

1 Big data and single cell transcriptomics: implications for ontological representation

2

3 Brian D. Aebermann¹, Mark Novotny¹, Trygve Bakken², Jeremy A. Miller², Alexander D. Diehl³,
4 David Osumi-Sutherland⁴, Roger S. Lasken¹, Ed S. Lein², Richard H. Scheuermann^{1,5}

5

6 ¹J. Craig Venter Institute, La Jolla, CA, USA

7 ²Allen Institute for Brain Science, Seattle, WA, USA

8 ³Department of Biomedical Informatics, University at Buffalo, Buffalo, NY, USA

9 ⁴European Molecular Biology Laboratory - European Bioinformatics Institute, Wellcome Trust
10 Genome Campus, Hinxton, United Kingdom

11 ⁵Department of Pathology, University of California San Diego, La Jolla, CA, USA

12

13 **Abstract** (250 words) – Cells are fundamental functional units of multicellular organisms, with
14 different cell types playing distinct physiological roles in the body. The recent advent of single
15 cell transcriptional profiling using RNA sequencing is producing “big data”, enabling the
16 identification of novel human cell types at an unprecedented rate. In this review, we summarize
17 recent work characterizing cell types in the human central nervous and immune systems using
18 single cell and single nuclei RNA sequencing, and discuss the implications that these discoveries
19 are having on the representation of cell types in the reference Cell Ontology (CL). We propose a
20 method based on random forest machine learning for identifying sets of necessary and sufficient
21 marker genes that can be used to assemble consistent and reproducible cell type definitions for

22 incorporation into the CL. The representation of defined cell type classes and their relationships
23 in the CL using this strategy will make the cell type classes findable, accessible, interoperable, and
24 reusable (FAIR), allowing the CL to serve as a reference knowledgebase of information about the
25 role that distinct cellular phenotypes play in human health and disease.

26

27 **Introduction** – Cells are probably the most important fundamental functional units of
28 multicellular organisms, since different cell types play different physiological roles in the body.
29 Although every cell of an individual organism has essentially the same genome structure, different
30 cells realize diverse functions due to differences in their *expressed* genome. In many cases,
31 abnormalities in gene expression form the physical basis of disease dispositions. Thus,
32 understanding and representing normal and abnormal cellular phenotypes can lead to the
33 development of biomarkers for diagnosing disease and the identification of critical targets for
34 therapeutic interventions.

35 Previous approaches used to characterize cell phenotypes have several drawbacks that limited their
36 ability to comprehensively identify the cellular complexity of human tissues. Transcriptional
37 profiling of bulk cell sample mixtures by microarray or RNA sequencing can simultaneously
38 assess gene expression levels and proportions of abundant known cell types, but precludes
39 identification of novel cell types and obscures the contributions of rare cell subsets to the gene
40 expression patterns present in the bulk samples. Flow cytometry provides phenotype information
41 at the single cell level, but is limited by the number of discrete markers that can be assessed, and
42 relies on prior knowledge of marker expression patterns. The recent establishment of methods for
43 single cell transcriptional profiling (1, 2) is revolutionizing our ability to understand complex cell

44 mixtures, avoiding the averaging phenomenon inherent in the analysis of bulk cell mixtures and
45 providing for an unbiased assessment of phenotypic markers within the expressed genome.

46 In order to compare experimental results and other information about cell types, a standard
47 reference nomenclature that includes consistent cell type names and definitions is required. The
48 Cell Ontology (CL) is a biomedical ontology developed to provide this standard reference
49 nomenclature for *in vivo* cell types in humans and major model organisms (3). However, the
50 advent of high-content single cell transcriptomics for cell type characterization has resulted in a
51 number of challenges for their representation in the CL (discussed in Bakken 2017 (4)). In this
52 paper, we review some of the recent discoveries that have resulted from the application of single
53 cell transcriptomics to human samples, and propose a strategy for defining cell types within the
54 CL based on the identification of necessary and sufficient marker genes, to support interoperable
55 and reproducible research.

56

57 **Application to the human brain** - Initial progress in neuronal cell type discovery by single cell
58 RNA sequencing (scRNAseq) focused on mouse cerebral, visual, and somatosensory cortices (5,
59 6, 7, 8, 9). More recently, technological advances, including RNAseq using single nuclei
60 (snRNAseq) instead of single cells (10, 11, 12), have extended these investigations into human
61 neuronal cell type discovery (13, 14). Comprehensive reviews of these recent advances have been
62 reported recently (15, 16).

63 Initial efforts toward human neuronal cell type discovery focused on identifying broad lineages.
64 Pollen *et al.* profiled 65 neuronal cells into six categories - neural progenitor cells, radial glia,
65 newborn neurons, inhibitory interneurons, and maturing neurons (17), while Darmanis *et al.*
66 sequenced 466 cells, also identifying six broad, but distinct, categories - oligodendrocytes,

67 astrocytes, microglia, endothelial cells, oligodendrocyte precursor cells (OPCs), and neurons (18).
68 Darmanis et al. further subtyped the adult neurons into 2 excitatory and 5 inhibitory types. More
69 recent single *nuclei* RNAseq investigations are attempting more comprehensive cell typing. Lake
70 *et al.* sampled 3,227 nuclei from 6 Brodmann areas, from which the neurons were classified into
71 8 excitatory and 8 inhibitory subtypes (13). Similarly, Boldog *et al.* sampled 769 nuclei from layer
72 1 of the Middle Temporal Gyrus (MTG) and identified 11 distinct inhibitory cell types (14).
73 Comparing results between these studies has been challenging given the different areas and layers
74 of cortex sampled. Many of the studies leveraged classical cell type markers derived from the
75 mouse scRNAseq literature. For example, SNAP25 expression was used to broadly define
76 neuronal cells, while GAD1 expression defined inhibitory interneurons. Additional classical
77 markers have then been used to subdivide the excitatory and inhibitory classes, such as CUX2 or
78 VIP respectively; however, these markers individually are still not specific enough to define
79 discrete cell type classes at the level of granularity revealed by clustering of the sc/snRNAseq data.
80 In fact, there has been surprisingly limited overlap in gene sets specific for individual cell type
81 clusters between studies, as the genes found in each study appear to be sensitive to both the context
82 and methodology used. For example, Lake *et al.* found that cluster In1 had CNR1 (Table S5 in
83 reference 13) as the highest ranked marker, while Boldog *et al.* found 7 distinct inhibitory types
84 that expressed this marker (Figure 3 in reference 14). Without a standardized methodology for
85 determining the necessary and sufficient marker genes and a corresponding marker gene reference
86 database, comparison of newly-identified cell types to those reported in previous studies requires
87 a complete reprocessing of the data.

88

89 **Application to the human immune system** - Single cell transcriptomic analysis has also been
90 applied to study the functional cell type diversity of the human immune system (reviewed in
91 Stubbington 2017 (19)). Bjorklund et al. used scRNAseq to explore the subtype diversity of
92 CD127⁺ innate lymphoid cells isolated from human tonsil, providing an in-depth transcriptional
93 characterization of the three major subtypes – ILC1, ILC2, and ILC3, and three additional subtypes
94 within the ILC3 class, by comparing their single cell transcriptional profiles (20).

95 Two recent studies explored the subtype diversity of dendritic cells in human blood. In addition
96 to identifying two conventional dendritic cell subtypes (cDC1 and cDC2) and one plasmacytoid
97 dendritic cell subtype, See et al. identified several subtypes that appear to correspond to precursor
98 cells, including one early uncommitted CD123⁺ pre-DC subset and two CD45RA⁺CD123^{lo}
99 lineage-committed subsets (pre-cDC1 and pre-cDC2), using cell sorting, scRNAseq, and *in vitro*
100 differentiation assays (21). Villani et al. used fluorescence-activated cell sorting and scRNAseq
101 to delineate six different dendritic cell subtypes (DC1 – 6) and four different monocyte subtypes
102 (Mono1 – 4), and showed that these different subtypes, which were defined based on their
103 transcriptional profiles, exhibited different functional capabilities for allogeneic T cell stimulation
104 and for cytokine production following TLR agonist stimulation (22).

105 Two recent studies have explored the phenotypes of immune cells infiltrating tumor specimens
106 using scRNAseq. In melanoma, Tirosh et al. found that the non-malignant tumor
107 microenvironment was composed of T cell, B cell, NK cell, endothelial cell, macrophage and
108 cancer-associated fibroblast (CAF) subsets (23). In contrast to the distinct transcriptional
109 phenotypes of the malignant component across individual melanoma specimens, common features
110 could be observed in the non-malignant components, with important therapeutic implications.
111 Expression of multiple complement factors by CAFs correlated with the extent of T cell

112 infiltration. T cells with activation-independent exhaustion profiles, characterized by expression
113 of co-inhibitory receptors (e.g. PD1 and TIM3), could be distinguished from cytotoxic T cell
114 profiles. Potential biomarkers that distinguish between exhausted and cytotoxic T cells could aid
115 in selecting patients for immune checkpoint blockade. In hepatocellular carcinoma, Zheng et al.
116 found clonal enrichment of both regulatory T cells and exhausted CD8 T cells using scRNAseq
117 and T cell receptor repertoire analysis (24). The diagnostic and prognostic significance of these
118 findings remain to be explored.

119 While these studies illustrate the power of single cell genomics in identifying important functional
120 cell subtypes, they also illuminate a major challenge in comparing the results from different
121 studies, due to the lack of a consistent, reusable approach for naming, defining, and comparing
122 new cell types being identified by these high content phenotyping technologies. For example, in
123 the two studies focused on the identification of dendritic cell subtypes, it is unclear if the cDC1
124 and cDC2 subtypes identified by See et al. correspond to the DC1 and DC2 subtypes identified by
125 Villani et al. Indeed, the only way to make this determination would be to perform a *de novo*
126 comparative analysis of the transcriptional profiles from both studies. For these studies to truly
127 comply with the newly emerging FAIR principles of open data (25), a robust reproducible strategy
128 for defining and representing new cell types will be essential to support their broad interoperability.

129

130 **Ontological representation** - Biomedical ontologies, as promoted by the Open Biomedical
131 Ontology (OBO) Foundry (26), provide a framework to name and define the types, properties and
132 relationships of entities in the biomedical domain. The Cell Ontology (CL) was established in
133 2005 to provide a standard reference nomenclature for *in vivo* cell types, including those observed
134 in specific developmental stages in humans and different model organisms (3). The semantic

135 hierarchy of CL is mainly constructed using two core relations – *is_a* and *develops_from*. Masci
136 *et al.* proposed a major revision to the CL using dendritic cells as the driving biological use case
137 in which the expression of specific marker proteins on the cell surface (e.g. receptor proteins) or
138 internally (e.g. transcription factors) would be used as the main *differentia* for the asserted
139 hierarchy (27). Diehl *et al.* applied this approach first to cell types of the hematopoietic system
140 and then later to the full CL (28, 29, 30). As of December 2017, the CL contained 2199 cell type
141 classes, with 583 classes within the hematopoietic cell branch alone.

142 We recently discussed the challenges faced by the CL in the era of high-throughput, high-content
143 single cell phenotyping technologies, including sc/snRNAseq (4). One of the key
144 recommendations was to establish a standard strategy for defining cell type classes that combine
145 three essential components:

- 146 • the minimum set of ***necessary and sufficient marker genes*** selectively expressed by the
147 cell type,
- 148 • a ***parent cell class*** in the Cell Ontology, and
- 149 • a ***specimen source description*** (anatomic structure + species).

150 In order to identify the set of necessary and sufficient marker genes from an sc/snRNAseq
151 experiment, we have developed a method – NSforest – that utilizes a random forest of decision
152 trees machine learning approach.

153 To illustrate how this approach can produce standard cell type definitions, we have applied the
154 method to a transcriptomic dataset derived from single nuclei isolated from the middle temporal
155 gyrus, cortical layer 1 of a post-mortem human brain specimen (Figure 1a in reference 14).
156 Transcriptional profiles obtained from RNA sequencing of a collection of single sorted nuclei was
157 used to identify 16 discrete cell types using an iterative data clustering approach. Based on the

158 expression of the previously characterized marker genes SNAP25 and GAD1 for broad classes, 11
159 inhibitory interneurons, 1 excitatory neuron and 4 glial cell type clusters were identified.

160 In the first step (Figure 1b), NSforest takes the gene expression data matrix of single nuclei with
161 their cell type cluster membership as input, and develops a classification model for each cell type
162 cluster by comparing each Cluster X versus all non-Cluster X profiles using the random forest
163 algorithm (31). In addition to the classification model itself, NSforest produces a ranked list of
164 features (genes) that are most informative for distinguishing between Cluster X and all of the other
165 clusters.

166 In the second step, NSforest constructs single decision trees using first the top gene, then the top
167 two genes, top three genes, etc., until a stable tree topology and optimal classification accuracy is
168 achieved. The minimum number of genes necessary to obtain this stable classification result
169 corresponds to the necessary and sufficient set of marker genes defining each cell type cluster
170 within this experimental context.

171 The expression of the complete set of marker genes obtained from applying NSforest to the single
172 nuclei dataset is illustrated in Figure 2. In most cases, the expression of three marker genes is
173 sufficient to define a cell type cluster, with a range of one to five necessary and sufficient marker
174 genes per cluster. Glial cell subtypes appear to be more distinct from each other, requiring
175 relatively few genes to sufficiently define the cell type. In contrast, neuronal subtypes appear to
176 be more similar, requiring more genes to achieve specificity. In some cases, a combination of both
177 positive and negative expression optimally defines a cell type cluster.

178 For one of the inhibitory interneuron cell types defined in this study (i5), we were able to connect
179 the distinct transcriptional profile with a previous cell type defined based on its unique cellular
180 morphology – the Rosehip cell (14). This then allows us to construct an ontological representation

181 that includes both a colloquial name, an alternative name, and a definition combining the necessary
182 and sufficient marker genes, a CL parent cell class, and specimen source information, as follows:

- 183 • Colloquial name – *rosehip neuron*
- 184 • Alternative name - *KIT-expressing MTG cortical layer 1 GABAergic interneuron,*
185 *human*
- 186 • Definition - *A human middle temporal gyrus cortical layer 1 GABAergic interneuron*
187 *that selectively expresses KIT, NTNG1, and POU6F2 mRNAs*

188 A complete set of cell type names and definitions for all cell type clusters identified in this
189 experiment is provided in Table 1.

190 These informal textual definitions can then be converted into formal ontological definitions,
191 represented in OWL as equivalent classes, using a set of logical axioms that combine assertions
192 about the parent cell class (interneuron), anatomic locations of the neuron cell body (soma),
193 functional capacity of the cell type (gamma-aminobutyric acid secretion), and marker gene
194 expression (expresses some KIT) requirements (Figure 3). Using semantic reasoners, these logical
195 axioms can then be used to infer novel characteristics, e.g. SubClass Of ‘cerebral cortex
196 GABAergic interneuron’.

197 The challenge remains of ensuring that these cell type definitions, whose necessary and sufficient
198 conditions are derived from analysis of data from one particular methodology (scRNAseq), are
199 compatible with both existing cell type classes in the CL and cell types defined using alternative
200 experimental methods and data analysis approaches. Working with CL developers, we are now
201 establishing an extension ontology module containing provisional definitions for novel cell types
202 that we and other research groups will contribute. Ontological reasoners will be used to link these
203 cell types to more general classes in the CL proper, structure them into an extended hierarchy, and

204 determine when separate research groups have defined similar or identical cell types. CL
205 developers will review these provisional cell types periodically to determine when multiple lines
206 of evidence provide sufficient support to promote particular cell type classes to the CL proper. In
207 this way we will ensure the integrity of the CL reference, while still allowing for the rapid
208 expansion of its content to accommodate cell types defined via these new technologies.

209

210 **Conclusions** – The application of high-throughput/high-content cytometry and single cell genomic
211 techniques is producing an explosion in the number of distinct cellular phenotypes being identified
212 in human specimens. For biomedical ontologies to stay relevant, it will be critical for ontology
213 developers to establish procedures for the processing and incorporation of representations derived
214 from these data-intensive technologies into reference ontologies in a timely fashion. The
215 representation of defined cell types and their relationships in the CL will serve as a reference
216 knowledgebase to support interoperability of information about the role of cellular phenotypes in
217 human health and disease.

218

219 **Acknowledgements**

220 This work was supported by the Allen Institute for Brain Science, the JCVI Innovation Fund, the
221 U.S. National Institutes of Health R21-AI122100 and U19-AI118626, and the California Institute
222 for Regenerative Medicine GC1R-06673-B. We thank Nik Schork, Jamison McCorrison, Pratap
223 Venepally, Lindsay Cowell, Bjoern Peters, and Sirarat Sarntivijai for helpful discussion.

224

225 **Literature Cited**

- 226 1. Tang, F., Barbacioru, C., Wang, Y., Nordman, E., Lee, C., Xu, N., Wang, X., Bodeau, J.,
227 Tuch, B.B., Siddiqui, A., et al. (2009) mRNA-Seq whole transcriptome analysis of a
228 single cell. *Nat. Methods*, 6, 377–382.
- 229 2. Tang, F., Barbacioru, C., Bao, S., Lee, C., Nordman, E., Wang, X., Lao, K., and Surani,
230 M.A. (2010) Tracing the derivation of embryonic stem cells from the inner cell mass by
231 single-cell RNA-Seq analysis. *Cell Stem Cell*, 6, 468–478.
- 232 3. Bard J., Rhee S.Y., Ashburner, M. (2005) An ontology for cell types. *Genome Biol.*, 6,
233 R21.
- 234 4. Bakken, T., Cowell, L., Aevermann, B.D., Novotny, M., Hodge, R, Miller, J.A., Lee, A.,
235 Chang, I., McCarrison, J., Pulendran, B., et al. (2017) Cell type discovery and
236 representation in the era of high-content single cell phenotyping. *BMC Bioinformatics*,
237 18 (Suppl. 17), 559.
- 238 5. Zeisel, A., Muñoz-Manchado, A.B., Codeluppi, S., Lönnerberg, P., La Manno, G., Juréus
239 A., Marques, S., Munguba, H., He, L., Betsholtz, C., et al. (2015) Brain structure. Cell
240 types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science*, 347,
241 1138–1142.
- 242 6. Macosko, E.Z., Basu, A., Satija, R., Nemesh, J., Shekhar, K., Goldman, M., Tirosh, I.,
243 Bialas, A.R., Kamitaki, N., Martersteck, E.M., et al. (2015) Highly parallel genome-wide
244 expression profiling of individual cells using nanoliter droplets. *Cell*, 161, 202–1214.
- 245 7. Li, C-L., Li, K-C., Wu, D., Chen, Y., Luo, H., Zhao, J-R., Wang, S-S., Sun, M-M., Lu, Y-
246 J., Zhong, Y-Q., et al. (2015) Somatosensory neuron types identified by high-coverage
247 single-cell RNA-sequencing and functional heterogeneity. *Cell Res.*, 26, 83-102.

- 248 8. Shekhar, K., Lapan, S.W., Whitney, I.E., Tran, N.M., Macosko, E.Z., Kowalczyk, M.,
249 Adiconis, X., Levin, J.Z., Nemesh, J., Goldman, M., et al. (2016) Comprehensive
250 classification of retinal bipolar neurons by single-cell transcriptomics. *Cell*, 166, 1308–
251 1323.e30.
- 252 9. Tasic, B., Menon, V., Nguyen, T.N., Kim, T.K., Jarsky, T., Yao, Z., Levi, B., Gray, L.T.,
253 Sorensen, S.A., Dolbeare, T., et al. (2016) Adult mouse cortical cell taxonomy revealed by
254 single cell transcriptomics. *Nat. Neurosci.*, 19, 335–46.
- 255 10. Grindberg, R.V., Yee-Greenbaum, J.L., McConnell, M.J., Novotny, M., O'Shaughnessy,
256 A.L., Lambert, G.M., Araúzo-Bravo, M.J., Lee, J., Fishman, M., Robbins, G.E., Lin, X., et
257 al. (2013) RNA-sequencing from single nuclei. *Proc. Natl. Acad. Sci. USA*, 110, 19802–
258 19807.
- 259 11. Krishnaswami, S.R., Grindberg, R.V., Novotny, M., Venepally, P., Lacar, B., Bhutani, K.,
260 Linker, S.B., Pham, S., Erwin, J.A., Miller, J.A., et al. (2016) Using single nuclei for RNA-
261 seq to capture the transcriptome of postmortem neurons. *Nat. Protoc.*, 11, 499–524.
- 262 12. Lacar, B., Linker, S.B., Jaeger, B.N., Krishnaswami, S., Barron, J., Kelder, M., Parylak, S.,
263 Paquola, A., Venepally, P., Novotny, M., et al. (2016) Nuclear RNA-seq of single neurons
264 reveals molecular signatures of activation. *Nat. Commun.*, 7, 11022.
- 265 13. Lake, B.B., Ai, R., Kaeser, G.E., Salathia, N.S., Yung, Y.C., Liu, R., Wildberg, A., Gao,
266 D., Fung, H.L., Chen, S., et al. (2016) Neuronal subtypes and diversity revealed by single-
267 nucleus RNA sequencing of the human brain. *Science*, 352, 1586–1590.
- 268 14. Boldog, E., Bakken, T, Hodge, R.D., Novotny, M, Aevermann, B.D., Baka, J., Borde, S.,
269 Close, J.L., Diez-Fuertes, F., Ding, S.L., et al. (2017) Transcriptomic and

- 270 morphophysiological evidence for a specialized human cortical GABAergic cell type.
271 preprint bioRxiv: <https://t.co/v53HzGEe3V>.
- 272 15. Johnson M.B., Walsh, C.A., (2017) Cerebral cortical neuron diversity and development at
273 the single-cell resolution. *Curr. Opin. Neurobiol.*, 42, 9-16.
- 274 16. Lein, E.S., Belgard, T.G., Hawrylycz, M., Molnár, Z. (2017) Transcriptomic Perspectives
275 on Neocortical Structure, Development, Evolution, and Disease. *Annu. Rev. Neurosci.*,
276 40, 629-652.
- 277 17. Pollen, A.A., Nowakowski, T.J., Shuga, J., Wang, X., Leyrat, A.A., et al. (2014) Low-
278 coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated
279 signaling pathways in developing cerebral cortex. *Nat. Biotechnol.*, 32, 1053–1058
- 280 18. Darmanis, S., Sloan, S.A., Zhang, Y., Enge, M., Caneda, C., et al. (2015) A survey of
281 human brain transcriptome diversity at the single cell level. *PNAS*, 112, 7285–7290.
- 282 19. Stubbington, M.J.T., Rozenblatt-Rosen, O., Regev, A., Teichmann, S.A. (2017) Single-cell
283 transcriptomics to explore the immune system in health and disease. *Science*, 358, 58-63.
- 284 20. Björklund, Å.K., Forkel, M., Picelli, S., Konya, V., Theorell, J., Friberg, D., Sandberg, R.,
285 Mjösberg, J. (2016) The heterogeneity of human CD127(+) innate lymphoid cells revealed
286 by single-cell RNA sequencing. *Nat. Immunol.*, 4, 451-460.
- 287 21. See, P., Dutertre, C.A., Chen, J., Günther, P., McGovern, N., Irac, S.E., Gunawan, M.,
288 Beyer, M., Händler, K., Duan, K., et al. (2017) Mapping the human DC lineage through
289 the integration of high-dimensional techniques. *Science*, 356.
- 290 22. Villani, A.C., Satija, R., Reynolds, G., Sarkizova, S., Shekhar, K., Fletcher, J., Griesbeck,
291 M., Butler, A., Zheng, S., Lazo, (2017) Single-cell RNA-seq reveals new types of human
292 blood dendritic cells, monocytes, and progenitors. *Science*, 356.

- 293 23. Tirosh, I., Izar, B., Prakadan, S.M., Wadsworth, M.H. 2nd, Treacy, D., Trombetta, J.J.,
294 Rotem, A., Rodman, C., Lian, C., Murphy, G., et al. (2016) Dissecting the multicellular
295 ecosystem of metastatic melanoma by single-cell RNA-seq. *Science*, 352, 189-196.
- 296 24. Zheng, C., Zheng, L., Yoo, J.K., Guo, H., Zhang, Y., Guo, X., Kang, B., Hu, R., Huang,
297 J.Y., Zhang, Q., et al. (2017) Landscape of Infiltrating T Cells in Liver Cancer Revealed
298 by Single-Cell Sequencing. *Cell.*, 169, 1342-1356.e16.
- 299 25. Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A.,
300 Blomberg, N., Boiten, J.W., da Silva Santos, L.B., Bourne, P.E., Bouwman, J., Brookes,
301 A.J., et al. (2016) The FAIR Guiding Principles for scientific data management and
302 stewardship. *Sci. Data*, 3, 160018.
- 303 26. Smith, B., Ashburner, M., Rosse, C., Bard, J., Bug, W., Ceusters, W., Goldberg, L.J.,
304 Eilbeck, K., Ireland, A., Mungall, C.J., et al. (2007) The OBO Foundry: coordinated
305 evolution of ontologies to support biomedical data integration. *Nat. Biotechnol.*, 11, 1251-
306 1255.
- 307 27. Masci, A.M., Arighi, C.N., Diehl, A.D., Lieberman, A.E., Mungall, C., Scheuermann,
308 R.H., Smith, B., Cowell, L.G. (2009) An improved ontological representation of dendritic
309 cells as a paradigm for all cell types. *BMC Bioinformatics*. 10, 70.
- 310 28. Diehl, A.D., Augustine, A.D., Blake, J.A., Cowell, L.G., Gold, E.S., Gondré-Lewis, T.A.,
311 Masci, A.M., Meehan, T.F., Morel, P.A., Nijnik, A., et al. (2011) Hematopoietic cell types:
312 prototype for a revised cell ontology. *J. Biomed. Inform.*, 1, 75-79.
- 313 29. Meehan, T.F., Masci, A.M., Abdulla, A., Cowell, L.G., Blake, J.A., Mungall, C.J., Diehl,
314 A.D. (2011) Logical development of the cell ontology. *BMC Bioinformatics*, 12, 6.

- 315 30. Diehl, A.D., Meehan, T.F., Bradford, Y.M., Brush, M.H., Dahdul, W.M., Dougall, D.S.,
316 He, Y., Osumi-Sutherland D, Ruttenberg A, Sarntivijai S, et al. (2016) The Cell Ontology
317 2016: enhanced content, modularization, and ontology interoperability. *J. Biomed.*
318 *Semantics.* 7, 44.
- 319 31. Berthold M., Cebron N., Dill F., Gabriel, T., Kotter, T., Meinl, T., Ohl, P., Thiel, K.,
320 Wiswedel, B., et al. (2008) KNIME: The Konstanz Information Miner. In: Preisach C,
321 Burkhardt H, Schmidt-Thieme L, Decker R, editors. (eds). *Data Analysis, Machine*
322 *Learning and Applications.* Springer Berlin Heidelberg, Chapter 38, 319–326.

323

324

325

326 **Figure legends**

327 Figure 1. Identification of necessary and sufficient marker genes using NSforest – a) A typical
328 single cell/single nuclei RNA sequencing workflow in which a tissue specimen is obtained, single
329 cells/nuclei isolated by fluorescence activated cell sorting, amplified cDNA quantified by
330 sequencing, and cell types identified by clustering the resultant transcriptional profiles. b) The
331 NSforest approach takes a data matrix of expression values (e.g. transcripts per million reads) of
332 genes (rows) in single cell/nuclei samples (columns) grouped by cell type cluster membership. In
333 the first step, the expression levels of genes are used as features in the random forest machine
334 learning procedure to train classification models comparing single cell/nuclei expression data in
335 one cell type cluster against single cell/nuclei expression data in all other clusters, for every cell
336 type cluster separately, using the Random Forest Learner in KNIME v3.1.2. Each cell type cluster

337 classification model is constructed from one hundred thousand trees using Information Gain Ratio
338 as the splitting criteria, where each decision tree is generated using the default bagging parameters
339 - the square root of the number of features and a bootstrap of samples equal to the training set size.
340 For each cell type cluster classification model, the method outputs usage statistics, including how
341 often each gene is used as a branching criterion and the number of times it was a candidate across
342 all random decision trees. By summing the frequency of use when a candidate across the first three
343 branching levels, the list of genes can be ranked by their usefulness in distinguishing one cell type
344 clusters from the other clusters. In the second step, single decision trees are constructed using the
345 first gene from the ranked list, the first two genes, the first three genes, etc. Each individual tree
346 is then assessed for classification accuracy and tree topology using the training data. Given the
347 objective of determining the necessary and sufficient marker genes, we apply additional criteria in
348 scoring the trees - we restrict each gene to being used in only one branch per tree, and find the
349 optimal classification for the target cluster only, rather than the overall classification score. The
350 addition of genes from the ranked list is stopped when an optimal classification or stable tree
351 topology is achieved. The minimum number of genes used to produce this optimal result
352 corresponds to the set of necessary and sufficient marker genes required to define the cell type
353 cluster.

354

355 Figure 2. Marker gene expression patterns in single nuclei grouped by cluster – A heatmap of
356 expression levels for the necessary and sufficient marker genes identified for all 16 clusters across
357 all single nuclei grouped by cell type cluster is shown, including 1 excitatory (e1), 11 inhibitory
358 (i1 – i11), and 4 glial (g1 – g4) cell type clusters. In total, 49 markers genes were selected as being

359 necessary and sufficient to distinguish these 16 different cell type clusters from cortical layer 1 of
360 the human brain middle temporal gyrus region.

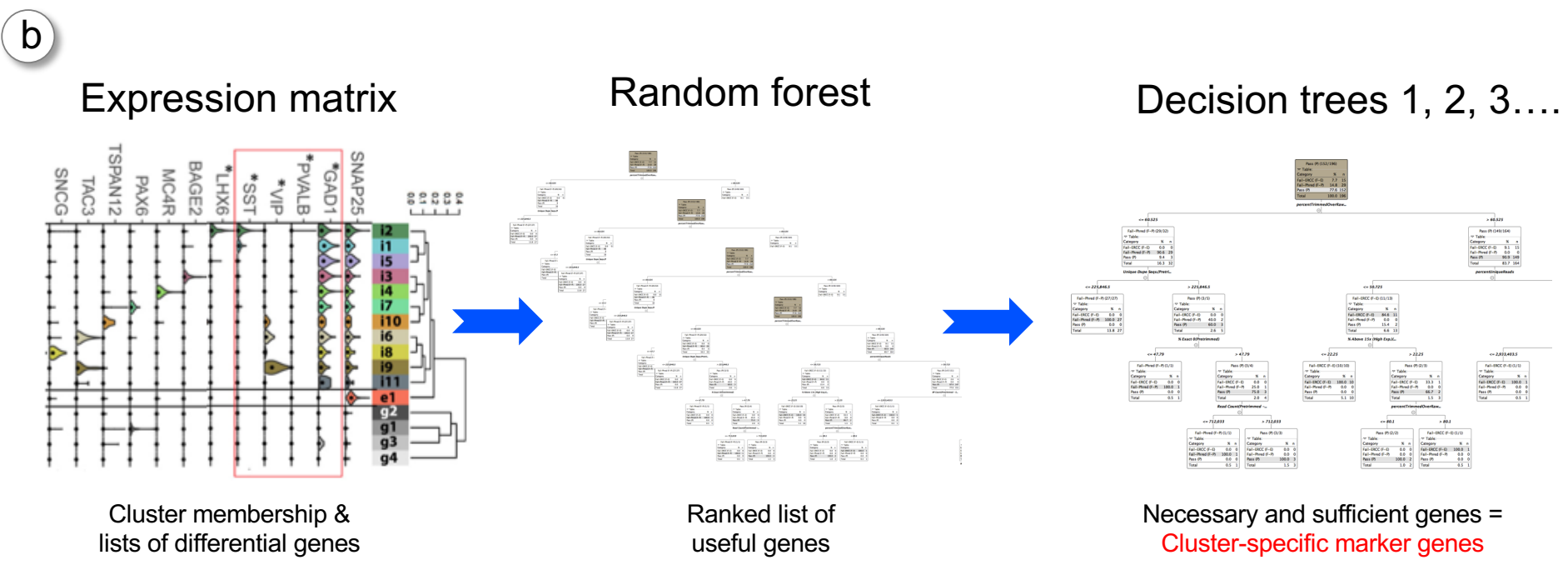
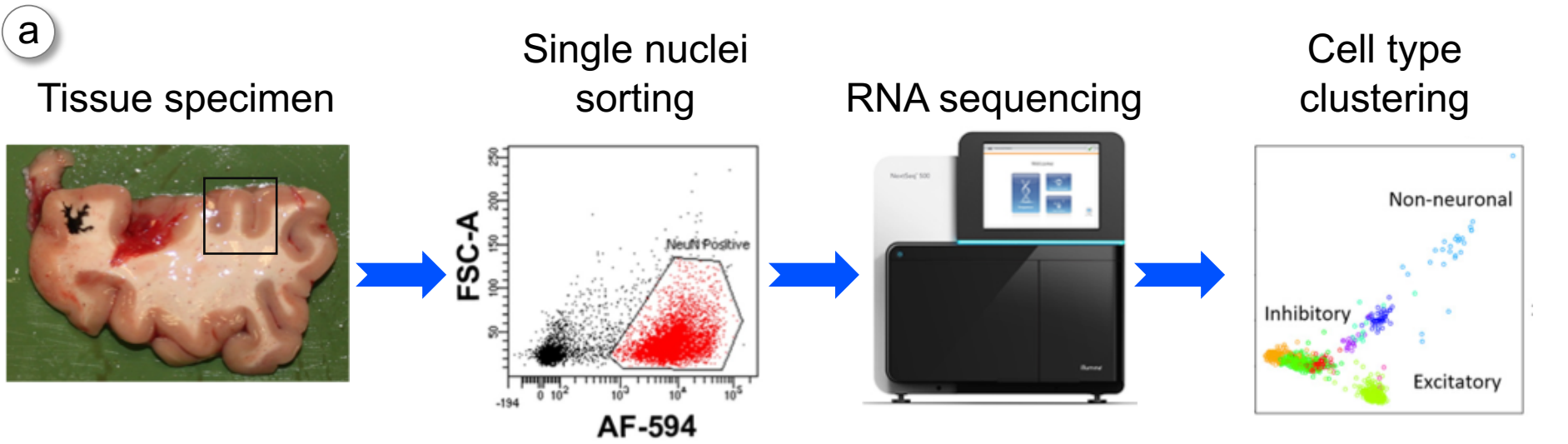
361
362 Figure 3. Formal rosehip neuron definition using logical axioms – A set of logical axioms about
363 the anatomic local of the cell body (soma), the functional capacity, and the necessary and sufficient
364 marker gene expressions are combined to construct an equivalent class cell type definition for the
365 rosehip neuron interneuron cluster – i4 (see Boldog 2017 (14) for more information about how this
366 cell type was characterized).

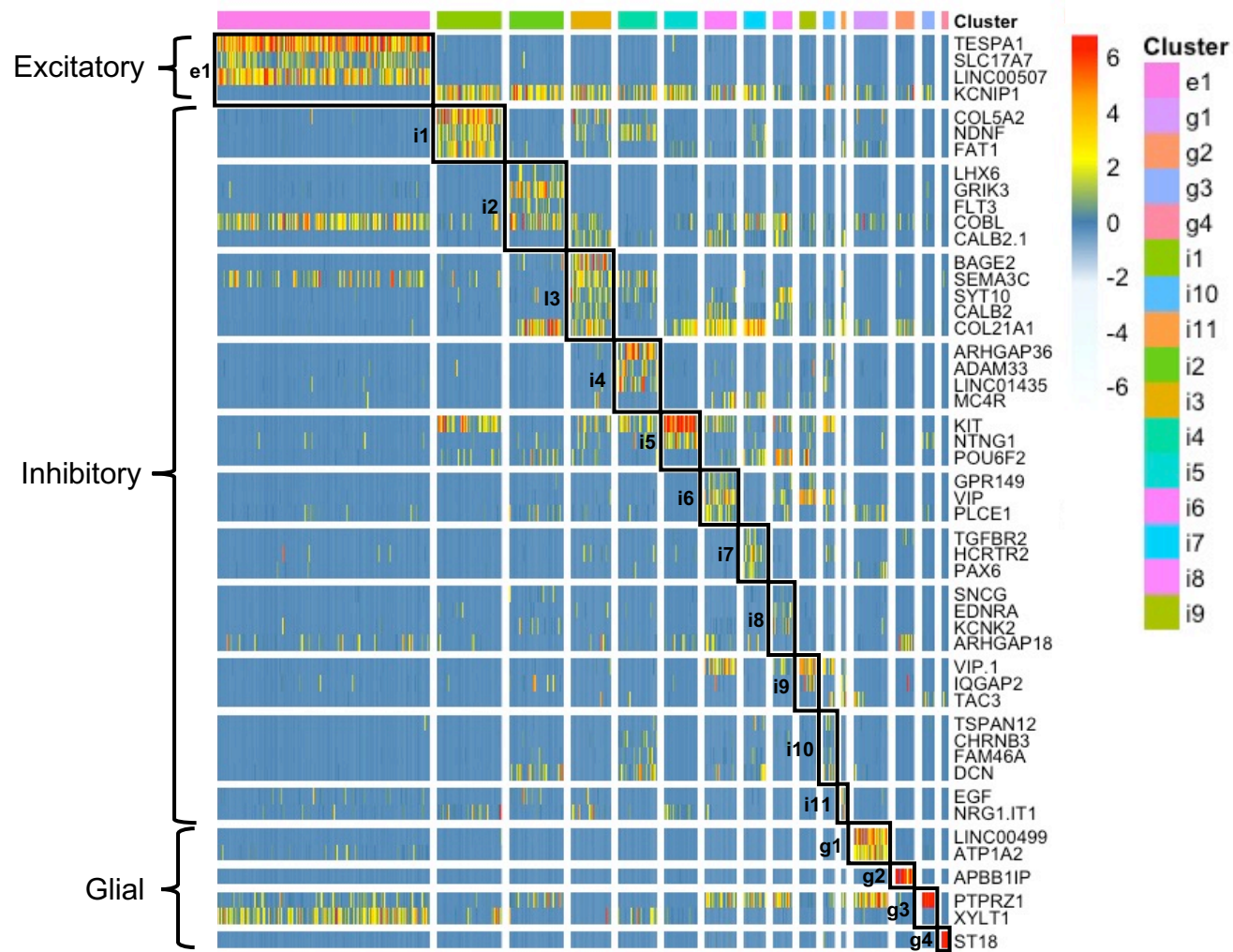
367

368 **Table 1. Cell types identified in cortical layer 1 of the human middle temporal gyrus**

Cluster ID	Cell Type Name	Cell Type Definition
e1	<i>TESPA1-expressing MTG cortical layer 1 excitatory neuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 excitatory neuron that selectively expresses TESPA1, LINC00507, and SLC17A7 mRNAs, and lacks expression of KCNIP1 mRNA</i>
i1	<i>COL5A2-expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 2 GABAergic interneuron that selectively expresses COL5A2 and NDNF and FAT1 mRNAs</i>
i2	<i>LHX6-expressing MTG cortical layer 2 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 2 GABAergic interneuron that selectively expresses LHX6, GRIK3, and FLT3, while of lacking expression of COBL and CALB2 mRNAs</i>
i3	<i>BAGE2 expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses BAGE2 and SEMA3C and SYT10 and CALB2 and COL21A1 mRNAs</i>
i4	<i>ARHGAP36 expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses ARHGAP36 and ADAM33 and LINC01435 and MC4R mRNAs</i>
i5	<i>KIT-expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses KIT and NTNG1 and POU6F2 mRNAs</i>
i6	<i>GPR149-expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses GPR149 and VIP and PLCE1 mRNAs</i>
i7	<i>TGFBR2 -expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses TGFBR2 and HCRTR2 and PAX6 mRNAs</i>
i8	<i>SNCG-expressing MTG cortical</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic</i>

	<i>layer 1 interneuron, human</i>	<i>interneuron that selectively expresses SNCG and EDNRA and KCNK2 and ARHGAP18 mRNAs</i>
i9	<i>VIP-expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses VIP and IQGAP2 and TAC3 mRNAs</i>
i10	<i>TSPAN12-expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses TSPAN12 and CHRN3 and FAM46A and DCN mRNAs</i>
i11	<i>EGF-expressing MTG cortical layer 1 interneuron, human</i>	<i>A human middle temporal gyrus cortical layer 1 GABAergic interneuron that selectively expresses EGF and NRG1-IT1 mRNAs</i>
g1	<i>Linc00499-expression MTG cortical layer 1 glial cell, human</i>	<i>A human middle temporal gyrus cortical layer 1 glial cell that selectively expresses Linc00499 and ATP1A2 mRNAs</i>
g2	<i>APBB1IP-expressing MTG cortical layer 1 glial cell, human</i>	<i>A human middle temporal gyrus cortical layer 1 glial cell that selectively expresses APBB1IP mRNAs</i>
g3	<i>PTPRZ1-expressing MTG cortical layer 1 glial cell, human</i>	<i>A human middle temporal gyrus cortical layer 1 glial cell that selectively expresses PTPRZ1 and XYLT1 mRNAs</i>
g4	<i>ST18 expressing MTG cortical layer 1 glial cell, human</i>	<i>A human middle temporal gyrus cortical layer 1 glial cell that selectively expresses ST18 mRNAs</i>





Annotations 

[rdfs:label](#) [language: en]

rosehip neuron

Description: 'rosehip neuron'

Equivalent To 

● interneuron

and ('has soma location' some 'cortical layer I')

and ('has soma location' some 'middle temporal gyrus')

and ('capable of' some 'gamma-aminobutyric acid secretion, neurotransmission')

and (expresses some KIT)

and (expresses some NTNG1)

and (expresses some POU6F2)

SubClass Of 

● 'synapsed to' some 'layer III pyramidal neuron'

☰ 'cerebral cortex GABAergic interneuron'