

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

3 **Authors:** Brian Aeversmann¹, Yun Zhang¹, Mark Novotny¹, Trygve Bakken², Jeremy
4 Miller², Rebecca Hodge², Boudewijn Lelieveldt^{3,4}, Ed Lein², Richard H.
5 Scheuermann^{1,4,5}

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

11 **Abstract**

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

29

30 Introduction

31 Cells are the fundamental functional units of life. In multicellular organisms, different
32 cell types play different physiological roles in the body. The identity and function of a
33 cell - the cell phenotype - is dictated by the subset of genes/proteins expressed in that
34 cell at any given point in time. Abnormalities in this expressed genome are disorders
35 that form the physical basis of disease (1) Thus, understanding normal and abnormal
36 cellular phenotypes is key for diagnosing disease and for identifying therapeutic targets.

37 Single cell transcriptomic technologies that measure cell phenotypes using single
38 cell/single nucleus RNA sequencing (scRNA-seq) are revolutionizing cellular biology.
39 The expression profiles produced by these technologies can be used to define cell
40 types and their states. Numerous atlas projects designed to provide a comprehensive
41 enumeration of normal cell types and states are currently underway, including the
42 Human Cell Atlas (2), California Institute for Regenerative Medicine (CIRM) (3-5),
43 LungMAP (6), Pancreas atlas (7), Heart atlas (8), and NIH Brain initiative (9). By
44 leveraging these atlases of normal cell types defined from healthy patients as
45 references, the role of expression deviations in disease are being investigated (10-12).

46 These projects rely primarily upon two scRNA-seq single cell technologies: Droplet
47 based technologies (13) or FACS sorting followed by Smart-Seq library preparation
48 (14). In the standard data processing pipeline, the raw sequencing reads are processed
49 using reference-based alignment, transcript reconstruction, and expression level
50 estimation. From the expression matrix, the typical downstream analysis workflows
51 produce a set of gene expression data clusters representing cells/nuclei with similar
52 expression patterns. We interpret these distinct transcriptional profiles to represent
53 distinct cell phenotypes, which include canonical cell types and distinct cell states, that
54 have achieved a state of equilibrium. This is in contrast to transitional trajectories in
55 developing tissues and populations in the process of responding to perturbations that
56 can show gradual expression pattern changes between discrete cell phenotypes (15-
57 17).

58 Despite the incredible promise of single cell transcriptomic analysis for identifying and
59 quantifying known and novel cell types, the cell type clusters and their transcriptional

60 phenotypes are not being formalized in a standardized way to ensure dissemination is
61 in accordance with FAIR principles (18). One approach for formalizing knowledge
62 representation and dissemination is to use the semantic framework provided by
63 biomedical ontologies. For cell phenotypes defined by single cell transcriptomics, the
64 Cell Ontology (CL) is an established biomedical ontology already applying ontological
65 methodologies that could be used to address FAIR-compliant cell phenotype
66 dissemination (19-22). With the rapid expansion in both datasets and cell phenotypes
67 being defined using scRNA-seq, the challenge will be to make the generation of these
68 semantic knowledge representations scalable.

69 To develop this scalable dissemination solution, we propose to define cell type
70 phenotypes based on the minimum combination of necessary and sufficient features
71 that capture cell type identity and uniquely characterize a discrete cell phenotype. In
72 the case of cell types identified by scRNA-seq experiments, these features correspond
73 to the combination of differentially expressed marker genes unique to a given gene
74 expression cluster. Historically, marker gene expression, especially at the protein level,
75 has been an essential tool to connect cell type identity with defining cell type functional
76 characteristics. For example, the classical markers CD19 and CD3 have been used
77 extensively to differentiate between B cells and T cells, while within the T cell population
78 CD3 and CD4 and CD8 are used to further separate helper and cytotoxic types (23). In
79 neurology, SLC17A7 and GAD1 are well-known markers for excitatory (glutamatergic)
80 and inhibitory (GABAergic) neuron types, respectively (24).

81 In the case of scRNA-seq expression clusters, an optimal cell type marker would be a
82 cellular feature or unique combination of features that provides high sensitivity and high
83 specificity for cell type classification. The ideal scenario would be to have a marker gene
84 that is expressed at high levels in all individual cells of a given cell type and not
85 expressed at all in any cell of any other cell type. We refer to this phenomenon as a
86 binary expression pattern. Finding markers with this binary expression pattern can be
87 quite useful for downstream experimental validation and investigation using
88 technologies such as multiplex fluorescence in situ hybridization (mFISH), quantitative
89 PCR, or flow cytometry. However, in complex tissues this ideal scenario is rarely
90 observed. Candidate marker genes are often expressed at high levels in the target

91 cluster and lower levels in off-target clusters. We refer to these markers as quantitative
92 markers as their discriminatory power is derived from specific expression level cutoffs,
93 which are dependent on the sensitivity of the assay being performed. Or a single binary
94 marker may be expressed in multiple related cell types.

95 Though similar in concept, determining markers from cell type clusters is different from
96 differential expression analysis (DE) in that DE analysis evaluates each gene for
97 expression level variation between groups, whereas marker genes are tested for their
98 classification power. The most common scRNA-seq analysis tools - Seurat (16) and
99 Scanpy (17) – handle gene selection in a similar fashion. After cluster analysis, genes
100 are evaluated by comparing expression in cells in a target cluster versus expression in
101 all other cells using, for example, the Wilcoxon Rank Sum test, which produces a gene
102 list that can be ranked by adjusted p-value. However, while the best marker genes for
103 discriminating and defining a cell type cluster are often found among the differentially
104 expressed genes, their utility in defining a cell type is not apparent from either their p-
105 value rank nor their fold difference in expression. Furthermore, DE analysis tests genes
106 individually, while defining cell types often requires the combined contribution of sets of
107 marker genes.

108 Here we describe Necessary and Sufficient Forest (NS-Forest) version 2.0, which
109 leverages the non-linear attributes of random forest feature selection and Binary
110 Expression Score ranking to discover marker gene combinations that can be used to
111 both define cell type phenotypes and in downstream biological investigations. The initial
112 version of NS-forest was based on a simple approach in which feature selection by
113 Random Forest machine learning was used to discover potential marker genes (21).
114 Here we describe user community-driven improvements upon this original methodology.
115 NS-Forest version 2.0 is available at <https://github.com/JCVenterInstitute/NSForest>
116 under an open source MIT license.

117 Results

118 User driven development of NS-Forest

119 NS-Forest version 2.0 was developed in close collaboration with the brain cell user
120 community. The primary goal was to optimize the NS-Forest method in order to discover
121 marker genes that can both uniquely define the cell type phenotype and aid in their
122 downstream experimental investigation. In order to accomplish this, several major
123 changes were made to NS-Forest version 1 (**Table 1**). First, negative markers were
124 removed by implementing a positive expression level filter. A negative marker is not
125 expressed in the target cluster while having expression in off-target clusters. These
126 markers are not optimal for many downstream assays or definitional purposes. These
127 genes are now filtered out by applying a cluster median expression threshold. The
128 default setting is zero; however, this can be changed to enrich for genes at varying
129 expression levels (**Figure 1C**).

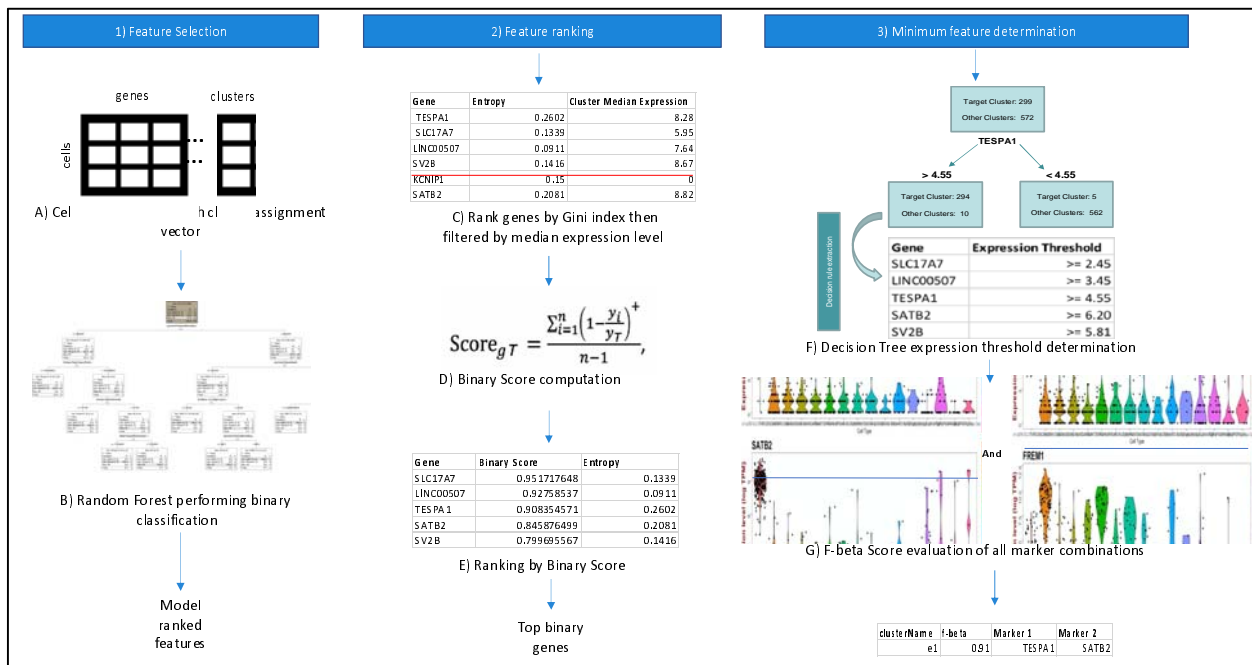
130 **Table 1:** Major changes between NS-Forest version 1.3 and version 2.0

| Workflow step | NS-Forest v1.3 | NS-Forest v2.0 |
|--|--|--|
| Feature Selection (Figure1A/B) | Random Forest (RF) driven feature selection | No changes |
| Feature Ranking (Figure1C) | Used RF ranking | Filtering of negative markers |
| Feature Ranking (Figure1D/E) | NA | Binary score reranking |
| Minimum feature determination (Figure1F) | Gene expression criteria determined by median cluster expression | Gene expression criteria determined by Decision Tree split |
| Minimum feature determination (Figure1G) | Calculate F-score | Calculate Fbeta-score |

131
132 Next, the way genes are ranked after Random Forest (RF) selection was refined. Genes
133 selected by RF have an expression level split point between target and off-target
134 clusters. Often the genes selected discriminate based on a specific value of expression
135 resulting in quantitative expression markers. As will be demonstrated below, these
136 quantitative markers are good for classification but are less useful in many downstream
137 biological assays. To address this issue, we optimized version 2.0 for the selection of
138 binary expression markers. Binary expression markers are characterized by having
139 expression within the target cell type while being expressed at low or negligible levels in
140 other cell types. We accomplished this by developing a new Binary Expression Score
141 metric with subsequent re-ranking based on this score after random forest feature
142 determination (**Figure 1D/E**).

143 Lastly, the marker gene evaluation framework was redesigned. In early NS-Forest
144 versions, the top ranked genes were evaluated by unweighted F-1 score in an additive

145 fashion. First, each of the top genes were tested individually and the best gene then
 146 removed from the list. Next all pairings with the previously determined gene were tested
 147 to find an improvement over the previous individual gene F-score. This additive process
 148 continued until the F-score plateaued or the selected number of top rank genes were all
 149 tested. In the new version of NS-Forest, all combinations of the selected top ranked
 150 genes are tested by weighted F-beta score. The F-beta score contains a weighting
 151 term, beta, that allows for emphasizing either precision or recall. By weighting toward
 152 precision, zero inflation (drop-out) can be controlled, which is a known technical artifact
 153 with scRNA-seq data. These adjustments result in better final marker gene
 154 combinations given the known limitations of scRNA-seq analysis (**Figure 1F/G**).



155

156 **Figure 1:** NS-Forest version 2.0 workflow. The method begins with a cell by gene matrix
 157 and cluster assignments. These are used to generate binary classification models using
 158 Random Forest. Features are extracted from the model and ranked by Gini index. Top
 159 features are filtered by expression level before being ranked by binary expression.
 160 Decision point expression level cutoffs are then determined for the most binary features
 161 and F-beta score used as an objective function to evaluate the discriminatory power of
 162 all combinations.

163

164 Performance Testing of Binary Scoring Approach

165 Simulation testing of the NS-Forest Binary Expression Score was performed to evaluate
166 re-ranking behavior. First, anticipated marker gene expression patterns were
167 themselves ranked by order of hypothetical preference (**Figure 2A**). The highest
168 preference was given to a marker gene that shows a binary expression pattern and is
169 only expressed in the target cluster (**Figure 2A(a)/2B**). Next, preference is given to a
170 marker gene that shows binary expression and is only expressed in the target cluster
171 and a limited number of off-target clusters (**Figure 2A(b)**). This is followed by
172 quantitative markers which have a high expression in the target cluster and lower
173 expression in off-target clusters (**Figure 2A(c)/2C**) or high expression in the target
174 cluster and a limited number of off-target clusters (**Figure 2A(d)**). The least preferred
175 pattern is when the marker is expressed at only slightly different levels between the
176 target and off-target clusters (**Figure 2A(e)/2D**).

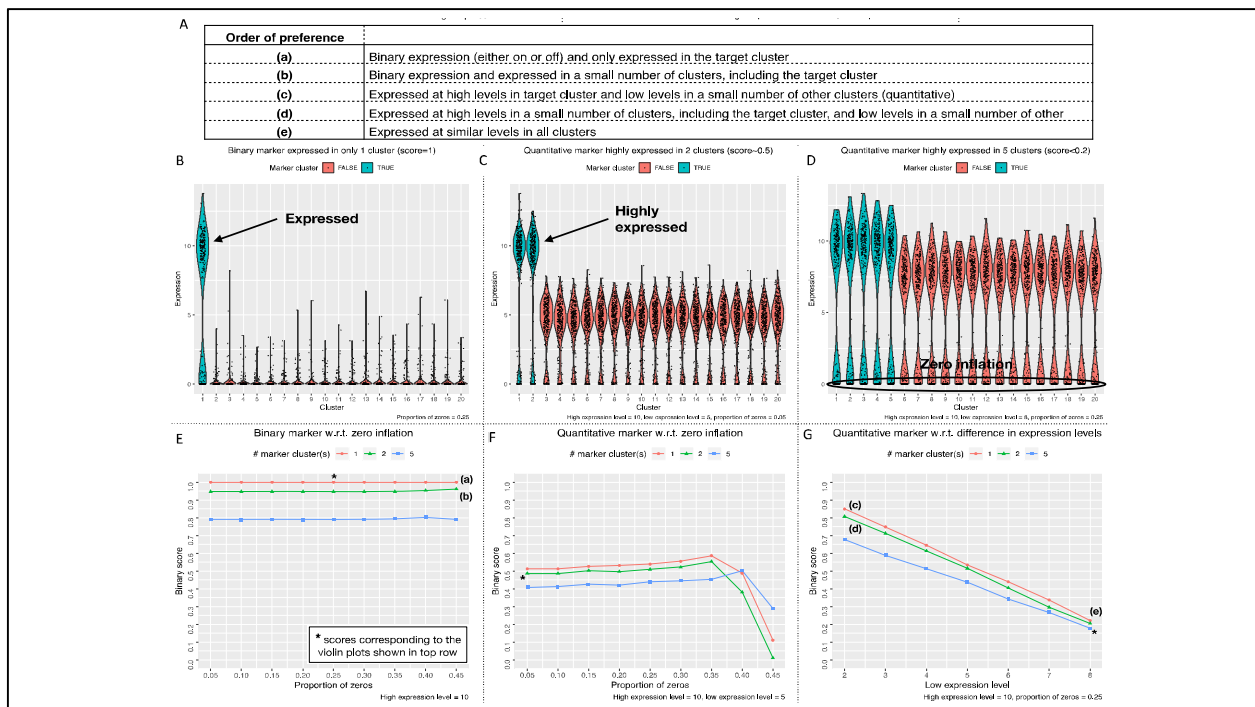
177 Simulations varying the binary expression pattern and level of zero inflation (**Figure 2E**)
178 were then performed. First, the ideal scenario of binary expression, as described above,
179 produced a simulated Binary Expression Score of 1.0 (**Figure 2E** red). When the
180 candidate marker gene was expressed in one (**Figure 2E** green) or four (**Figure 2E**
181 blue) off-target clusters, the Binary Expression Score decreased to 0.95 and 0.90,
182 respectively. In addition, these scores were robust to high zero inflation proportions,
183 demonstrating no decrease in Binary Expression Score up to 45% zero values.

184 Next, quantitative marker expression patterns was added to the simulation (**Figures 2F**
185 **& G**) by varying the number of off-target clusters with high expression levels and adding
186 moderate expression to the remaining off-target clusters. In all cases in which
187 quantitative difference in expression are simulated, the Binary Expression Scores are
188 substantially reduced (**Figures 2F**). In the best case, where only the target cluster has
189 high expression and the off-target clusters have moderate expression, the Binary
190 Expression Score was 0.52. Further Binary Expression Score reductions are found
191 when the high expression levels are present in additional off-target clusters. Adjusting
192 the level of zero inflation for these scenarios showed that these Binary Expression

193 Scores were also robust to increasing zero inflation levels until they drop dramatically
 194 after 35% zero values.

195 Finally, simulations were performed to again test how a high-expressing marker is
 196 effected by the addition of 1 or 4 high expressing off-target clusters together with
 197 increasing expression levels in the remaining off-target clusters from low (2) to high (8)
 198 expression (**Figures 2G**). With the remaining off-target clusters held at low expression
 199 levels, these three scenarios returned high Binary Expression Scores [0.7-0.85], but
 200 these Binary Expression Scores quickly decreased with increasing levels of off-target
 201 expression. For example, when the off-target expression level was set to 6, all three
 202 high-expressing off-target scenarios returned Binary Expression Scores below 0.5. In
 203 the worst case, where the candidate marker has relatively high expression in all off-
 204 target clusters, the Binary Expression Score was less than 0.2.

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



210 **Figure 2:** Performance testing of Binary Expression Score. In panel A) Possible marker
211 gene expression patterns were ranked by order of preference. Violin plots showing
212 three different expression scenarios: panel B) binary expression only in the target
213 cluster, panel C) quantitative expression with high expression in the target cluster and
214 one other cluster and large differences in expression in the other off-target clusters,
215 panel D) quantitative expression with high expression in the target cluster and four other
216 cluster, small differences in expression in the other off-target clusters, and higher levels
217 of zero inflation. Below, line graphs show the full range of tested simulations from which
218 the above example violin plots are taken. For panels E-G there were three defined test
219 cases: the red where there is one cluster with high expression of the marker gene,
220 green where there are two clusters with high expression of the marker gene, and lastly
221 blue where 5 clusters have high expression of the marker gene. Panel E) gives
222 simulation testing of the Binary Expression Score increasing the proportion of zeros
223 while maintaining off target expression at zero. Panel F) off-target clusters were given
224 moderate levels of expression while the proportion of zeros was increased. In panel G)
225 expression levels were varied in all off-target clusters from low (2) to high expression
226 (8).

227 Marker Gene Comparison Between NS-Forest Versions

228 To evaluate the differences in results between NS-Forest v1.3 and v2.0, we analyzed
229 marker genes selected for cell type clusters generated from single nuclei transcriptomes
230 prepared from all layers (1-6) of the human middle temporal gyrus (MTG) obtained from
231 postmortem and surgically resected samples (**sTables 1-3**). For this dataset, three
232 broad classes of cells were identified: excitatory neurons (10,708 cells), inhibitory
233 neurons (4,297 cells), and non-neuronal cells (923 cells). These nuclei were clustered
234 iteratively by first clustering into the larger groups, followed by subsequent re-clustering
235 within each group until 75 putative cell types were found (25). From left to right of the
236 hierarchical clustering of clusters shown at the top of both heatmaps there are 46
237 inhibitory, 23 excitatory, and 6 non-neuronal types (**Figure 3**). Subsequent figures
238 investigating these cell type clusters are ordered by these taxonomic relationship. In
239 Figure 3, the 155 marker genes determined by NS-Forest version 1.3 and the 157

240 marker genes determined by version 2.0 are the rows while the clusters are columns
241 where the row normalized expression level is reflected in the color gradient of high in
242 red to low in blue/white. The diagonal corresponds to the marker gene set for each
243 cluster type. From the heatmaps it is clear that the binary expression has dramatically
244 improved between NS-Forest version 1.3 and version 2.0 as the diagonal contains more
245 genes with red or bright yellow levels of expression and off-diagonal expression levels
246 are more blue (closer to 0).

247 Given the intention of the Binary Expression Score ranking step to preferentially find
248 marker genes with binary expression, there are tradeoffs in both the number of genes
249 required and the classification power when compared to markers ranked by importance
250 from the random forest model. In general, NS-Forest version 2.0 requires more unique
251 genes for a given dataset. In the case of the full MTG dataset the increase is marginal
252 requiring only two additional unique genes (155 vs 157 genes); however, a larger
253 difference in the number of marker genes required has been observed for other
254 datasets (data not shown). Furthermore, the genes that have a high Binary Expression
255 Score are usually not the same genes that were ranked highest by Information Gain
256 within the random forest model. This suggests that in terms of pure classification the
257 markers determined by v2.0 might be expected to underperform. To directly compare
258 the F-scores between these two versions of NS-Forest, an additional analysis was run
259 setting the beta weight of the F-score to 1 in version 2.0 thereby making it directly
260 comparable to version 1.3. As expected, the mean F-score for version 1.3 (0.62) was
261 slightly higher than the mean F-score for version 2.0 (0.58); however, the average
262 Binary Expression Score for the version 1.3 markers was significantly lower at 0.72
263 versus 0.94 for version 2.0. (For cluster-by-cluster correlations of F-scores and Binary
264 Expression Scores see **sFigure 2**).

265 Within the major branches of the taxonomy, major subclasses are labeled by important
266 neurological markers such as VIP, SST, and PVALB within the inhibitory subclass, and
267 RORB and FEZF2 within the excitatory subclass. Binary marker genes for the cell types
268 within these subclass branches of the taxonomy can be more difficult to determine,
269 especially when there are many closely related types. For example, both the SST and
270 PVALB contain a number of closely related cell types. When looking at the expression

271 of the marker genes, we can see that between these two major subclasses there is little
272 expression overlap; however, within each subclass there are number of closely-related
273 cell types that show overlapping expression, for example the cell types circled in red
274 within the SST subclass or the types circled in blue within the PVALB subclass (**Figure**
275 **3C**). These cell types tend to have lower F-beta scores and lower marker Binary
276 Expression Scores.

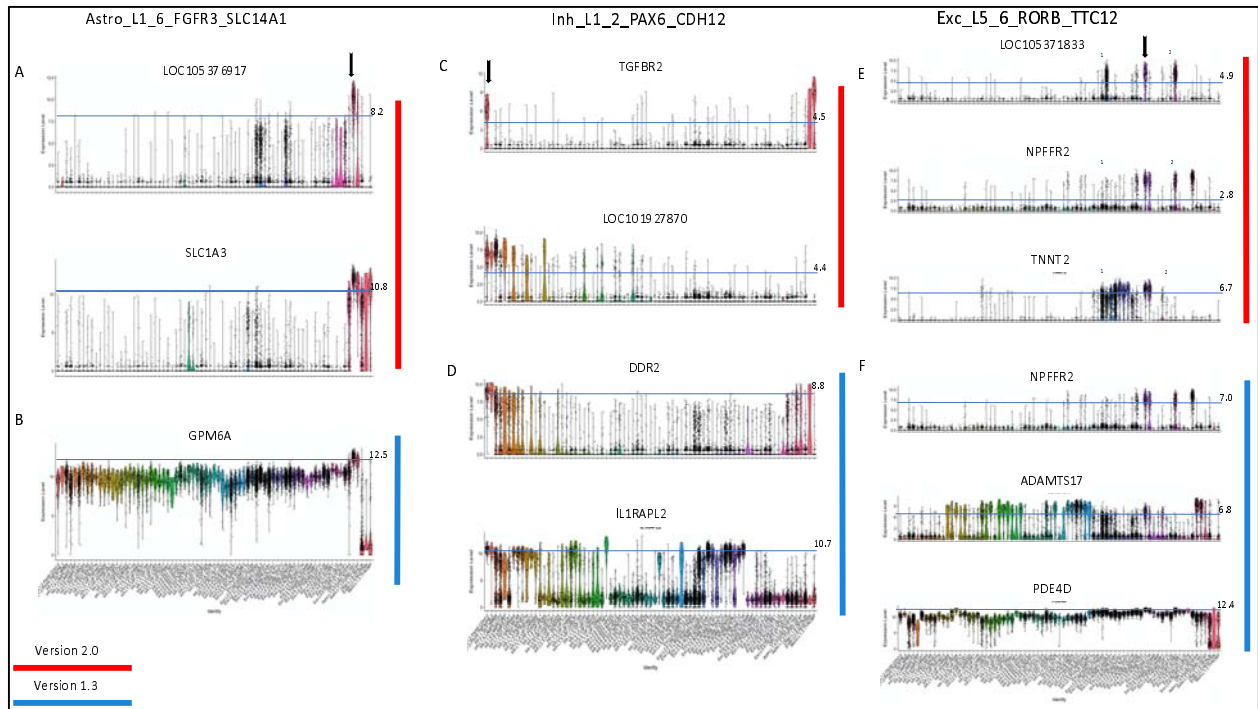
277 Looking in more detail at the properties of the marker genes selected for individual cell
278 types, we can clearly see the differences between NS-Forest Version 1.3 and 2.0. The
279 expression patterns for the astrocyte cell type Astro_L1_6_FGFR3_SLC14A1 show
280 these differences in the marker genes selected by the two NS-Forest version in more
281 detail (**Figure 4A/B**). NS-Forest v1.3 selects a single marker gene to best discriminate
282 this cluster, while v2.0 selects two. NS-Forest v1.3 selects only the GPM6A gene which
283 performs well at classifying this cell type along a quantitative boundary at the high log2
284 expression level of 12.5, but also shows intermediate expression in several off-target
285 clusters centered around 10 (**Figure 4A**). Consequently, this quantitative marker is
286 good for classification only when this small window of expression difference is
287 discernible. In contrast, version 2.0 selects LOC105376917 and SLC1A3, both of which
288 have binary expression patterns across clusters (**Figure 4B**). LOC105376917 is highly
289 expressed in only the target cluster and one additional closely-related off-target cluster.
290 Adding SLC1A3 further improves classification by removing cells from this off-target
291 cluster.

292 In the case of the inhibitory neuron Inh_L1_2_PAX6_CDH12 both v1.3 and v2.0 select
293 two marker genes; however, their characteristics are very different (**Figure 4C/D**). NS-
294 Forest v1.3 again found markers that classified along quantitative boundaries. DDR2 is
295 expressed in all the related clusters in the taxonomy and in some glial clusters at the far
296 end of the taxonomy. The addition of IL1RAPL2 removes the glial clusters and improves
297 the classification; however, IL1RAPL2 is another example of a quantitative marker as it
298 separates the target cluster from the related cluster by narrow differences in expression.
299 NS-Forest v2.0 selected two highly binary markers: TGFBR2, which is very specific to
300 only two clusters, the target cluster and a non-neuronal type at the other end of the

301 taxonomy. The addition of the LOC101927870 gene eliminates cells in the non-neuronal
302 cluster to refine the classification.

303 Lastly, the excitatory neuron *Exc_L5_6_RORB_TTC12* required three markers by both
304 NS-Forest versions to optimize the classification (**Figure 4E/F**). Again, as previously
305 described, NS-Forest v1.3 determined genes that used a quantitative boundary for
306 classification while NS-Forest v2.0 discovered binary markers. A more detailed look at
307 these binary markers provides a clear demonstration of the combinatorics employed by
308 NS-Forest v2.0. Within the target cluster, demarcated by the arrow, all three markers
309 have high expression; however, the off-target excitatory clusters marked as 1 and 2 also
310 express some but not all these markers. By leveraging the combinatorics of the three
311 marker combination, highly discriminative solution is obtained. Gene LOC105371833 is
312 the most binary marker; however, it has high expression in a number of off-target cells
313 in clusters 1 and 2. The addition of the *NPFFR2* gene removes most of the false
314 positives in cluster 1, while adding the *TNNT2* gene removes the false positives from
315 cluster 2. Together this combination of three marker genes discriminates
316 *Exc_L5_6_RORB_TTC12* from other excitatory cell types.

317 These results show that while adding the Binary Expression Score criteria does slightly
318 decrease the overall classification power of the markers selected, it dramatically
319 increases the binary expression pattern making the markers more useful for
320 downstream applications.



326

327 **Figure 4:** Violin plots for marker gene expression for a selection of cell type clusters
328 representative of the three major classes in the taxonomy: glial cells, inhibitory neurons,
329 and excitatory neurons. Panels A, C, and E (red) give markers determined by NS-Forest
330 v2.0 while panels B, D, and F (blue) give markers from NS-Forest v1.3. Expression is
331 given in log₂ scale. Expression thresholds are demarcated by light blue lines and values
332 are given on the right. Thresholds for NS-Forest v2.0 were determined by decision tree
333 split points, while NS-Forest v1.3 were fixed at the cluster median expression for that
334 gene.

335

336 Characterization of NS-Forest v2.0 Markers

337 Overall, the results from NS-Forest v2.0 reflect the high quality of the data and
338 clustering analysis as NS-Forest is a supervised machine learning method and is reliant
339 on the quality of the clustering results. The median number of markers required for
340 optimal classification was 2, with only two clusters needing 4 markers, producing a
341 mean F-beta score of 0.69. Overall, the 75 clusters required 157 unique genes to
342 achieve optimal classification. Occasionally, marker genes are shared between clusters,

343 with eleven genes that were not unique [MOXD1, MME, LOC101928196, SULF1,
344 NPFFR2, LINC01583, TAC1, COL15A1, LOC401478, CPED1, TAC3].

345 Out of the 157 NS-Forest v2.0 marker genes, 37 (24%) were long non-coding RNAs
346 (lncRNAs) or uncharacterized loci (LOCs). Non-coding RNAs have been previously
347 found to be prevalent when analyzing RNA-seq data from single neuronal cells or nuclei
348 and surprisingly these non-coding RNAs had higher specificity as markers when
349 compared to coding genes (27). In particular, lncRNAs are known to show cell line
350 specific expression (28). In contrast, little is known about the LOC genes. These genes
351 are particularly intriguing as they are highly specific to individual cell types and are likely
352 important for their function. As such, they represent areas of unknown biology
353 discovered by scRNA-seq and NS-Forest machine learning that warrant further
354 investigation.

355 For the characterized marker genes, the most enriched annotations both by adjusted p-
356 value and number of genes involved are for signaling (signal peptide, Signal,
357 GO:0007218~neuropeptide) and extracellular matrix (Glycoprotein, Extracellular matrix,
358 GO:0005615~extracellular space, GO:0005578~proteinaceous extracellular matrix,
359 GO:0030198~extracellular matrix organization, GO:0005576~extracellular region,
360 GO:0031012~extracellular matrix), including neuropeptide, GO:0007218~neuropeptide
361 signaling pathway, and calcium (**sTable 4**). There are fewer genes annotated with these
362 specific functions as neurology is a rapidly expanding field; however, many other genes
363 assessed here are known signaling peptides in other contexts and would benefit from
364 further characterization in neurological context. Taken together, these results suggests
365 that specific signaling pathways and extracellular signaling molecules are key to
366 neuronal cell type identity.

367 Comparison with other Marker Gene Selection Approaches

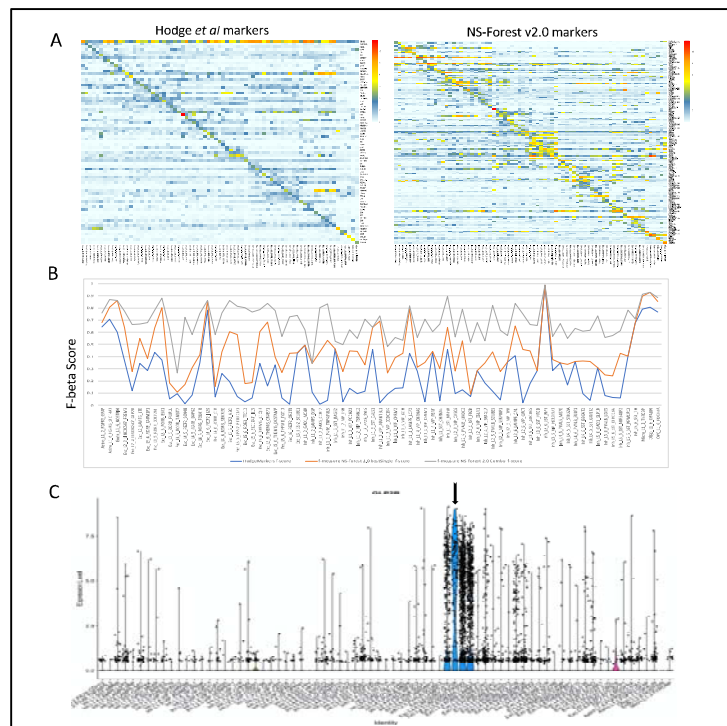
368 To understand how the NS-Forest marker genes compare to previously published
369 markers for the human middle temporal gyrus (MTG), we compared the NS-Forest
370 markers to those determined in Hodge *et al* (25) using a different binary expression
371 approach for use in cell cluster naming. In addition to a broad marker determined by the

372 taxonomy and prior knowledge (such as GAD1 or SST), a single marker per cell type
373 cluster was assigned. In total, sixteen of the seventy-five Hodge markers overlapped
374 with the NS-Forest markers [BAGE2, GGH, CASC6, NPY, HPGD, STK32A, ADGRG6,
375 TH, MEPE, PENK, CARM1P1, TWIST2, IL26, SULF1, ADAMTSL1, PDGFRA]. These
376 sixteen were spread across the taxonomy, representing cell type clusters from all three
377 major cell type lineages. Unscaled heatmaps of mean gene expression per cluster for
378 both the Hodge and NS-Forest marker sets (**Figure 5A**) demonstrate that both are
379 characterized by binary expression patterns, having a higher expression along the
380 diagonal versus off-diagonal; however, the Hodge markers have an overall lower mean
381 expression level of 4.8 log₂ CPM in comparison with the mean expression for the NS-
382 Forest markers of 7.0 log₂ CPM.

383 One major difference between these two approaches is that the Hodge marker set
384 contains a single marker per cluster to label a distinct cluster phenotype while NS-
385 Forest selects combinations of markers that optimize classification power. By running
386 the Hodge markers through NS-forest v2.0, we estimated F-beta scores for the single
387 Hodge markers in order to compare their classification power to the best single NS-
388 Forest markers, and the NS-Forest combination of markers (**Figure 5B**). Overall, the
389 trend lines show that the F-beta scores for single markers, both blue and orange lines,
390 follow a similar trajectory with some clusters being more difficult to classify than others,
391 i.e., having lower F-beta scores. However, the NS-Forest combination of markers,
392 shown in grey, demonstrate that combinations of markers yield a uniformly higher power
393 of discrimination over a single marker, regardless of how the single best marker is
394 chosen.

395 When evaluating the F-beta scores for the Hodge markers, it became clear that many
396 had elevated false positive rates. To directly compare the two sets of markers, we
397 computed the false discovery rate (FDR= FP/FP+TP) for each cell type and averaged
398 across the entire set. The Hodge markers had an average FDR of 0.7 versus 0.14 for
399 the NS-Forest markers. GLP2R, which is a marker for Exc_L2_4_LINC00507_GLP2R,
400 offers a good visual example (**Figure 5C**). This gene expressed in the target cluster but
401 also the nearest cell types within the LINC00507 group. NS-Forest also has difficulty

402 finding markers for this cell cluster phenotype, requiring 3 markers in total; however, in
403 combination these markers help reduced the FDR rate from 0.89 to 0.11.



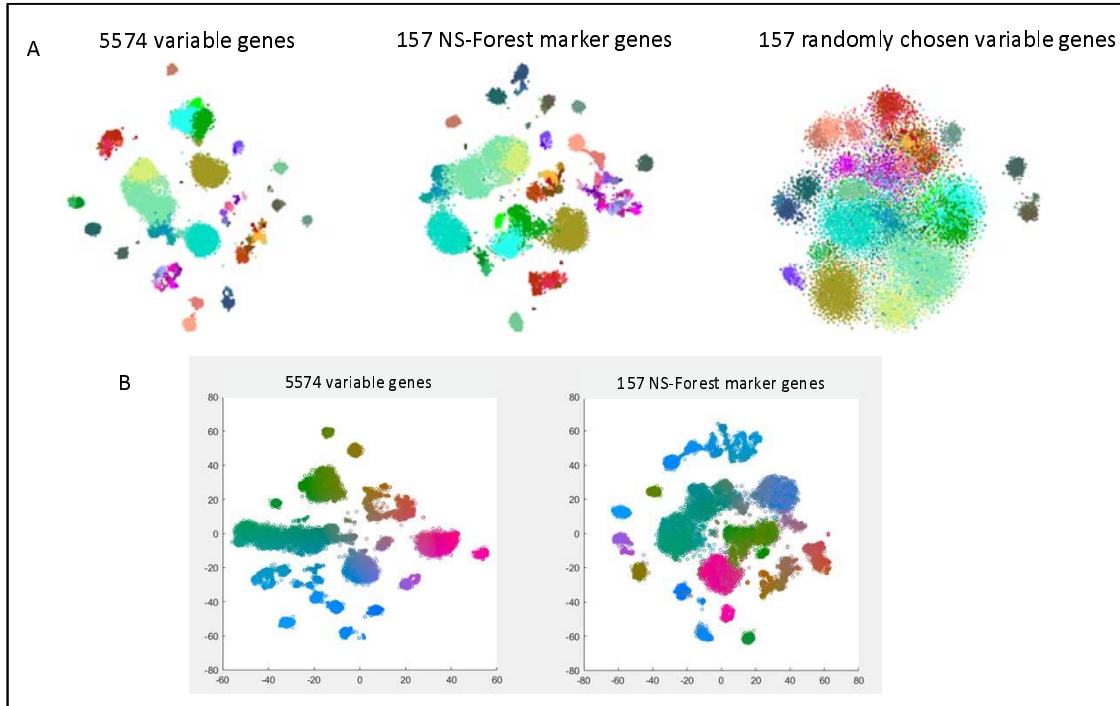
404
405 **Figure 5:** Comparison of Hodge *et al* markers to NS-Forest v2.0 for the full MTG data.
406 Panel A gives an unscaled heatmap where the rows are the mean expression per gene
407 and the columns are clusters. Panel B) Gives the F-beta scores for the single Hodge
408 marker, the best NS-Forest single marker, and the combination of markers found by
409 NS-Forest. C) an example violin plot of a binary expression pattern selected for by the
410 method used by Hodge *et al* for cluster Exc_L2_4_LINC00507_GLP2R.

411
412 **Validation of Human MTG NS-Forest v2.0 Markers**

413 At current, the ground truth for the neuron types and their marker genes in human MTG
414 taxonomy is not available as it is an active area of investigation. Consequently, a true
415 biological validation of the marker genes is not possible. As an alternative, we asked the
416 question, do the minimum set of marker genes selected by NS-Forest capture the
417 underlying structure of cell type identity reflected in the entire expressed transcriptome?
418 To do this, we generated tSNE plots using the complete 5574 variable genes used for

419 the original MTG clustering, the minimum set of 157 NS-Forest v2.0 marker genes, and
420 157 genes randomly selected from the complete variable genes list. These embeddings
421 were then painted using the cell type assignments from the MTG taxonomy. From the
422 tSNE plots it is clear that the NS-Forest markers closely recapitulate the clustering
423 structure of the complete variable genes set, much better than the randomly selected
424 genes (**Figure 6A**). For example, in the bottom of the complete variable gene tSNE
425 there is a light salmon and dark salmon colored group of clusters, and these two
426 clusters are preserved in the right hand side of the NS-Forest marker tSNE, whereas in
427 the tSNE from the randomly selected variable genes these two clusters spread out and
428 a third brown cluster is now merged with light salmon cluster. Examples like this can be
429 seen throughout the three embeddings. A more quantitative analysis of these tSNE
430 embeddings using the Nearest-Neighbor Preservation metric shows that both the
431 precision and recall are higher using the 157 NS-Forest markers compared with 50
432 sampling of 157 genes randomly selected from the variable gene set (**sFigure 3**).

433 In addition, the local embedding structures within a given tSNE cluster also appear to be
434 well preserved (**Figure 6B**). The complete variable gene tSNE was painted using a
435 color gradient based on the coordinate positioning. This yields a visual way of
436 comparing where individual nuclei are located within the full tSNE embedding versus
437 other tSNE embeddings. The NS-Forest marker tSNE was then painted using the
438 colors derived from the full tSNE gradient. The fact that the same color gradients are
439 observed in the NS-Forest embedding demonstrates that the positional gradients, and
440 thus the nuclei-to-nuclei relationships, in the NS-Forest embedding closely reflect the
441 positional gradients in the complete full tSNE embedding. For example, in the full tSNE
442 there is a long cluster of nuclei beginning on the left in green that extends toward the
443 middle moving into bluish green and ending with a purplish blue. This same cluster, with
444 the same color gradient is preserved within the center left cluster of the NS-Forest
445 tSNE.



446

447 **Figure 6:** Validation of the 157 NS-Forest v2.0 Middle Temporal Gyrus (MTG) marker
448 genes. In panel A) we have tSNE plots for the full DE list of 5574 genes, the 157 NS-
449 Forest markers, and 157 genes randomly selected from the variable gene list. In panel
450 B) left is a tSNE generated from the full variable gene list of 5574 genes colored by
451 coordinate position while right is a tSNE generated using the 157 NS-Forest markers
452 then painted by nuclei according to the color scheme established in the full tSNE on the
453 left.

454 Discussion

455 Here we describe NS-Forest version 2.0. Development was driven by user community
456 requirements for data derived cell type phenotype definitions that are testable in future
457 investigations. To this end, a number of changes were made after the random forest
458 feature selection. In earlier version of NS-Forest, negative markers were occasionally
459 found. These are marker genes that are expressed in most off-target clusters but not
460 the target cluster. Given that testing for a something that is not expressed is
461 methodologically difficult, it was decided to avoid this category of markers. By
462 implementing a median expression level cutoff greater than zero for the target cluster,

463 we removed all possible negative marker genes. In addition, this cutoff also defines one
464 basic characteristic of a NS-Forest Marker: they are required to be expressed in greater
465 than half of the cells within the cell type cluster.

466 NS-Forest v1.3 contained simple random forest feature selection approach that
467 discovered quantitative markers that were good for classification but generally
468 problematic for further biological investigation. This limitation of random forest feature
469 selection may be shared with other machine learning methods. Consequently, a ranking
470 step was incorporated to select for markers with binary expression patterns. Simulation
471 testing performed on this Binary Expression Score ranking step demonstrated that it
472 selected for marker genes with binary expression patterns and accurately ranked them
473 according to level of binary expression. As a result, NS-Forest v2.0 demonstrated clear
474 improvement in the enrichment for binary expression patterns but at a small cost to the
475 overall classification power and number of marker genes necessary. Consequently, If
476 the user requires classification with less requirement for downstream investigation, then
477 we would recommend using NS-Forest v1.3; however, in all other cases NS-Forest v2.0
478 is recommended. Both versions are available as official releases at the github
479 repository.

480 Beyond their use for defining and investigating cell type phenotypes, necessary and
481 sufficient marker genes also offer a dimensionality reduction with limited loss of fidelity
482 to the originally clustering solution. This dimensionality reduction offers a feasible way of
483 representing the clustering solution with a minimal amount of information which is ideal
484 for data dissemination. These marker genes can then be used to generate a reference
485 knowledgebase for cell types, in effect generating an expression barcode of marker
486 genes for a given cell phenotype.

487 As mentioned above, NS-Forest markers are optimized for downstream experimental
488 investigation. There are a number of assays for which known markers could facilitate
489 biological investigation such qPCR and the burgeoning field of spatial transcriptomics
490 based on multiplex FISH. To date a number of projects have used NS-Forest markers
491 for these purposes. For example, qPCR probes based on NS-Forest markers were
492 made to detect genes in scRNA-seq libraries from myeloid dendritic cells (mDCs) FACS

493 sorted from peripheral blood in patients treated with the Hepatitis B vaccine (30,
494 publication in preparation). In a similar fashion, gene probes were designed based on
495 NS-Forest markers for cell type detection using a number of spatial transcriptomic
496 technologies. These technologies aim to resolve the location within the tissue of cell
497 types derived from scRNA-seq generated taxonomies (31).

498 Another possible application of NS-Forest is to utilize selected gene sets of particular
499 interest as input to produce marker gene sets designed to capture specific cell type
500 properties. For example, the input of gene sets composed of transcription factors could
501 reveal master regulators of developmental programs (32). Input gene sets composed of
502 neuropeptides and neurotransmitters could be used to shed new light on the specific
503 signaling properties of different neuronal cell subsets (33). Input gene sets composed
504 of cell surface markers could be used to identify markers for use in FACS sorting.

505 As the number of experiments performed and datasets made publicly available
506 dramatically increase, the greater biological community is left with the monumental task
507 of integrating these data into a consensus of canonical cell types. With cell phenotypes
508 defined by NS-Forest marker genes, we can move ahead with the creation of a
509 dissemination framework that defines ontological classes based upon these molecular
510 markers as the necessary and sufficient criteria in an axiomatic semantic representation
511 compliant with FAIR principles. Ontological representation has numerous advantages
512 over simple vocabularies, including the structuring of knowledge in a computationally
513 readable format so that findings from many experiments can be easily accessible and
514 “reasoning” can be performed to ensure the consistency of the representation as the
515 knowledge rapidly grows. These provisional instances of “cell type clusters” defined by
516 NS-Forest markers can form the basis for the instantiation of an ontology class that can
517 be in the future adopted into the official Cell Ontology (CL). Progress is already
518 underway in developing programmatic and scalable methods to handle the amount of
519 single cell data being generated. This ontological representation can address several
520 pressing needs of the wider biological research community. Producing an easy, visually
521 accessible overview of the results of many single cell experiments in a traversable
522 structure while preserving the hierarchical relationships inherent in a taxonomy of cells.
523 In addition, this ontology will provide a platform for integration with other data modalities

524 such as cell morphology, electrophysiology, cell-cell interactions. A provisional cell
525 ontology (pCL) generated in this manner for Middle Temporal Gyrus and primary motor
526 cortex is available for exploration at <https://bioportal.bioontology.org/ontologies/PCL> .

527 Methods

528 NS- Forest version 2.0

529 Initial Feature Selection: The NS-Forest version 2.0 workflow (**Figure 1a-b**) begins with
530 a cell-by-gene expression matrix, with an additional column containing cluster
531 membership labels, produced by any expression data clustering method applied to
532 single cell/nucleus RNA sequencing (scRNA-seq) datasets. This cluster-labelled
533 expression matrix is then used to generate Random Forest classification models
534 distinguishing each target cluster from all other clusters (binary classification) using
535 RandomForestClassifier scikit. RandomForestClassifier hyperparameters were left at
536 default except that the number of trees was set at 10,000 to give sufficient coverage of
537 the sample and gene expression feature space; necessary coverage for a given feature
538 space is estimated as the square root of the number of samples (~10,000 cells) times
539 the square root of the number of features (~10,000 genes). From the resulting Random
540 Forest model, the average Gini Impurity value is used to initially rank genes based on
541 their feature importance.

542 Feature Re-ranking Based on Positive Binary Expression: Re-ranking the features after
543 initial Random Forest selection begins with positive expression filtering (**Figure 1c**). By
544 default, genes with a median cluster expression of 0 are removed in order to exclude
545 genes that are not expressed in the relevant cluster, which we refer to as negative
546 markers, or show high zero inflation. This parameter is tunable and can be adjusted
547 according to the desired positive expression level.

548 Next, genes are re-ranked to enrich for genes with binary expression patterns (**Figure**
549 **1d**). A “Binary Expression Score” was developed to select for genes that show all-or-
550 none expression patterns, with expression in the target cluster and as few other cell

551 type clusters as possible. The Binary Expression Score is calculated for each gene in
552 the initial Random Forest feature list according to the equation:

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

553

554 where y_i is the median gene expression level for each cluster i , y_T is the median
555 expression in the target cluster, and n is the number of clusters. This results in a Binary
556 Expression Score in the range of 0 – 1, with a Binary Expression Score of 1 being the
557 ideal case where the gene is only expressed in the target cluster (**Figure 1e**).

558 Minimum Feature Combination Determination: After the top genes are re-ranked based
559 on positive binary expression, they are then tested for their classification power
560 individually and in combination. First, the top M genes (6 genes by default) are used to
561 generate individual decision trees to determine the optimal expression level cut-off
562 value for each gene (**Figure 1F**). The maximum leaf nodes parameter is set at two,
563 thereby ensuring a single split point per tree. From these trees, the optimal gene
564 expression threshold at the split point is extracted.

565 To evaluate the discriminative power of a given combination of candidate marker genes,
566 we use the F-beta score as an objective function. The F-score is the harmonic mean of
567 precision and recall providing equal weight for these two classification measures. The
568 F-beta score includes a beta term that allows for the weighting of the function towards
569 either precision (beta < 1) or recall (beta > 1) (**Figure 1G**). The beta for the analysis
570 described here was estimated empirically at 0.5 (**Supplemental Figure 1**).

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

575 Simulation Testing of the Binary Expression Score

576 Simulation studies were conducted to investigate the properties of the Binary
577 Expression Score weighting using a three-component mixture model to reflect the zero-

578 inflation technical artifact and the background and positive expression signals in real
579 data distributions. Denoting X as the gene expression value, our simulated data follow a
580 mixture distribution:

$$581 \quad P(X = x) = \pi_1 \cdot \delta_0(x) + \pi_2 \cdot f_{\text{Gamma}}(x) + \pi_3 \cdot f_{\text{Normal}}(x)$$

582 Where $\delta_0(x)$ is the probability density function of the degenerate distribution at 0 for the zero-
583 inflation technical artifact, $f_{\text{Gamma}}(x)$ is the probability density function of a Gamma
584 distribution (with hyperparameters α and β) for low level background expression from off-target
585 cells or on-target cells with low expression, and $f_{\text{Normal}}(x)$ is the probability density function of a
586 Normal distribution (with hyperparameters μ and σ) for positive expression
587 signals; parameters π_1 , π_2 and π_3 are the corresponding mixture weights for each component
588 such that $\pi_1, \pi_2, \pi_3 > 0$ and $\pi_1 + \pi_2 + \pi_3 = 1$. In our simulations, we generated 20 clusters with
589 300 cells in each cluster. We designed cases where the simulated gene is expressed at high
590 levels in 1, 2, or 5 clusters. Both binary and quantitative markers were simulated for on-target
591 and off-target clusters by setting different parameters and hyperparameters in the mixture
592 model.

593 snRNA-seq Data

594 The scRNA-seq data evaluated here were obtained from the Allen Institute for Brain
595 Science (<https://portal.brain-map.org/atlas-and-data/rnaseq>). Experimental design,
596 including tissue sampling and data processing, can be found in Hodge *et al.* (23). In
597 brief, layers 1-6 of the human Middle Temporal Gyrus (MTG) were vibratome sectioned,
598 nuclei were extracted and labelled for NeuN expression. Nuclei were then FACS sorted
599 and libraries were generated using the Smart-Seq4 and Nextera XT chemistries. Data
600 processing and clustering were then performed as detailed in (22).

601 NS-Forest v2.0 was run using the cluster assignments given in Hodge *et al.* (23). Cells
602 not assigned to a cluster were removed from the analysis. CPM expression values were
603 $\log_2(x+1)$ transformed and genes with a sum of zero median expression across all
604 clusters were removed. After filtering, 15,928 cells and 13,946 genes remained. Given
605 the size of the input matrix, we increased the number of trees in the random forest
606 model from the default of ten thousand to fifty thousand.

607 Marker Validation

608 In order to demonstrate the preservation of the cell type clustering characteristics using
609 NS-Forest marker genes, tSNE embeddings were generated using Cytosplore. The
610 original clustering solution is represented by an embedding generated from the 5574
611 variable genes used for the iterative clustering originally performed (22, 23). Additional
612 embeddings were made using the combined set of 157 marker genes for all cell type
613 clusters determined by NS-Forest version 2.0, and 50 sets of 157 genes chosen at
614 random from the original 5574 genes.

615 Data Access

616 All data used herein is publicly available.

617

618 Acknowledgements

619

620 This work was supported by the Allen Institute for Brain Science, the JCVI Innovation
621 Fund, the U.S. National Institutes of Health (R21-AI122100 and U19-AI118626), the
622 California Institute for Regenerative Medicine (GC1R-06673-B), the Wellcome
623 Trust 208379/Z/17/Z, and from the Chan Zuckerberg Initiative DAF, an advised fund of
624 the Silicon Valley Community Foundation (2018-182730).

625

626 Disclosure Declaration

627 None to make.

628 References

- 629 1. Scheuermann RH, Ceusters W, Smith B. 2009. Toward an Ontological Treatment
630 of Disease and Diagnosis **Summit on Translational Bioinformatics** 2009:116-
631 120. PMID: 21347182; PMCID: PMC3041577.
632
- 633 2. Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E, Bodenmiller B,
634 Campbell P, Carninci P, Clatworthy M, Clevers H, et al. 2017. The Human Cell
635 Atlas. *Elife*. 5;6:e27041. doi: 10.7554/eLife.27041. PMID: 29206104; PMCID:
636 PMC5762154.
637
- 638 3. Darmanis S, Sloan SA, Zhang Y, Enge M, Caneda C, Shuer LM, Hayden
639 Gephart MG, Barres BA, Quake SR. 2015. A survey of human brain
640 transcriptome diversity at the single cell level. *Proc Natl Acad Sci U S A*. Jun
641 9;112(23):7285-90. PMID: 26060301; PMC: PMC4466750
642
- 643 4. Enge M, Arda HE, Mignardi M, Beausang J, Bottino R, Kim SK, Quake SR. 2017.
644 Single-Cell Analysis of Human Pancreas Reveals Transcriptional Signatures of
645 Aging and Somatic Mutation Patterns. *Cell*. Oct 5;171(2):321-330.e14. PMID:
646 28965763; PMC: PMC6047899
647
- 648 5. Nowakowski TJ, Bhaduri A, Pollen AA, Alvarado B, Mostajo-Radji MA, Di Lullo E,
649 Haeussler M, Sandoval-Espinosa C, Liu SJ, Velmeshev D, et al. 2017.
650 Spatiotemporal gene expression trajectories reveal developmental hierarchies of
651 the human cortex. *Science*. Dec 8;358(6368):1318-1323. PMID: 29217575;
652 PMC: PMC5991609
653
- 654 6. Schiller HB, Montoro DT, Simon LM, Rawlins EL, Meyer KB, Strunz M, Vieira
655 Braga FA, Timens W, Koppelman GH, Budinger GRS, et al. 2019. The Human
656 Lung Cell Atlas: A High-Resolution Reference Map of the Human Lung in Health
657 and Disease. *Am J Respir Cell Mol Biol*. Jul;61(1):31-41. doi:
658 10.1165/rcmb.2018-0416TR. PMID: 30995076; PMCID: PMC6604220.
659
- 660 7. Muraro MJ, Dharmadhikari G, Grün D, Groen N, Dielen T, Jansen E, van Gorp L,
661 Engelse MA, Carlotti F, de Koning EJ, van Oudenaarden A. 2016. A Single-Cell
662 Transcriptome Atlas of the Human Pancreas. *Cell Syst*. Oct 26;3(4):385-394.e3.
663 doi: 10.1016/j.cels.2016.09.002. Epub 2016 Sep 29. PMID: 27693023; PMCID:
664 PMC5092539.
665
- 666 8. Asp M, Giacomello S, Larsson L, Wu C, Fürth D, Qian X, Wärdell E, Custodio J,
667 Reimegård J, Salmén F, et al. 2019. A Spatiotemporal Organ-Wide Gene

- 668 Expression and Cell Atlas of the Developing Human Heart. *Cell*. Dec
669 12;179(7):1647-1660.e19. doi: 10.1016/j.cell.2019.11.025. PMID: 31835037.
670
- 671 9. Mott MC, Gordon JA, Koroshetz WJ. 2018. The NIH BRAIN Initiative: Advancing
672 neurotechnologies, integrating disciplines. *PLoS Biol*. Nov 26;16(11):e3000066.
673 doi: 10.1371/journal.pbio.3000066. PMID: 30475794; PMCID: PMC6283590.
674
- 675 10. Levitin HM, Yuan J, Sims PA. 2018. Single-Cell Transcriptomic Analysis of
676 Tumor Heterogeneity. *Trends Cancer*. Apr;4(4):264-268. doi:
677 10.1016/j.trecan.2018.02.003. Epub 2018 Mar 9. PMID: 29606308; PMCID:
678 PMC5993208.
679
- 680 11. Al-Dalahmah O, Sosunov AA, Shaik A, Ofori K, Liu Y, Vonsattel JP, Adorjan I,
681 Menon V, Goldman JE. 2020. Single-nucleus RNA-seq identifies Huntington
682 disease astrocyte states. *Acta Neuropathol Commun*. Feb 18;8(1):19. doi:
683 10.1186/s40478-020-0880-6. PMID: 32070434; PMCID: PMC7029580.
684
- 685 12. Chaudhry F, Isherwood J, Bawa T, Patel D, Gurdziel K, Lanfear DE, Ruden DM,
686 Levy PD. 2019. Single-Cell RNA Sequencing of the Cardiovascular System: New
687 Looks for Old Diseases. *Front Cardiovasc Med*. 2019;6:173. Dec 10.
688 doi:10.3389/fcvm.2019.00173
689
- 690 13. Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, Ziraldo SB,
691 Wheeler TD, McDermott GP, Zhu J, et al. 2017. Massively parallel digital
692 transcriptional profiling of single cells. *Nat Commun* 8, 14049.
693 <https://doi.org/10.1038/ncomms14049>
694
- 695 14. Picelli S, Björklund ÅK, Faridani OR, Sagasser S, Winberg G, Sandberg R. 2013.
696 Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat*
697 *Methods*;10(11):1096-1098. doi:10.1038/nmeth.2639
698
- 699 15. Krishnaswami SR, Grindberg RV, Novotny M, Venepally P, Lacar B, Bhutani K,
700 Linker SB, Pham S, Erwin JA, Miller JA, et al. 2016. Using single nuclei for RNA-
701 seq to capture the transcriptome of postmortem neurons. *Nat*
702 *Protoc*;11(3):499-524. doi:10.1038/nprot.2016.015
703
- 704 16. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM 3rd, Hao
705 Y, Stoeckius M, Smibert P, Satija R. 2019. Comprehensive Integration of Single-
Cell Data. *Cell*, **177**, 1888-1902. doi: [10.1016/j.cell.2019.05.031](https://doi.org/10.1016/j.cell.2019.05.031).

- 706 17. Wolf F, Angerer P, Theis F. 2018. SCANPY: large-scale single-cell gene
707 expression data analysis. *Genome Biol* **19**, 15. [https://doi.org/10.1186/s13059-](https://doi.org/10.1186/s13059-017-1382-0)
708 [017-1382-0](https://doi.org/10.1186/s13059-017-1382-0)
709
- 710 18. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A,
711 Blomberg N, Boiten JW, da Silva Santos LB, Bourne PE, et al. 2016. The FAIR
712 Guiding Principles for scientific data management and stewardship. *Sci Data* **3**,
713 160018. <https://doi.org/10.1038/sdata.2016.18>
714
- 715 19. Bard J, Rhee SY, Ashburner M. 2005. An ontology for cell types. *Genome Biol.*;
716 **6**(2):R21. doi:10.1186/gb-2005-6-2-r21
717
- 718 20. Diehl AD, Augustine AD, Blake JA, Cowell LG, Gold ES, Gondré-Lewis TA,
719 Masci AM, Meehan TF, Morel PA, Nijnik A, et al. 2011. Hematopoietic cell types:
720 prototype for a revised cell ontology. *J Biomed Inform.*; **44**(1):75-79.
721 doi:10.1016/j.jbi.2010.01.006
722
- 723 21. Meehan TF, Masci AM, Abdulla A, Cowell LG, Blake JA, Mungall CJ, Diehl AD.
724 2011. Logical development of the cell ontology. *BMC Bioinformatics*. 2011;**12**:6.
725 doi:10.1186/1471-2105-12-6.
726
- 727 22. Bakken T, Cowell L, Aevermann BD, Novotny M, Hodge R, Miller JA, Lee A,
728 Chang I, McCorrison J, Pulendran B, et al. 2017. Cell type discovery and
729 representation in the era of high-content single cell phenotyping. *BMC*
730 *Bioinformatics*. Dec 21;**18**(Suppl 17):559. doi: 10.1186/s12859-017-1977-1.
731 PMID: 29322913; PMCID: PMC5763450.
732
- 733 23. Kuby J, Kindt TJ, Goldsby RA, Osborne BA. 2007. *Kuby Immunology*. San
734 Francisco: W.H. Freeman. ISBN 1-4292-0211-4.
735
- 736 24. Tasic B, Yao Z, Graybuck LT, Smith KA, Nguyen TN, Bertagnolli D, Goldy J,
737 Garren E, Economo MN, Viswanathan S, et al. 2018. Shared and distinct
738 transcriptomic cell types across neocortical areas. *Nature*.
739 **2018**;563(7729):72-78. doi:10.1038/s41586-018-0654-5
740
- 741 25. Hodge RD, Bakken TE, Miller JA, Smith KA, Barkan ER, Graybuck LT, Close JL,
742 Long B, Johansen N, Penn O, et al. 2019. Conserved cell types with divergent
743 features in human versus mouse cortex. *Nature*. 2019;**573**(7772):61-68.
744 doi:10.1038/s41586-019-1506-7
745

- 746 26. Aevermann BD, Novotny M, Bakken T, Miller JA, Diehl AD, Osumi-Sutherland D,
747 Lasken RS, Lein ES, Scheuermann RH. 2018. Cell type discovery using single-
748 cell transcriptomics: implications for ontological representation. *Hum Mol Genet.*
749 2018;27(R1):R40-R47. doi:10.1093/hmg/ddy100
750
- 751 27. Bakken TE, Hodge RD, Miller JA, Yao Z, Nguyen TN, Aevermann B, Barkan E,
752 Bertagnolli D, Casper T, Dee N, et al. 2018. Single-nucleus and single-cell
753 transcriptomes compared in matched cortical cell types. *PLoS One.*
754 2018;13(12):e0209648. doi:10.1371/journal.pone.0209648
755
- 756 28. Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A,
757 Lagarde J, Lin W, Schlesinger F, et al. 2012. Landscape of transcription in
758 human cells. *Nature.* 2012; 489: 101–8.
759
- 760 29. Huang DW, Sherman BT, Lempicki RA. 2009. Bioinformatics enrichment tools:
761 paths toward the comprehensive functional analysis of large gene lists. *Nucleic*
762 *Acids Res.* 2009;37(1):1-13.
763
- 764 30. Scheuermann RH, Novotny M, Aevermann B, Ben-Othman R, Liu A,
765 Sadarangani M, Kollmann T. Differential abundance of mDC subsets predict
766 response to Hepatitis B vaccination. *J Immunol* May 1, 2018, 200 (1 Supplement)
767 180.1.
768
- 769 31. Perkel JM. 2019. Starfish enterprise: finding RNA patterns in single cells. *Nature.*
770 2019;572(7770):549-551. doi:10.1038/d41586-019-02477-9
771
- 772 32. Cui Y, Zheng Y, Liu X, Yan L, Fan X, Yong J, Hu Y, Dong J, Li Q, Wu X, et al.
773 2019. Single-Cell Transcriptome Analysis Maps the Developmental Track of the
774 Human Heart. *Cell Rep.* 2019;26(7):1934-1950.e5.
775 doi:10.1016/j.celrep.2019.01.079
776
- 777 33. Smith SJ, Sümbül U, Graybuck LT, Collman F, Seshamani S, Gala R, Gliko O,
778 Elabbady L, Miller JA, Bakken TE, et al. 2019. Single-cell transcriptomic
779 evidence for dense intracortical neuropeptide networks. *Elife.* 2019;8:e47889.
780 doi:10.7554/eLife.47889
781
- 782 34. Moffitt JR, Hao J, Wang G, Chen KH, Babcock HP, Zhuang X. 2016. High-
783 throughput single-cell gene-expression profiling with multiplexed error-robust
784 fluorescence in situ hybridization. *Proc Natl Acad Sci U S A.*
785 2016;113(39):11046-11051. doi:10.1073/pnas.1612826113

787 **Supplemental Figures (legends only)**

788 **sFigure 1:** Empirical determination of the beta weighting parameter for the F-beta
789 score. All 75 clusters from full Middle Temporal Gyrus data were used to estimate the
790 average true positive (TP), false positive (FP), false negative (FN), and true negative
791 (TN), and number of markers for all clusters at beta values of 0.01, 0.5, 1, 1.5 and 2.

792 **sFigure 2:** Correlation between F-scores and average Binary Expression Scores per
793 cluster for v1.3 and v2.0. The F-score per cluster was computed using a beta=1 for NS-
794 Forest v2.0 to make both versions comparable.

795 **sFigure 3:** Quantitative assessment of Nearest-Neighbor Preservation metric (NNP, by
796 Venna et al. and IM). In brief, this is computed as follows: for each data point, the K-
797 Nearest-Neighborhood (KNN) in the high-dimensional space is compared with the KNN
798 in the reduced-dimensional space. Average precision/recall curves are generated by
799 taking into account high-dimensional neighborhoods of increasing size up to $K_{max} = 50$.
800 The True-Positive number is the intersection between high-dimensional and the low
801 dimensional neighborhood based on 157 selected genes. The precision is computed as
802 TP/K and the recall as TP/K_{max} . In A) red curve: NNP curves for the random-forest
803 selected 157 genes, while blue curves: NNP curves for 50 random gene sets of the
804 same size (selected from the full 5574 high-variance gene set). In B) green curve is the
805 tSNE generated from the complete 5574 variable genes, the red curve: NNP curves for
806 the random-forest selected 157 genes, while blue curves: NNP curves for 50 random
807 gene sets of the same size (selected from the full 5574 high-variance gene set)

808

809 sTable 1: Complete NS-Forest results for v1.3

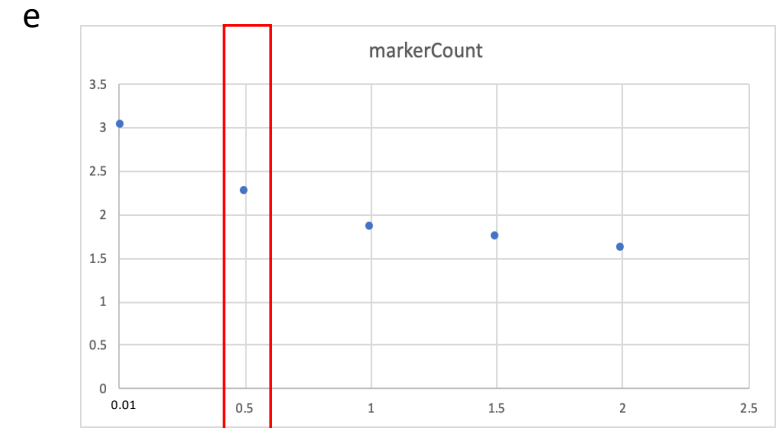
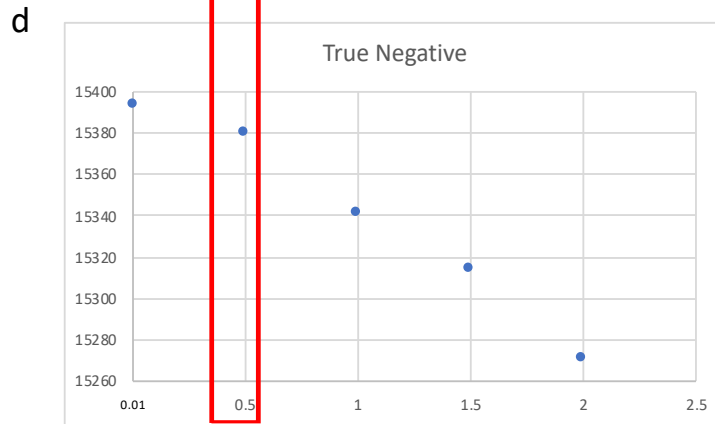
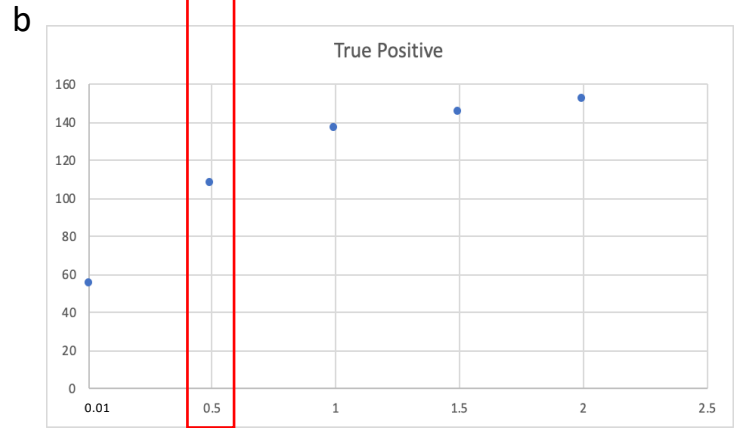
810 sTable 2: Complete NS-Forest results for v2.0

811 sTable 3: Supplemental ranked binary results from NS-Forest v2.0

812 sTable 4: Enrichment of annotations

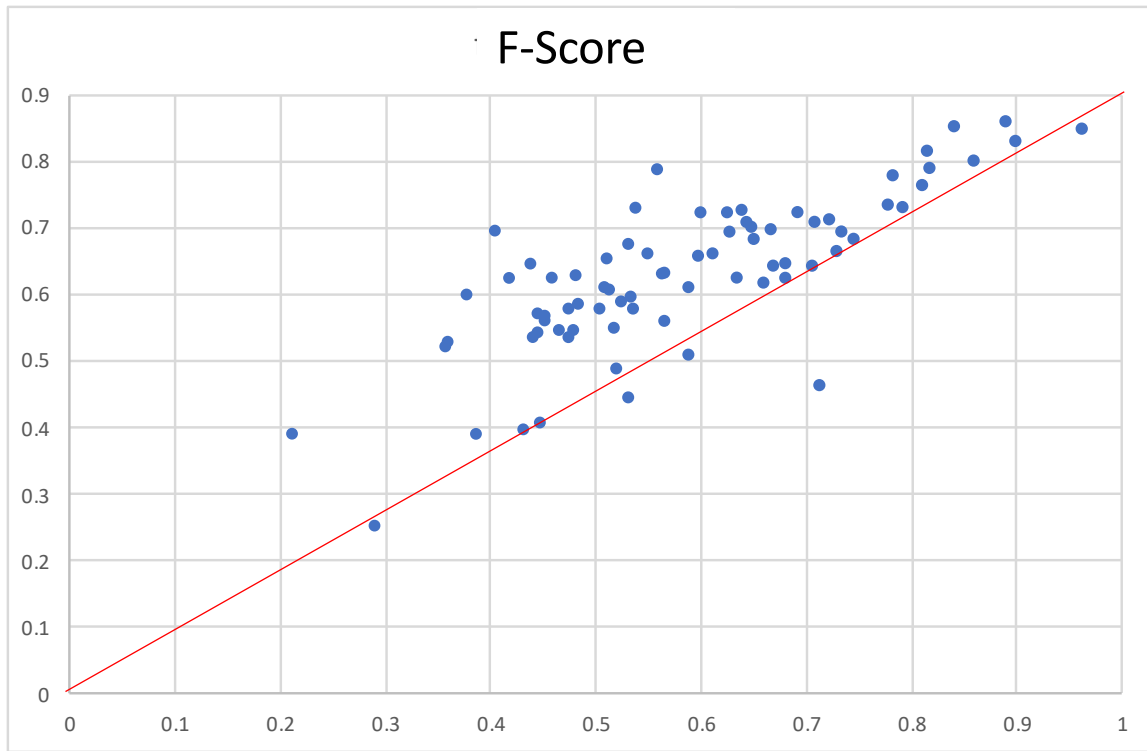
813

of events



Beta score

NS-FOREST V1.3



NS-FOREST V2.0

