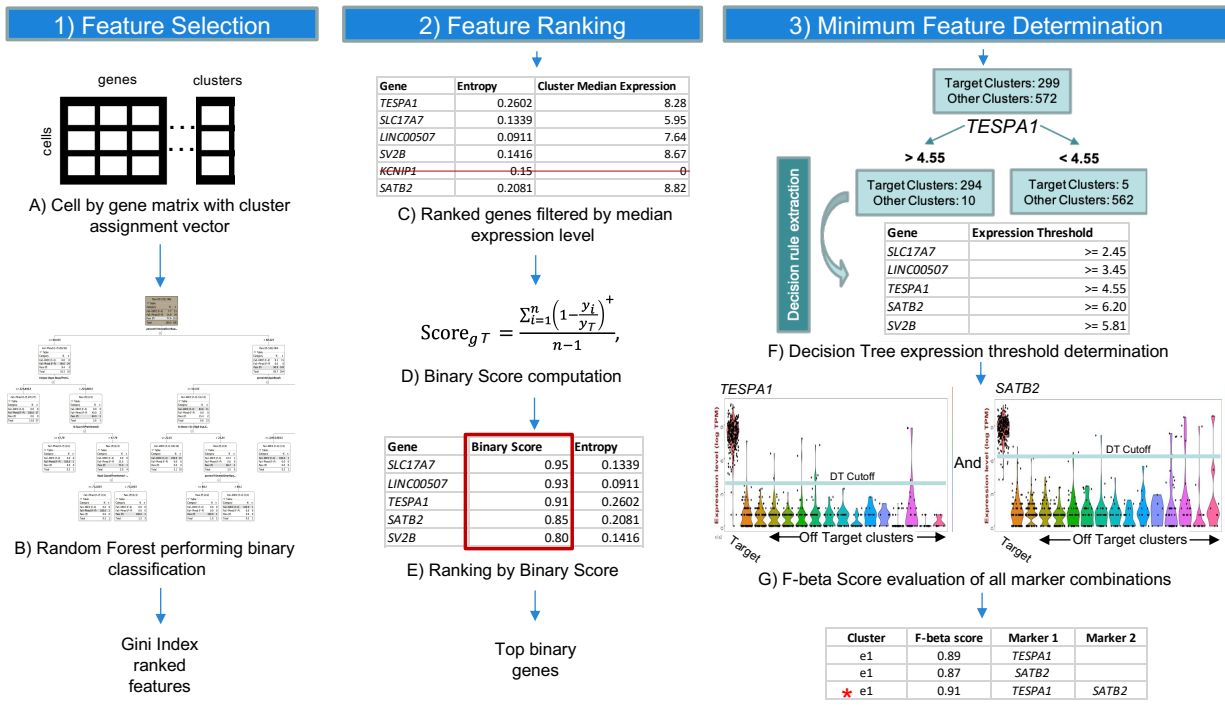
102 Here we describe Necessary and Sufficient Forest (NS-Forest) version 2.0, which  
103 improves on the simple approach to feature selection implemented in the initial version of  
104 NS-Forest (Aevermann et al. 2018). By leveraging the non-linear attributes of random  
105 forest feature selection, NS-Forest v2.0 identifies optimal combination of markers for  
106 classification while simultaneously enriching for genes with binary expression patterns.

107

## 108 Results

### 109 User driven development of NS-Forest

110 NS-Forest v2.0 was developed in close collaboration with the neuroscience user  
111 community. The primary goal was to further optimize the NS-Forest method in order to  
112 discover marker genes that can be better used for both unique cell type definition and  
113 downstream experimental investigation (**Figure 1**). In order to accomplish this, several  
114 major changes were made to NS-Forest v1.3 (**Table 1**). First, negative markers were  
115 removed by implementing a positive expression level filter (**Figure 1C**). A negative marker  
116 is defined as a gene that is not expressed in the target cluster while having expression in  
117 off-target clusters. These markers are not optimal for many downstream assays or  
118 definitional purposes. These genes are now filtered out by applying a cluster median  
119 expression threshold, with a default setting of zero.



120

121 **Figure 1: NS-Forest version 2.0 workflow** -The method begins with a cell-by-gene  
 122 expression matrix with cluster assignments for each cell (A). This clustered expression  
 123 matrix is used to generate binary classification models for each cell cluster using the  
 124 random forest machine learning method. Features are extracted from the model and  
 125 ranked by Gini Index (B). Top features are filtered by expression level to remove negative  
 126 markers (C) before being re-ranked by Binary Expression Score (D-E). Decision branch  
 127 expression level cutoffs are derived from decision tree analysis for the most binary  
 128 features (F) and F-beta score used as an objective function to evaluate the discriminatory  
 129 power of all permutations of selected markers (G).

130 **Table 1: Major changes between NS-Forest v1.3 and v2.0**

Workflow step	NS-Forest v1.3	NS-Forest v2.0
Feature Selection (Figure1A/B)	Random Forest selection of candidate features	No change
Feature Filtering (Figure1C)	None	Filtering of negative markers
Feature ranking (Figure 1D/E)	Gini index only	Gini index and Binary Expression Score reranking

Expression threshold determination (Figure1F)	Thresholds determined by median cluster expression	Thresholds determined by decision tree analysis
Minimum feature determination (Figure1G)	F1-score optimization by stepwise addition of ranked genes	F1 beta-score of all permutations of top ranked genes

131

132 Next, the way genes are ranked after random forest selection was refined. Genes  
 133 selected by random forest have an expression level threshold that is optimized to  
 134 distinguish between target and off-target clusters. Often the genes selected discriminate  
 135 based on a specific expression value resulting in quantitative expression markers. While  
 136 these quantitative markers may be good for classification, they are less useful in many  
 137 downstream biological assays. To address this issue, we modified NS-Forest v2.0 to  
 138 enrich for selection of binary expression markers. Binary expression markers are  
 139 characterized by having expression within the target cell type while being expressed at  
 140 low or negligible levels in other cell types. We accomplished this by developing a new  
 141 Binary Expression Score metric with subsequent re-ranking of the candidate markers  
 142 produced by random forest feature selection based on this score (**Figure 1D/E**).

143 Lastly, the marker gene evaluation framework was redesigned. In the initial NS-Forest  
 144 version, top ranked genes were evaluated using an unweighted F1 score in an additive  
 145 fashion. Candidate gene produced by random forest were ranked by unweighted F1 and  
 146 the top gene selected. Next the second ranked gene was added to the top ranked gene  
 147 to determine if an improvement in the F-score was obtained. This stepwise additive  
 148 process continued until the F-score plateaued or the selected number of top rank genes  
 149 were all tested.

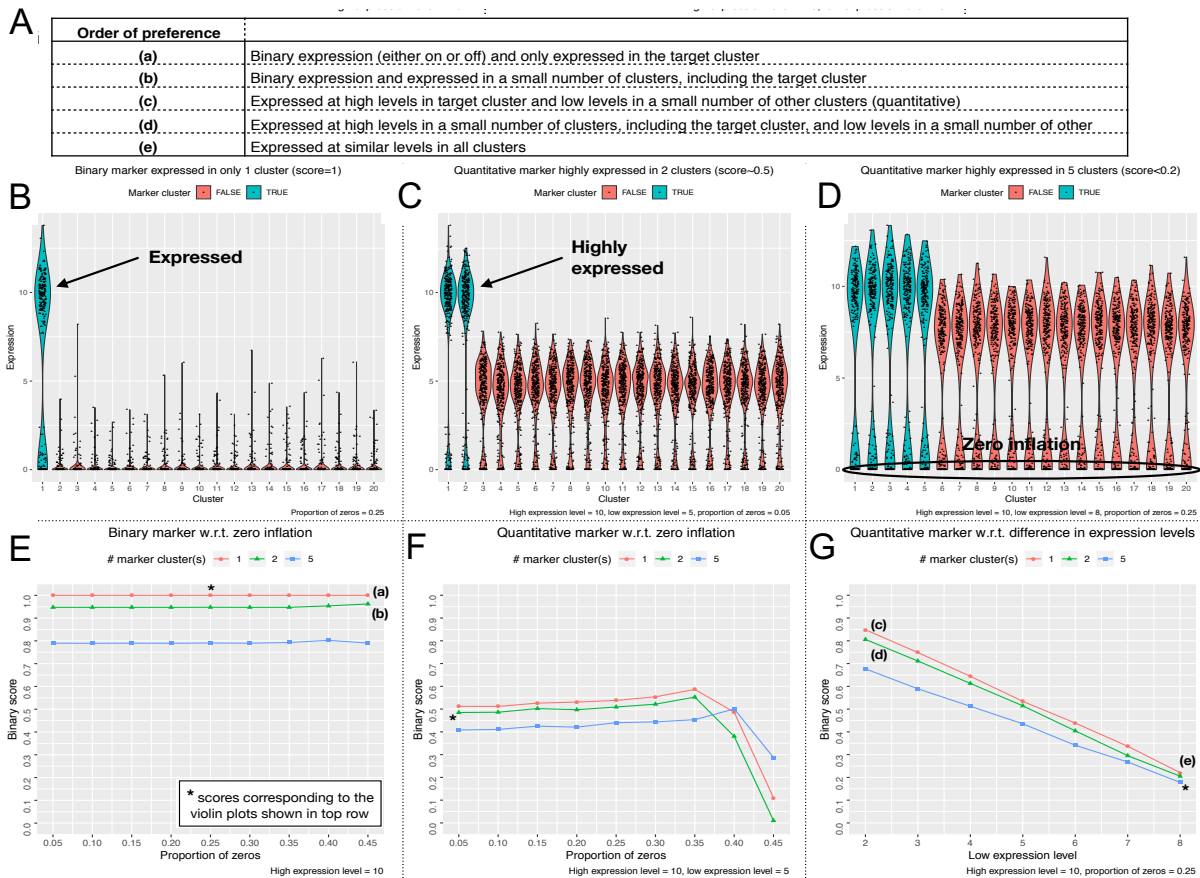
150 In NS-Forest v2.0, all permutations of the selected top ranked genes are tested and their  
 151 performance assessed using the weighted F-beta score. The F-beta score contains a  
 152 weighting term, beta, that allows for emphasizing either precision or recall. By weighting  
 153 for precision (the contributions of false positives) versus recall (the contributions of false  
 154 negatives), we limit the impact of zero inflation (a.k.a. drop-out), a known technical artifact  
 155 with scRNA-seq data, on marker gene assessment. In addition, by testing all  
 156 permutations of candidate marker genes, local optima resulting from gene ranking can be

157 avoided. These adjustments result in better final marker gene combinations given the  
158 known limitations of scRNA-seq analysis (**Figure 1F/G**).

159

## 160 Performance Testing of the Binary Expression Score Approach

161 Simulation testing of the NS-Forest Binary Expression Score was performed to evaluate  
162 the impact of different data characteristics on re-ranking behavior. First, anticipated  
163 marker gene expression patterns were themselves ranked by order of theoretical  
164 preference (**Figure 2A**). The highest preference was given to a marker gene that shows  
165 a binary expression pattern and is only expressed in the target cluster (**Figure 2A(a)/2B**).  
166 The next highest preference is given to a marker gene that shows binary expression and  
167 is only expressed in the target cluster and a limited number of off-target clusters (**Figure**  
168 **2A(b)**). This is followed by quantitative markers which have high expression in the target  
169 cluster and lower expression in off-target clusters (**Figure 2A(c)/2C**) or high expression  
170 in the target cluster and a limited number of off-target clusters (**Figure 2A(d)**). The least  
171 preferred pattern is when the marker is expressed at only slightly different levels between  
172 the target and off-target clusters (**Figure 2A(e)/2D**). The Binary Expression Score  
173 developed (see Methods section) was designed to quantify this order of expression  
174 pattern preference with a range of 0 (least desirable) to 1 (most desirable).



175

176 **Figure 2: Performance testing of Binary Expression Score** – Gene expression data  
 177 was simulated as described in the Methods section for different expression scenarios. A)  
 178 Possible marker gene expression patterns were ranked by order of preference. Panels B  
 179 – D show violin plots for three different expression scenarios: B) binary expression only  
 180 in the target cluster, C) quantitative expression with high expression in the target cluster  
 181 and one other cluster and large differences in expression in the other off-target clusters,  
 182 and D) quantitative expression with high expression in the target cluster and four other  
 183 cluster, small differences in expression in the other off-target clusters, and higher levels  
 184 of zero inflation. Panels E-G show line graphs of the full range of tested simulations from  
 185 three defined test cases: one cluster with high expression of the marker gene (red), two  
 186 clusters with high expression of the marker gene (green), and five clusters with high  
 187 expression of the marker gene (blue). E) Proportion of zeros was increased while  
 188 maintaining off target expression at zero. F) Off-target clusters were given moderate

189 levels of expression while the proportion of zeros was increased. G) Expression levels  
190 were varied in all off-target clusters from low (2) to high expression (8).

191

192 Simulations varying the binary expression pattern and level of zero inflation (**Figure 2E**)  
193 were then generated to test the performance of the Binary Expression Score developed.  
194 First, the ideal scenario of binary expression only in the target cluster produced a  
195 simulated Binary Expression Score of 1 (**Figure 2E** red). When the candidate marker  
196 gene was expressed in one (**Figure 2E** green) or four (**Figure 2E** blue) off-target clusters,  
197 the Binary Expression Score decreased to 0.95 and 0.80, respectively. These scores  
198 were robust to high zero inflation proportions, demonstrating no decrease in Binary  
199 Expression Score up to 45% zero values.

200 Next, quantitative marker expression patterns were added to the simulation (**Figures 2F**  
201 **& G**) by varying the number of off-target clusters with high expression levels and adding  
202 moderate expression to other off-target clusters. In all cases in which quantitative  
203 differences in expression were simulated, the Binary Expression Scores was reduced  
204 accordingly (**Figures 2F**). In the best case, where only the target cluster had high  
205 expression and the off-target clusters have moderate expression, the Binary Expression  
206 Score was 0.52. Further Binary Expression Score reductions were found when the high  
207 expression levels are present in additional off-target clusters. Adjusting the level of zero  
208 inflation for these scenarios showed that these Binary Expression Scores were also  
209 robust to increasing zero inflation levels until they dropped dramatically above 35% zero  
210 values.

211 Finally, simulations were performed to again test how a high-expressing marker is  
212 affected by the addition of 1 or 4 high expressing off-target clusters together with  
213 increasing expression levels in the remaining off-target clusters from low (2) to high (8)  
214 expression (**Figures 2G**). With the remaining off-target clusters held at low expression  
215 levels, these three scenarios returned high Binary Expression Scores [0.7-0.85], but these  
216 Binary Expression Scores quickly decreased with increasing levels of off-target  
217 expression. For example, when the off-target expression level was set to 6, all three high-  
218 expressing off-target scenarios returned Binary Expression Scores below 0.5. In the worst

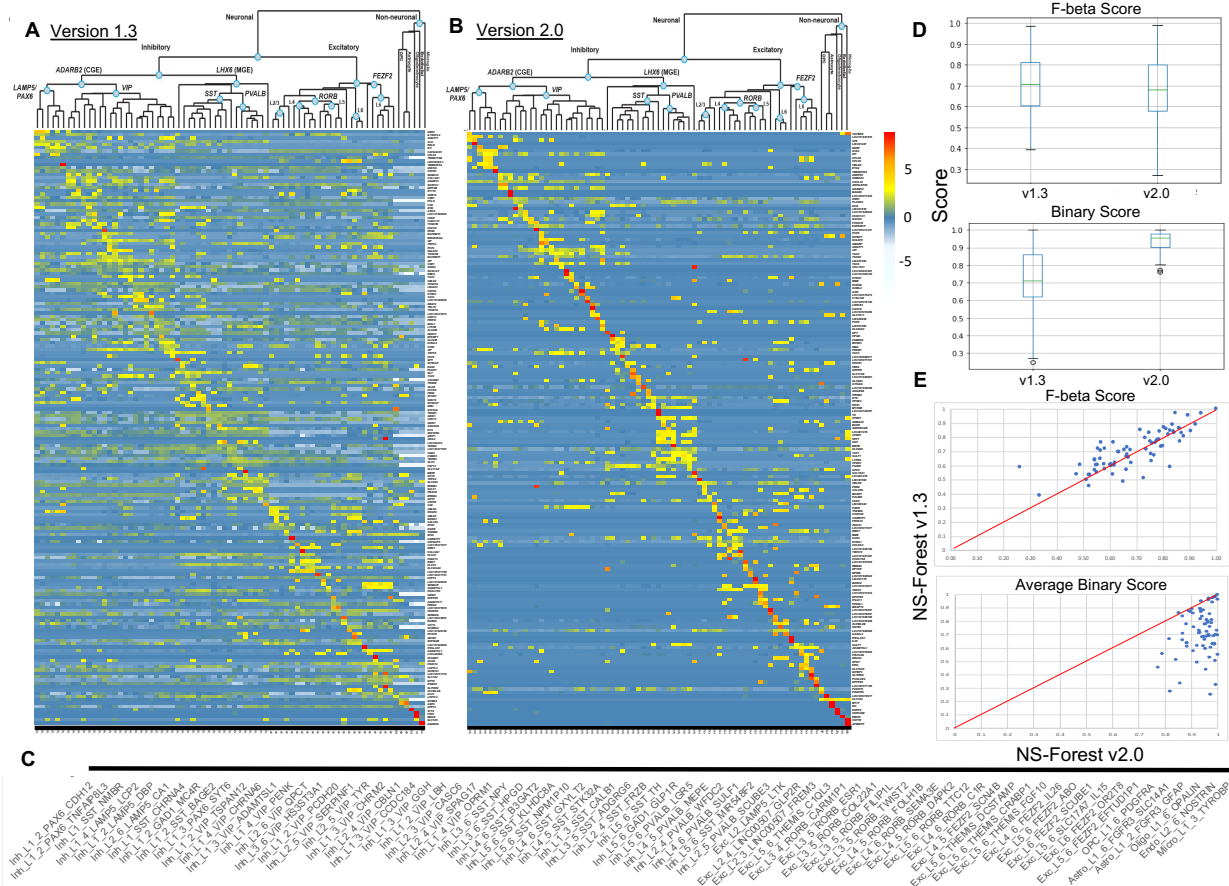
219 case, where the candidate marker had relatively high expression in all off-target clusters,  
220 the Binary Expression Score was less than 0.2.

221 These simulations demonstrate that the Binary Expression Score value produced by the  
222 algorithm recapitulates the preferred expression pattern ranking order (**Figure 2A**). In all  
223 simulations tested, the Binary Expression Scores decreased with the addition of marker  
224 expression in off-target clusters and were robust to zero inflation.

225

### 226 Marker Gene Comparison Between NS-Forest Versions

227 To evaluate the differences in results between NS-Forest v1.3 and v2.0, we analyzed  
228 marker genes selected for cell type clusters generated from single nuclei transcriptomes  
229 prepared from all cortical layers (1-6) of the human middle temporal gyrus (MTG) obtained  
230 from postmortem and surgically resected samples. For this dataset, three broad classes  
231 of cells were initially identified: excitatory neurons (10,708 nuclei), inhibitory neurons  
232 (4,297 nuclei), and non-neuronal cells (923 nuclei). The median depth of sequencing was  
233 2.6 +/- 0.5 million reads per nucleus, with a median gene detection of 9046 for neurons  
234 and 6432 for non-neuronal cells. These nuclei were clustered iteratively by first clustering  
235 into the larger groups, followed by subsequent re-clustering within each group until 75  
236 putative cell types were found (see Hodge et al. 2019 for more details on the dataset and  
237 the iterative clustering methodology). From left to right of the hierarchical clustering of  
238 clusters shown at the top of both heatmaps, there are 46 inhibitory, 23 excitatory, and 6  
239 non-neuronal cell types identified (**Figure 3**). Subsequent figures investigating these cell  
240 type clusters are ordered by these taxonomic relationships (**Figure 3C**).



241

242 **Figure 3: Comparing NS-Forest v1.3 and v2.0 marker gene sets - Heatmaps of NS-**

243 **Forest v1.3 (A) and v2.0 (B) markers from human middle temporal gyrus. The taxonomy**

244 **along the top of each heatmap is based upon the hierarchical clustering result described**

245 **in (Bakken et al. 2018; Hodge et al. 2019). Expression values are  $\log_2$  CPM cluster**

246 **medians normalized by row. The colors correspond to the normalized median expression**

247 **level for the marker gene (rows) for a given cell type cluster (columns), with high**

248 **expression (greater than five) in red and low expression (zero to negative five) in**

249 **blue/white. C) A blowup of the cell type labels corresponding to the heat map columns in**

250 **parts A and B. D) Box plots of F-beta and Binary Scores produced by NS-Forest v1.3 and**

251 **v2.0 for all 75 cell type marker gene combinations. E) Correlation of F-beta and Binary**

252 **Scores between NS-Forest v1.3 and v2.0.**

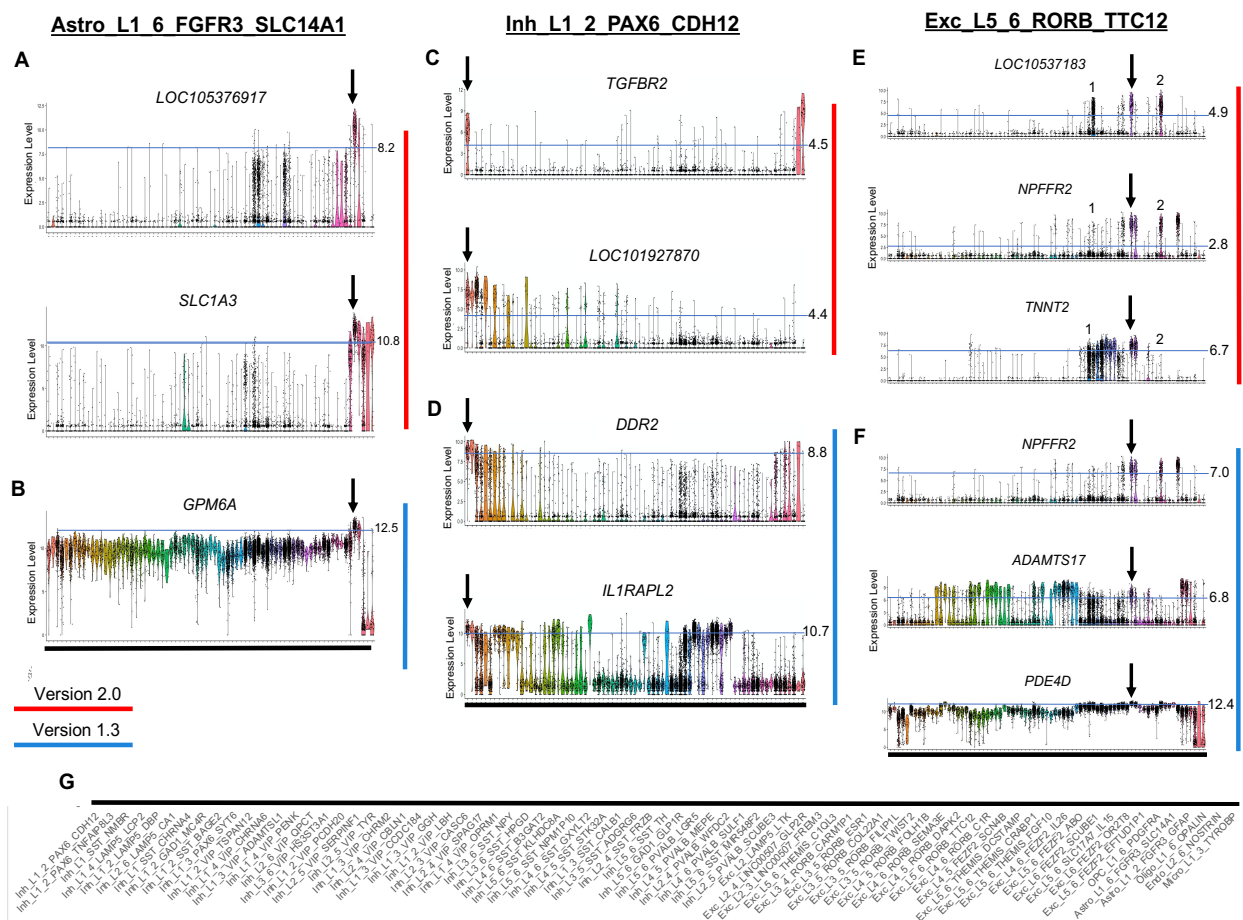
253

254 In total, 155 and 157 marker genes to optimally distinguish between these 75 different  
255 cell type classes were identified by NS-Forest v1.3 and v2.0, respectively (**Supplemental**  
256 **Tables S1-3**). Of these two unique sets of markers, 51 were common to both sets (~30%).  
257 For each method, the average number of markers per cell type was just above 2 (2.4 and  
258 2.3 respectively). This trend of cell types requiring a combination of an average of 2-3  
259 markers has been seen in other datasets and other tissue types (Aevermann et al. 2018,  
260 Aevermann et al. 2021), with cell types requiring only one marker reflecting very distinct  
261 types, such as the non-neuronal types found in this brain region. From the heatmaps  
262 (**Figure 3A/B**) it is clear that the selection of genes with binary expression patterns has  
263 dramatically improved between NS-Forest v1.3 and v2.0. The diagonal for NS-Forest v2.0  
264 contains more genes with high expression levels and, importantly, the off-diagonal  
265 expression levels are closer to zero, which demonstrates binary marker expression on a  
266 global level. Cluster median expression for markers genes, are provided in Supplemental  
267 Tables S4 and S5.

268 Given the objective of the Binary Expression Score ranking step to preferentially find  
269 marker genes with binary expression, there are tradeoffs in both the number of genes  
270 required and the classification power when compared to markers ranked strictly by  
271 importance from the random forest model in NS-Forest v1. In general, NS-Forest v2.0  
272 requires more unique genes for a given dataset. In the case of the full MTG dataset, the  
273 increase is marginal, requiring only two additional unique genes (155 vs 157 genes);  
274 similar differences in the number of marker genes required have been observed in other  
275 datasets (Aevermann et al. 2021). Furthermore, the genes that have a high Binary  
276 Expression Score are usually not the same genes that were ranked highest by Gini Index  
277 in the random forest models. This suggests that, in terms of pure classification, the  
278 markers identified by v2.0 might be expected to underperform as compare to NS-Forest  
279 v1.3. To directly compare the F scores between these two versions of NS-Forest, an  
280 additional analysis was run setting the beta weight of the F score to 0.5 in v1.3 thereby  
281 making it directly comparable to v2.0. As expected, the median F-beta Score for v2.0  
282 (0.68) was slightly lower than for v1.3 (0.71) (**Figure 3D**) and also slightly lower on a  
283 cluster-by-cluster basis (**Figure 3E**). However, the Binary Expression Scores for the v1.3  
284 markers were significantly lower – mean of 0.72 for v1.3 versus 0.94 for v2.0

285 **(Figure 3D&E)**. These results show that while adding the Binary Expression Score criteria  
286 does slightly decrease the overall classification power of the markers selected, it  
287 dramatically increases the binary expression pattern, making the markers more useful for  
288 many downstream experimental applications.

289 To demonstrate more clearly the differences between markers determined either by NS-  
290 Forest v1.3 or NS-Forest v2.0, we looked at one cell type cluster from each major group  
291 (non-neuronal, inhibitory neuron, and excitatory neuron) in the taxonomy **(Figure 4)**. For  
292 clarity, cluster labels are given along the bottom **(Figure 4G)**. The expression patterns for  
293 the astrocyte cell type Astro\_L1\_6\_FGFR3\_SLC14A1 illustrates the differences in marker  
294 gene characteristics broadly **(Figure 4A/B)**. NS-Forest v1.3 selects a single marker gene  
295 to best discriminate this cluster, whereas v2.0 selects two. NS-Forest v1.3 selects only  
296 the *GPM6A* gene which performs well at classifying this cell type along a quantitative  
297 boundary at the high  $\log_2$  expression level of 12.5, but also shows intermediate  
298 expression centered around 10 in many off-target clusters **(Figure 4A)**. Consequently,  
299 this quantitative marker is good for classification only when this small window of  
300 expression difference is discernible. In contrast, version 2.0 selects *LOC105376917* and  
301 *SLC1A3*, both of which have binary expression patterns across clusters **(Figure 4B)**.  
302 *LOC105376917* is highly expressed only in the target cluster and one additional closely  
303 related off-target cluster. Adding *SLC1A3* further improves classification by discarding  
304 cells from this off-target cluster.



305

306 **Figure 4: Marker gene expression for representative cell type clusters of the three**  
 307 **major taxonomy classes: non-neuronal, inhibitory neurons, and excitatory neurons**  
 308 **- Panels A, C, and E (red) show markers determined by NS-Forest v2.0; panels B, D, and**  
 309 **F (blue) show markers from NS-Forest v1.3. Expression levels violin plots are log<sub>2</sub> CPMs**  
 310 **with cell types enumerated along the x-axis in taxonomic order. Expression thresholds**  
 311 **are demarcated by light blue lines and cutoff values are given on the right. Thresholds for**  
 312 **NS-Forest v2.0 were determined by decision tree split points, while NS-Forest v1.3 were**  
 313 **fixed for a given gene at the expression level where 75% of cells had expression within**  
 314 **target cluster. G) Taxonomy ordered labels corresponding to the x-axis of all violin plots.**

315

316 In the case of the inhibitory neuron Inh\_L1\_2\_PAX6\_CDH12, both v1.3 and v2.0 select  
 317 two marker genes; however, their characteristics are very different (**Figure 4C/D**). NS-

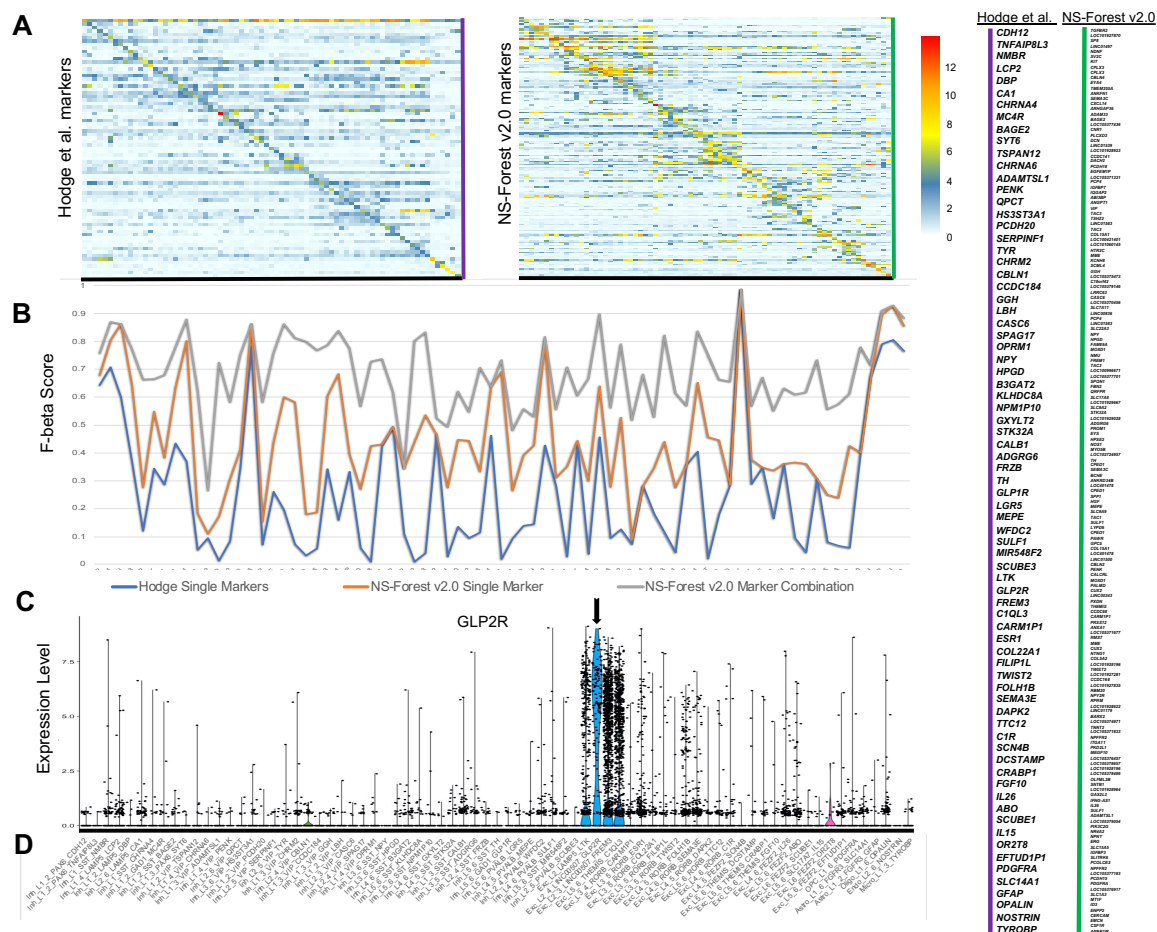
318 Forest v1.3 again found markers that classified along quantitative boundaries. *DDR2* is  
319 expressed in all the related clusters in the taxonomy and in some glial clusters at the far  
320 end of the taxonomy. The addition of *IL1RAPL2* removes the glial clusters and improves  
321 the classification; however, *IL1RAPL2* is another example of a quantitative marker as it  
322 separates the target cluster from the related cluster by narrow differences in expression.  
323 NS-Forest v2.0 selected two highly binary markers: *TGFBR2*, which is very specific to  
324 only two clusters, the target cluster and a non-neuronal type at the other end of the  
325 taxonomy. The addition of the *LOC101927870* gene eliminates cells in the non-neuronal  
326 cluster to refine the classification.

327 Lastly, the excitatory neuron Exc\_L5\_6\_RORB\_TTC12 required three markers by both NS-  
328 Forest versions to optimize classification (**Figure 4E/F**). Again NS-Forest v1.3 identified  
329 genes that used a quantitative boundary for classification whereas NS-Forest v2.0  
330 discovered binary markers. A more detailed look at these binary markers provides a clear  
331 demonstration of the combinatorics captured by NS-Forest v2.0. Within the target cluster,  
332 demarcated by the arrow, all three markers have high expression; however, the off-target  
333 excitatory clusters marked as 1 and 2 also express some but not all these markers. By  
334 leveraging the combinatorics of the three-marker combination, a highly discriminative  
335 solution is obtained. Gene *LOC105371833* is the most binary marker; however, it has  
336 high expression in a number of off-target cells in clusters 1 and 2. The addition of the  
337 *NPFFR2* gene removes most of the false positives in cluster 1, while adding the *TNNT2*  
338 gene removes the false positives from cluster 2. Together this combination of three  
339 marker genes discriminates Exc\_L5\_6\_RORB\_TTC12 from other excitatory cell types.

#### 340 Comparison with Previous MTG Marker Genes

341 To understand how the NS-Forest marker genes compare to previously published  
342 markers for the human middle temporal gyrus (MTG), we compared the NS-Forest  
343 markers to those reported in Hodge et al 2019 using a different binary expression  
344 approach used for cell cluster naming. In addition to a broad marker determined by the  
345 taxonomy and prior knowledge (such as *GAD1* or *SST*), a single marker gene per cell  
346 type cluster was assigned in Hodge *et al.* Sixteen of the seventy-five Hodge markers  
347 overlapped with the NS-Forest markers [*BAGE2*, *GGH*, *CASC6*, *NPY*, *HPGD*, *STK32A*,

348 *ADGRG6, TH, MEPE, PENK, CARM1P1, TWIST2, IL26, SULF1, ADAMTSL1, PDGFRA*.  
 349 These sixteen were spread across the taxonomy, representing cell type clusters from all  
 350 three major cell type lineages. Unscaled heatmaps of mean gene expression per cluster  
 351 for both the Hodge and NS-Forest marker sets (**Figure 5A**) demonstrate that both are  
 352 characterized by largely binary expression patterns, having a higher expression along the  
 353 diagonal versus off-diagonal. However, the Hodge markers have an overall lower mean  
 354 expression level of 4.8 log<sub>2</sub> CPM in comparison with the mean expression for the NS-  
 355 Forest markers of 7.0 log<sub>2</sub> CPM.



356

357 **Figure 5: Comparison of Hodge et al. markers with NS-Forest v2.0 for the full MTG**  
 358 - A) Unscaled heatmap for both sets of markers where the values are the mean  
 359 expression per gene. B) F-beta scores (y-axis) for the single Hodge marker gene (blue),  
 360 the best NS-Forest single marker gene (orange), and the combination of marker genes  
 361 found by NS-Forest (grey).

362 selected by the method used by Hodge et al for cluster Exc\_L2\_4\_LINC00507\_GLP2R  
363 with expression given as  $\log_2$  CPMs. For all panels, cell type clusters are listed along the  
364 x-axis in taxonomic order. D) Taxonomy ordered labels corresponding to the x-axis of the  
365 heatmaps in panel A and also the violin plot in panel C.

366

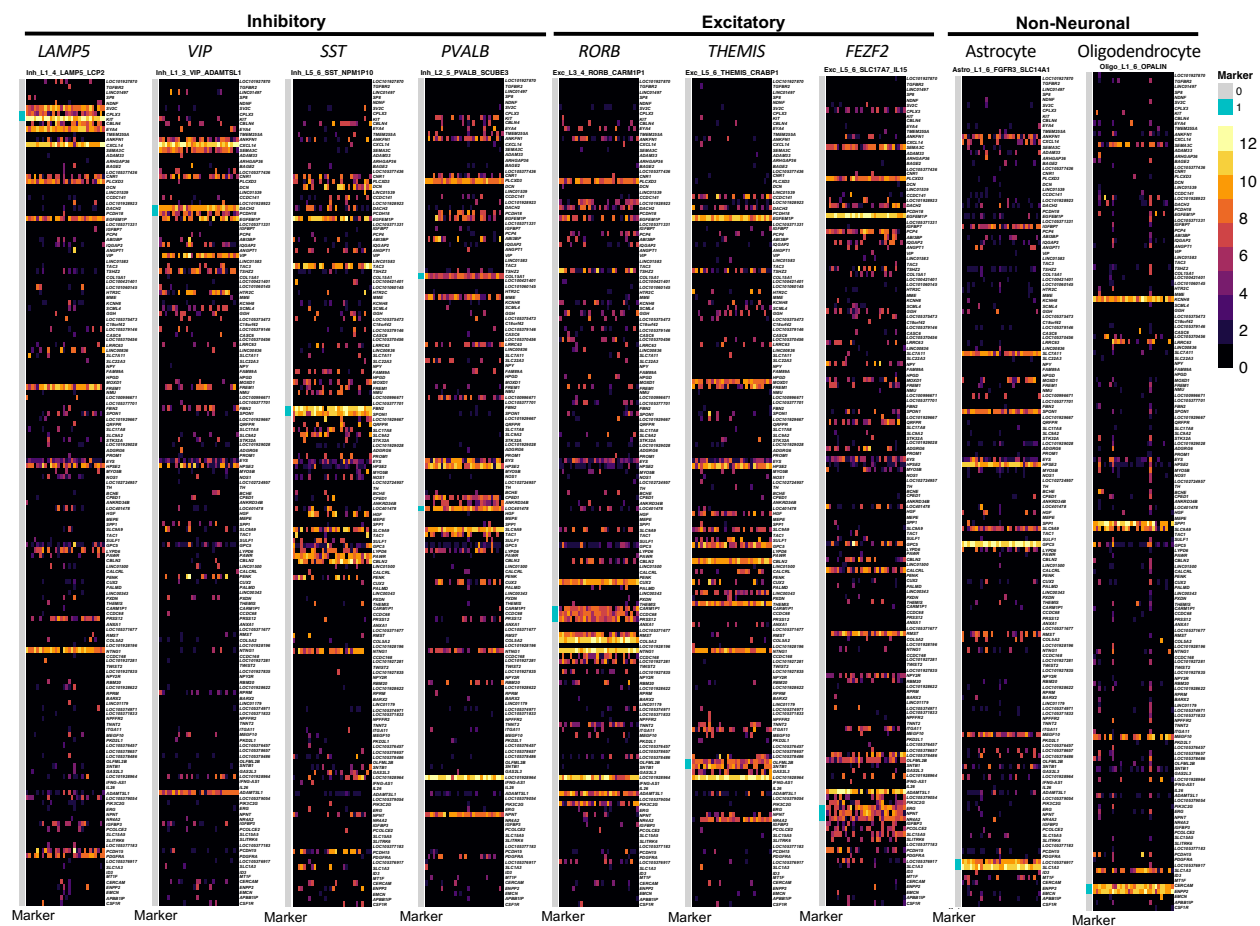
367 One major difference between these two approaches is that the Hodge marker set  
368 contains a single marker per cluster, selected to label a distinct cluster phenotype,  
369 whereas NS-Forest selects combinations of markers that optimize classification power.  
370 By running the Hodge markers through NS-forest v2.0, we estimated F-beta scores for  
371 the single Hodge markers in order to compare their classification power to the best single  
372 NS-Forest markers and the NS-Forest marker combinations (**Figure 5B**). Overall, the  
373 trend lines show that the F-beta scores for single markers (blue and orange lines) follow  
374 a similar trajectory with some clusters being more difficult to classify than others, i.e.,  
375 having lower F-beta scores. However, the NS-Forest marker combinations (grey line)  
376 provide a uniformly higher power of discrimination over either single marker, regardless  
377 of how the single best marker is chosen.

378 When evaluating the F-beta scores for the Hodge markers, it became clear that many  
379 had elevated false positive rates. To directly compare the two sets of markers, we  
380 computed the false discovery rate ( $FDR = FP / (FP + TP)$ ) for each cell type and averaged  
381 across the entire set. The Hodge markers had an average FDR of 0.7 versus 0.14 for the  
382 NS-Forest markers. *GLP2R*, which is a marker for Exc\_L2\_4\_LINC00507\_GLP2R, offers  
383 a good visual example (**Figure 5C**). This gene is expressed in the target cluster but also  
384 the nearest cell types within the *LINC00507* group. NS-Forest also has difficulty finding  
385 markers for this cell cluster phenotype, requiring 3 markers in total; however, in  
386 combination these markers help reduced the FDR rate from 0.89 to 0.11. For clarity,  
387 cluster labels for the x-axis are given along the bottom (**Figure 5D**).

## 388 NS-Forest Markers as Cell Type Barcode

389 From the complete panel of NS-Forest marker genes for a given dataset, it is possible to  
390 generate a “transcriptional barcode” for each cell type. As an illustration, barcodes

391 randomly selected nuclei from 9 different cell types representing each major subclass in  
 392 the taxonomy with the 157 NS-Forest v2.0 markers displayed as rows and individual  
 393 nuclei as columns is shown (**Figure 6**). The markers that are specific for the given cluster  
 394 are demarcated in pink within the blue bar along the left side of the barcode. The distinct  
 395 pattern of these transcriptional barcodes are clearly apparent and include not only distinct  
 396 expression of the specific marker genes in the target cells but also variable but distinct  
 397 expression patterns of marker genes from other clusters. Barcodes for all cell types within  
 398 the human MTG are provided in **Supplemental Figures S1-S9**. These barcodes can be  
 399 used as a clear visualization of a given cell type within the context of its dataset or  
 400 projected onto new datasets to demonstrate cell type similarity.

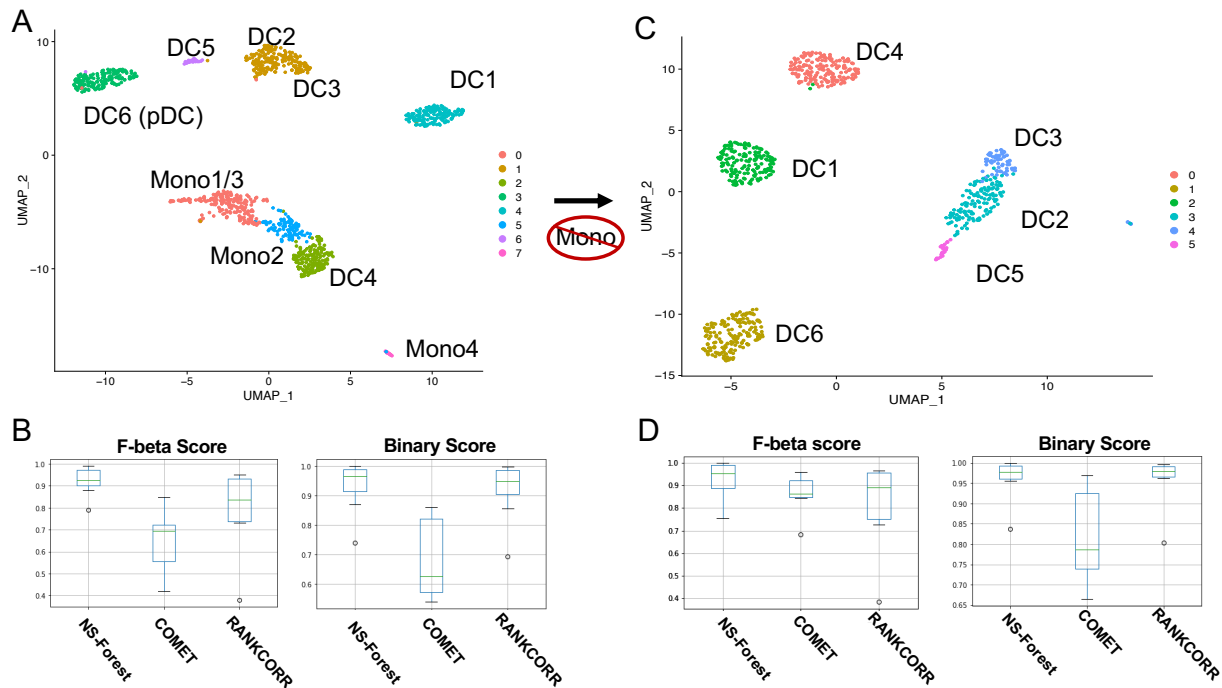


405 astrocytes, and oligodendrocytes. Each cell type is represented by 30 individual cells  
406 selected at random (columns) with the heatmap color-coded by  $\log_2$  CPM expression  
407 values for each marker gene (rows).

#### 408 Comparison with other Marker Gene Selection Approaches

409 In order to assess the performance of NS-Forest v2.0 for marker gene selection, we  
410 compared it to two other marker gene selection tools - COMET (Delaney et al. 2019) and  
411 RANKCORR (Vargo et al. 2020). These three tools were evaluated using an independent  
412 monocyte/dendritic cell dataset produced by Villani et al, 2017. For the 8 clusters  
413 produced by reprocessing these data (**Figure 7A**), NS-Forest v2.0, COMET, and  
414 RANKCORR required 17, 16, and 28 markers, respectively, to produce optimal  
415 classification results. When comparing F-beta and Binary Scores (**Figure 7B**), NS-Forest  
416 v2.0 was found to outperform both COMET and RANKCORR. While there was a  
417 significant overlap of 10 genes between NS-Forest v2.0 and RANKCORR, none were  
418 shared with COMET. Given the overlap, it is not surprising to find that both the F-beta  
419 and Binary Scores are close between NS-Forest v2.0 and RANKCORR. The median F-  
420 beta Scores were [0.92 > 0.84 > 0.69] for NS-Forest v2.0 > RANKCORR > COMET, and  
421 the median Binary Scores were [0.97 > 0.95 > 0.62] for NS-Forest v2.0 > RANKCORR >  
422 COMET.

423



424

425 **Figure 7: Results from marker gene set determination for monocyte and dendritic**  
 426 **cell types described in Villani et al. 2017 - A)** Louvain clustering result for all monocytes  
 427 and dendritic cell types with labels indicating the cell types defined in Villani 2017. In  
 428 comparison with the original result derived from iterative clustering, monocyte 1 and 3  
 429 and dendritic type DC2 and DC3 have been merged in this clustering result. B) Box plots  
 430 showing F-beta Scores and Binary Scores for markers determined for the clusters in  
 431 panel A by NS-Forest v2.0, COMET, and RANKCORR. C) Louvain clustering results of  
 432 dendritic cells only with labels indicating the cell type defined in Villani 2017. D) Box plots  
 433 showing F-beta Scores and Binary Scores for markers determined for the cell type  
 434 clusters in panel C by NS-Forest v2.0, COMET, and RANKCORR.

435

436 Clustering was also performed after removal of the monocytes, which resulted in 6  
 437 clusters corresponding to the DC1-DC6 types as characterized in the original study  
 438 (**Figure 7C**). For the 6 clusters produced by reprocessing these data, NS-Forest v2.0,  
 439 COMET, and RANKCORR required 9, 12, and 19 markers, respectively, to produce optimal  
 440 classification results. Three markers were shared by all methods and 4 markers shared

441 between NS-Forest v2.0 and COMET and between NS-Forest v2.0 and RANKCORR.  
442 Again NS-Forest v2.0 outperformed both COMET and RANKCORR, but the F-beta score  
443 results were more comparable with these clusters. The median F-beta Scores were [0.95  
444 > 0.89 > 0.86] for NS-Forest v2.0 > RANKCORR > COMET, and the average Binary  
445 Scores were [0.979 > 0.978 > 0.78] for NS-Forest v2.0 > RANKCORR > COMET.

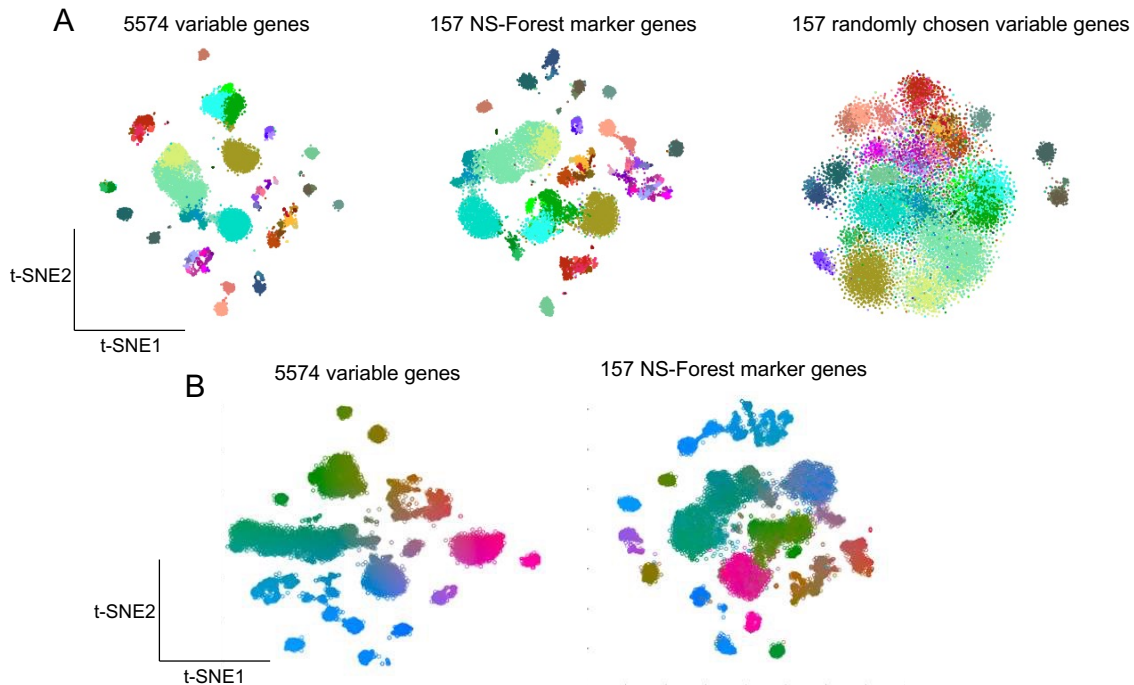
446 In general, these results show that NS-Forest v2.0 outperforms these other marker gene  
447 identification methods. However, it should be noted that these other methods were not  
448 designed for the purpose of selecting the minimum set of marker genes. COMET may  
449 have underperformed because its XL-mHG framework, that uses the X and L parameters,  
450 optimizes for the true positives and false positives, whereas NS-Forest v2.0 and its  
451 associated metrics are more focused on false positives only. And while RANKCORR  
452 performance is comparable to NS-Forest v2.0, it required substantially more marker  
453 genes for optimal performance. Thus, NS-Forest v2.0 appears to be optimal for the  
454 specific use case of finding the minimal set of markers for maximal classification  
455 accuracy. Marker genes identified by each method are given in **Supplemental Fig S10**.

456 Of the three methods, NS-Forest v2.0 was the most time intensive. Both COMBAT and  
457 RANKCORR completed in under 3 minutes when analyzing both the monocyte/dendritic  
458 cell and dendritic cells only dataset, whereas NS-Forest v2.0 took 45 minutes at the  
459 default settings. To investigate the performance of NS-Forest v2.0 further, the dendritic  
460 cell dataset was run while varying the parameters (**Supplemental Fig S11**). While the  
461 number of trees used during the random forest modeling step did not have a significant  
462 impact on the run time, the number of genes tested for all permutations was the limiting  
463 step. We found that running 1-4 genes resulted in run times under 4 minutes, using 5  
464 genes increased this to 7 minutes, and using the default of 6 genes resulted in a 45-  
465 minute runtime. Looking at the resulting F-beta scores for each run, we can see clear  
466 improvement in marker determination up to the 5 gene selection and only a slight  
467 improvement thereafter. Consequently, it may be advisable to change the default to 5  
468 genes in situations where time is a limiting factor.

469

## 470 Validation of Human MTG NS-Forest v2.0 Markers

471 The ground truth for neuron types and their marker genes in human MTG is not known  
472 as this is currently an active area of investigation. Consequently, a true biological  
473 validation of the marker genes is not possible. As an alternative, we asked the question,  
474 do the minimum set of marker genes selected by NS-Forest capture the underlying  
475 diversity of cell type identity reflected in the entire expressed transcriptome? To do this,  
476 we generated t-SNE plots using the complete set of 5574 variable genes used for the  
477 original MTG clustering, the minimum set of 157 NS-Forest v2.0 marker genes, and sets  
478 of 157 genes randomly selected from the complete variable genes list. These embeddings  
479 were then painted using the cell type assignments from the MTG taxonomy. From the t-  
480 SNE plots, it is clear that the NS-Forest markers closely recapitulate the clustering and  
481 embedding structure of the complete variable genes set, much better than the randomly  
482 selected genes (**Figure 8A**). For example, in the bottom of the complete variable genes  
483 t-SNE there are light salmon- and dark salmon-colored groups of clusters; these two  
484 clusters are nicely preserved in the right-hand side of the NS-Forest marker t-SNE,  
485 whereas in the t-SNE from the randomly selected variable genes these two clusters are  
486 spread out and a third brown cluster is now merged with light salmon cluster. Examples  
487 like this can be seen throughout the three embeddings. A more quantitative analysis of  
488 these t-SNE embeddings using the Nearest-Neighbor Preservation metric showed that  
489 both the precision and recall are higher using the 157 NS-Forest markers compared with  
490 50 sampling of 157 genes randomly selected from the variable gene set (**Supplemental**  
491 **Figure S12**).



492

493 **Figure 8: Validation of NS-Forest v2.0 MTG marker genes** - A) t-SNE plots generated  
 494 using the full 5574 variable gene list, the 157 NS-Forest v2.0 markers, and 157 genes  
 495 randomly selected from the variable genes list painted by taxonomic assignment. B) t-  
 496 SNE map generated from the full 5574 variable gene list was painted by CIELAB color  
 497 space using coordinate position for each nuclei (left). t-SNE map generated using the 157  
 498 NS-Forest markers was then painted according to the CIELAB color space established in  
 499 the complete variable genes t-SNE (right).

500

501 In addition, the local embedding structures as reflected by expression gradients within a  
 502 given t-SNE cluster, also appear to be well preserved (**Figure 8B**). The complete variable  
 503 genes t-SNE map was painted using coordinate positioning. This yields a visual way of  
 504 comparing where individual nuclei are located within the full t-SNE embedding versus  
 505 other t-SNE embeddings. The NS-Forest marker t-SNE was then painted using the colors  
 506 derived from the complete variable genes t-SNE. The fact that the same color gradients  
 507 are observed in the NS-Forest embedding demonstrates that the positional gradients, and  
 508 thus the nuclei-to-nuclei relationships, in the NS-Forest embedding closely reflect the  
 509 positional gradients in the complete variable genes t-SNE embedding. For example, in

510 the full t-SNE there is a long cluster of nuclei beginning on the left in green that extends  
511 toward the middle, moving into bluish green, and ending in purplish blue. This same  
512 cluster, with the same color gradient is preserved within the center left cluster of the NS-  
513 Forest t-SNE.

514

## 515 Characterization of NS-Forest v2.0 Markers

516 Overall, the results from NS-Forest v2.0 reflect the high quality of the data and clustering  
517 analysis; as a supervised machine learning method, NS-Forest v2.0 is reliant on the  
518 quality of the clustering results. The median number of markers required for optimal  
519 classification was 2, with only two clusters needing 4 markers, producing a mean F-beta  
520 score of 0.69. Overall, the 75 clusters required 157 unique genes to achieve optimal  
521 classification. Occasionally, marker genes are shared between clusters, with eleven  
522 genes that were not unique (*MOXD1*, *MME*, *LOC101928196*, *SULF1*, *NPFFR2*,  
523 *LINC01583*, *TAC1*, *COL15A1*, *LOC401478*, *CPED1*, *TAC3*).

524 Out of the 157 NS-Forest v2.0 marker genes, 37 (24%) were long non-coding RNAs  
525 (lncRNAs) or uncharacterized loci (LOCs). Non-coding RNAs have been previously found  
526 to be prevalent when analyzing RNA-seq data from single neuronal cells or nuclei and,  
527 surprisingly, these non-coding RNAs had higher specificity as markers when compared  
528 to coding genes (Bakken et al. 2018). In particular, lncRNAs are known to show cell line-  
529 specific expression (Djebali et al. 2012). In contrast, little is known about the LOC genes.  
530 These genes are particularly intriguing as they are highly specific to individual cell types  
531 and are likely important for their function. As such, they represent areas of unknown  
532 biology discovered by scRNA-seq and NS-Forest machine learning that warrant further  
533 investigation.

534 For the characterized marker genes, the most enriched annotations both by adjusted p-  
535 value and number of genes involved are for signaling (signal peptide, signal, secreted),  
536 including neuropeptide signaling (GO:0007218~neuropeptide) and calcium, and  
537 extracellular matrix (glycoprotein, extracellular matrix, GO:0005615~extracellular space,  
538 GO:0005578~proteinaceous extracellular matrix, GO:0030198~extracellular matrix

539 organization, GO:0005576~extracellular region, GO:0031012~extracellular matrix), and  
540 calcium (**Supplemental Table S6**). There are fewer genes annotated with specific  
541 neurological functions in the marker gene list as molecular neuroscience is a relatively  
542 nascent field. However, many of genes assessed here are known signaling peptides in  
543 other contexts and would benefit from further characterization in a neurological context.  
544 Taken together, these results suggest that specific signaling pathways and extracellular  
545 signaling molecules are key to neuronal cell type identity.

546

## 547 Discussion

548 Here we describe the development and performance of NS-Forest version 2.0, a method  
549 for the identification of cell type-specific gene expression markers from scRNA-seq data.  
550 Development was driven by user community requirements for data-driven cell type  
551 definitions that are testable in future investigations. To this end, a number of changes  
552 were made after the random forest feature selection step. In earlier versions of NS-Forest,  
553 negative markers were occasionally found. These are marker genes that are expressed  
554 in many off-target clusters but not the target cluster. Given that experimental testing for a  
555 gene that is not expressed is methodologically difficult, NS-Forest v2.0 was designed to  
556 avoid this category of markers. By implementing a median expression level cutoff greater  
557 than zero for the target cluster, all possible negative marker genes were removed. In  
558 addition, this cutoff also defines another core characteristic of NS-Forest Markers:  
559 selected marker genes must be expressed in greater than half of the individual cells within  
560 the cell type cluster.

561 In addition to negative markers, the standard random forest feature selection approach  
562 used in early NS-Forest versions discovered quantitative markers that were good for  
563 classification but problematic for further biological investigation. This limitation of random  
564 forest feature selection could be shared with other machine learning methods.  
565 Consequently, a ranking step to select marker genes with binary expression patterns was  
566 incorporated. Simulation testing performed to assess this Binary Expression Score  
567 ranking step demonstrated that marker genes with binary expression patterns were

568 preferentially selected and accurately ranked according to the levels of binary expression.  
569 As a result, NS-Forest v2.0 demonstrated clear improvement in the enrichment for binary  
570 expression patterns, with a nominal impact on the overall classification power and number  
571 of marker genes necessary. Consequently, if a user prefers the highest level of  
572 classification accuracy without the practical constraint imposed by many types of  
573 downstream investigations, NS-Forest v1.3 might be preferred. But if binary expression  
574 for downstream application is important, NF-Forest v2.0 would be the best choice. Both  
575 versions are available as official Github releases.

576 Beyond their use for defining and investigating cell types, necessary and sufficient marker  
577 genes also offer a dimensionality reduction with limited loss of fidelity to the originally  
578 clustering solution. This dimensionality reduction offers a feasible way of representing the  
579 clustering solution with a minimal amount of information, which is ideal for data  
580 dissemination. These marker genes can then be used to generate a reference  
581 knowledgebase for cell types, generating expression barcodes that can be used to  
582 identify these cell types within new datasets. Indeed, NS-Forest marker genes have been  
583 used to facilitate reference cell type matching in the FR-Match algorithm (Zhang et al.  
584 2020).

585 As mentioned above, NS-Forest markers are optimized for downstream experimental  
586 investigation. There are a number of assays for which known markers could facilitate  
587 biological investigation, such as qPCR and the burgeoning field of spatial transcriptomics  
588 based on multiplex FISH. To date, a number of projects have used NS-Forest markers  
589 for these purposes. For example, qPCR probes based on NS-Forest markers were made  
590 to detect genes in scRNA-seq libraries from myeloid dendritic cells (mDCs) FACS sorted  
591 from peripheral blood in patients treated with the Hepatitis B vaccine (Aevermann et al.  
592 2021). In a similar fashion, gene probes were designed based on NS-Forest markers for  
593 cell type detection using a number of spatial transcriptomic technologies. These  
594 technologies aim to resolve the location of cell types derived from scRNA-seq generated  
595 taxonomies within intact tissue specimens (Perkel 2019).

596 Another possible application of NS-Forest is to utilize selected gene sets of particular  
597 interest as input to produce marker gene sets designed to capture specific cell type

598 properties. For example, the input of gene sets composed of transcription factors could  
599 reveal master regulators of developmental programs (Cui et al. 2019). Input gene sets  
600 composed of neuropeptides and neurotransmitters could be used to shed new light on  
601 the specific signaling properties of different neuronal cell subsets (Smith et al. 2019). Input  
602 gene sets composed of cell surface markers could be used to identify markers for use in  
603 fluorescence-activated cell sorting.

604 As the number of experiments performed and datasets made publicly available  
605 dramatically increase, the greater biological community is left with the monumental task  
606 of integrating these data into a consensus of canonical cell types. With cell types defined  
607 by NS-Forest marker genes, we can move ahead with the creation of a dissemination  
608 framework that defines ontological classes based upon these molecular markers as the  
609 necessary and sufficient criteria in an axiomatic semantic representation compliant with  
610 FAIR principles. Ontological representations have numerous advantages over simple  
611 vocabularies, including the structuring of knowledge in a computationally readable format  
612 so that findings from many experiments can be easily accessible and “reasoning” can be  
613 performed to ensure the consistency of the representation as the knowledge rapidly  
614 grows. These instances of “cell type clusters” defined by NS-Forest markers can form the  
615 basis for the instantiation of an ontology class for adoption into the official Cell Ontology  
616 (CL). Progress is already underway in developing programmatic and scalable methods to  
617 handle the volumes of single cell data being generated. This ontological representation  
618 can address several pressing needs of the wider biological research community,  
619 producing an easy, visually accessible overview of the results of many single cell  
620 experiments in a traversable structure while preserving the hierarchical relationships  
621 inherent in a taxonomy of cell types. In addition, this ontology will provide a platform for  
622 integration with other data modalities, such as cell morphology, electrophysiology, and  
623 cell-cell interactions. A Provisional Cell Ontology (pCL) generated in this manner for  
624 human middle temporal gyrus and human, mouse, and marmoset primary motor cortex  
625 is available for exploration at <https://bioportal.bioontology.org/ontologies/PCL> .

626 Development of NS-Forest is ongoing; a number of functionalities are planned for near  
627 term release. One major update to NS-Forest v2.0 will be to add the option to run marker  
628 determination within a hierarchical framework, e.g., to determine markers for a series of

629 cluster labels that reflect a relational structure such as a taxonomy dendrogram. Another  
630 key aspect will be to include cross-validation or some other methodology to estimate the  
631 reliability of a given marker gene for a given cell type cluster. On a broader level,  
632 incorporating NS-Forest into the library of easily available Scanpy plugins is a high  
633 priority. Lastly, we will be increasing the number of output reports to facilitate the  
634 generation of ontological type artifacts, including OWL and RDF representations.

635

## 636 Methods

### 637 NS- Forest version 2.0

638 **Initial Feature Selection:** The NS-Forest v2.0 workflow (**Figure 1A/B**) begins with a cell-  
639 by-gene expression matrix, with an additional column containing cluster membership  
640 labels, produced by any expression data clustering method applied to single cell/nucleus  
641 RNA sequencing (scRNA-seq) datasets. This cluster-labelled expression matrix is then  
642 used to generate random forest classification models distinguishing each target cluster  
643 from all other clusters (binary classification) using RandomForestClassifier scikit.  
644 RandomForestClassifier hyperparameters were left at default except that the number of  
645 trees was set at 10,000 to give sufficient coverage of the sample and gene expression  
646 feature space; necessary coverage for a given feature space is estimated as the square  
647 root of the number of samples (~10,000 cells) times the square root of the number of  
648 features (~10,000 genes). From the resulting random forest model, the average Gini  
649 Index value is used to initially rank genes based on their feature importance. The output  
650 from the random forest model is a ranked list of all the input features from most informative  
651 to least informative.

652 **Feature Re-ranking Based on Positive Binary Expression:** Re-ranking the features  
653 after initial random forest ranking begins with selecting the top 15 genes ranked by Gini  
654 index. It is critical to limit the number of genes before reranking by binary expression as  
655 the Binary Expression Score does not necessarily correlate to their importance in the

656 classification context. As such, increasing the number of genes for reranking would  
 657 potentially lower the overall classification power. Positive expression filtering (**Figure 1C**)  
 658 is then performed by removing genes with a median cluster expression of 0 in order to  
 659 exclude genes that are not expressed in the relevant cluster, which we refer to as negative  
 660 markers, or show high zero inflation. The “Median\_Expression\_Level” parameter, default  
 661 value of 0, is tunable and can be adjusted according to the dataset.

662 Next, genes are re-ranked to enrich for genes with binary expression patterns (**Figure**  
 663 **1D**). A “Binary Expression Score” was developed to enrich for genes that show all-or-  
 664 none expression patterns, with expression in the target cluster and as few other cell type  
 665 clusters as possible. The Binary Expression Score is calculated for each gene in the initial  
 666 random forest feature list according to the equation:

$$\text{Score}_{gT} = \frac{\sum_{i=1}^n \left(1 - \frac{y_i}{y_T}\right)^+}{n-1},$$

667  
 668 where  $y_i$  is the median gene expression level for each cluster  $i$ ,  $y_T$  is the median expression  
 669 in the target cluster, and  $n$  is the number of clusters while  $(\cdot)^+$  denotes the non-negative  
 670 value of a real number. This results in a Binary Expression Score in the range of 0 – 1, with  
 671 a Binary Expression Score of 1 being the ideal case where the gene is only expressed in  
 672 the target cluster (**Figure 1E**). The final list of 15 genes is ranked first on the Binary  
 673 Expression Score and then by the Gini Index value. This guarantees that any genes with  
 674 Binary Expression Score ties are ranked by informativeness rather than lexicographically.  
 675

676 **Estimation of Expression Thresholds for Evaluation:** After the top genes are re-  
 677 ranked based on positive binary expression, they are then tested for their classification  
 678 power individually and in combination. First the top  $M$  genes, a tunable parameter  
 679 “Genes\_to\_testing”, set to 6 genes by default, are used to generate individual decision  
 680 trees to determine the optimal expression level cut-off value for each gene (**Figure 1F**).  
 681 The maximum leaf nodes parameter is set at two, thereby ensuring a single split point per  
 682 tree. From these trees, the optimal gene expression threshold at the split point is  
 683 extracted.

684 **Minimum Feature Combination Determination:** To evaluate the discriminative power  
 685 of a given combination of candidate marker genes, we use the F-beta score as an  
 686 objective function:

$$F_{\beta} = (1 + \beta^2) \frac{\text{precision} \cdot \text{recall}}{\beta^2 \cdot \text{precision} + \text{recall}}$$

687  
 688 The F-score is the harmonic mean of precision and recall providing equal weight for these  
 689 two classification measures. The F-beta score includes a beta term that allows for the  
 690 weighting of the function towards either precision (beta < 1) or recall (beta > 1) (**Figure**  
 691 **1G**). The beta for the analysis described here was estimated empirically at 0.5. In brief,  
 692 the empirical selection of 0.5 was based on a balance of the average values for the  
 693 confusion matrix across all cell type clusters while varying the beta parameter. At a beta  
 694 of 0.5, there was an optimum reached in the confusion matrix while averaging  
 695 approximately 2 markers per cell type cluster (**Supplemental Fig S13**). This parameter  
 696 should be evaluated for each dataset, as it adjusts for the amount of zero inflation within  
 697 the data. Here we are analyzing Smart-seq data which is known to have comparatively  
 698 lower zero inflation versus droplet-based methodologies.

699 Finally, all permutations of the top ranked genes (6 genes by default) are then evaluated  
 700 at the expression levels determined earlier by decision tree analysis. The F-beta scores  
 701 for all permutations are written to a complete results file and the gene feature combination  
 702 producing the best F-beta score result selected per cluster.

### 703 Simulation Testing of the Binary Expression Score

704 Simulation studies were conducted to investigate the properties of the Binary Expression  
 705 Score weighting using a three-component mixture model to reflect the zero-inflation  
 706 technical artifact and the background and positive expression signals in real data  
 707 distributions. Denoting  $X$  as the gene expression value, our simulated data follow a  
 708 mixture distribution:

$$709 \quad P(X = x) = \pi_1 \cdot \delta_0(x) + \pi_2 \cdot f_{\text{Gamma}}(x) + \pi_3 \cdot f_{\text{Normal}}(x)$$

710 Where  $\delta_0(x)$  is the probability density function of the degenerate distribution at 0 for the  
711 zero-inflation technical artifact,  $f_{\text{Gamma}}(x)$  is the probability density function of a Gamma  
712 distribution (with hyperparameters  $\alpha$  and  $\beta$ ) for low level background expression from off-  
713 target cells or on-target cells with low expression, and  $f_{\text{Normal}}(x)$  is the probability density  
714 function of a Normal distribution (with hyperparameters  $\mu$  and  $\sigma^2$ ) for positive expression  
715 signals; parameters  $\pi_1$ ,  $\pi_2$  and  $\pi_3$  are the corresponding mixture weights for each  
716 component such that  $\pi_1, \pi_2, \pi_3 > 0$  and  $\pi_1 + \pi_2 + \pi_3 = 1$ . In our simulations, we  
717 generated 20 clusters with 300 cells in each cluster. We designed cases where the  
718 simulated gene is expressed at high levels in 1, 2, or 5 clusters. Both binary and  
719 quantitative markers were simulated for on-target and off-target clusters by setting  
720 different parameters and hyperparameters in the mixture model.

## 721 scRNA-seq Data

722 The scRNA-seq data evaluated here were obtained from the Allen Institute for Brain  
723 Science (<https://portal.brain-map.org/atlasses-and-data/rnaseq>). Experimental design,  
724 including tissue sampling and data processing, can be found in Hodge et al. 2019. In brief,  
725 layers 1-6 of the human middle temporal gyrus (MTG) were vibratome sectioned, nuclei  
726 were extracted and labelled for NeuN expression. Nuclei were then FACS sorted and  
727 libraries were generated using the Smart-Seq4 and Nextera XT chemistries. Data  
728 processing and clustering were then performed as detailed in (Bakken et al. 2018).

729 NS-Forest v2.0 was run using the cluster assignments given in Hodge et al 2019. Nuclei  
730 not assigned to a cluster were removed from the analysis. CPM expression values were  
731  $\log_2(x+1)$  transformed and genes with a sum of zero median expression across all  
732 clusters were removed. After filtering, 15,928 nuclei and 13,946 genes remained. Given  
733 the size of the input matrix, we increased the number of trees in the random forest model  
734 from the default of ten thousand to fifty thousand.

## 735 Marker Validation

736 In order to demonstrate the preservation of the cell type clustering characteristics using  
737 NS-Forest marker genes, t-SNE embeddings were generated using Cytosplore (Höllt et

738 al. 2016; van Unen et al. 2017). The original clustering solution is represented by an  
739 embedding generated from the 5574 variable genes used for the iterative clustering  
740 originally performed (Hodge et al. 2019). Additional embeddings were made using the  
741 combined set of 157 marker genes for all cell type clusters determined by NS-Forest v2.0,  
742 and a set of 157 genes chosen at random from the original 5574 genes. Figures were  
743 generated using two different painting strategies. The first painted cells based upon the  
744 cluster assignment given in the taxonomy. The second painted using a CIELAB color  
745 space on the coordinate positioning giving a visual way of comparing the relative location  
746 of individual nuclei between the full t-SNE embedding and other t-SNE embeddings.

747 In addition, a more quantitative analysis of these t-SNE embeddings using the Nearest-  
748 Neighbor Preservation metric was performed. In brief, this is computed as follows: for  
749 each data point, the K-Nearest-Neighborhood (KNN) in the high-dimensional space is  
750 compared with the KNN in the reduced-dimensional space. Average precision/recall  
751 curves are generated by taking into account high-dimensional neighborhoods of  
752 increasing size up to  $K_{max} = 50$ . The True-Positive (TP) number is the intersection  
753 between the high-dimensional and low-dimensional neighborhoods based on the 157  
754 selected genes. The precision is computed as  $TP/K$  and the recall as  $TP/K_{max}$  (Venna  
755 et al. 2010; Ingram et al. 2015; Pezzotti et al. 2020).

## 756 Comparison to Other Marker Gene Methodologies

757 Comparisons of marker gene methodologies was performed using the monocyte/dendritic  
758 cell dataset detailed in Villani et al, 2017. This dataset was chosen because it is well  
759 characterized in the associated publication and offers a range of defined cell types that  
760 vary in their difficulty to classify. Raw data was obtained from GEO  
761 (<https://www.ncbi.nlm.nih.gov/geo/>) using accession GSE94820 and then processed  
762 using a standard Seurat analysis (Stuart et al. 2019) in two ways: first the entire dataset  
763 was processed and clustered, and second the monocytes were removed followed by  
764 processing and clustering of the dendritic cell populations only. These analyses were  
765 independent and not iterative. For both analyses, cells were filtered that had less than  
766 1000 genes and the top 2500 variable genes were selected. The complete dataset had a

767 total of 1103 cells while the dendritic cell dataset had 750 cells. After processing, the  
768 resulting datasets were analyzed by Louvain clustering and visualized by UMAP  
769 embedding.

770 Clustering assignments and expression matrices containing the top 10,000 variable  
771 genes were used to perform marker determination using NS-Forest v2.0, COMET  
772 (Delaney et al. 2019), and RANKCORR (Vargo et al. 2020). All three methods were run  
773 using default parameters with COMET being run using <http://www.cometsc.com/comet>  
774 web submission. To compare the resulting marker gene sets, NS-Forest v2.0 was used  
775 to compute the Binary Score and F-beta score for all results.

776 To benchmark NS-Forest v2.0, the dendritic only dataset was used to estimate the  
777 computations time. Two different parameters were tested: the number of trees used in  
778 the random forest model generation and the number of top genes for which all  
779 permutation were tested.

780

## 781 Software Availability

782 NS-Forest version 2.0 is available at <https://github.com/JCVenterInstitute/NSForest>  
783 under an open-source MIT license. Source code is also available with this manuscript  
784 labeled “NS\_Forest\_v2.ipynb”. Protocol is available at [protocols.io](https://protocols.io):  
785 [dx.doi.org/10.17504/protocols.io.un7evhn](https://doi.org/10.17504/protocols.io.un7evhn).

786

787

## 788 Acknowledgements

789 This work was supported by the U.S. National Institutes of Health (R21-AI122100 and  
790 U19-AI118626), the California Institute for Regenerative Medicine (GC1R-06673-B), the  
791 Wellcome Trust 208379/Z/17/Z, the Chan Zuckerberg Initiative DAF, an advised fund of  
792 the Silicon Valley Community Foundation (2018-182730), the NWO Gravitation 2019  
793 grant: BRAINSCAPES: A Roadmap from Neurogenetics to Neurobiology (NWO:  
794 024.004.012) and NWO TTW project 3DOMICS (NWO: 17126)

795

## 796 Disclosure Declaration

797 No disclosures to declare.

798

799

800 **References**

- 801 Aevermann BD, Novotny M, Bakken T, Miller JA, Diehl AD, Osumi-Sutherland D, Lasken RS, Lein  
802 ES, Scheuermann RH. 2018. Cell type discovery using single-cell transcriptomics: implications  
803 for ontological representation. *Hum Mol Genet.* 2018;27(R1):R40-R47. doi:10.1093/hmg/ddy100  
804
- 805 Aevermann BD, Shannon CP, Novotny M, Ben-Othman R, Cai B, Zhang Y, Ye JC, Kobor MS,  
806 Gladish N, Lee A, Blimke TM, Hancock RE, Llibre A, Duffy D, Koff WC, Sadarangani M, Tebbut  
807 SJ, Kollmann TR, Scheuermann RH. "Machine learning-based single cell and integrative analysis  
808 reveals that baseline mDC predisposition predicts protective Hepatitis B vaccine response"  
809 medRxiv 2021.02.22.21251864; doi: <https://doi.org/10.1101/2021.02.22.21251864>  
810
- 811 Al-Dalahmah O, Sosunov AA, Shaik A, Ofori K, Liu Y, Vonsattel JP, Adorjan I, Menon V, Goldman  
812 JE. 2020. Single-nucleus RNA-seq identifies Huntington disease astrocyte states. *Acta*  
813 *Neuropathol Commun.* Feb 18;8(1):19. doi: 10.1186/s40478-020-0880-6. PMID: 32070434;  
814 PMCID: PMC7029580.  
815
- 816 Asp M, Giacomello S, Larsson L, Wu C, Fürth D, Qian X, Wärdell E, Custodio J, Reimegård J,  
817 Salmén F, et al. 2019. A Spatiotemporal Organ-Wide Gene Expression and Cell Atlas of the  
818 Developing Human Heart. *Cell.* Dec 12;179(7):1647-1660.e19. doi: 10.1016/j.cell.2019.11.025.  
819 PMID: 31835037.  
820
- 821 Bakken T, Cowell L, Aevermann BD, Novotny M, Hodge R, Miller JA, Lee A, Chang I, McCarrison  
822 J, Pulendran B, et al. 2017. Cell type discovery and representation in the era of high-content  
823 single cell phenotyping. *BMC Bioinformatics.* Dec 21;18(Suppl 17):559. doi: 10.1186/s12859-017-  
824 1977-1. PMID: 29322913; PMCID: PMC5763450.  
825
- 826 Bakken TE, Hodge RD, Miller JA, Yao Z, Nguyen TN, Aevermann B, Barkan E, Bertagnolli D,  
827 Casper T, Dee N, et al. 2018. Single-nucleus and single-cell transcriptomes compared in matched  
828 cortical cell types. *PLoS One.* 2018;13(12):e0209648. doi:10.1371/journal.pone.0209648  
829
- 830 Bard J, Rhee SY, Ashburner M. 2005. An ontology for cell types. *Genome Biol.*; 6(2):R21.  
831 doi:10.1186/gb-2005-6-2-r21  
832
- 833 Chaudhry F, Isherwood J, Bawa T, Patel D, Gurdziel K, Lanfear DE, Ruden DM, Levy PD. 2019.  
834 Single-Cell RNA Sequencing of the Cardiovascular System: New Looks for Old Diseases. *Front*  
835 *Cardiovasc Med.* 2019;6:173. Dec 10. doi:10.3389/fcvm.2019.00173  
836
- 837 Cui Y, Zheng Y, Liu X, Yan L, Fan X, Yong J, Hu Y, Dong J, Li Q, Wu X, et al. 2019. Single-Cell  
838 Transcriptome Analysis Maps the Developmental Track of the Human Heart. *Cell Rep.*  
839 2019;26(7):1934-1950.e5. doi:10.1016/j.celrep.2019.01.079  
840

- 841 Darmanis S, Sloan SA, Zhang Y, Enge M, Caneda C, Shuer LM, Hayden Gephart MG, Barres  
842 BA, Quake SR. 2015. A survey of human brain transcriptome diversity at the single cell level. *Proc*  
843 *Natl Acad Sci U S A.* Jun 9;112(23):7285-90. PMID: 26060301; PMC: PMC4466750  
844
- 845 Delaney C, Schnell A, Cammarata LV, Yao-Smith A, Regev A, Kuchroo VK, Singer M.  
846 Combinatorial prediction of marker panels from single-cell transcriptomic data. *Mol Syst Biol.* 2019  
847 Oct;15(10):e9005. doi: 10.15252/msb.20199005. PMID: 31657111; PMCID: PMC6811728.  
848
- 849 Diehl AD, Augustine AD, Blake JA, Cowell LG, Gold ES, Gondré-Lewis TA, Masci AM, Meehan  
850 TF, Morel PA, Nijnik A, et al. 2011. Hematopoietic cell types: prototype for a revised cell ontology.  
851 *J Biomed Inform.*; 44(1):75-79. doi:10.1016/j.jbi.2010.01.006  
852
- 853 Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W,  
854 Schlesinger F, et al. 2012. Landscape of transcription in human cells. *Nature.* 2012; 489: 101–8.  
855
- 856 Enge M, Arda HE, Mignardi M, Beausang J, Bottino R, Kim SK, Quake SR. 2017. Single-Cell  
857 Analysis of Human Pancreas Reveals Transcriptional Signatures of Aging and Somatic Mutation  
858 Patterns. *Cell.* Oct 5;171(2):321-330.e14. PMID: 28965763; PMC: PMC6047899
- 859 Höllt T, Pezzotti N, van Unen V, Koning F, Eisemann E, Lelieveldt B, Vilanova A.  
860 2016. Cytosplore: Interactive Immune Cell Phenotyping for Large Single-Cell Datasets. *Computer*  
861 *Graphics Forum (Proceedings of EuroVis),* 35(3): pp. 171—180.
- 862 Hodge RD, Bakken TE, Miller JA, Smith KA, Barkan ER, Graybuck LT, Close JL, Long B,  
863 Johansen N, Penn O, et al. 2019. Conserved cell types with divergent features in human versus  
864 mouse cortex. *Nature.* 2019;573(7772):61-68. doi:10.1038/s41586-019-1506-7  
865
- 866 Huang DW, Sherman BT, Lempicki RA. 2009. Bioinformatics enrichment tools: paths toward the  
867 comprehensive functional analysis of large gene lists. *Nucleic Acids Res.* 2009;37(1):1-13.  
868
- 869 Ingram S, Munzner T. 2015. Dimensionality reduction for documents with nearest neighbor  
870 queries. *Neurocomputing* 150 (2015), 557–569. 8  
871
- 872 Krishnaswami SR, Grindberg RV, Novotny M, Venepally P, Lacar B, Bhutani K, Linker SB, Pham  
873 S, Erwin JA, Miller JA, et al. 2016. Using single nuclei for RNA-seq to capture the transcriptome  
874 of postmortem neurons. *Nat Protoc*;11(3):499-524. doi:10.1038/nprot.2016.015  
875
- 876 Kuby J, Kindt TJ, Goldsby RA, Osborne BA. 2007. *Kuby Immunology.* San Francisco: W.H.  
877 Freeman. ISBN 1-4292-0211-4.  
878
- 879 Levitin HM, Yuan J, Sims PA. 2018. Single-Cell Transcriptomic Analysis of Tumor Heterogeneity.  
880 *Trends Cancer.* Apr;4(4):264-268. doi: 10.1016/j.trecan.2018.02.003. Epub 2018 Mar 9. PMID:  
881 29606308; PMCID: PMC5993208.  
882

- 883 Meehan TF, Masci AM, Abdulla A, Cowell LG, Blake JA, Mungall CJ, Diehl AD. 2011. Logical  
884 development of the cell ontology. *BMC Bioinformatics*. 2011;12:6. doi:10.1186/1471-2105-12-6.  
885
- 886 Moffitt JR, Hao J, Wang G, Chen KH, Babcock HP, Zhuang X. 2016. High-throughput single-cell  
887 gene-expression profiling with multiplexed error-robust fluorescence in situ hybridization. *Proc*  
888 *Natl Acad Sci U S A*. 2016;113(39):11046-11051. doi:10.1073/pnas.1612826113  
889
- 890 Mott MC, Gordon JA, Koroshetz WJ. 2018. The NIH BRAIN Initiative: Advancing  
891 neurotechnologies, integrating disciplines. *PLoS Biol*. Nov 26;16(11):e3000066. doi:  
892 10.1371/journal.pbio.3000066. PMID: 30475794; PMCID: PMC6283590.  
893
- 894 Muraro MJ, Dharmadhikari G, Grün D, Groen N, Dielen T, Jansen E, van Gurp L, Engelse MA,  
895 Carlotti F, de Koning EJ, van Oudenaarden A. 2016. A Single-Cell Transcriptome Atlas of the  
896 Human Pancreas. *Cell Syst*. Oct 26;3(4):385-394.e3. doi: 10.1016/j.cels.2016.09.002. Epub 2016  
897 Sep 29. PMID: 27693023; PMCID: PMC5092539.  
898
- 899 Nowakowski TJ, Bhaduri A, Pollen AA, Alvarado B, Mostajo-Radji MA, Di Lullo E, Haeussler M,  
900 Sandoval-Espinosa C, Liu SJ, Velmeshev D, et al. 2017. Spatiotemporal gene expression  
901 trajectories reveal developmental hierarchies of the human cortex. *Science*. Dec  
902 8;358(6368):1318-1323. PMID: 29217575; PMC: PMC5991609  
903
- 904 Perkel JM. 2019. Starfish enterprise: finding RNA patterns in single cells. *Nature*.  
905 2019;572(7770):549-551. doi:10.1038/d41586-019-02477-9  
906
- 907 Pezzotti N, Thijssen J, Mordvintsev A, Höllt T, Lew BV, Lelieveldt BPF, Eisemann E, Vilanova  
908 A. 2020. GPGPU Linear Complexity t-SNE Optimization. *IEEE Transactions on Visualization and*  
909 *Computer Graphics*, vol. 26, no. 1, pp. 1172-1181, Jan. 2020, doi: 10.1109/TVCG.2019.2934307.  
910
- 911 Picelli S, Björklund ÅK, Faridani OR, Sagasser S, Winberg G, Sandberg R. 2013. Smart-seq2 for  
912 sensitive full-length transcriptome profiling in single cells. *Nat Methods*;10(11):1096-1098.  
913 doi:10.1038/nmeth.2639  
914
- 915 Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E, Bodenmiller B, Campbell P,  
916 Carninci P, Clatworthy M, Clevers H, et al. 2017. The Human Cell Atlas. *Elife*. 5;6:e27041. doi:  
917 10.7554/eLife.27041. PMID: 29206104; PMCID: PMC5762154.  
918
- 919 Scheuermann RH, Ceusters W, Smith B. 2009. Toward an Ontological Treatment of Disease and  
920 Diagnosis Summit on Translational Bioinformatics 2009:116-120. PMID: 21347182; PMCID:  
921 PMC3041577.  
922
- 923 Scheuermann RH, Novotny M, Aevermann B, Ben-Othman R, Liu A, Sadarangani M, Kollmann  
924 T. Differential abundance of mDC subsets predict response to Hepatitis B vaccination. *J Immunol*  
925 May 1, 2018, 200 (1 Supplement) 180.1.  
926

927 Schiller HB, Montoro DT, Simon LM, Rawlins EL, Meyer KB, Strunz M, Vieira Braga FA, Timens  
928 W, Koppelman GH, Budinger GRS, et al. 2019. The Human Lung Cell Atlas: A High-Resolution  
929 Reference Map of the Human Lung in Health and Disease. *Am J Respir Cell Mol Biol.*  
930 Jul;61(1):31-41. doi: 10.1165/rcmb.2018-0416TR. PMID: 30995076; PMCID: PMC6604220.  
931

932 Smith SJ, Sümbül U, Graybuck LT, Collman F, Seshamani S, Gala R, Gliko O, Elabbady L, Miller  
933 JA, Bakken TE, et al. 2019. Single-cell transcriptomic evidence for dense intracortical  
934 neuropeptide networks. *Elife.* 2019;8:e47889. doi:10.7554/eLife.47889  
935

936 Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM 3rd, Hao Y, Stoeckius M,  
937 Smibert P, Satija R. 2019. Comprehensive Integration of Single-Cell Data. *Cell*, **177**, 1888-1902.  
938 doi: 10.1016/j.cell.2019.05.031.  
939

940 Tasic B, Yao Z, Graybuck LT, Smith KA, Nguyen TN, Bertagnolli D, Goldy J, Garren E, Economo  
941 MN, Viswanathan S, et al. 2018. Shared and distinct transcriptomic cell types across neocortical  
942 areas. *Nature.* 2018;563(7729):72-78. doi:10.1038/s41586-018-0654-5  
943

944 van Unen V, Höllt T, Pezzotti N, Li N, Reinders MJT, Eisemann E, Koning F, Vilanova A, Lelieveldt  
945 BPF. 2017. Visual analysis of mass cytometry data by hierarchical stochastic neighbour  
946 embedding reveals rare cell types. *Nat Commun.* 2017 Nov 23;8(1):1740. doi: 10.1038/s41467-  
947 017-01689-9. PMID: 29170529; PMCID: PMC5700955.  
948

949 Vargo AHS, Gilbert AC. A rank-based marker selection method for high throughput scRNA-seq  
950 data. *BMC Bioinformatics.* 2020 Oct 23;21(1):477. doi: 10.1186/s12859-020-03641-z. PMID:  
951 33097004; PMCID: PMC7585212.  
952

953 Venna J, Peltonen J., Nybo K, Aidos H, Kaski S. 2010. Information retrieval perspective to  
954 nonlinear dimensionality reduction for data visualization. *The Journal of Machine Learning*  
955 *Research* 11 (2010), 451–490. 8  
956

957 Villani AC, Satija R, Reynolds G, Sarkizova S, Shekhar K, Fletcher J, Griesbeck M, Butler A,  
958 Zheng S, Lazo S, Jardine L, Dixon D, Stephenson E, Nilsson E, Grundberg I, McDonald D, Filby  
959 A, Li W, De Jager PL, Rozenblatt-Rosen O, Lane AA, Haniffa M, Regev A, Hacohen N. Single-  
960 cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors.  
961 *Science.* 2017 Apr 21;356(6335):eaah4573. doi: 10.1126/science.aah4573. PMID: 28428369;  
962 PMCID: PMC5775029.  
963

964 Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, Blomberg N, Boiten  
965 JW, da Silva Santos LB, Bourne PE, et al. 2016. The FAIR Guiding Principles for scientific data  
966 management and stewardship. *Sci Data* 3, 160018. <https://doi.org/10.1038/sdata.2016.18>  
967

968 Wolf F, Angerer P, Theis F. 2018. SCANPY: large-scale single-cell gene expression data  
969 analysis. *Genome Biol* **19**, 15. <https://doi.org/10.1186/s13059-017-1382-0>  
970

971 Zhang Y, Aevermann BD, Bakken TE, Miller JA, Hodge RD, Lein ES, Scheuermann RH. FR-  
972 Match: robust matching of cell type clusters from single cell RNA sequencing data using the  
973 Friedman-Rafsky non-parametric test. *Brief Bioinform.* 2020 Nov 30:bbaa339. doi:  
974 10.1093/bib/bbaa339. Epub ahead of print. PMID: 33249453.  
975  
976 Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, Ziraldo SB, Wheeler TD,  
977 McDermott GP, Zhu J, et al. 2017. Massively parallel digital transcriptional profiling of single cells.  
978 *Nat Commun* 8, 14049. <https://doi.org/10.1038/ncomms1404>  
979