**Title Page**

**Title:** FR-Match: Robust matching of cell type clusters from single cell RNA sequencing data using the Friedman-Rafsky non-parametric test

**Authors:** Yun Zhang[1], Brian D. Aevermann[1], Trygve E. Bakken[2], Jeremy A. Miller[2], Rebecca D. Hodge[2], Ed S. Lein[2], Richard H. Scheuermann[1,3,4]

**Affiliations:** [1]J. Craig Venter Institute, La Jolla, CA, USA; [2]Allen Institute for Brain Science, Seattle, WA, USA; [3]Department of Pathology, University of California San Diego, La Jolla, CA, USA; [4]Division of Vaccine Discovery, La Jolla Institute for Immunology, La Jolla, CA, USA

**Corresponding Author:** Richard H. Scheuermann, 858-200-1876, rscheuermann@jcvi.org

## Abstract

Single cell/nucleus RNA sequencing (scRNAseq) is emerging as an essential tool to unravel the phenotypic heterogeneity of cells in complex biological systems. While computational methods for scRNAseq cell type clustering have advanced, the ability to integrate datasets to identify common and novel cell types across experiments remains a challenge. Here, we introduce a cluster-to-cluster cell type matching method – FR-Match – that utilizes supervised feature selection for dimensionality reduction and incorporates shared information among cells to determine whether two cell type clusters share the same underlying multivariate gene expression distribution. FR-Match is benchmarked with existing cell-to-cell and cell-to-cluster cell type matching methods using both simulated and real scRNAseq data. FR-Match proved to be a stringent method that produced fewer erroneous matches of distinct cell subtypes and had the unique ability to identify novel cell phenotypes in new datasets. *In silico* validation demonstrated that the proposed workflow is the only self-contained algorithm that was robust to increasing numbers of true negatives (i.e. non-represented cell types). FR-Match was applied to two human brain scRNAseq datasets sampled from cortical layer 1 and full thickness middle temporal gyrus. When mapping cell types identified in specimens isolated from these overlapping human brain regions, FR-Match precisely recapitulated the laminar characteristics of matched cell type clusters, reflecting their distinct neuroanatomical distributions. An R package and Shiny application are provided at https://github.com/JCVenterInstitute/FRmatch for users to interactively explore and match scRNAseq cell type clusters with complementary visualization tools.

**Keywords:** single cell RNA sequencing, data integration, feature selection, cell types, cellular neuroscience, non-parametric test

## 1  Introduction

2    Global collaborations, including the Human Cell Atlas [1] and the NIH BRAIN Initiative [2], are

3    making rapid advances in the application of single cell/nucleus RNA sequencing (scRNAseq) to

4    characterize the transcriptional profiles of cells in healthy and diseased tissues as the basis for

5    understanding fundamental cellular processes and for diagnosing, monitoring, and treating

6    human diseases. The standard workflow for processing and analysis of scRNAseq data

7    includes steps for quality control to remove poor quality data based on quality metrics [3-5],

8    sequence alignment to reference genomes/transcriptomes [6-8], and transcript assembly and

9    quantification [8, 9] to produce a gene expression profile (transcriptome) for each individual cell.

10    In most cases, these expression profiles are then clustered [10-13] to group together cells with

11    similar gene expression phenotypes, representing either discrete cell types or distinct cell states.

12    Once these cell phenotype clusters are defined, it is also useful to identify sensitive and specific

13    marker genes for each cell phenotype cluster that could be used as targets for quantitative PCR,

14    probes for *in situ* hybridization assays, and other purposes (e.g. semantic cell type

15    representation where biomarkers can be used for defining cell types based on their necessary

16    and sufficient characteristics [14, 15]).

17    A major challenge emerging from the broad application of these scRNAseq technologies is the

18    ability to compare transcriptional profiles across studies. In some cases, basic normalization [16,

19    17] or batch correction [18, 19] methods have been used to combine multiple scRNAseq

20    datasets with limited success. Recently, several computational methods have been developed

21    to address this challenge more comprehensively [20-25]. General steps in these methods

22    include feature selection/dimensionality reduction and quantitative learning for matching. Scmap

23    [20] is a method that performs cell-to-cell (scmapCell) and cell-to-cluster (scmapCluster)

24    matchings. The feature selection step is unsupervised and based on a combination of

3

25    expression levels and dropout rates, pooling genes from all clusters in the reference dataset.

26    Matching is based on agreement of nearest neighbor searching using multiple similarity

27    measures. Seurat (Version 3) [21, 22] provides a cell-to-cell matching method within its suite of

28    scRNAseq analysis tools. Feature selection is unsupervised and selects highly variable features

29    in the reference dataset to define the high-dimensional space. Both query and reference cells

30    are aligned in a search space projected by PCA-based dimensionality reduction and canonical

31    correlation analysis, to transfer cluster labels through "anchors".  Among many others [23-25],

32    these methods have focused on individual cell level strategies when comparing a query dataset

33    to a reference dataset, not relying on clustering results to guide supervised feature selection or

34    cluster-level matching.


35    Here, we present a supervised cell phenotype matching strategy, called FR-Match, for cluster-

36    to-cluster cell transcriptome integration across scRNAseq experiments. Utilizing *a priori* learned

37    cluster labels and computationally- or experimentally-derived marker genes, FR-Match uses the

38    Friedman-Rafsky statistical test [26, 27] (FR test) to learn the multivariate distributional

39    concordance between query and reference data clusters in a graphical model. In this

40    manuscript, we first illustrate the matching properties of FR test in this scRNAseq adaptation

41    using thorough simulation and validation studies in comparison with other popular matching

42    methods. We then use FR-Match to match brain cell types defined in the full thickness of human

43    middle temporal gyrus (MTG) neocortex with cell types defined in a Layer 1 dissection of MTG

44    using public datasets from the Cell Types Database of the Allen Brain Map (www.brain-

45    map.org). We also report the cell types that are consistently matched between the two brain

46    regions using multiple matching methods. An R-based implementation, user guide, and Shiny

47    application for FR-Match are available in the open-source GitHub repository:

48    https://github.com/JCVenterInstitute/FRmatch.


49    **Results**

4

### *FR-Match: cluster-to-cluster mapping of cell type clusters*

51  FR-Match, is a novel application of the Friedman-Rafsky test [26, 27], a non-parametric

52  statistical test for multivariate data comparison, tailored for single cell clustering results. FR-

53  Match takes clustered gene expression matrices from query and reference experiments and

54  returns the FR statistic with p-value as evidence that the query and reference cell clusters being

55  compared are matched or not, i.e. they share a common gene expression phenotype. The

56  general steps of FR-Match (Figure 1a) include: i) select informative marker genes using, for

57  example, the NS-Forest marker gene selection algorithm [14]; ii) construct minimum spanning

58  trees for each pair of query and reference clusters (different colors); iii) remove all edges that

59  connect a node from the query cluster with a node from the reference cluster, and iv) calculate

60  FR statistics and p-values by counting the number of subgraphs remaining in the minimum

61  spanning tree plots. Intuitively, the larger the FR statistic, the stronger the evidence that the cell

62  clusters being compared represent the same cell transcriptional phenotype.

63                                      [Figure 1 here]

### *Supervised marker gene selection provides unique cell type clusters "barcodes"*

65  We adopted the NS-Forest algorithm [14] v2.0 (https://github.com/JCVenterInstitute/NSForest)

66  to select informative marker genes for a given cell type cluster. Applying NS-Forest feature

67  selection to the cortical Layer 1 and full thickness MTG datasets produced a collection of 34 and

68  157 marker genes that, in combination, can distinguish the 16 cortical Layer 1 [28] and 75 full

69  MTG [29] cell type clusters, respectively. These markers include well known neuronal marker

70  genes like *SATB2*, *LHX6*, *VIP*, *NDNF*, *NTNG1*, etc. (Supplementary Figure 1). The selected

71  marker genes display on-off binary expression patterns producing, in combination, a unique

72  gene expression "barcode" for each cell cluster (Figure 1b). In addition to producing marker

73  genes for each of the individual cell type clusters, this composite barcode serves as an effective

5

74     dimensionality reduction strategy that captures gene features that are informative for every cell

75     type cluster. The collection of informative marker genes effectively creates an essential

76     subspace that reflects the composite cell cluster phenotype structure in the single cell gene

77     expression data. Thus, supervised feature selection by NS-Forest was used as the

78     dimensionality reduction step for the FR-Match method in this study. Although NS-Forest was

79     used for marker gene selection here, FR-Match is compatible with any feature

80     selection/dimensionality reduction approach that selects informative cluster classification

81     features.

82     ***Matching performance in cross-validation and simulation studies***

83     To assess the performance of FR-Match in comparison with other matching methods, we

84     generated cross-validation datasets utilizing the cortical Layer 1 data and its known 15 cell type

85     clusters for validation studies (excluding the smallest cluster in the original studies with too few

86     cells). Matching was performed using six implementations of the three core methods: FR-Match

87     (using NS-Forest genes), FR-Match incorporating p-value adjustment (FR-Match adj.), scmap

88     (scmapCluster) with default gene selection (500 genes based on dropout proportions), scmap

89     with NS-Forest marker genes (scmap+NSF), scmap with extended NS-Forest marker genes

90     (scmap+NSF.ext) (see Methods section), and Seurat with default gene selection (top 2000

91     highly variable genes). (Seurat with NS-Forest marker genes was not reported since the results

92     were similar to the results obtained using default marker genes.)

93     **<u>Cross-validation assessment of 1-to-1 positive matches</u>**

94     In the two-fold cross-validation study, half of the cells serve as the query dataset and the other

95     half as the reference dataset. Exactly one 1-to-1 true positive match should be identified for

96     each cluster. Figure 2a displays the average matching rate over the cross-validation iterations,

97     where true positives are expected to lay along the diagonal. Four implementations, FR-Match,

98      FR-Match adj., scmap+NSF.ext, and Seurat had excellent performance with 0.93~1 true

99      positive rates (TPR) calculated as the grand average of the diagonal entries. Scmap using its

100     default gene selection approach performed sub-optimally, especially for glial cell types. This is

101     likely due to the fact that informative marker genes for these cell types were not selected using

102     the dropout rate-based feature selection criterion (Supplementary Figure 2). However, using

103     NS-Forest marker genes (scmap+NSF) instead of its default genes resulted in a significant

104     improvement in scmap performance, suggesting that supervised feature selection is

105     advantageous for cell type matching in general. FR-Match implementations had median

106     matching accuracies approaching 0.98 and above, while the next tier performers,

107     scmap+NSF.ext and Seurat, had median accuracies around 0.95 (Figure 2b). Sensitivity and

108     specificity metrics further break down the accuracy measure and indicate the balance between

109     the diagonal (true positive, a.k.a. sensitivity) and off-diagonal (true negative, a.k.a. specificity)

110     matching performance. FR-Match after p-value adjustment is the only algorithm that identified

111     all positive matches. Most methods had very high specificities, whereas FR-Match adj. had

112     somewhat lower specificity due to slightly more false positives.

113                                        [Figure 2 here]

114     **Cross-validation assessment of 1-to-0 negative matches**

115     Leave-*K*-cluster-out cross-validation was used to test the performance of these methods under

116     circumstances where one or more cell phenotypes is missing from the reference datasets, i.e. a

117     situation where a novel cell type has been discovered. The left-out cluster(s) should have 1-to-0

118     match(s) and should be unassigned. While FR-Match implementations clearly identified the left-

119     out cluster as unassigned, other methods produced inappropriate matching when query cell

120     types were missing from the reference dataset (Figure 3). Figure 3a shows results for when the

121     i5 cluster was left out; Supplementary Figures 3-8 show results for when other cluster were left-

7

122    out in turn. Both FR-Match implementations easily identified the true negative match and

123    correctly labeled the query i5 cluster as unassigned. Other methods partially or primarily mis-

124    matched the query cluster (i5) to a similar yet distinct cluster (i1), as seen in the UMAP

125    embedding where the query i5 nuclei are nearest neighbors to the reference i1 nuclei

126    (Supplementary Figure 9). The accuracy measure for leave-1-cluster-out cross-validation again

127    suggests that the FR-Match method is the best performer with median accuracies approaching

128    0.99 (Figure 3b). Furthermore, as we removed more and more reference clusters, the FR-Match

129    method showed robust precision-recall curve that consistently outperformed default

130    implementations of scmap and Seurat in ROC analysis (Figure 3c). Seurat's curve deteriorated

131    because its current implementation lacks an option for unassigned matches; therefore, all cells

132    in the query dataset were forced to map somewhere in the reference dataset. Interestingly,

133    scmap implementations with NS-Forest selected features also had robust precision-recall

134    curves with respect to the increasing number of true negatives.

135                                          [Figure 3 here]

136    The leave-*K*-cluster-out cross-validation has important implications for the capability of each

137    matching method to detect novel cell types in new data sets that are not present in the

138    reference datasets when integrating single cell experiments. In this important use case, the FR-

139    Match method exhibits desirable properties for novel cell phenotype discovery.

140    **Simulation of under- and over-partitioning during upstream clustering**

141    Accurate cell type determination from scRNAseq analysis is dependent on accurate partitioning

142    of the cellular transcriptomes into clusters based on their similarity.  Existing neuroscientific

143    knowledge [28] suggests that the 15 cortical Layer 1 cell clusters are the current "optimal"

144    clustering of the human brain upper cortical layer scRNAseq data. By combining and splitting

145    these optimal cell type clusters, we simulated under- and over-partitioning scenarios of the

8

146    upstream clustering analysis. Figure 4a summarizes five cluster partitions ranging from 3 to 18

147    clusters with F-measure scores indicating the classification power of partition-specific marker

148    genes. The "Top nodes" under-partitioning combines clusters into the three top-level broad cell

149    type classes: inhibitory neurons, excitatory neurons, and non-neuronal cells, producing well

150    known GABAergic, glutamatergic, and neuroglia markers with high F-measure score. The "Mid

151    nodes" under-partitioning combines three groups of closely related GABAergic clusters – i1 + i5,

152    i3 + i4, and i6 + i8 + i9 – resulting in 11 clusters. Over-partitioning of either one (e1) or three (i1,

153    i2, and i3) clusters was performed by running k-means clustering with k = 2 independently for

154    each cluster to simulate real over-partitioning scenarios.


155                                    [Figure 4 here]


156    It is important to note that over- and under-partitioning will also have an effect on the gene

157    selection step; it would be predicted that marker gene selection algorithms would have difficulty

158    finding maker genes specific for over-partitioned clusters, which would be reflected in the drop

159    in F-measure scores. Indeed, particularly low F-measure scores may be a good indication of

160    cluster over-partitioning. Figure 4b describes the expected effects on marker gene identification

161    and FR-Match performance after p-value adjustment when clusters are under-, optimally-, and

162    over-partitioned. The types of marker genes that would be selected with different reference

163    cluster partitioning scenarios would impact their ability to effectively drive cluster matching.


164    Supplementary Figures 10-15 show the matching results of all considered matching methods in

165    various partitioning scenarios. The FR-Match and Seurat methods showed good quality and

166    expected matching results in most partitioning scenarios; scmap had the same problem with the

167    unmatched glial clusters. Seurat showed excellent performance when reference clusters were

168    under-partitioned, but poor performance when query clusters were under-partitioned. Overall,

169    the FR-Match method had stable matching performance in the cluster partitioning simulations.

9

170    Indeed, 1-to-many and many-to-1 matching results using FR-Match could possibly indicate

171    under- or over-partitioning of the upstream clustering step in scRNAseq data analysis.

172    **Simulation of scenarios in which imperfect marker genes are included**

173    Though we recommend using the NS-Forest algorithm to select the minimum set of informative

174    marker genes, users may also want to use their own feature list as the input to FR-Match. There

175    may be other cases where non-informative marker genes have been included. In order to

176    assess the performance of FR-Match with respect to less than ideal marker gene lists, we use

177    simulation to evaluate the matching performance in two scenarios: i) when there are non-

178    informative (i.e. noisy) genes in the features selected, and ii) when some informative marker

179    genes are missing from the feature list with or without non-informative genes. Throughout this

180    simulation study, the FR-Match adj. implementation was used.

181    To simulate scenario (i), we used the 32 NS-Forest marker genes associated with the 15 cell

182    types in the Layer 1 data, together with randomly selected genes from the 16,497 available

183    genes in the dataset. In this scenario, the barcoding pattern of the informative marker genes

184    were preserved, whereas the random genes showed more noisy and non-specific expression

185    patterns in the "barcode" plots (Supplementary Figure 16a). In the simulations, we increased the

186    number of extra genes added from 1 to 15; FR-Match was very robust to noisy genes in each

187    simulated case with true positive rate staying close to 1 (Supplementary Figure 16b). Other

188    performance measures – accuracy, sensitivity (true positive rate), and specificity (true negative

189    rate) – all stayed well-above 0.9, suggesting that the overall performance of FR-Match was

190    stable and robust, even when the marker gene list contained up to 30% non-informative genes

191    (15 extra genes) (Supplementary Figure 16c). Increasing the number of non-informative genes

192    may slightly impact the specificity due to more false positives (off-diagonal intensities in

10

193    Supplementary Figure 16b) and therefore leads to the slight downward trend of the overall

194    accuracy.

195    For simulation scenario (ii), we generated two subcases to illustrate the impact of interfering

196    with different combinations of marker genes on the matching performance. In the first subcase,

197    we removed marker genes for three very distinct cell types: an excitatory cell type (e1), a glial

198    cell type (g1), and an inhibitory cell type (i1); and used the remaining NS-Forest marker genes

199    to match all cell types in the Layer 1 dataset. Surprisingly, each cell type was matched correctly

200    most of the time with an overall true positive rate of 0.98 (Supplementary Figure 17a). We also

201    replaced the removed marker genes with the same number of random genes; the matching

202    performance was also very good, and the impact of the changes in the marker gene list was

203    insignificant (Supplementary Figure 17a). In the second subcase, we considered

204    removing/replacing the marker genes for two related inhibitory cell types: i1 and i2. Without

205    marker genes that distinguish these similar cell types, FR-Match matched the i1 and i2 cell

206    types to each other (i.e. a many-to-many match) while maintaining the distinction from other cell

207    types with informative classification markers (Supplementary Figure 17b). The "barcode" plots

208    for i1 and i2 became generally non-selective with random expression of some other inhibitory

209    markers in the background (Supplementary Figure 17c). Such indistinct "barcode" plots may be

210    an effecting warning for many-to-many matches.  The absence of good classification markers is

211    most harmful to specificity (due to false positives), while sensitivity (true positive rate) remains

212    high (Supplementary Figure 17d).

213    In summary, as long as informative marker genes with good classification power are selected,

214    FR-Match is robust to other non-informative genes included in the feature list. Many-to-many

215    matching results by FR-Match may be a good indicator of the absence of informative marker

216    genes between the mis-matched cell types.

217    *Cell type mapping between cortical Layer 1 and full MTG*

218    We next extended the validation testing to a more realistic real-world scenario where a new

219    dataset has been generated in the same tissue region using slightly different experimental and

220    computational platforms. We tested FR-Match with p-value adjustment using two single nucleus

221    RNA sequencing datasets from overlapping human brain regions – the single apical layer of the

222    MTG cerebral cortex (cortical Layer 1), in which 16 discrete cell types were identified [28], and

223    the full laminar depth of the MTG cerebral cortex, in which 75 distinct cell types were identified

224    [29]. We selected NS-Forest combinatorial marker genes separately for each dataset. The

225    marker gene sets may contain overlapping genes for some cell types, e.g. *CUX2* is a useful

226    marker gene for more than one layer 2-3 cell types in combination with other marker genes;

227    classification power of these combinatorial marker genes are evaluated in detail in another

228    study [30].

229    Matching results were assessed from two perspectives: i) agreement with prior knowledge such

230    as layer metadata from the design of these experiments [28, 29], and ii) agreement with other

231    matching methods. Since these datasets targeted the same cortical region with overlapping

232    laminar sampling, we expect that matching algorithm should find 1-to-1 matches of each cell

233    types in the cortical Layer 1 data to one cell type in layers 1-2 from the full MTG data. The final

234    matching results were concluded from two matching directions: Layer 1 query to MTG reference

235    with MTG markers, and MTG query to Layer 1 reference with Layer 1 markers. The two-way

236    matching approach was applied to all comparable matching algorithms.

237    **FR-Match uniquely maps cell types reflecting the overlapping anatomic regions**

238    Using FR-Match, we mapped each of the 13 Layer 1 cell types uniquely to one MTG cell type

239    (Figure 5a), i.e. 1-to-1 two-way matches. These matches precisely reflect the overlapping

240    anatomic regions in these two independent experiments in that the matched MTG cell types all

12

241    have an "L1" layer indicator in their nomenclature. The one exception for the Layer 1 e1 cluster

242    likely reflects the incidental capture of upper cortical layer 2 excitatory neurons in the original

243    Layer 1 experiment [28]. And while most of the *SST* cell subtypes are located in deeper cortical

244    layers, FR-Match specifically selected the small number of L1 *SST* clusters as top matches. The

245    same was true for *VIP* and *LAMP5* cell subtypes. The minimum spanning tree plots produced by

246    FR-Match provide a clear visualization of matched and unmatched cell clusters (Figure 5b).

247                                        [Figure 5 here]

248    To validate further, we compared the matching results to the hierarchical taxonomy of MTG cell

249    types [29], which reflects cell type relatedness (left side of Figure 5a). First, the block of one-

250    way matches in Box A precisely corresponds to a specific sub-clade of *VIP*-expressing cells with

251    close lineage relationships, suggesting that one-way FR-Match results are evidence of closely

252    related cell types. Second, FR-Match correctly identified excitatory neurons that were

253    incidentally captured from upper Layer 2 in the cortical Layer 1 experiment in Box B,

254    corresponding to L2/3 excitatory neurons in the full MTG dataset. Third, Box C suggests under-

255    partitioning of the Layer 1 astrocyte cluster as multiple two-way matches were found for the

256    same cluster.

257    Directional one-way matching results are shown in Supplementary Figure 18. Though different

258    matching patterns are observed from each direction, they reflect the fact that these datasets are

259    measuring different cell types. There are some cases where the difference might be due to the

260    cell complexity in the datasets, e.g. the *VIP* or *SST* types, and this might be leading to the

261    dynamic range and skewness of p-value distributions for each query cluster.

262    **Cell type mapping using other existing approaches**

13

263     In mapping cell types between cortical Layer 1 and the full MTG, both FR-Match and Seurat

264     produced similar unique two-way matches (Figure 6). Examining all matching results and all

265     matching algorithms, FR-Match produced the most "conservative" mapping of cell types. The

266     other matching algorithms produced matching results that had more sparsely-distributed *VIP*

267     types (Box A), and were not laminar specific (Box B). Among all approaches, glial cell types

268     were mapped somewhat differently (Box C), probably due to their overall lower sampling and

269     distinct phenotypes compared to the majority of GABAergic and glutamatergic neurons.


270                                    [Figure 6 here]


271     FR-Match shows three advantages over the alternative methods. First, by using supervised

272     feature selection for each cell type, major and minor cell populations are equally represented in

273     the reduced-dimensional space for cell type matching. This strategy would also benefit other

274     matching methods with sub-optimal feature selection/dimensionality reduction. Second, FR-

275     Match clearly excludes the matching of cell types that are only present in one of the datasets.

276     Third, FR-Match allows one-to-multiple and unassigned matches, which allows for detecting

277     potential cluster partitioning issues and the discovery of novel cell types.


278     The other existing cell-level matching approaches naturally provide the probabilistic cluster-level

279     matching of cell types as the percentage of matched cells in query cluster (Supplementary

280     Figures 19-22); a deterministic cluster-level match would depend on the selection of an *ad-hoc*

281     cutoff of the probabilistic matching. Thus, deterministic cell type mapping or discovery of novel

282     cell types would be difficult as i) individual cells may be alike in the same broad cell class even if

283     the specific cell type may not be present in the reference dataset, and ii) the probabilistic cutoff

284     may be subjective. Therefore, both scmap and Seurat identified many more non-specific one-

285     way matches than FR-Match, which uses an objective p-value cutoff.

286    Combining all results, we finally report 15 high-confidence ensemble matches between Layer 1

287    and full MTG cell types in Supplementary Table 1.

**The effects of alternative gene selection and cell clustering methods on matching**

**performance**

290    To further elucidate the impact of alternative gene selection or cell clustering choices on cluster

291    matching, we performed the following analyses.

292    In the two brain datasets, cell types are defined and characterized by a domain knowledge-

293    guided iterative clustering [13] and transcriptomically-derived markers [28, 29]. The

294    nomenclature used to describe these cell types consists of the broad cell class (inhibitory,

295    excitatory, and glial cells), layering information (for the MTG dataset), one marker gene for the

296    subclass node in the taxonomy tree (e.g. *VIP*, *SST*, etc.), and one marker gene for the leaf node

297    cluster. For example, the "Inh_L1_2_PAX6_CDH12" from the MTG dataset means the inhibitory

298    neurons located in layer 1-2 within the *PAX6*-subclass/subbranch expressing *CDH12*. The leaf

299    node marker genes are preferentially selected by a binary scoring scheme [29] different from

300    the one used by NS-Forest. Thus, the "cell type naming genes" provide an alternative

301    informative marker gene set.

302    To assess matching performance using a different set of informative marker genes, we replaced

303    the NS-Forest marker genes by these cell type naming genes for both datasets, followed by the

304    same matching approaches. 26 and 87 naming genes were defined for the Layer 1 and full

305    MTG datasets, respectively, out of which, 9 and 18 genes are in common between the naming

306    genes and the NS-Forest marker genes, respectively. Using cell type naming genes, FR-Match,

307    scmap, and Seurat all performed slightly differently with less ideal matching patterns

308    (Supplementary Figure 23). Overall fewer matches were identified; and the identified matches

309    were less specific (i.e. mapping to neighboring cell types). This is probably because using only

15

310     one leaf node marker gene may not be enough to fully capture the differences between those

311     closely related leaf node cell types. From these matching results, we may conclude that NS-

312     Forest selects better sets of informative markers than the other approach in this example, which

313     has an impact on all three matching methods; less optimal feature selection will negatively

314     impact matching regardless of the matching methods.

315     In another analysis, we compared the matching performance of FR-Match, scmap, and Seurat

316     with respect to a different clustering method. The community detection Louvain method [10] is

317     one of the most commonly used clustering methods for scRNAseq analysis. We applied Louvain

318     clustering (implemented in the Seurat R package, with resolution = 1) to the full MTG dataset,

319     which resulted in 26 reasonably segregated clusters in the UMAP low-dimensional embedding

320     space (Supplementary Figure 24a). Matching results with the Louvain clusters are shown in

321     Supplementary Figure 24b. FR-Match produced similar matching results regardless of the

322     clustering methods: each Layer 1 cluster is strongly matched (two-way match) to some Louvain

323     cluster of the full MTG dataset. Many-to-one and one-to-many matches are observed since the

324     generic Louvain method appears to have under-partitioned the data in comparison with the

325     original expert-curated iterative clusters, which agrees with the matching patterns we observed

326     in our simulations. Matching by scmap and Seurat with the Louvain clusters shows the same

327     problems as with the original clusters, i.e. excessive unassigned matches (scmap), and non-

328     specific matches of the Layer 1 excitatory cluster (scmap and Seurat). Using different clustering

329     methods will lead to different matching results depending on the clustering quality. As long as

330     the clusters are reasonably good, FR-Match is able to detect high quality matches regardless of

331     the clustering methods.

332     **Cell type matching using batch integration**

333    To date, there are more than 10 methods that have been proposed to correct the batch effects

334    of scRNAseq data; among them, Harmony [31], LIGER [32], and Seurat 3 [21] are the

335    recommended algorithms for batch integration [33]. Only Seurat is an end-to-end pipeline that

336    inputs multiple scRNAseq data batches and outputs cell-to-cell alignment between batches. By

337    summarizing the cell-level batch integration with prior cluster memberships of the cells, we

338    compared the performance of Seurat for cell type matching with FR-Match in previous

339    subsections. In this subsection, we implemented a workaround for Harmony and LIGER to

340    transfer the batch integration outputs to produce putative cell type matches.

341    We applied Harmony (Supplementary Figure 25-26) and LIGER (Supplementary Figure 27-28)

342    individually to integrate the Layer 1 and MTG datasets; both methods showed effective "batch-

343    effect" removal in the UMAP (Supplementary Figure 25b-c) or tSNE (Supplementary Figure

344    27b-c) low-dimensional embedding. For both Harmony and LIGER, the outputs from the

345    algorithms are the integrated cells in some dimensionally reduced spaces; joint clustering can

346    then be conducted on the integrated data spaces (Supplementary Figure 25d, Supplementary

347    Figure 27d); and cell type matching can be inferred from the "river" plots (Supplementary Figure

348    26a, Supplementary Figure 28a) between the input batches through the common joint clusters.

349    We transferred the river plot to a one-to-one correspondent cell type matching heatmap, with

350    each match indicating there exists a path between the two cell types in the river plot. Note that

351    the heatmap is non-directional for a given set of edges of the river plot. Through such a

352    workaround, we obtained cell type matching results for Harmony (Supplementary Figure 26b)

353    and LIGER (Supplementary Figure 28b) in a similar format as FR-Match. It is clear that the

354    batch integration approaches produce matches in blocks (i.e. many-to-many matches), and do

355    not effectively yield the specific matches within these blocks if multiple related cell subtypes are

356    presented. These batch integration methods were not originally designed for the task of cell type

357    integration; therefore, it is not surprising that they produce sub-optimal results.

17

358 **Discussion**

359 FR-Match offers a cluster-level approach for mapping cell phenotypes identified in scRNAseq

360 experiments. It extends the current cell-level matching algorithms by: i) borrowing information

361 from all the cells in the same cluster using a statistical test that provides both probabilistic

362 matching in p-values and objective p-value thresholds for deterministic matching, and ii)

363 providing simple visualization of cell type data clouds in the minimum spanning tree graphical

364 representation. Matching results of FR-Match are relatively conservative yielding highly specific

365 matches, which can confirm cell type equivalence, lead to novel cell type discovery, and

366 diagnose upstream clustering problems. Among many other scRNAseq data integration

367 strategies, this approach combines informative feature selection and cluster-level integration of

368 the NS-Forest and FR-Match software suites, producing intuitive results with high interpretability,

369 including useful intermediate results such as binary marker genes and minimum spanning tree

370 graphs for users to monitor and gain meaningful insights from the mapping solutions.

371 Based on the computational and statistical investigation of both simulated and real datasets, we

372 conclude that: i) the FR-Match and Seurat methods show excellent performance in mapping

373 neuronal and glial cell types using snRNAseq data from human brain; and ii) supervised feature

374 selection, such as the NS-Forest algorithm, appears to produce excellent marker gene

375 combinations that can be used as an effective feature selection/dimensionality reduction

376 technique for cell type mapping with multiple methods, including FR-Match and scmap. Scmap

377 is a consensus method that requires at least two of the three association metrics – cosine

378 similarity, Pearson and Spearman correlations – to be in agreement as the last step to

379 determine a match, thus the comparative analysis results of the matching methods reported

380 here may also serve as a reference guide for matching performance using those association

381 metrics.

18

382    One of the biggest challenges in scRNAseq alignment at the moment seems to be the proper

383    assignment of cells from a cell type found in only one dataset. These cells are often matched to

384    a closely related cell type in a second dataset. In this regard, FR-Match appears to be superior

385    in being able to determine which cell types from two datasets are *not* matched, for novel cell

386    type discovery.

387    For all compared methods in this study, it's interesting to note that under-partitioning the query

388    clusters leads to degraded performance, except if the reference clusters are also under-

389    partitioned. This suggests that a useful strategy would be to map to reference types in a

390    hierarchical manner by first mapping to broad classes of references types and then moving

391    down the tree to finer types until ambiguous matches appear. The negative effect of under-

392    partitioned clusters also applies to the nested classes of heterogeneous cell types.

393    Automated cell type integration of independent scRNAseq datasets remains challenging.

394    Creating an unbiased, high-resolution and comprehensive cell type reference would be a critical

395    task for the whole single cell research community. Consensus mapping schemes that survey

396    both cell-level and cluster-level matchings will be useful for establishing such a reference data

397    atlas. We believe that final mapping of the brain cell types agreed upon by the type of bi-

398    directionally and multi-level matchings reported here represents the best-practice for

399    computational cell type mapping, requiring minimal expert intervention.

400    Single cell evaluation is a fast-evolving field. Although not fully explored here, we expect FR-

401    Match to be applicable to cross-platform, cross-specimen, cross-anatomy, and cross-species

402    matching of scRNAseq clustered data. The effect of dropouts and the dynamic range of single

403    cell sequencing data from protocols other than the Smart-seq [34] protocol stand out as key

404    challenges to be overcome. To address these challenges, we are now developing add-on

405    features to the core FR-Match algorithm, including imputation techniques [35] for the relatively

406      high dropout rates in 10X Genomics droplet-based protocols [36], and moment-based

407      normalization options [37] for the discrete and dispersed values produced in single cell spatial

408      *in-situ* hybridization protocols [38-40]. Preliminary results of mapping Smart-seq cell clusters to

409      10X cell clusters suggest that FR-Match will be useful for cross-platform cell type matching

410      when appropriate dropout imputation and data normalization upstream steps are included in the

411      computational pipeline (data not shown). While these emerging technologies will produce more

412      complicated data integration challenges, the adaptation of methods like FR-Match are poised to

413      play an essential role in the broad integration of scRNAseq cell phenotyping experiments.

414      **Methods**

415      ***The cell type matching problem***

416      Consider two single cell RNA sequencing experiments – one query/new experiment and one

417      reference experiment. A cell-by-gene expression matrix for each experiment is obtained by

418      standard scRNAseq data processing and analysis workflows, including quality control, reference

419      alignment, sequence assembly, and transcript quantification. Cell cluster labels are also

420      obtained from clustering analysis using, for example, the community detection Louvain

421      algorithm [10], and/or other domain specific knowledge. These cell clusters represent

422      transcriptionally-distinct cellular phenotypes within each experiment. The cell type matching

423      problem is whether a pair of query and reference cell clusters identified in related but

424      independent experiments are instances of the same or different transcriptionally-defined cell

425      phenotypes.

426      We propose a computational solution to the cell type matching problem – FR-Match – an

427      adaptation the Friedman-Rafsky statistical test for scRNAseq data, which takes two input

428      datasets (query and reference) each with a gene expression matrix and cell cluster membership

429      labels (Figure 1a).  Importantly, FR-Match uses a set of informative marker genes that

20

430    characterize the reference cell type clusters. Dimensionality reduction is done by imposing the

431    same set of marker genes on the query dataset, to select the most informative features shared

432    with the reference dataset. For each pair of cross-dataset clusters, we perform cluster-to-cluster

433    matching via the Friedman-Rafsky statistical test. As a result, FR-Match outputs the following

434    types of match (format: query-to-reference): 1-to-0 or unassigned (indicative of a novel cell type),

435    1-to-1 (indicative of a uniquely matched cell type), 1-to-many (indicative of an under-partitioned

436    query cluster or over-partitioned reference cluster), many-to-1 (indicative of an over-partitioned

437    query cluster or an under-partitioned reference cluster).

438    ***Necessary and sufficient marker gene identification by random forest***

439    In order to perform dimensionality reduction, random forest machine learning as implemented in

440    the NS-Forest algorithm [14, 15, 30] (v2.0 at https://github.com/JCVenterInstitute/NSForest)

441    was used to select necessary and sufficient marker genes for each reference cell type cluster.

442    NS-Forest includes steps for: i) feature selection, ii) feature ranking, and iii) minimum feature

443    determination. Let $X$ be an $n \times m$ dimensional cell-by-gene matrix, where $n$ is the number of

444    cells and $m$ is the number of genes. Let $y$ be an $n \times 1$ vector of cluster labels. In step (i),

445    random forest models, with 10,000 decision trees each, are built for input data $X$ and each

446    cluster label in $y$ under a binary classification scheme. From each random forest model, the

447    average information gain based on the Gini index for each gene is extracted, which is then used

448    as a measure of feature importance to rank the gene features. In step (ii), for the top 15 ranked

449    genes, a binary expression score for gene $g$ in cluster $k$ is calculated as

450    $$\text{Score}_{g,k} = \frac{\Sigma_{k'=1}^{K}\left(1 - \frac{med_{g,k'}}{med_{g,k}}\right)^{+}}{K-1},$$

451    where $med_{g,k}$ is the median expression level of gene $g$ in cluster $k$, $K$ is the total number of

452    clusters, and $(\cdot)^{+}$ defines the non-negative value of the equation. The binary expression score

453    ranges from 0 to 1, where 1 indicates absolute binaryness, i.e. the gene exclusively expressed

454    in the target cluster and not at all in non-target clusters. In step (iii), the top 6 genes from step (ii)

455    are selected and all combinations are evaluated by the F-beta score. F-beta is an F-measure

456    weighted by $\beta$ such that

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

458    $\beta = 0.5$ was set to weight precision more than recall, which compensates the effect of false

459    negatives dropouts due to technical artifacts in scRNAseq experiments. The output from step (iii)

460    is a minimum set of marker genes for each cell type cluster (usually $1 - 4$), whose expression in

461    combination is sufficient to discriminate the target cell type cluster from the rest of the cells. In

462    addition to the minimum set of NS-Forest marker genes, the algorithm also provides an

463    extended list of binary marker genes as a supplementary output from step (ii), which may

464    achieve higher discriminative power under certain circumstances. The top 15 NS-Forest genes

465    for each cell type formed an NS-Forest extended gene list as an alternative feature selection

466    option for matching algorithms. For a more detailed discussion of the choice of the number of

467    top genes used in NS-Forest v2.0, see Aevermann et al. [30]

468    ***Friedman-Rafsky test***

469    The Friedman-Rafsky (FR) test [26] is a multivariate generalization of the non-parametric two-

470    sample comparison problem. This classical statistical test is distribution free. Consider two

471    general distributions $F_X$ and $F_Y$ for samples $(x_1, \cdots, x_m)$ and $(y_1, \cdots, y_n)$ in a $k$-dimensional space,

472    respectively. (In the context of FR-Match, the $x$'s and $y$'s denote the expression profiles of each

473    cell in the query and reference clusters; $m$ and $n$ are the number of cells in each cluster; $k$ is

474    determined by the number of informative marker genes from the reference dataset). Under the

475    hypothesis testing framework, the original FR test is designed for testing

22

476 $$H_0: F_X = F_Y \qquad \text{versus} \qquad H_1: F_X \neq F_Y,$$

477 in which the null hypothesis states that the cells from both query and reference clusters are from

478 the same transcriptional distribution; the alternative hypothesis states that the two cell

479 populations are from different transcriptional distributions. Thus, the cell type matching problem

480 becomes a statistical test to detect comparisons for which $H_0$ is true.

481 The underlying model of the FR test is a graphical model based on the minimum spanning tree

482 of pooled samples (Figure 1a). In the multi-dimensional informative marker gene space, cells

483 from different clusters (indicated by colors) are pooled and form a mixture of data points. A

484 complete graph can be constructed, which connects all cells to each other and uses the edge

485 length to preserve the pairwise Euclidean distance between cells in the original space. Next, the

486 complete graph is trimmed to a tree graph that connects all cells with the minimum total length

487 of edges, i.e. the minimum spanning tree. Edges that connect cells of different clusters are then

488 removed and the number of disjoint subtrees is counted. Intuitively, if there are a large number

489 of subtrees, it implies that the pooled cells are closely interspersed and therefore more likely to

490 be from the same multivariate gene expression distribution.

491 Formally, let $R$ be the total number of subtrees – "multivariate runs" in the FR test framework,

492 with mean $E(R)$ and variance $\mathrm{Var}(R)$ directly derived from graph theory. The FR statistic is

493 defined as

494 $$W = \frac{R - E(R)}{\mathrm{Var}(R)^{1/2}}.$$

495 Friedman and Rafsky showed that the asymptotic distribution of $W$ follows a standard normal

496 distribution for large sample sizes:

497 $\qquad W \sim N(0,1)$ as $m, n \rightarrow \infty$ with $m/n$ bounded away from 0 and $\infty$. $\qquad$ (1)

23

498     For the hypothesis testing purpose, $H_0$ is rejected for small values of $W$, i.e. p-value is one-

499     sided such that $p = \Pr(W \leq w)$. Note that, in the cell type matching application, we determine a

500     *match* if $p > 0.05$, but other p-value thresholds could also be used.

501     ***FR-Match method***

502     Extending from the classical statistical test, FR-Match is a novel application of FR test to

503     approach the cell type matching problem with scRNAseq data. The full FR-Match algorithm not

504     only implements the basic testing procedure, but also adapts modifications for specific issues

505     pertaining to the scRNAseq application. A major issue is that two cell clusters to be compared

506     may have very different cluster sizes, such as a dozen cells versus hundreds of cells

507     (Supplementary Figure 29). The unbalanced cluster sizes will often cause two problems: i)

508     unstable statistical power as the ratio of cluster sizes deviates from the asymptotic condition,

509     and ii) exponentially long computational time needed for constructing minimum spanning tree for

510     large number of cells. To address these problems, an iterative subsampling scheme was

511     implemented, which repeatedly performs sampling without replacement of $S$ cells, or all cells if

512     $S$ > cluster size, from each cell cluster for $B$ times. Default values of $S$ and $B$ are 10 and 1000,

513     respectively, but are tunable. The median p-value of all iterations is outputted. Other

514     modifications include filtering small clusters with less than $C$ cells each, and p-value adjustment

515     for multiple hypothesis testing correction. Empirically, $C = 10$ was chosen for defining a cell type

516     cluster with high confidence since it appeared to provide enough cell instances to be

517     representative. It is suggested to set $S = C$, but it is not a necessary condition for the algorithm.

518     A disproportionate ratio of $S$ to $C$ would adversely affect the underlying statistical assumptions

519     due to the unmet asymptotic condition in Equation (1).

520     As an alternative to the asymptotic theory, permutation testing is a widely-accepted practical

521     choice for approximating the null distribution of the FR statistic in a hypothesis testing

24

522    framework [41]. We designed a simple technical simulation to compare the statistical properties

523    of the FR test, FR permutation test, and FR test with subsampling scheme, with respect to the

524    major pragmatic concern of imbalanced cluster sizes that specifically pertains to the cell type

525    matching problem. We generated multivariate data from a Multivariate Normal (MVN)

526    distribution ($k = 40$ dimensions). Random samples were drawn from $(x_1, \cdots, x_m) \sim MVN(\mu =$

527    $0, \Sigma = I)$ and $(y_1, \cdots, y_n) \sim MVN(\mu = 0 + \delta, \Sigma = I)$, where $I$ is the identity matrix. Under the null,

528    $\delta = 0$, i.e. no location difference between the $x$- and $y$-samples; under the alternative, we set

529    $\delta = 0.4$ for moderate location shift in their distributions. To simulate the imbalanced cluster sizes,

530    we fixed one cluster size $m = 10$ and varied the other cluster size $n = 10, 20, 100, 200$. The ROC

531    analysis (Supplementary Figure 30) confirm that the permutation test is a very good

532    approximation of the FR test based on asymptotic theory; however, both tests show

533    deteriorating ROC curves when the sample sizes were very imbalanced ($n = 200$, blue curve).

534    In contrast, FR test with subsampling shows the most ideal property – better ROC curve and

535    larger AUC value – as sample size (i.e. cluster size in this context) increases. Therefore, the

536    iteratively subsampling scheme was adopted in the FR-Match algorithm.

537    Though the subsampling parameter $S$ was initially chosen based on practical considerations, we

538    also provide more simulation results for guiding the choice of $S$ here. Based on the same

539    simulation design as above, we evaluated the AUC values for FR subsampling tests with

540    $S = 10, 20, 30$, and benchmarked with the FR test (Supplementary Figure 31). When both input

541    cluster sizes $m$ and $n$ vary from 10 to 200, the FR subsampling test with $S = 10$ outperforms all

542    other choices with the FR test showing the highest AUC values in all simulated cases with $m$

543    and $n$. This is potentially due to the expectation that the choice of $S$ should embrace the right

544    balance between gathering enough samples to represent the whole cluster and avoiding local

545    structures in the cluster (i.e. large subtrees of the same color in an MST). We believe this might

546    be related to the "effective" dimensionality of the data space characterized by $\Sigma$ and other

25

547     distributional properties, which will be an interesting topic for future statistical research. In this

548     manuscript, the choice of $S$ is supported by empirical evidence; readers should use their own

549     judgement on the choice of $S$ for their own datasets.

550     In the Layer 1 and full MTG matching analysis reported in this manuscript, tunable parameters

551     were set at the default values described above. When a sequence of FR-Match p-values were

552     computed for each pair of Layer 1 cell type and MTG cell type, Benjamini & Yekutieli [42] p-

553     value adjustment was applied for multiple hypothesis testing correction before the final

554     determination of a cell type match.

555     ***Determining cluster-level match for the cell-level matching methods***

556     In comparison with other popular matching methods, a voting rule was adopted after obtaining

557     the cell-level matching results from algorithms scmap (cell-to-cluster) and Seurat (cell-to-cell).

558     Scmap provides a map: query cell → reference cluster. We calculate the % of reference cluster

559     labels grouped by the query cell labels, and thereby obtain a quantitative measure ranging from

560     0 to 1 that indicates the probability of being the same cell type between the query and reference

561     cell clusters. Similarly, the Seurat alignment is extended to query cell → reference cell →

562     reference cluster, and calculate the cluster-to-cluster matching measure in the same way. For a

563     specific query cluster, its cluster-level match is determined by the votes of its member cells for

564     their mapped reference cluster labels. An ad-hoc threshold at 30% was used for defining a

565     deterministic match, which accounts for both the detection of a substantial proportion of query

566     cells matched to one reference cluster and the possibility that some query clusters might be

567     matched to multiple reference clusters. If the 30%-criterion is not met, then the query cluster is

568     defined as unassigned in the matching results. The cluster-level matching results may change

569     depending on the ad-hoc threshold used. For example, if changing the threshold to 40%, Seurat

570     would identify the same set of two-way matches, but with three fewer one-way matches

26

571     (Supplementary Figure 32). A data-driven decision on such a threshold can be guided by the

572     distribution of % of matched cells in Supplementary Figures 19-22.

573     ***Cross-validation and simulation design***

574     Data generation for the cross-validation and simulation studies were from the cortical Layer 1

575     data with 15 cell clusters [28] (excluding one cluster, i11, with too few cells). All cross-validation

576     designs  were two-fold by evenly splitting data into training and testing in proportion to the

577     original cluster sizes. All cross-validations were repeated 20 times each design.

578     Real data-guided simulations were used to mimic under-/over-partitioned scenarios (Figure 4).

579     "Top nodes" under-partitions are cells merged into three broad classes: GABAergic inhibitory

580     neurons, glutamatergic excitatory neurons, and neuroglial cells. "Mid nodes" under-partitions are

581     cells merged into similar inhibitory neurons according to the constellation diagram of cluster

582     network from the original study [28]; for the purpose of simulation, i1 and i5, i3 and i4, and i6, i8,

583     and i9 were merged. For over-partitions, large cell clusters were split by running k-means

584     clustering with k = 2 independently for each over-partitioned cluster. "Split e1" divided the

585     excitatory cluster into two sub-clusters of sizes 180 and 119 cells, resulting in 16 (= 15 + 1)

586     over-partitioned clusters. "Split i1, i2, i3" divided each of the inhibitory clusters into two sub-

587     clusters of sizes 56 and 34, 39 and 38, 32 and 24 cells, respectively, resulting in 18 (= 15 + 3)

588     over-partitioned clusters in total. NS-Forest marker genes were identified for each of the

589     simulated datasets. Matching performances of the under-/over-partitioned datasets were

590     evaluated through two-fold cross-validation repeated 20 times.

591     **Data availability**

592     Two published single-nucleus RNA-seq datasets from the Allen Institute of Brain Science of

593     human brain were used: i) cortical Layer 1 of middle temporal gyrus (MTG) [28] and ii) full

594    thickness MTG [29] (https://portal.brain-map.org/atlases-and-data/rnaseq/human-mtg-smart-

595    seq). The Layer 1 dataset contains expression data from 871 intact nuclei that form 16 cell type

596    clusters, including four non-neuronal type clusters, one excitatory neuron type cluster, and 11

597    inhibitory neuron type clusters. The MTG dataset contains filtered expression data from 15,603

598    nuclei that form 75 cell type clusters, subdivided into six non-neuronal type clusters, 24

599    excitatory neuron type clusters, and 45 inhibitory neuron type clusters. These cell type clusters

600    are regarded as transcriptionally distinct cell types with nomenclature asserted after iterative

601    clustering analysis [13]. Gene-level read count values were preprocessed to log-CPM (counts

602    per million) values for all nuclei.

603    The same high level data processing steps were used for both datasets, although the details

604    varied slightly:

605    1.  Whole postmortem brain specimens or neurosurgical tissue samples were collected from

606        adult male and female donors with 'control' condition (i.e. non-disease).

607    2.  Nuclei were isolated from microdissected tissue pieces to avoid damage to neurons [43],

608        and single nuclei were sorted using FACS instruments. The gating strategy included

609        doublet detection gates and gates on neuronal marker NeuN signal.

610    3.  RNA sequencing was performed using the SMART-Seq platform and multiplex library

611        preparation.

612    4.  STAR alignment of raw reads to human genome sequence, and sequence quantification

613        using standard Bioconductor packages were performed. Gene expression levels were

614        reported as counts per million (CPM) of exon and intron reads.

615    5.  Nuclei passing quality control criteria were included for clustering analysis.

616    6.  Iterative clustering procedure based on community detection were performed to group

617        nuclei into transcriptomic cell types [13]. Dropouts were accounted for while selecting

618        differentially expressed genes, and PCA was used for dimensionality reduction.

28

619    7. Clusters identified as donor-specific were flagged as outliers, and manually inspected for

620        cluster-level QC before exclusion.

621 **Key Points**

622    • Feature selection plays a key role in scRNAseq data integration of cell type clusters;

623      using supervised feature selection instead of approaches based on dropout rates

624      significantly improves the performance of existing cell type matching methods, e.g.

625      'scmap'.

626    • The random forest-based 'NS-Forest' marker gene selection algorithm is an effective

627      dimensionality reduction tool that produces an informative set of necessary and sufficient

628      genes for characterizing reference cell types.

629    • The cluster-level cell type matching method 'FR-Match', which builds upon a non-

630      parametric multivariate statistical test, shows robustness against missing reference cell

631      types, i.e. novel query cell types.

632    • FR-Match precisely matched common cell types from two independent scRNAseq

633      experiments that reflect the laminar characteristics of the two anatomically overlapping

634      brain regions.

635    • FR-Match software provides barcode plots and minimum spanning tree graphs for the

636      query and reference cell type clusters, which are user-friendly visualization tools for

637      insightful data exploration of scRNAseq data clusters.

638 **Funding**

**Author contributions**

Y.Z. and R.H.S. designed the study, conceived the statistical model, and wrote the manuscript. Y.Z. and B.D.A. developed the software suites. Y.Z. and B.D.A applied the software to real data analysis. Y.Z., B.D.A, T.E.B., J.A.M., and R.H.S. interpreted the real data analysis. R.D.H., T.E.B., J.A.M., R.H.S., and E.S.L. performed the single nucleus RNA sequencing experiments used. R.H.S. and E.S.L. supervised the work.

**Yun Zhang** is a Staff Scientist and Biostatistician in the Informatics Department at the J. Craig Venter Institute.

**Brian D. Aevermann** is Senior Bioinformatics Analyst in the Informatics Department at the J. Craig Venter Institute.

**Trygve E. Bakken** is an Assistant Investigator at the Allen Institute for Brain Science.

**Jeremy A. Miller** is a Senior Scientist at the Allen Institute for Brain Science.

**Rebecca D. Hodge** is an Assistant Investigator at the Allen Institute for Brain Science.

**Ed S. Lein** is a Senior Investigator at the Allen Institute for Brain Science.

**Richard H. Scheuermann** is a Professor and Director in the Informatics Department at the J. Craig Venter Institute.

**Figure Legends**

**Figure 1. FR-Match schematic and marker gene "barcodes". (a)** FR-Match cluster-to-cluster matching schematic diagram. Input data: query/new and reference datasets, each with cell-by-gene expression matrix and cell cluster membership labels. Step I: dimensionality reduction by selecting expression data of reference cell type marker genes from the query dataset. Here, we

30

664    use the NS-Forest marker genes selected for the reference cell types. Step II: Cluster-to-cluster

665    matching through the Friedman-Rafsky (FR) test. From left to right: multivariate data points of

666    cell transcriptional profiles (colored by cell cluster labels) in a reduced dimensional (reference

667    marker gene expression) space; construct a complete graph by connecting each pair of vertices

668    (i.e. cells); find the minimum spanning tree that connects all vertices with minimal summed edge

669    lengths; remove the edges that connect vertices from different clusters; count the number of

670    disjoint subgraphs, termed "multivariate runs" and denoted as $R$; calculate the FR statistic $W$,

671    which has asymptotically a standard normal distribution. **(b)** "Barcodes" of the cortical Layer 1

672    NS-Forest marker genes in four Layer 1 clusters. Heatmaps show marker gene expression

673    levels of 30 randomly selected cells in each cell cluster. The "Marker" column indicates if the

674    gene is a marker gene of the cluster or not (1=yes, 0=no).

675

676    **Figure 2. Cross-validation results.** Two-fold cross-validation were repeated 20 times on the

677    cortical Layer 1 data with all clusters. Training (reference) and testing (query) data were evenly

678    split in proportion to the cluster sizes. Cluster-level matching results for the cell-level matching

679    methods were summarized as the most mapped cluster labels beyond a defined threshold (see

680    Methods section). Matching output: 1 if a match; 0 otherwise. If a query cluster is not matched to

681    any reference cluster, then it is unassigned. **(a)** Heatmaps show the average matching result for

682    each matching method. True positive rate (TPR) is calculated as the average of the diagonal

683    matching rates, i.e. true positives. **(b)** Median, interquartile range, and full range of accuracy,

684    sensitivity, and specificity of all cluster-matching results in cross-validation for each matching

685    method is shown.

686

687    **Figure 3. Leave-*K*-cluster-out cross-validation results.** The same cross-validation settings

688    as in Figure 2 were used. After data split, $K \geq 1$ reference clusters were held-out to simulate the

689    situation in which the query dataset contains one or more novel cell type clusters. **(a)** Heatmaps

690    show the average matching result for each matching method when the i5 "rosehip" cluster was

691    left out. **(b)** Accuracy, sensitivity, and specificity of the leave-1-cluster-out cross-validation

692    performance for each matching method is shown. Each cluster was left out in turn, and

693    performance was evaluated across all turns. **(c)** Precision-Recall Curves of the leave-*K*-cluster-

694    out cross-validation performance for $K = 1, 3, 5,$ and 7 are shown and Area-Under-the-Curves

695    (AUC) statistics are calculated. Performance was evaluated across 20 iterations of randomly

696    selected *K* clusters. Curves for the FR-Match with and without p-value adjustment have the

697    same shape since the adjustment preserves the order of p-values. Note that the Seurat

698    package by default does not provide for unassigned cells/clusters as a direct output.

699

700    **Figure 4. Design of the under-, optimally-, and over-partitioned cluster simulations and**

701    **their matching properties. (a)** A schematic of simulating cluster partitions. The optimal

702    partitioning produced nodes where cells were consistently co-clustered across 100 bootstrap

703    iterations for clustering and curated by domain expert knowledge [13, 28]. Connectivity (edge

704    width) between nodes are measured by the number of intermediate cells/nuclei shared by

705    similar nodes. Two under-partition scenarios, "Mid nodes" and "Top nodes", were simulated by

706    merging similar/hierarchically-connected nodes (e.g. i1 + i5 clusters and all inhibitory clusters,

707    respectively). Two over-partition scenarios, split e1 and split i1, i2, and i3, were simulated by

708    splitting those large size clusters by k-means clustering with k = 2. Median F-measure of the

709    NS-Forest marker genes for each partition are reported in the table. **(b)** FR-Match properties

710    and expected marker gene types with respect to under-, optimally-, and over-partitioned

711    reference and query cluster scenarios, summarized from the simulation results (Supplementary

712     Figure 11). Green blocks in the table are cases with high true positive rate (TPR); red blocks are

713     warning cases with low TPR.

714

715     **Figure 5. FR-Match results for cell type matching between the cortical Layer 1 and full**

716     **MTG datasets. (a)** Two-way matching results are shown in three colors: red indicates that a

717     pair of clusters are matched in both directions (Layer 1 query to MTG reference with MTG

718     markers, and MTG query to Layer 1 reference with Layer 1 markers); yellow indicates that a pair

719     of clusters are matched in only one direction; and blue indicates that a pair of clusters are not

720     matched. The hierarchical taxonomy of the full MTG clusters is from the original study [29]. FR-

721     Match produced 13 unique, and two non-unique two-way matches between the two datasets.

722     Box A shows densely located one-way matches in the subclade of *VIP*-expressing clusters. Box

723     B shows incidentally captured cells from upper cortical Layer 2 mixed in the Layer 1 e1 cluster.

724     Box C shows the non-unique two-way matches of astrocyte clusters. **(b)** Examples of matched

725     and unmatched minimum spanning tree plots from the FR-Match graphical tool. Top row:

726     examples of two-way matched inhibitory clusters. Middle row: examples of two-way matched

727     non-neuronal clusters. Bottom row: examples of unmatched excitatory clusters from different

728     layers. Legend: cluster name (cluster size).

729

730     **Figure 6. Cell type matching results between the cortical Layer 1 and full MTG datasets**

731     **using other matching methods.** Two-way cluster-level matching results for the cell-level

732     matching methods were summarized as the most mapped cluster labels beyond a defined

733     threshold (see Methods section). Box A shows matches in the *VIP*-expressing subclade. Box B

734     shows matches spanning multiple layers among the MTG clusters. Box C shows matches of
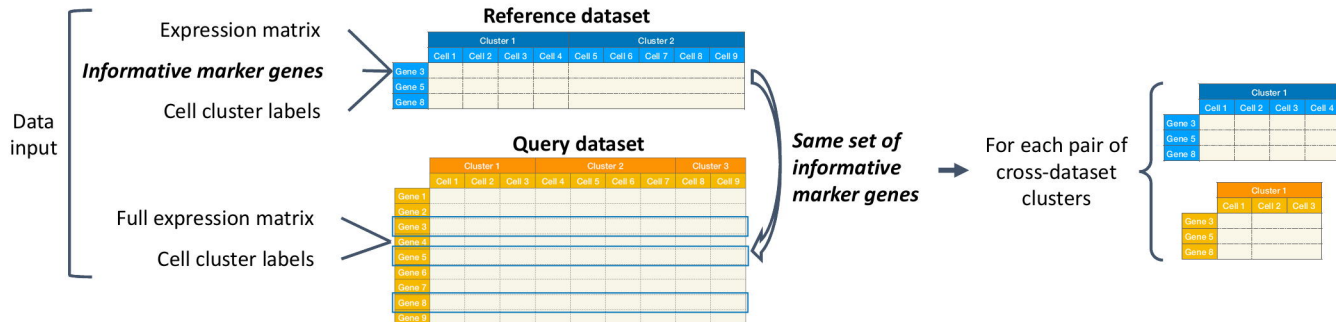
735     glial clusters.

33

**References**

1.  Regev, A., et al., *The Human Cell Atlas.* Elife, 2017. **6**.
2.  *The impact of the NIH BRAIN Initiative.* Nat Methods, 2018. **15**(11): p. 839.
3.  Aevermann, B., et al. *Production of a preliminary quality control pipeline for single nuclei Rna-Seq and its application in the analysis of cell type diversity of post-mortem human brain neocortex.* in *PACIFIC SYMPOSIUM ON BIOCOMPUTING 2017.* 2017. World Scientific.
4.  Ilicic, T., et al., *Classification of low quality cells from single-cell RNA-seq data.* Genome Biol, 2016. **17**: p. 29.
5.  Islam, S., et al., *Characterization of the single-cell transcriptional landscape by highly multiplex RNA-seq.* Genome Res, 2011. **21**(7): p. 1160-7.
6.  Kim, D., et al., *Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype.* Nature biotechnology, 2019. **37**(8): p. 907-915.
7.  Dobin, A., et al., *STAR: ultrafast universal RNA-seq aligner.* Bioinformatics, 2013. **29**(1): p. 15-21.
8.  Li, H. and R. Durbin, *Fast and accurate short read alignment with Burrows-Wheeler transform.* Bioinformatics, 2009. **25**(14): p. 1754-60.
9.  Pertea, M., et al., *StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.* Nature biotechnology, 2015. **33**(3): p. 290.
10. Blondel, V.D., et al., *Fast unfolding of communities in large networks.* Journal of statistical mechanics: theory and experiment, 2008. **2008**(10): p. P10008.
11. Wolf, F.A., P. Angerer, and F.J. Theis, *SCANPY: large-scale single-cell gene expression data analysis.* Genome Biol, 2018. **19**(1): p. 15.
12. Kiselev, V.Y., et al., *SC3: consensus clustering of single-cell RNA-seq data.* Nature methods, 2017. **14**(5): p. 483.
13. Bakken, T.E., et al., *Single-nucleus and single-cell transcriptomes compared in matched cortical cell types.* PloS one, 2018. **13**(12): p. e0209648.
14. Aevermann, B.D., et al., *Cell type discovery using single-cell transcriptomics: implications for ontological representation.* Hum Mol Genet, 2018. **27**(R1): p. R40-R47.
15. Bakken, T., et al., *Cell type discovery and representation in the era of high-content single cell phenotyping.* BMC Bioinformatics, 2017. **18**(Suppl 17): p. 559.
16. Lun, A.T., K. Bach, and J.C. Marioni, *Pooling across cells to normalize single-cell RNA sequencing data with many zero counts.* Genome Biol, 2016. **17**: p. 75.
17. Bacher, R., et al., *SCnorm: robust normalization of single-cell RNA-seq data.* Nat Methods, 2017. **14**(6): p. 584-586.
18. Haghverdi, L., et al., *Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors.* Nat Biotechnol, 2018. **36**(5): p. 421-427.
19. Polanski, K., et al., *BBKNN: fast batch alignment of single cell transcriptomes.* Bioinformatics, 2020. **36**(3): p. 964-965.
20. Kiselev, V.Y., A. Yiu, and M. Hemberg, *scmap: projection of single-cell RNA-seq data across data sets.* Nat Methods, 2018. **15**(5): p. 359-362.
21. Stuart, T., et al., *Comprehensive Integration of Single-Cell Data.* Cell, 2019.
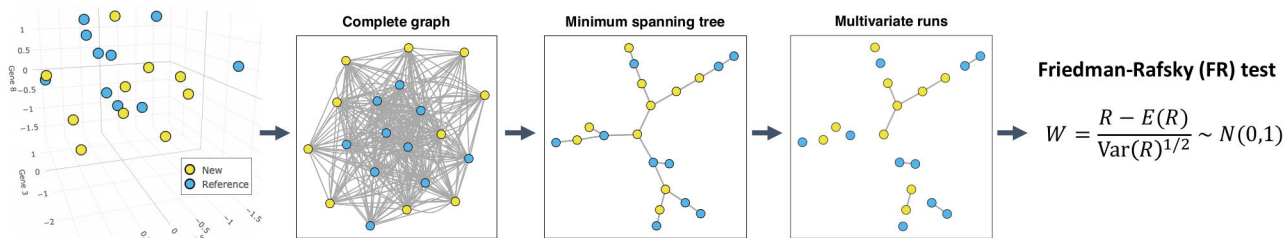
778     22.     Butler, A., et al., *Integrating single-cell transcriptomic data across different conditions,*
779             *technologies, and species.* Nat Biotechnol, 2018. **36**(5): p. 411-420.
780     23.     Hie, B., B. Bryson, and B. Berger, *Efficient integration of heterogeneous single-cell*
781             *transcriptomes using Scanorama.* Nat Biotechnol, 2019. **37**(6): p. 685-691.
782     24.     Lin, Y., et al., *scMerge leverages factor analysis, stable expression, and pseudoreplication*
783             *to merge multiple single-cell RNA-seq datasets.* Proc Natl Acad Sci U S A, 2019. **116**(20):
784             p. 9775-9784.
785     25.     Johansen, N. and G. Quon, *scAlign: a tool for alignment, integration, and rare cell*
786             *identification from scRNA-seq data.* Genome Biol, 2019. **20**(1): p. 166.
787     26.     Friedman, J.H. and L.C. Rafsky, *Multivariate generalizations of the Wald-Wolfowitz and*
788             *Smirnov two-sample tests.* The Annals of Statistics, 1979: p. 697-717.
789     27.     Hsiao, C., et al., *Mapping cell populations in flow cytometry data for cross-sample*
790             *comparison using the Friedman-Rafsky test statistic as a distance measure.* Cytometry A,
791             2016. **89**(1): p. 71-88.
792     28.     Boldog, E., et al., *Transcriptomic and morphophysiological evidence for a specialized*
793             *human cortical GABAergic cell type.* Nat Neurosci, 2018. **21**(9): p. 1185-1195.
794     29.     Hodge, R.D., et al., *Conserved cell types with divergent features in human versus mouse*
795             *cortex.* Nature, 2019. **573**(7772): p. 61-68.
796     30.     Aevermann, B., et al., *NS-Forest: A machine learning method for the objective*
797             *identification of minimum marker gene combinations for cell type determination from*
798             *single cell RNA sequencing.* bioRxiv, 2020.
799     31.     Korsunsky, I., et al., *Fast, sensitive and accurate integration of single-cell data with*
800             *Harmony.* Nature Methods, 2019. **16**(12): p. 1289-1296.
801     32.     Welch, J.D., et al., *Single-cell multi-omic integration compares and contrasts features of*
802             *brain cell identity.* Cell, 2019. **177**(7): p. 1873-1887. e17.
803     33.     Tran, H.T.N., et al., *A benchmark of batch-effect correction methods for single-cell RNA*
804             *sequencing data.* Genome Biology, 2020. **21**(1): p. 12.
805     34.     Picelli, S., et al., *Full-length RNA-seq from single cells using Smart-seq2.* Nat Protoc, 2014.
806             **9**(1): p. 171-81.
807     35.     Zhang, L. and S. Zhang, *Comparison of computational methods for imputing single-cell*
808             *RNA-sequencing data.* IEEE/ACM Trans Comput Biol Bioinform, 2018.
809     36.     Zheng, G.X., et al., *Massively parallel digital transcriptional profiling of single cells.* Nat
810             Commun, 2017. **8**: p. 14049.
811     37.     Vallejos, C.A., et al., *Normalizing single-cell RNA sequencing data: challenges and*
812             *opportunities.* Nat Methods, 2017. **14**(6): p. 565-571.
813     38.     Moffitt, J.R., et al., *High-throughput single-cell gene-expression profiling with*
814             *multiplexed error-robust fluorescence in situ hybridization.* Proc Natl Acad Sci U S A,
815             2016. **113**(39): p. 11046-51.
816     39.     Shah, S., et al., *In Situ Transcription Profiling of Single Cells Reveals Spatial Organization*
817             *of Cells in the Mouse Hippocampus.* Neuron, 2016. **92**(2): p. 342-357.
818     40.     Perkel, J.M., *Starfish enterprise: finding RNA patterns in single cells.* Nature, 2019.
819             **572**(7770): p. 549-551.
820     41.     Holmes, S. and W. Huber, *Modern statistics for modern biology.* 2018: Cambridge
821             University Press.

822   42.    Benjamini, Y. and D. Yekutieli, *The control of the false discovery rate in multiple testing*
823          *under dependency.* Annals of statistics, 2001: p. 1165-1188.
824   43.    Krishnaswami, S.R., et al., *Using single nuclei for RNA-seq to capture the transcriptome*
825          *of postmortem neurons.* Nature Protocols, 2016. **11**(3): p. 499-524.

826

**a** I. Dimensionality reduction

**II. Cluster-to-cluster matching**

Friedman-Rafsky (FR) test

$$W = \frac{R - E(R)}{\mathrm{Var}(R)^{1/2}} \sim N(0,1)$$

**b**

Layer1.e1_e299_SLC17A7_L5b_Cdh13

Layer1.g1_g48_GLI3_Astro_Gja1

Layer1.i1_i90_COL5A2_Ndnf_Car4

Layer1.i5_i47_TRPC3_Ndnf_Car4

**a**

FR-Match (TPR = 0.94)   FR-Match adj. (TPR = 1)   scmap (TPR = 0.54)

scmap+NSF (TPR = 0.7)   scmap+NSF.ext (TPR = 0.96)   Seurat (TPR = 0.93)

**b**

Accuracy   Sensitivity   Specificity

FR-Match, FR-Match adj., scmap, scmap+NSF, scmap+NSF.ext, Seurat

**a**

FR-Match (leave–i5–out)  FR-Match adj. (leave–i5–out)  scmap (leave–i5–out)

scmap+NSF (leave–i5–out)  scmap+NSF.ext (leave–i5–out)  Seurat (leave–i5–out)

**b**

Accuracy

Sensitivity

Specificity

**c**

Leave–1–cluster–out

PRAUC
FR-Match: 0.99
scmap: 0.78
scmap+NSF: 0.99
scmap+NSF.ext: 1
Seurat: 0.99

Leave–3–cluster–out

PRAUC
FR-Match: 0.98
scmap: 0.77
scmap+NSF: 0.99
scmap+NSF.ext: 0.98
Seurat: 0.97

Leave–5–cluster–out

PRAUC
FR-Match: 0.98
scmap: 0.74
scmap+NSF: 0.98
scmap+NSF.ext: 0.98
Seurat: 0.94

Leave–7–cluster–out

PRAUC
FR-Match: 0.95
scmap: 0.71
scmap+NSF: 0.97
scmap+NSF.ext: 0.93
Seurat: 0.9

**a**

| Cluster partitioning | | # clusters | F-measure |
|---|---|---|---|
| Under-partition | **Top nodes** | 3 | 0.959 |
| | **Mid nodes** | 11 | 0.929 |
| Optimal-partition | **15 clusters** | 15 | 0.874 |
| Over-partition | **Split e1** | 16 | 0.864 |
| | **Split i1, i2, i3** | 18 | 0.818 |

**b**

| | | Query cluster | | |
|---|---|---|---|---|
| | Partition<br>*Marker type (F-measure)* | Under-partition | Optimal-partition | Over-partition |
| **Reference cluster** | Under-partition<br>*Common marker (higher)* | 1-to-1<br>(TPR=0.98) | Many-to-1<br>(TPR=0.94) | Many-to-1<br>(TPR=0.95) |
| | Optimal-partition<br>*Specific markers (high)* | 1-to-many or missing<br>(TPR=0.72) | 1-to-1<br>(TPR=1) | Many-to-1<br>(TPR=0.99) |
| | Over-partition<br>*Noisy markers (low)* | 1-to-many or missing<br>(TPR=0.78) | 1-to-many<br>(TPR=0.99) | Many-to-many<br>(TPR=0.94) |

**a** FR-match

Row labels (top to bottom):
- MTG.Inh_L1_2_PAX6_CDH12
- MTG.Inh_L1_1_PAX6_TNFAIP8L3
- MTG.Inh_L1_1_SST_NMBR
- MTG.Inh_L1_4_LAMP5_LCP2
- MTG.Inh_L1_2_LAMP5_DBP
- MTG.Inh_L2_6_LAMP5_CA1
- MTG.Inh_L1_1_SST_CHRNA4
- MTG.Inh_L1_2_GAD1_MC4R
- MTG.Inh_L1_2_SST_BAGE2
- MTG.Inh_L1_3_PAX6_SYT6
- MTG.Inh_L1_2_VIP_TSPAN12
- MTG.Inh_L1_4_VIP_CHRNA6
- MTG.Inh_L1_3_VIP_ADAMTSL1
- MTG.Inh_L1_4_VIP_PENK
- MTG.Inh_L2_6_VIP_QPCT
- MTG.Inh_L3_6_VIP_HS3ST3A1
- MTG.Inh_L1_2_VIP_PCDH20
- MTG.Inh_L2_5_VIP_SERPINF1
- MTG.Inh_L2_5_VIP_TYR
- MTG.Inh_L1_3_VIP_CHRM2
- MTG.Inh_L2_4_VIP_CBLN1
- MTG.Inh_L1_3_VIP_CCDC184
- MTG.Inh_L1_3_VIP_GGH
- MTG.Inh_L1_2_VIP_LBH
- MTG.Inh_L2_3_VIP_CASC6
- MTG.Inh_L2_4_VIP_SPAG17
- MTG.Inh_L4_4_VIP_OPRM1
- MTG.Inh_L3_6_SST_NPY
- MTG.Inh_L3_6_SST_HPGD
- MTG.Inh_L4_6_SST_B3GAT2
- MTG.Inh_L5_6_SST_KLHDC8A
- MTG.Inh_L5_6_SST_NPM1P10
- MTG.Inh_L4_6_SST_GXYLT2
- MTG.Inh_L4_5_SST_STK32A
- MTG.Inh_L1_3_SST_CALB1
- MTG.Inh_L3_5_SST_ADGRG6
- MTG.Inh_L2_4_SST_FRZB
- MTG.Inh_L5_6_SST_TH
- MTG.Inh_L5_6_GAD1_GLP1R
- MTG.Inh_L5_6_PVALB_LGR5
- MTG.Inh_L4_5_PVALB_MEPE
- MTG.Inh_L2_4_PVALB_WFDC2
- MTG.Inh_L4_6_PVALB_SULF1
- MTG.Inh_L5_6_SST_MIR548F2
- MTG.Inh_L2_5_PVALB_SCUBE3
- MTG.Exc_L2_LAMP5_LTK
- MTG.Exc_L2_4_LINC00507_GLP2R
- MTG.Exc_L2_3_LINC00507_FREM3
- MTG.Exc_L5_6_THEMIS_C1QL3
- MTG.Exc_L3_4_RORB_CARM1P1
- MTG.Exc_L3_5_RORB_ESR1
- MTG.Exc_L3_5_RORB_COL22A1
- MTG.Exc_L3_5_RORB_FILIP1L
- MTG.Exc_L3_5_RORB_TWIST2
- MTG.Exc_L4_5_RORB_FOLH1B
- MTG.Exc_L4_5_RORB_SEMA3E
- MTG.Exc_L4_6_RORB_DAPK2
- MTG.Exc_L4_6_RORB_TTC12
- MTG.Exc_L4_6_RORB_C1R
- MTG.Exc_L4_5_FEZF2_SCN4B
- MTG.Exc_L5_6_THEMIS_DCSTAMP
- MTG.Exc_L5_6_THEMIS_CRABP1
- MTG.Exc_L5_6_THEMIS_FGF10
- MTG.Exc_L4_6_FEZF2_IL26
- MTG.Exc_L5_6_FEZF2_ABO
- MTG.Exc_L6_FEZF2_SCUBE1
- MTG.Exc_L5_6_SLC17A7_IL15
- MTG.Exc_L6_FEZF2_OR2T8
- MTG.Exc_L5_6_FEZF2_EFTUD1P1
- MTG.OPC_L1_6_PDGFRA
- MTG.Astro_L1_6_FGFR3_SLC14A1
- MTG.Astro_L1_2_FGFR3_GFAP
- MTG.Oligo_L1_6_OPALIN
- MTG.Micro_L1_3_TYROBP
- unassigned

Legend:
- Two-way Match (red)
- One-way match (yellow)
- No match (blue)

**b**

query.i7_i31_CLMP_Ndnf_Cxcl14
- query (31)
- Inh_L1_2_PAX6_CDH12 (90)

query.Inh_L1_2_PAX6_CDH12
- query (90)
- i7_i31_CLMP_Ndnf_Cxcl14 (31)

query.g2_g27_APBB1IP_Micro_Ctss
- query (27)
- Micro_L1_3_TYROBP (63)

query.Micro_L1_3_TYROBP
- query (63)
- g2_g27_APBB1IP_Micro_Ctss (27)

query.e1_e299_SLC17A7_L5b_Cdh13
- query (299)
- Exc_L4_6_RORB_C1R (160)

query.Exc_L4_6_RORB_C1R
- query (160)
- e1_e299_SLC17A7_L5b_Cdh13 (299)