

1 Identification of genetic markers for cortical areas using a Random 2 Forest classification routine and the Allen Mouse Brain Atlas

3
4 Natalie Weed, Trygve Bakken, Nile Graddis, Nathan Gouwens,
5 Daniel Millman, Michael Hawrylycz, Jack Waters*

6
7 Allen Institute for Brain Science

8 615 Westlake Ave N

9 Seattle WA 98109

10
11
12 *email: jackw@alleninstitute.org.

13 14 Abstract

15 The mammalian neocortex is subdivided into a series of ‘cortical areas’ that are functionally and
16 anatomically distinct, and are often distinguished in brain sections using histochemical stains and other
17 markers of protein expression. We searched the Allen Mouse Brain Atlas, a database of gene expression,
18 for novel markers of cortical areas. We employed a random forest algorithm to screen for genes that
19 change expression at area borders. We found novel genetic markers for 19 of 39 areas and provide code
20 that quickly and efficiently searches the Allen Mouse Brain Atlas.

21 Introduction

22 The mammalian neocortex is classified into a series of anatomically and functionally distinct regions or
23 ‘cortical areas’ (Brodmann, 1909; Glasser *et al.*, 2016). Areas are often identified using histochemical
24 stains and antibodies to visualize differences in protein expression across cortex. Examples include
25 cytochrome oxidase histochemistry and antibodies against m2 muscarinic receptors (Wang, Sporns &
26 Burkhalter, 2012). Furthermore, global expression signatures of cortical areas have been identified in
27 human (Hawrylycz *et al.*, 2012), rhesus monkey (Bernard *et al.*, 2012) and mouse (Hawrylycz *et al.*,
28 2010), but few genes have been identified with distinct transitions between adjacent areas. We
29 reasoned that there may be genetic markers of cortical areas that have not been identified and that we
30 might identify additional markers by screening the Allen Mouse Brain Atlas, a database containing in situ
31 hybridization information for thousands of genes (Lein *et al.*, 2007). We developed numerical tools to
32 screen the many thousands of images in the database, using a random forest algorithm to identify
33 changes in gene expression at the boundaries of cortical areas defined in the Allen Mouse Brain
34 Reference Atlas (Kuan *et al.*, 2015). We found novel genetic markers for several areas. In addition, we
35 provide code that searches the Allen Mouse Brain Atlas quickly and efficiently for differences in gene
36 expression between cortical areas. With only minor modification, our code could be adapted to search
37 for genes that mark other brain regions, including subcortical nuclei.

38

39 Methods and Results

40 Our aim was to locate changes in gene expression between cortical regions in the mouse. From the Allen
41 Mouse Brain Atlas, we took coronal in situ hybridization (ISH) data resampled to a canonical 3D
42 reference space and overlaid the borders of cortical regions from the Allen Mouse Brain Reference Atlas.
43 To identify genes with differential expression along these boundaries, we used a Random Forest
44 algorithm.

45

46 Horizontal Projections

47 We obtained ISH data for 4345 genes from the Allen Mouse Brain Atlas (brain-map.org/api/index.html).
48 ISH data were of coronal sections (Figure 1A). However, the perspective that best captures most borders
49 delineating cortical areas while eliminating excess information is the horizontal plane. To obtain a
50 horizontal plane perspective from coronal sections, we created two projections for each gene: a ‘top
51 projection’ and a ‘flat map projection’. Each projection was created in three steps, with the first two

52 steps being common to both projections. Firstly, we isolated cortical fluorescence and eliminated
53 fluorescence from subcortical structures by applying a mask derived from the Allen SDK (2015)
54 structure_tree class (Figure 1B and C). Secondly, we created a maximum intensity surface projection: for
55 each pixel on the cortical surface, we projected the fluorescence in the underlying tissue along a line
56 perpendicular to the pial surface of cortex. One might think of this first step as creating a curved sheet
57 of fluorescence intensity values at the surface of cortex. Finally, we projected these surface values to the
58 horizontal plane, creating a top projection (Figure 1D) or we ‘unfurled’ the curved cortical sheet to
59 create a flat map (Figure 1G). The flat map was particularly valuable in the study lateral cortical regions,
60 which are under-represented in top projections.

61 All ISH data in the Allen Mouse Brain Atlas are spatially registered to the Allen Mouse Brain
62 Reference Atlas ([http://help.brain-](http://help.brain-map.org/display/mousebrain/Documentation?preview=/2818169/8454277/MouseCCF.pdf)
63 [map.org/display/mousebrain/Documentation?preview=/2818169/8454277/MouseCCF.pdf](http://help.brain-map.org/display/mousebrain/Documentation?preview=/2818169/8454277/MouseCCF.pdf)). Hence all
64 data utilized are inherently co-aligned with the Allen Mouse Brain Reference Atlas and the locations of
65 brain areas can be readily superimposed on the ISH results. To locate cortical regions in the top
66 projection and create a cortical area map, we extracted the corresponding cortical area masks using the
67 structure_tree class and projected these masks to the horizontal plane, as described for ISH projections.
68 Simplification of three-dimensional data into two dimensions allowed for fast quantitative analysis as
69 well as easy visualization of expression patterns.

70

71 Random Forest Algorithm

72 When examining the ISH results, two limitations became apparent. Firstly, there are gaps in some data
73 sets, with missing data manifest as dark pixels in coronal images or dark medial-lateral bands in the top
74 projections (Figure 2A). Secondly, there is pronounced section-to-section variability in mean
75 fluorescence that appears as coronal banding or ‘stripes’ in top projections (Figure 2B). Together these
76 two effects often result in variation in pixel values, independent of variation due to differential gene
77 expression. These data properties complicate the comparison of fluorescence along the anterior-
78 posterior axis and, thereby, the comparison of expression between cortical regions. Rather than attempt
79 to mitigate these issues directly, we trained a Random Forest algorithm to classify pixels as either inside
80 or outside each cortical region, essentially learning the variance in the data.

81 We examined 39 cortical regions from the Allen Mouse Brain Reference Atlas for potential gene
82 markers. Each search involved comparison of one cortical region to expression patterns of all genes,
83 imputed as independent variables to the Random Forest algorithm. Random forest was implemented in

84 Python using the scikit-learn package (Pedregosa *et al.*, 2011). Nodes were determined by Gini Index
85 criteria $\sum_{k=1}^K (p_{mk} (1 - p_{mk}))$. Each random forest consisted of 100 decision trees. Random state was
86 initialized at 0. The importance of each variable was also determined by Gini Index criteria – reduction
87 in Gini Index each time a split occurred was attributed to the variable, and that variable-associated
88 reduction was divided by total reduction in Gini Index across the entire random forest to return the
89 variable importance value. Total variable importance across all genes summed to one.

90 For each cortical region, three outputs from the random forest were produced and analyzed: (1) a
91 confusion matrix, indicating the success rate of the classification algorithm; (2) the list of all 4345 genes,
92 ranked in decreasing order of variable importance, where importance is a pseudo-measure of the
93 expression predictive power across the cortical border; and (3) the importance values, one for each
94 gene. The Random Forest was trained to distinguish pixels within a cortical region from pixels in the
95 surrounding area outside of the cortical region. The inputs to the Random Forest algorithm were the
96 gene expression fluorescence intensity values of all 4345 genes for each pixel and the corresponding
97 labels for each pixel as cortical region or surrounding region. Surrounding pixels were identified by
98 dilating the region mask by 30 iterations using SciPy ndimage package in Python, translating to roughly
99 30 pixels in distance in each direction. Pixels were split into training and test sets, with 100 randomly
100 selected pixels held out as a test set and the remaining pixels forming the training set. This represented
101 less than 1% of total pixels classified for each cortical area. Data was divided using the scikit-learn
102 model_selection package in Python. Hence the training array input into the Random Forest algorithm
103 consisted of a 2D array of dimensions 4345 by $N - 100$, where 4345 is the number of genes and N is the
104 number of pixels within the dilated mask, and each cell in the array corresponding to a luminance value
105 of the pixel. A second array of dimensions 1 by $N - 100$ indicated the binary labels, inside or outside the
106 cortical region (Figure 2D). After training, performance of the algorithm was tested on the held-out
107 pixels (array dimensions 4345 by 100) for which the binary classification was withheld. Withheld pixels
108 were randomly selected, creating a test set that was representative of the cortical area: balanced inside
109 and outside the cortical region, and varying in distance from the cortical area boundary. Hence the
110 algorithm returned the cross-validated binary classification for 100 withheld pixels, which was compared
111 to known classification and used to plot a confusion matrix (Figure 2E), summarizing performance of the
112 Random Forest. The displayed confusion matrix is averaged over all folds for the specified cortical
113 region.

114 Results for primary somatosensory barrel field are illustrated in Figure 2E-G. The model correctly
115 classified 52 of 54 test pixels within the barrel field and 44 of 46 test pixels outside barrel field, resulting

116 in a combined model accuracy of 96% (Figure 2E). Most genes exhibited low variable importance (Figure
117 2F). We ranked genes by their random forest variable importance values. The gene with rank 1 exhibited
118 a distinct change along the border (Figure 2G). The gene with rank 10 exhibited a subtler change and the
119 gene at rank 100 exhibited no obvious change along the border (Figure 2G). Hence the Random Forest
120 algorithm accurately classified most pixels and, via a ranked list of genes, identified a short list of genes
121 that might act as putative genetic markers of the cortical region.

122

123 Genetic markers of cortical areas

124 To identify genetic markers for each cortical area, we manually inspected the top projections or flat
125 maps for the top 10 genes, as determined by the Random Forest results. Adequate information for
126 classification of barrel field was included in the 10 highest importance genes since running our analysis
127 with only these 10 genes as inputs conserved prediction accuracy at 96%. Of the 45 cortical regions
128 tested, we identified potential genetic markers for 19 (Table 1). 6 cortical areas were determined too
129 small to reliably examine, resulting in 39 regions to explore. Of the six markers identified by Hawrylycz *et*
130 *al.* (2010), three were extracted by our method (*Man1a*, somatomotor; *Rorb*, somatosensory; *Scnn1a*,
131 ventral retrosplenial), one included an area that was not explored (*Smoc1*, gustatory), one was in a
132 region where many other strong markers were identified (*Rreb1*, retrosplenial), and one was not
133 identified (*Hap1*, ectohippocampal). Selected potential genetic markers indicated relatively high sensitivity for
134 area marking, low specificity due to the point selection process, bilateral expression, and entire cortical
135 area contrast.

136 Examples of expression patterns are provided in Figure 3. For primary somatosensory cortex barrel
137 field, we identified *Rspo1* as a strong candidate gene (Figure 3A). Expression of *Rspo1* is relatively high in
138 the barrel field, moderate through somatosensory areas, and low in motor cortex. There were multiple
139 markers for motor cortex, including *Wnt7b* (Figure 3B), but we found no compelling markers for primary
140 or secondary motor cortex. *Rorb* was also identified as a potential marker, specifically for primary
141 sensory cortices (Figure 3C). This provided an additional positive control that our method was robust
142 and effective, as *Rorb* is an established marker for primary sensory areas (Hawrylycz *et al.*, 2010; Zhuang
143 *et al.*, 2017). *Cdh24* marked primary auditory cortex (Figure 3E). We found multiple genes that labeled
144 all or subregions of retrosplenial cortex. For example, *Tmem215* marked dorsal retrosplenial cortex (and
145 primary somatosensory cortex) and *Npsr1* marked all of retrosplenial cortex (Figure 3D, F). In flat maps,
146 *Serpinf1* was identified as a marker of the frontal pole (Figure 3G) and temporal association cortex was
147 marked by *Lifr* (Figure 3H). Notably, some of these genes exhibit mediolateral stripes in the top

148 projection, indicating that our method is robust to the missing data and expression-independent
149 variability in signal.

150

151 Cellular Basis of Potential Genetic Markers

152 Our Random Forest searches, applied to ISH data, identified genes that marked cortical area borders,
153 but provided no insight into the cellular basis of the genetic markers. Does, for example, a change in
154 gene expression result from an abrupt change in the density of a cell type with unique gene expression;
155 or might the border result from a change in gene expression by a cell type that straddles the border?
156 Gene expression in primary visual cortex and anterior lateral motor cortex has been studied using single-
157 cell RNA sequencing (Tasic *et al*, 2018). From this transcriptomic data set and associated analysis tools
158 (<https://github.com/AllenInstitute/scrattch.vis>), we examined cell types that express the marker genes
159 identified in this study by top view projections (Figure 4). Markers were expressed in many different cell
160 types. Several genes (*Adcy8*, *Bmp3*, *Cacnb3*, *Npsr1*, *Vgf*, *Zmat4*) were expressed mostly in neurons and
161 not non-neuronal cells, suggesting that the border-related change in expression of these genes was
162 neuronal. Some genes were expressed mostly in a cell sub-population, suggesting that there is likely a
163 border-related change in the density of these cells or of their expression of one gene. For example,
164 *Rspo1*, *Serpinf1* and *Man1a* are expressed in layer 4 excitatory neurons, vascular and leptomenigeal
165 cells (VLMC) and macrophages, respectively. Unsurprisingly, our results are consistent with changes in
166 gene expression marking cortical areas arising from changes in cell density in some instances and from
167 changes in gene expression within a cell population in other instances. Importantly, both instances
168 appear to have been detected by our Random Forest analysis of ISH data.

169

170 Discussion

171 We used a Random Forest algorithm to identify a short list of potential gene markers from thousands of
172 candidate genes, applying this approach to 39 cortical regions in the mouse. Our results identified 44
173 putative markers, marking 19 of the explored regions.

174 The spatial resolution and number of genes in the database places limits on the conclusions we can
175 draw. Firstly, the voxel size of the ISH quantification in the database is 200 μm . Once missing and
176 variable data is considered, the maximum accuracy we can hope to achieve is on the order of hundreds
177 of micrometers, resulting in an imperfect match between area borders and gene expression.

178 Subsequent experiments such as immunostaining for the genes we have identified would be necessary
179 to confirm our results and to assess the accuracy with which each gene marks borders. Furthermore, the

180 database includes coronal ISH images for 4345 genes. It may be that genes not sampled here mark some
181 of the 25 cortical regions for which we were unable to identify markers. Repeating our analysis on a
182 larger data set, should one become available, might identify further markers.

183 Alternative methods include gene identification by direct comparison of expression difference along
184 the cortical border. However, pooling of more pixels than those available solely along borders was
185 necessary to overcome high luminance variability across pixels and coronal sections. For this reason, and
186 the difficulty of direct quantification of variance in our data set, we decided to pursue random forest
187 classification as our selected model. Variable importance may be inaccurately skewed towards higher
188 sampled variables or continuous data types, and thus unusable; however, because our predictor
189 variables exhibit identical scale of measurement and data type, importance rank can be taken as
190 unbiased (Strobl *et al.*, 2007). Random forest uses a bootstrapped subset of variables at each splitting
191 node when building decision trees. By accumulating many splits on previously subdivided pixels, genes
192 are evaluated at subregions of the cortical area. Given this property, we find that occasionally genes
193 with relatively high variable importance display marking of a single border rather than the entire cortical
194 area. However, if a gene exhibits clear marking of all cortical borders, it is shown with higher variable
195 importance than an alternate gene expression pattern marking only a single border. Random forest is an
196 accurate, computationally efficient, and easily interpretable method of classification. This was important
197 as many of our data sets, especially for larger cortical areas like somatosensory areas, reached sizes of
198 almost 250,000 pixels, evaluated at 4345 genes. Each output predicted took less than a minute to
199 compile the data set, run computations, and produce outputs on a desktop computer. By maximizing
200 concurrent computation across all available cores, the time required to run is minimized while not
201 sacrificing predictive power of our model, as exhibited by the high accuracy of the random forest.

202 By dilating the cortical area mask a small amount instead of comparing the area of interest to the
203 entire cortex, we allowed for differential expression of the gene in more distant parts of cortex. This is
204 by design, as expression far from the desired cortical region does not impact the ability of the gene to
205 mark the border. However, potential uniquely expressed genes are still a subset of those that can be
206 identified with our method, and our method could be readily modified to solely identify uniquely
207 expressed genes. Similarly, the method could be easily extended to investigate laminar differences or
208 expression patterns in subcortical structures as masks for cortical layers and for subcortical structures
209 are included in the Allen Mouse Brain Reference Atlas.

210 References

- 211 Bernard A, Lubbers LS, Tanis KQ, Luo R, Podtelezchnikov AA, Finney EM, McWhorter MM, Serikawa K,
212 Lemon T, Morgan R, Copeland C, Smith K, Cullen V, Davis-Turak J, Lee CK, Sunkin SM, Loboda AP,
213 Levine DM, Stone DJ, Hawrylycz MJ, Roberts CJ, Jones AR, Geschwind DH, Lein ES (2012)
214 Transcriptional architecture of the primate neocortex. *Neuron* 73(6), 1083-1099.
- 215 Brodmann K (1909). *Vergleichende Lokalisationslehre der Grosshirnrinde in ihren Prinzipien dargestellt*
216 *auf Grund des Zellenbaues*. Leipzig: Barth JA.
- 217 Glasser MF, Coalson TS, Robinson EC, Hacker CD, Harwell J, Yacoub E, Ugurbil K, Andersson J, Beckmann
218 CF, Jenkinson M, Smith SM, Van Essen DC (2016) A multi-modal parcellation of human cerebral
219 cortex. *Nature* 536, 171-178.
- 220 Hawrylycz MJ, Lein ES, Guillozet-Bongaarts AL, Shen EH, Ng L, Miller JA, van de Lagemaat LN, Smith KA,
221 Ebbert A, Riley ZL, Abajian C, Beckmann CF, Bernard A, Bertagnolli D, Boe AF, Cartagena PM,
222 Chakravarty MM, Chapin M, Chong J, Dalley RA, Daly BD, Dang C, Datta S, Dee N, Dolbeare TA, Faber
223 V, Feng D, Fowler DR, Goldy J, Gregor BW, Haradon Z, Haynor DR, Hohmann JG, Horvath S, Howard
224 RE, Jeromin A, Jochim JM, Kinnunen M, Lau C, Lazarz ET, Lee C, Lemon TA, Li L, Li Y, Morris JA, Overly
225 CC, Parker PD, Parry SE, Reding M, Royall JJ, Schulkin J, Sequeira PD, Slaughterbeck CR, Smith SC,
226 Sodt AJ, Sunkin SM, Swanson BE, Vawter MP, Williams D, Wohnoutka P, Zielke HR, Geschwind DH,
227 Hof PR, Smith SM, Koch C, Grant SGN, Jones AR (2012) An anatomically comprehensive atlas of the
228 adult human brain transcriptome. *Nature* 489, 391-399.
- 229 Hawrylycz MJ, Bernard A, Lau C, Sunkin SM, Chakravarty MM, Lein ES, Jones AR, Ng L (2010) Areal
230 and laminar differentiation in the mouse neocortex using large scale gene expression data.
231 *Methods* 50(2), 113-121.
- 232 Kuan L, Li Y, Lau C, Feng D, Bernard A, Sunkin SM, Zeng H, Dang C, Hawrylycz M, Ng L, Li Y LC (2015)
233 Neuroinformatics of the Allen Mouse Brain Connectivity Atlas. *Methods* 73, 4–17.
- 234 Lein, E.S, Hawrylycz MJ, Ao N, Ayres M, Bensinger A, Bernard A, Boe AF, Boguski MS, Brockway KS,
235 Byrnes EJ, Chen L, Chen L, Chen TM, Chin MC, Chong J, Crook BE, Czaplinska, A, Dang CN, Datta S,
236 Dee NR, Desaki AL, Desta T, Diep E, Dolbeare TA, Donelan MJ, Dong HW, Dougherty JG, Duncan BJ,
237 Ebbert AJ, Eichele G, Estin LK, Faber C, Facer BA, Fields R, Rischer SR, Fliss TP, Frensley C, Gates SN,
238 Glattfelder KJ, Halverson KR, Hart MR, Hohmann JG, Howell MP, Jeung DP, Johnson RA, Karr PT,
239 Kawal R, Kidney JM, Knapik RH, Kuan CL, Lake JH, Laramée AR, Larsen KD, Lau C, Lemon TA, Liang AJ,
240 Liu Y, Luong LT, Michaels J, Morgan JJ, Morgan RJ, Mortrud MT, Mosqueda NF, Ng LL, Ng R, Orta GJ,
241 Overly CC, Pak TH, Parry SE, Pathak SD, Pearson OC, Puchalski RB, Riley ZL, Rockett HR, Rowland SA,

242 Royall JJ, Ruiz MJ, Sarno NR, Schaffnit K, Shapovalova NV, Sivasay T, Slaughterbeck CR, Smith SC,
243 Smith KA, Smith BI, Sodt AJ, Stewart NN, Stumpf KR, SUnkin SM, Sutram M, Tam A, Temmer CD,
244 Thaller C, Thompson CL, Varnam LR, Visel A, Whitlock RM, Wohnoutka PE, Wolkey CK, Wong VY,
245 Wood M, Yaylaoglu MB, Young RC, Youngstrom BL, Yuan XF, Zhang B, Zwingman TA, Jones AR (2007)
246 Genome-wide atlas of gene expression in the adult mouse brain. *Nature* 445, 168-176.

247 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss
248 R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E (2011)
249 Scikit-learn: Machine Learning in Python. *JMLR* 12, 2825-2830.

250 Strobl C, Boulestin AL, Zeileis A, Hothorn T (2007) Bias in random forest variable importance measures:
251 Illustrations, sources and a solution. *BMC Bioinformatics* 8, 25.

252 Tasic B, Yao Z, Graybiack LT, Smith KA, Nguyen TN, Bertagnolli D, Goldy J, Garren E, Economo MN,
253 Viswanathan S, Osnat P, Bakken T, Menon V, Miller J, Fong O, Hirokawa KE, Lathia K, Rimorin C, Tieu
254 M, Larsen R, Casper T, Barkan E, Kroll M, Parry S, Shapovalova NV, Hirschstein D, Pendergraft J,
255 Sullivan HA, Kim TK, Szafer A, Dee N, Groblewski P, Wickersham I, Cetin A, Harris JA, Levi BP, Sunkin
256 SM, Madisen L, Daigle TL, Looger L, Bernard A, Phillips J, Lein E, Hawrylycz M, Svoboda K, Jones AR,
257 Koch C, Zeng H (2018) Shared and distinct transcriptomic cell types across neocortical areas. *Nature*
258 563, 72-28.

259 Wang Q, Sporns O, Burkhalter A (2012) Network analysis for corticocortical connections reveals ventral
260 and dorsal processing streams in mouse visual cortex. *J Neurosci* 32(13), 4386-4399.

261 Zhuang J, Ng L, Williams D, Valley M, Li Y, Garrett M, Waters J (2017) An extended retinotopic map of
262 mouse cortex. *eLife* 6, e18372.

263
264

265 Acknowledgements

266 We thank the Allen Institute founder, Paul G. Allen, for his vision, encouragement and support.

267 Figure Legends

268 **Figure 1. Creation of top projections from coronal images of gene expression.**

269 (A) Image of a coronal section from the Allen Mouse Brain Atlas. Gene: *Rorb*. ISH intensity is normalized
270 to a range of 0 to 1. Color scale shown in panel C. (B) Binary cortex mask with value = 1 for cortical pixels
271 and value = 0 for subcortical pixels. (C) Product of the images in A and B, resulting in ISH intensity values
272 in cortical pixels and zero's in subcortical locations. (D) Schematic illustration of the projection process
273 used to generate top projections. (E) Top projection for *Rorb*. Dashed line indicates location of section in
274 panel A. (F) Cortical boundaries from the Allen Mouse Brain Reference Atlas, overlaid onto the gene
275 expression top view of panel E. (G) Schematic illustration of the projection process used to generate flat
276 map top projections. (H) Flat map projection for *Rorb*. Dashed line indicates location of section in panel
277 A. (I) Cortical boundaries from the Allen Mouse Brain Reference Atlas, overlaid onto the gene expression
278 top view of panel H.

279

280 **Figure 2. Random forest algorithm: method and example results.**

281 (A) A coronal image (left) and the top projection (right) for gene *Nvl*. Note the missing data (black
282 pixels). (B) Top view for gene *Adra1d*. Note the pronounced variation in density along the a-p axis. (C)
283 Binary mask for primary somatosensory cortex barrel field (SSp-bfd). Light gray inside SSp-bfd; darker
284 grey marks pixels in surrounding region. White lines: boundaries of cortical areas in Allen Mouse Brain
285 Reference Atlas. (D) Schematic illustration of arrays input into Random Forest algorithm. Columns
286 correspond to gene, rows to pixels in the top projection data set. Each value is an ISH luminance value.
287 Classification of pixel is taken from the reference mask (panel C). (E) Confusion matrix output from
288 Random Forest algorithm for SSp-bfd. 0 indicates point outside SSp-bfd, 1 indicates point inside SSp-bfd.
289 (F) Gene importance histogram. Importance values approximate a logarithmic distribution. (G) Examples
290 of genes that mark SSp-bfd, with overlaid Allen Mouse Brain Reference Atlas borders. *Nov* rank 1,
291 importance 0.022. *Hlf* rank 10, importance value 0.0081. *Stap2* rank 100, importance 0.0018.

292

293 **Figure 3. Examples of markers for cortical areas.**

294 (A) R-Spondin 1 labels primary somatosensory cortex barrel field. (B) Wnt Family Member 7B labels
295 primary motor cortex. (C) Retinoid-Related Orphan Receptor, Beta labels primary sensory areas. (D)
296 Transmembrane Protein 215 labels dorsal retrosplenial cortex. (E) Cadherin 24 labels primary auditory
297 cortex. (F) Neuropeptide S Receptor 1 labels retrosplenial cortex. (G) Serpin Family F Member 1 labels
298 frontal pole. (H) Leukemia Inhibitory Factor Receptor labels temporal association cortex. Panels A-F

299 provide examples of genes identified in top projections, panels G and H of genes identified in flat maps.
300 Cortical areas of interest are marked with cyan borders.

301

302 **Figure 4. Single Cell RNA-sequencing expression plot**

303 Log-transformed average expression of top-projection identified potential areal marker genes in mouse
304 cortical cells grouped into subtypes of three major cell classes: inhibitory neurons, excitatory neurons,
305 and non-neuronal cells. Expression data was measured by RNA-sequencing of single cells isolated from
306 Primary Visual Cortex (VISp) and Anterior Lateral Motor Cortex (ALM). Color corresponds to expression
307 value, with warmer colors indicating high expression and cooler colors indicating low expression. Max
308 value indicates the maximum expression per gene, measured in average CPM per cluster. CPM = counts
309 per million reads.

310

311 **Table 1. List of Potential Genetic Markers**

312 Potential cortical boundary genetic markers, listed by cortical region, as identified by Random Forest
313 variable importance classifier. All explored regions listed. Regions with no listed genes displayed no clear
314 potential genetic marker. Each gene was identified independently. Asterisks indicate regions explored
315 and genes identified with flat map projections.

figure 1

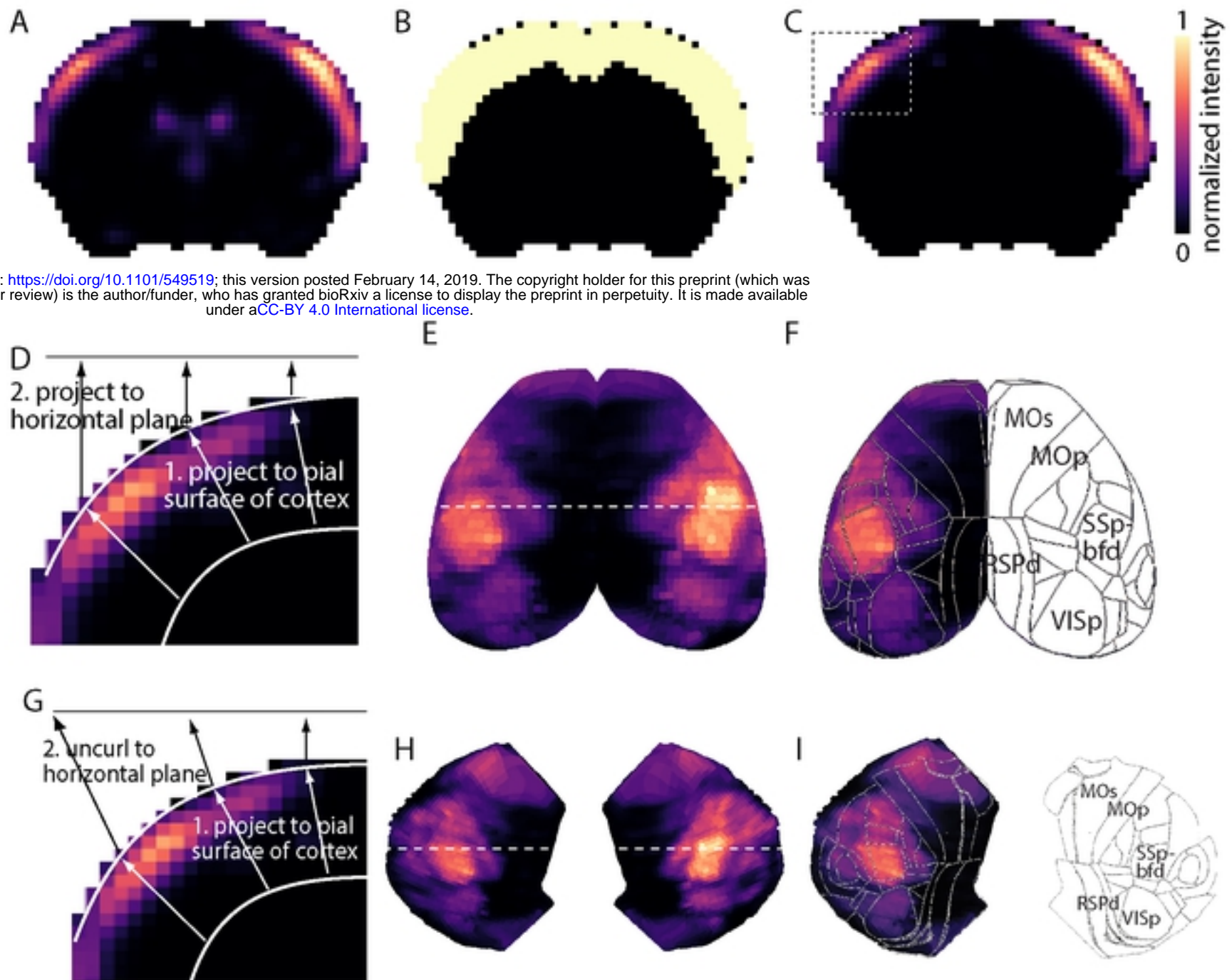


Figure 1

figure 4

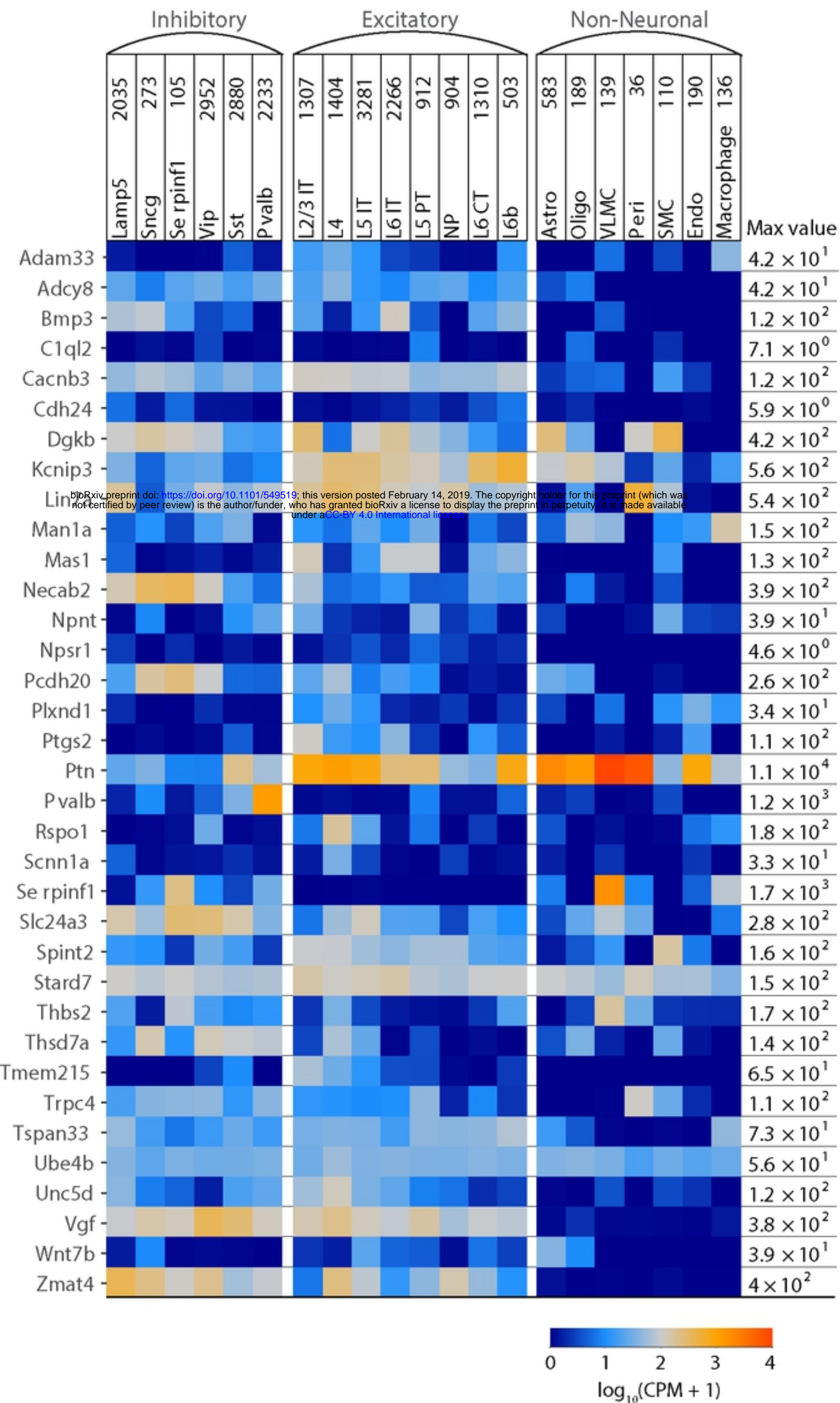


Figure 4

figure 2

doi: <https://doi.org/10.1101/549519>; this version posted February 14, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

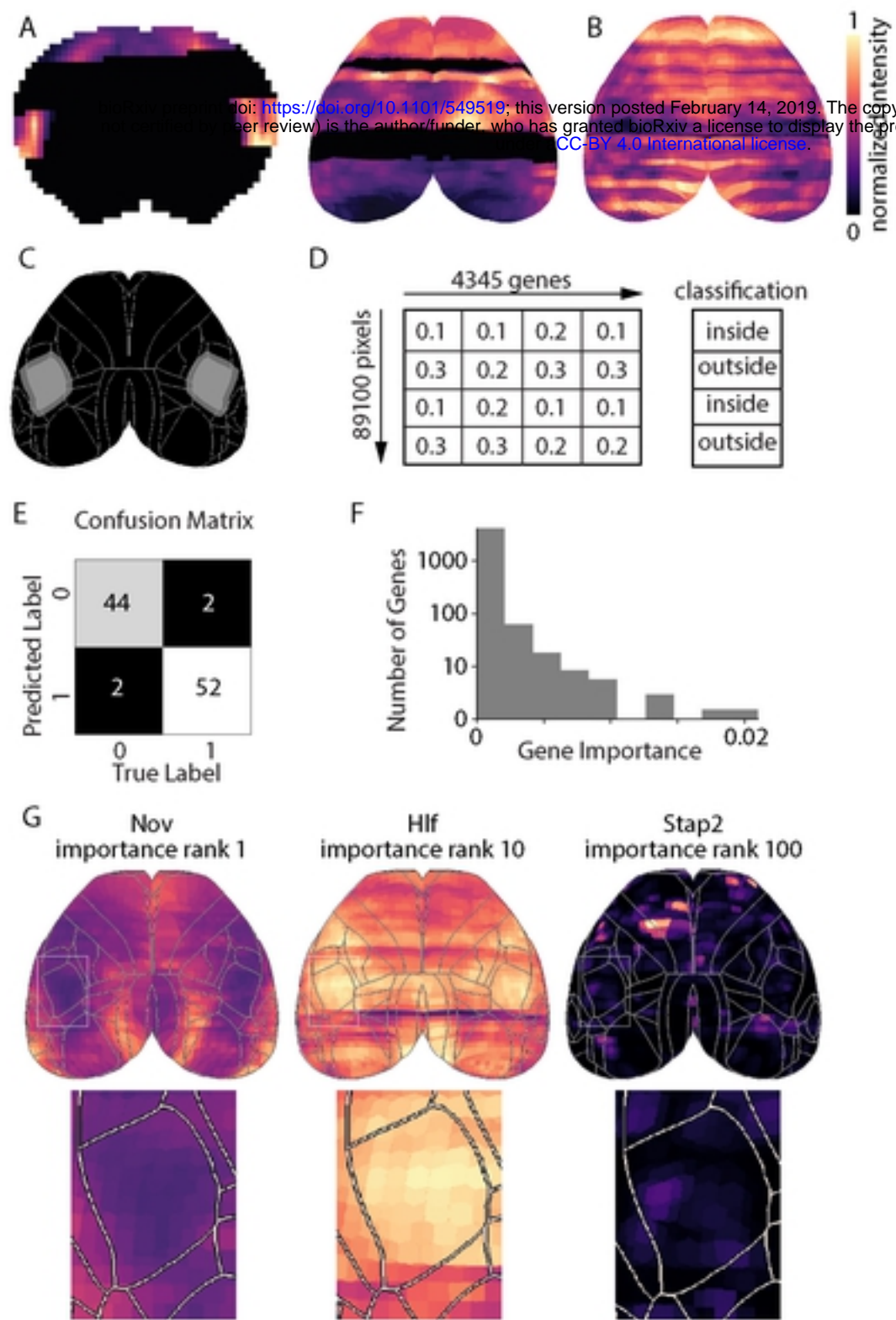


Figure 2

figure 3

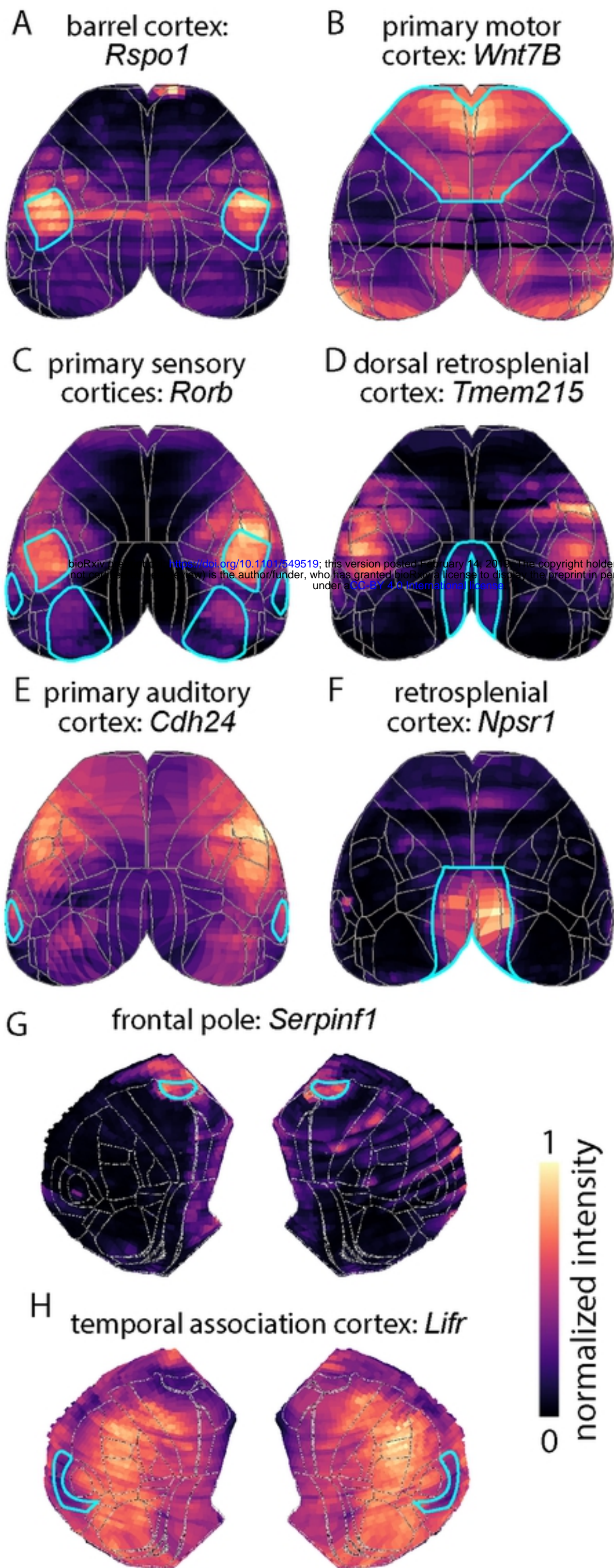


Figure 3

Table 1

Cortical Region	Potential Genetic Markers
Somatomotor	<i>Ube4b, Lin7a, Man1a, Wnt7b</i>
- Primary	
- Secondary	
Somatosensory	<i>Adam33, Vgf, Kcnip3, Rorb, Pcdh20</i>
- Primary	<i>Pvalb</i>
-- Primary, barrel	<i>Rspo1</i>
-- Primary, nose	
-- Primary, l. limb	
-- Primary, mouth	
-- Primary, u. limb	
-- Primary, trunk	<i>Trpc4</i>
-- Primary, unassigned	
Auditory	<i>Vspan5, *Coch</i>
- Primary*	<i>Unc5d, Zmat4, Rspo1, Cdh24, *Ptn, *Chn2</i>
- Dorsal*	<i>Dgkb</i>
- Posterior*	
- Ventral*	
Visual	<i>Ptgs2</i>
- Primary	<i>Slc24a3, Thbs2</i>
- Lateral	
- Anterolateral	
- Anteromedial	<i>Stard7</i>
- Posterolateral	
- Posteromedial	
- Postrhinal	<i>Bmp3, Plxnd1, Spint2</i>
- Laterointermediate	
- Rostrolateral	
Anterior Cingulate	
- Dorsal	
Retrosplenial	<i>Cacnb3, Npsr1, Mas1, C1ql2</i>
- Lateral agranular	
- Dorsal*	<i>Tmem215, Npnt</i>
- Ventral*	<i>Adcy8, Scnn1a, Sm1399, Necab2, Thsd7a, *Dpysl5</i>
Temporal Association*	<i>*Lifr</i>
Ectorhinal*	<i>*Kctd4</i>
Medial Orbital*	<i>*Dtx1</i>
Visceral*	
Prelimbic*	
Frontal Pole*	<i>*Serpinf1</i>

bioRxiv preprint doi: <https://doi.org/10.1101/549519>; this version posted February 14, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.