

1 **TITLE:**

2
3 **Single-cell epigenomic identification of inherited risk loci in Alzheimer's and**
4 **Parkinson's disease**

5
6 **AUTHOR LIST AND AFFILIATIONS:**

7
8 M. Ryan Corces^{1,2}, Anna Shcherbina^{3,4}, Soumya Kundu^{4,5}, Michael J. Gloudemans¹, Laure
9 Frésard¹, Jeffrey M. Granja^{2,4,6}, Bryan H. Louie^{1,2}, Shadi Shams^{2,4}, S. Tansu Bagdatli^{2,4}, Maxwell
10 R. Mumbach^{2,4}, Bosh Liu^{1,7}, Kathleen S. Montine¹, William J. Greenleaf^{2,4,8,9}, Anshul
11 Kundaje^{4,5}, Stephen B. Montgomery^{1,4}, Howard Y. Chang^{2,4,10,11,*}, Thomas J. Montine^{1,*}

12
13 ¹Department of Pathology, Stanford University School of Medicine, Stanford, CA, USA.

14 ²Center for Personal Dynamic Regulomes, Stanford University, Stanford, CA, USA.

15 ³Department of Biomedical Data Science, Stanford University School of Medicine, Stanford,
16 CA, USA.

17 ⁴Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA.

18 ⁵Department of Computer Science, Stanford University, Stanford, CA, USA.

19 ⁶Program in Biophysics, Stanford University, Stanford, CA, USA.

20 ⁷Department of Biology, Stanford University, Stanford, CA, USA.

21 ⁸Department of Applied Physics, Stanford University, Stanford, CA, USA.

22 ⁹Chan-Zuckerberg Biohub, San Francisco, CA, USA.

23 ¹⁰Program in Epithelial Biology, Stanford University, Stanford, CA, USA.

24 ¹¹Howard Hughes Medical Institute, Stanford University, Stanford, CA, USA.

25
26 *Correspondence should be addressed to T.J.M. (tmontine@stanford.edu) or H.Y.C.
27 (howchang@stanford.edu)

28
29 **Contact Information**

30 Thomas J. Montine, MD, PhD

31 Stanford University School of Medicine

32 Lane L235, 300 Pasteur Dr., Stanford, CA, 94305-5324

33 Email: tmontine@stanford.edu

34 Phone: 650-725-9352

35

36 Howard Y. Chang, MD, PhD

37 Stanford University School of Medicine

38 CCSR 2155c, 269 Campus Drive, Stanford, CA 94305-5168

39 Email: howchang@stanford.edu

40 Phone: 650-736-0306

41 **ABSTRACT**

42
43 Genome-wide association studies (GWAS) have identified thousands of variants associated with
44 disease phenotypes. However, the majority of these variants do not alter coding sequences, making
45 it difficult to assign their function. To this end, we present a multi-omic epigenetic atlas of the
46 adult human brain through profiling of the chromatin accessibility landscapes and three-
47 dimensional chromatin interactions of seven brain regions across a cohort of 39 cognitively healthy
48 individuals. Single-cell chromatin accessibility profiling of 70,631 cells from six of these brain
49 regions identifies 24 distinct cell clusters and 359,022 cell type-specific regulatory elements,
50 capturing the regulatory diversity of the adult brain. We develop a machine learning classifier to
51 integrate this multi-omic framework and predict dozens of functional single nucleotide
52 polymorphisms (SNPs), nominating gene and cellular targets for previously orphaned GWAS loci.
53 These predictions both inform well-studied disease-relevant genes, such as *BIN1* in microglia for
54 Alzheimer's disease (AD) and reveal novel gene-disease associations, such as *STAB1* in microglia
55 and *MAL* in oligodendrocytes for Parkinson's disease (PD). Moreover, we dissect the complex
56 inverted haplotype of the *MAPT* (encoding tau) PD risk locus, identifying ectopic enhancer-gene
57 contacts in neurons that increase *MAPT* expression and may mediate this disease association. This
58 work greatly expands our understanding of inherited variation in AD and PD and provides a
59 roadmap for the epigenomic dissection of noncoding regulatory variation in disease.

60
61 **INTRODUCTION**

62
63 Alzheimer's disease (AD) and Parkinson's disease (PD) affect ~50 and ~10 million individuals
64 world-wide, as two of the most common neurodegenerative disorders. Several large consortia have
65 assembled genome-wide association studies (GWAS) that associate genetic variants with clinical
66 diagnoses of probable AD dementia¹⁻⁴ or probable PD⁵⁻⁷, or with their characteristic pathologic
67 features. These efforts have led to the identification of dozens of potential risk loci for these
68 prevalent neurodegenerative diseases. One goal of these studies was to build more precise
69 molecular biomarkers of AD or PD, efforts that are beginning to yield encouraging results with
70 polygenic risk scores⁸. The other major goal was to gain deeper insight into the molecular
71 pathogenesis of disease and thereby inform novel therapeutic targets. Some of the risk loci contain
72 coding variants and so have credibility as putative disease mediators. However, most risk loci are
73 in noncoding regions and so it remains unclear if the nominated (often nearest) gene is the
74 functional disease-relevant gene, or if some other gene is involved⁹. Furthermore, even if the
75 nominated gene is a true positive, the noncoding risk locus might regulate additional genes. These
76 challenges remain a fundamental gap in interpreting the etiology of neurodegenerative diseases
77 and detecting high-confidence therapeutic targets.

78 To an extent not achieved in other organs, human brain function is closely coupled to region
79 and thus cellular composition. However, GWAS are agnostic to the regional and cellular
80 heterogeneity of the brain, making it difficult to *a priori* predict which brain regions or specific

81 cell types may mediate the phenotypic association. In addition, functional noncoding SNPs would
82 be predicted to exert their effects through alteration of gene expression via perturbation of
83 transcription factor binding and regulatory element function⁹. Moreover, such regulatory elements
84 are highly cell type-specific¹⁰. Thus, comprehensive nomination of putative functional noncoding
85 SNPs in the brain requires cataloging the regulatory elements that are active in every brain cell
86 type in the correct organismal and regional context. These critical data will illuminate the
87 functional significance of genetic risk loci in the molecular pathogenesis of common
88 neurodegenerative diseases.

89 Here, we have further expanded upon the current understanding of inherited variation in
90 neurodegenerative disease through implementation of a multi-omic framework that enables
91 accurate prediction of functional noncoding SNPs. This framework layers bulk Assay for
92 Transposase-accessible chromatin using sequencing (ATAC-seq)¹¹, single-cell ATAC-seq
93 (scATAC-seq)¹², and HiChIP enhancer connectome^{13,14} data over a machine learning classifier to
94 predict putative functional SNPs driving association with neurodegenerative diseases. Through
95 these efforts, we pinpoint putative target genes and cell types of several noncoding GWAS locus
96 in AD and PD, enabling the identification of putative driver polymorphisms regulating expression
97 of key disease-relevant genes and nominating novel gene-cell type associations. Moreover, our
98 integrative framework provides a roadmap for application of this data and technology to any
99 neurological disorder, thus enabling a more comprehensive understanding of the role of inherited
100 noncoding variation in disease.

101

102 RESULTS

103

104 Chromatin accessibility landscapes identify brain regional epigenomic heterogeneity

105 We profiled the chromatin accessibility landscapes of 7 brain regions across 39 cognitively healthy
106 individuals to deeply characterize the role of the noncoding genome in neurodegenerative diseases
107 (Supplementary Table 1). These brain regions include distinct isocortical regions [superior and
108 middle temporal gyri (SMTG, Brodmann areas 21 and 22), parietal lobe (PARL, Brodmann area
109 39), and middle frontal gyrus (MDFG, Brodmann area 9)], striatum at the level of the anterior
110 commissure [caudate nucleus (CAUD) and putamen (PTMN)], hippocampus (HIPP) at the level
111 of the lateral geniculate nucleus, and the substantia nigra (SUNI) at the level of the red nucleus
112 (Figure 1a). These regions were chosen to represent the diversity of brain functionality and cell
113 type composition, and to be the most relevant to prevalent neurodegenerative diseases. In total, we
114 generated 268 ATAC-seq libraries from 140 macrodissected brain samples, with technical
115 replicates for 128 of the 140 samples. From these 268 ATAC-seq libraries, we compiled a merged
116 set of 186,559 peaks reproducible across at least 30% of samples within a given brain region
117 (Figure 1b and Supplementary Table 2; see Methods). Dimensionality reduction via t-distributed
118 stochastic neighbor embedding (t-SNE) identified 4 distinct clusters of samples, grouped roughly
119 by the major brain region (isocortex, striatum, hippocampus, and substantia nigra; Figure 1c).
120 Similar groupings were observed in principal component analysis with nearly 40% of the variance

121 explaining the difference between striatal and non-striatal brain regions (Supplementary Fig 1a-
122 b). These samples showed no clustering based on covariates such as biological sex, post-mortem
123 interval, or *APOE* genotype (Supplementary Fig 1c-d and Supplementary Table 1). Originally, the
124 samples in this cohort were selected from two clinically similar but pathologically distinct research
125 participants: (i) cognitively normal individuals with no or low neuropathological features of AD,
126 or (ii) cognitively normal individuals with intermediate or high burden of neuropathological
127 features of AD^{15,16}. Comparison of these clinico-pathologically normal and clinically resilient
128 donor subgroups showed no statistically significant differences in bulk chromatin accessibility in
129 any of the brain regions profiled (Supplementary Fig. 1e). The variability across these donor
130 subgroups was minimal in comparison to the differences in chromatin accessibility observed
131 across different brain regions (Supplementary Fig. 1f). For this reason, these donor subgroups were
132 treated as a single group in the remainder of analyses.

133 Assessment of regional variation in chromatin accessibility through “feature binarization”
134 (see Methods) identified 28,077 peaks showing region-specific or multi-region-specific
135 accessibility (Figure 1d). For example, 14,628 and 1,734 peaks were identified with significantly
136 increased chromatin accessibility only in striatum or substantia nigra, respectively (Figure 1d).
137 These peak sets showed enrichment for key brain-related transcription factors (TFs) in the FOX,
138 NEUROD, and OLIG families, consistent with suspected brain-relevant enhancers and promoters
139 (Figure 1d). Moreover, some peaks within these sets were in the vicinity of key cell lineage-
140 defining genes such as the dopamine receptor D2 (*DRD2*) in striatal regions, iroquois homeobox
141 3 (*IRX3*) in the substantia nigra, and potassium voltage-gated channel modifier subfamily S
142 member 1 (*KCNS1*) in the isocortical regions (Figure 1e). Notably, while the hippocampus shares
143 many peaks with other regions, we identified only 29 peaks that showed significantly increased
144 chromatin accessibility specifically in this region. Taken together, these results indicate an
145 extensive degree of brain regional heterogeneity that is likely representative of the functional and
146 cellular diversity of the brain regions studied here.
147

148 **ATAC-seq refines interpretation of inherited risk variants in neurodegeneration**

149 Using this atlas of regional chromatin accessibility, we sought to identify functional noncoding
150 regulatory elements that may be impacted by disease-associated genetic variation identified
151 through genome-wide association studies. Approximately 90% of phenotype-associated GWAS
152 polymorphisms reside in noncoding DNA¹⁷, making it difficult to predict a putative functional
153 impact. Moreover, linkage disequilibrium (LD) makes it difficult to pinpoint a single causative
154 SNP when many other nearby SNPs are co-inherited. To resolve these complexities, we used a
155 multi-tiered approach to predict which GWAS SNPs may be functional. First, we identified a
156 compendium of SNPs that could be associated with either AD or PD (Supplementary Table 3, see
157 Methods). To do this, we identified (i) any SNPs passing genome-wide significance in recent
158 GWAS^{1-3,5-7}, (ii) any SNPs exhibiting colocalization of GWAS and eQTL signal, and (iii) any
159 SNPs in linkage disequilibrium with a SNP in the previous two categories. In total, this identified
160 9,741 SNPs including 3,245 unique SNPs across 44 loci associated with AD and 6,496 unique

161 SNPs across 86 loci associated with PD, with a single locus containing 34 SNPs appearing in both
162 diseases. We then performed LD score regression to identify brain regional enrichment of
163 neurodegeneration-related SNPs in noncoding regulatory regions. However, these regional
164 analyses showed minimal enrichment of GWAS SNPs in peak regions associated with any of the
165 brain regions profiled (Supplementary Fig. 2a-b). These results provide evidence against a possible
166 regional effect involving most cell types in a particular area of the brain, but leave open the
167 possibility of involvement of specific cell types in specific regions of the brain. Thus, we
168 hypothesized that a single-cell-based approach could provide more granularity in identifying the
169 precise cell types mediating disease-relevant genetic associations.
170

171 **Single-cell ATAC-seq captures regional and cell type-specific heterogeneity**

172 To test this hypothesis and to better understand brain-regional cell type-specific chromatin
173 accessibility landscapes, we performed single-cell chromatin accessibility profiling in 10 samples
174 spanning the isocortex (N=3), striatum (N=3), hippocampus (N=2), and substantia nigra (N=2)
175 (Supplementary Table 1). In total, we profiled chromatin accessibility in 70,631 individual cells
176 (Figure 2a) after stringent quality control filtration (Supplementary Fig. 2c and Supplementary
177 Table 4). Unbiased iterative clustering^{12,18} of these single cells identified 24 distinct clusters
178 (Figure 2a) which were assigned to known brain cell types based on gene activity scores (see
179 Methods) compiled from chromatin accessibility signal in the vicinity of key lineage-defining
180 genes^{18,19} (Figure 2b and Supplementary Fig. 2c). For example, chromatin accessibility at the
181 myelin associated glycoprotein (*MAG*) gene locus defined clusters corresponding to
182 oligodendrocytes while genes such as vesicular glutamate transporter 1 (*VGLUT1 / SLC17A7*) and
183 vesicular GABA transporter (*VGAT / SLC32A1*) defined excitatory and inhibitory neurons,
184 respectively (Figure 2b). Additionally, 13 of the 24 clusters showed regional specificity with some
185 clusters being made up almost entirely from a single brain region (Figure 2c and Supplementary
186 Table 4). This is most obvious for neuron, astrocyte, and oligodendrocyte precursor cell (OPC)
187 clusters which show clear region-specific differences in clustering (Supplementary Fig. 3a-b).
188 From this cluster-based perspective, we did not identify any clusters that were clearly segregated
189 by gender but the sample size used in this study was not powered to make such a determination
190 (Supplementary Fig. 3c). Cumulatively, we defined 8 distinct cell groupings and identified one
191 cluster (Cluster 18) as putative doublets that we excluded from downstream analyses (Figure 2a
192 and Supplementary Fig. 3d). These cell groupings varied largely in the total number of cells per
193 grouping (Supplementary Fig. 3e) and showed distinct donor and regional compositions
194 (Supplementary Fig. 3f-i).

195 Using these robustly defined clusters, we then called peaks of pseudo-bulk chromatin
196 accessibility to create a union set of 359,022 reproducible peaks (Supplementary Table 5). Overall,
197 89% of the bulk ATAC-seq peaks were overlapped by a peak called in the scATAC-seq data
198 (Figure 2d). Conversely, only 34% of the scATAC-seq peaks were overlapped by a peak from the
199 bulk ATAC-seq peak set (Figure 2d). This is consistent with the known difficulty in identifying
200 peaks in bulk data derived from cell types that comprise less than 20% of the total cells in the

201 tissue²⁰. These results highlight the utility of single-cell methods in situations where cell type-
202 specific peaks are difficult to identify from bulk tissues containing multiple distinct cell types at
203 varying frequencies.

204 This single-cell ATAC-seq-derived peak set enabled the identification of 221,062 highly
205 cell type-specific peaks (Figure 2e). These peaks, comprising more than 60% of all peaks identified
206 in our single-cell data, were selected to be specific to a single cell type or specifically shared across
207 up to three cell types using “feature binarization” (see Methods). For example, some peaks are
208 shared across the 3 different neuronal groups (excitatory, inhibitory, nigral) while others are shared
209 across astrocytes, OPCs, and oligodendrocytes (Figure 2e, Supplementary Table 6). However, the
210 majority of cell type-specific peaks are uniquely accessible in a single cell type; for example,
211 microglia show 45,196 peaks that are specifically accessible in microglia and not in any of the
212 other cell types profiled (Figure 2e). In total, more than 47% of the peaks called in our single-cell
213 ATAC-seq data are specific to a single cell type (Supplementary Table 6) with the vast majority
214 of these cell type-specific peaks remaining undetected in our bulk ATAC-seq analyses. To predict
215 which TFs may be responsible for establishing and maintaining these cell type-specific regulatory
216 programs, we performed motif enrichment analyses of peaks specific to each cell type (Figure 2f).
217 We identified many known drivers of cell type identity, such as motifs specific to SOX9 and
218 SOX10 in oligodendrocytes^{21,22}, or to ASCL1 in OPCs^{23,24}. Lastly, TF footprinting from our
219 scATAC-seq-derived cell type-specific chromatin accessibility data showed enrichment of binding
220 of key lineage defining TFs SPI1 and JUND in microglia and neurons, respectively (Figure 2g).
221 Overall, these results provide a reference map of chromatin accessibility in the adult brain at single-
222 cell resolution.

223

224 **Single-cell ATAC-seq provides reference cell populations for deconvolution of cell type- 225 specific signals in bulk data**

226 Using the cell type-specific signals present in our scATAC-seq data (Supplementary Fig. 4a), we
227 performed cell type deconvolution of our bulk ATAC-seq data using CIBERSORT²⁵
228 (Supplementary Table 7). Using our 8 cell type classification, we deconvolved the ATAC-seq
229 signal from all 140 samples profiled by bulk ATAC-seq in this study, finding clear and expected
230 patterns of cell type abundance such as a relative absence of excitatory neurons in the striatum
231 (Supplementary Fig. 4b). Similarly, deconvolution based on clusters shows expected patterns
232 including the mapping of signal from Cluster 14 (nigral astrocytes) specifically to samples from
233 the substantia nigra, and mapping of signal from Cluster 2 (striatal inhibitory neurons) specifically
234 to samples from the striatum (Supplementary Fig. 4c). By comparing the CIBERSORT prediction
235 to the observed “ground truth” in the scATAC-seq data for the 10 samples profiled here, we were
236 able to assess the performance of the cell type-specific and cluster-specific classifiers
237 (Supplementary Fig. 4d-e). As would be expected, the cell type-specific classifier showed better
238 performance than the cluster-specific classifier, largely due to over- or under-prediction of closely
239 related clusters, such as the oligodendrocytic Clusters 19-23, by the cluster-specific classifier
240 (Supplementary Fig. 4e). Application of the cell type-specific and cluster-specific classifiers to

241 each individual bulk ATAC-seq sample profiled above showed a striking degree of variability in
242 the bulk data based on predicted cell type abundance (Supplementary Fig. 4f-g). Such large
243 differences in cell type composition can hamper efforts to find differential features, further
244 supporting the use of single-cell approaches to understand complex tissues and disease states
245 where small disease-specific variation may be overshadowed by larger differences in cell type
246 composition across samples.
247

248 **Single-cell ATAC-seq identifies brain region-specific differences in glial cells**

249 Our dissection of the cell type-specific chromatin landscapes in adult brain identified clusters that
250 are both region- and cell type-specific such as Cluster 14 which is comprised almost exclusively
251 of astrocytes from the substantia nigra (Figure 2c and Supplementary Table 4). This observation
252 indicates that certain brain cell types may show region-specific variation. This phenomenon has
253 been very well described in neurons, with, for example, inhibitory neurons from the striatum
254 (largely medium spiny neurons) differing substantially from inhibitory neurons outside of the
255 striatum²⁶. Murine oligodendrocytes²⁷ and astrocytes²⁸ also show regional differences in
256 morphology, function, and gene expression. However, the brain-regional variation of glial cells in
257 humans remains less well understood. To address this, we grouped cells into one of the 8 broad
258 cell types defined above and created pseudo-bulk reference populations from the cumulative data
259 (see Methods). Using these region-cell type combinations, we calculated Pearson correlations for
260 all regions across a single cell type (Supplementary Fig. 5a). As expected, neuronal cell types
261 showed the most regional variation.

262 Glial cells, however, also showed substantial regional variation, with astrocytes showing
263 the most variation followed by OPCs (Supplementary Fig. 5a). Within astrocytes, the greatest
264 difference was found between the substantia nigra and the isocortex, indicating that the function
265 or composition of astrocytes may differ across these brain regions. Differential peak analysis
266 identified significant differences in chromatin accessibility near transcriptional regulators that may
267 help explain the observed regional astrocytic differences (Supplementary Fig. 5b and
268 Supplementary Table 8). In particular, nigral astrocytes showed significantly increased
269 accessibility at the forkhead box B1 (*FOXB1*), *IRX1*, *IRX2*, *IRX3*, and *IRX5* genes. Conversely,
270 isocortical astrocytes showed significantly increased accessibility at the *FOGX1*, zic family
271 member 2 (*ZIC2*), and *ZIC5* genes. These changes in chromatin accessibility would be expected
272 to correlate with similar changes in gene expression for the annotated genes. Moreover, the gene
273 activity scores of these genes are definitional for the region-cell subtypes with, for example,
274 *FOXB1* being active only in nigral astrocytes and *ZIC2* and *ZIC5* being active in all other astrocytes
275 (Supplementary Fig. 5c-d). Of particular interest, the observed FOX switch from *FOGX1* in
276 isocortical (and hippocampal/striatal) astrocytes to *FOXB1* in nigral astrocytes and the significant
277 changes in chromatin accessibility at the IRX genes represent a potential transcriptional lineage
278 control mechanism that could help to better understand region-specific functional differences in
279 these astrocytes. Notably, diencephalic brain regions such as the substantia nigra have previously
280 been shown to express *FOXB1*²⁹, *IRX1*³⁰, and *IRX3*³¹ during early brain development, thus

281 explaining part of this broad TF-based lineage control. These transcriptional regulators could be
282 exploited to drive differentiation programs to, for example, create regionally biased glial cells in
283 vitro.

284 In addition to controlling regional astrocytic identity, chromatin accessibility at *IRX* genes
285 was also found to differentiate nigral OPCs from isocortical OPCs (Supplementary Fig. 5d-e).
286 Similarly, *FOXG1* also showed significantly more accessibility in isocortical OPCs, echoing the
287 observations from astrocytes. Lastly, chromatin accessibility at the *PAX3* gene locus was
288 significantly higher in nigral OPCs compared to isocortical OPCs (Supplementary Fig. 5d-e).
289 Taken together, these results identify shared and disparate transcriptional regulatory programs that
290 likely control regional differences amongst astrocytes and OPCs in the substantia nigra and
291 isocortex.

292 Compared to astrocytes, oligodendrocytes and microglia showed less regional variation in
293 chromatin accessibility (Supplementary Fig. 5f-g). While a small number of genes showed highly
294 significant regional differences in oligodendrocytes (Supplementary Fig. 5h), very few genes
295 showed appreciable regional differences among microglia. As noted previously, the regional
296 differences observed in glial cells are a small fraction of the size and magnitude of regional
297 differences observed in neurons (Supplementary Fig. 5i-j), further emphasizing the importance of
298 single-cell approaches to study complex tissues.
299

300 **Single-cell ATAC-seq pinpoints the cellular targets of GWAS polymorphisms**

301 Having generated high-quality cell type-specific chromatin accessibility profiles using scATAC-
302 seq, we sought to refine our previous interpretation of GWAS polymorphisms. More specifically,
303 we aimed to use these data to predict which cell type(s) may be the functional targets of various
304 polymorphisms. When using peaks called in bulk ATAC-seq, we found that 78 LD-expanded
305 SNPs in AD and 186 LD-expanded SNPs in PD overlapped peak regions. Combining our bulk
306 ATAC-seq and scATAC-seq peak sets, we found that 438 SNPs in AD and 880 SNPs in PD
307 directly overlapped peak regions. This represents a 5-fold increase in the number of SNPs observed
308 to overlap peaks called from bulk ATAC-seq alone (Supplementary Table 3), illustrating the
309 importance of cell type-specific interrogation of noncoding regions to dissect GWAS
310 polymorphisms. Cell type-specific LD score regression using AD and PD GWAS results revealed
311 a significant increase in per-SNP heritability for AD in the microglia peak set, reinforcing previous
312 studies^{2,32,33} (Figure 3a and Supplementary Table 9). Similar analyses in PD showed no significant
313 enrichment in SNP heritability in any particular cell type, perhaps indicating that the cellular bases
314 of PD are more heterogeneous than AD (Figure 3a). Though not a focus of the current study, we
315 note that the data generated here can be used to inform the cellular ontogeny of any brain-related
316 GWAS. For example, we observe a striking enrichment of SNP heritability for schizophrenia,
317 neuroticism, and attention deficit hyperactivity disorder in excitatory and inhibitory neurons
318 (Figure 3a). We also confirmed that the heritability of GWAS SNPs from traits not directly related
319 to brain cell types, such as lean body mass, were not enriched in any of the tested brain cell types
320 and that cell types not expected to be involved in brain-related diseases show no enrichment of

321 SNP heritability for brain-related disease SNPs (Supplementary Fig. 6a). Thus, combination of our
322 scATAC-seq data with our curated list of disease-relevant SNPs enables prediction of the cellular
323 targets of each polymorphism.
324

325 **Three-dimensional chromatin landscapes nominate novel target genes of inherited risk 326 variants**

327 In addition to understanding the cell type-specific impacts of an individual polymorphism, we also
328 wanted to predict the gene(s) that may be the direct regulatory targets of a given noncoding
329 polymorphism. We reasoned that the vast majority of functional GWAS SNPs would reside in
330 noncoding sequences and therefore exert their effects through modulation of enhancer or promoter
331 activity. As such, we mapped the enhancer-centric three-dimensional (3D) chromatin architecture
332 in multiple brain regions using HiChIP for histone H3 lysine 27 acetylation (H3K27ac) which
333 marks active enhancers and promoters (Figure 3b and Supplementary Fig. 6b). In total, we
334 generated 3D interaction maps for 6 of the 7 regions profiled by ATAC-seq (putamen was excluded
335 given the high overlap with the caudate nucleus) with an average of 158 million valid interaction
336 pairs identified per region (Supplementary Fig. 6c). These maps led to the identification of 833,975
337 predicted 3D interactions across all brain regions profiled of which 331,730 (40%) were
338 reproducible in at least two brain regions (Supplementary Fig. 6d and Supplementary Table 10).
339 Of these loops, 29.2% had an ATAC-seq peak present in one anchor, 67.4% had an ATAC-seq
340 peak present in both anchors, and 3.4% did not overlap any ATAC-seq peaks identified in either
341 the bulk or scATAC-seq datasets (Supplementary Fig. 6e). Additionally, correlated variation of
342 chromatin accessibility in peaks across single cells has been shown to predict functional
343 interactions between regulatory elements^{19,34}. Using this co-accessibility framework, we predicted
344 regulatory interactions from our scATAC-seq data (Supplementary Fig. 6f), identifying 2,822,924
345 putative interactions between regions of chromatin accessibility (Supplementary Table 10). This
346 set of interactions showed only moderate overlap (~20%) with our HiChIP data, consistent with
347 the ability of this technique to identify cell type-specific regulatory interactions, whereas HiChIP
348 of bulk brain tissue is better suited for identification of more shared regulatory interactions
349 (Supplementary Fig. 6f). Together, these two techniques define a compendium of putative
350 regulatory interactions in the various brain regions studied here.

351 To predict which genes may be altered by noncoding GWAS polymorphisms, we first
352 classified GWAS loci according to whether their phenotypic association was likely mediated by
353 alterations in the coding or noncoding genome (Figure 3c). Across AD and PD, this identified 17
354 loci that harbored likely functional coding alterations, 68 loci that harbored likely functional
355 noncoding alterations, 9 loci that could be associated with putatively functional coding and
356 noncoding alterations, and 22 loci that did not harbor any SNPs in coding regions nor any SNPs
357 in regulatory regions identified in our chromatin accessibility data (Supplementary Table 3). These
358 “unknown” loci likely represent noncoding associations in cell types that were not adequately
359 represented in our analysis. From the original set of 9,741 disease-related SNPs, we identified 438
360 SNPs for AD and 880 SNPs for PD that overlapped peak regions of chromatin accessibility. Of

361 these SNPs, 395 and 531 were involved in a putative enhancer-promoter interaction identified in
362 our HiChIP or co-accessibility data for AD and PD, respectively (Supplementary Table 3).
363 Cumulatively, this enabled the identification of 433 and 516 genes putatively affected by the
364 activity of GWAS polymorphisms in AD and PD, respectively (Figure 3d-e). These gene sets are
365 enriched for biological processes known to be implicated in AD and PD including lipoprotein
366 particle clearance¹ (AD) and synaptic vesicle recycling³⁵ (PD) (Supplementary Fig. 6g-h).

367

368 **Machine learning predicts putative functional SNPs and identifies the molecular ontogeny
369 of disease associations**

370 To disentangle further the molecular underpinnings of AD and PD associations, we developed a
371 multi-omic approach to predict functional noncoding GWAS polymorphisms (Figure 4a and
372 Supplementary Fig. 7a). This approach is anchored in the use of a machine learning framework to
373 score the allelic effect of a SNP on chromatin accessibility. Using the gapped k -mer support vector
374 machine (gkm-SVM) framework³⁶, we trained models to learn the patterns and grammars of
375 chromatin accessibility using our scATAC-seq data (Figure 4b). Specifically, for each cluster (cell
376 type) identified from the scATAC-seq data, we provided 1000-bp sequences centered at all of the
377 peak regions from the cluster-specific pseudo-bulk ATAC-seq data and an equal number of GC-
378 matched non-accessible genomic sequences to a gkm-SVM classifier and trained it to predict
379 whether each sequence is accessible or not. The gkm-SVM models for all 24 scATAC-seq clusters
380 exhibited high prediction performance on held-out test sequences (Supplementary Fig. 7b-c),
381 across all folds of a 10-fold validation training paradigm (Supplementary Fig. 7d).

382 Next, we used three complementary approaches, GkmExplain³⁷, *in silico* mutagenesis³⁸,
383 and deltaSVM³⁹ to predict the allelic impact of 1677 candidate SNPs on chromatin accessibility in
384 each cluster by providing the sequences corresponding to both alleles of each SN to the models for
385 each of the 24 clusters. All three approaches showed high concordance of predicted allelic effects
386 across all candidate SNPs (Supplementary Fig. 7e). In total, among the 1677 SNPs that we scored,
387 we identified 44 high-confidence, and 41 moderate-confidence SNPs that the model predicts will
388 have a functional consequence on chromatin accessibility via identifiable TF binding sites.
389 Integration of these predictions with our colocalization, HiChIP, and scATAC-seq data sets
390 allowed for a comprehensive interrogation of the epigenetic effects of noncoding polymorphisms
391 in AD and PD (Figure 4a and Supplementary Table 3).

392 This multi-omic approach identifies two main categories of novel associations: established
393 disease-related genes where the precise causative SNP remains unknown, and novel genes
394 previously not implicated in disease pathogenesis. In each of these categories, our integrative
395 analysis implicates SNP-gene associations that are supported by (i) the presence of the SNP in an
396 ATAC-seq peak (Tier 3), (ii) a colocalization, HiChIP interaction, or co-accessibility correlation
397 linking the SNP to one or more genes (Tier 2), and in many cases (iii) orthogonal prediction of
398 SNP function via either allelic imbalance (Supplementary Fig. 7f), machine learning predictions,
399 or both (Tier 1) (Supplementary Fig. 7a). Allelic imbalance refers to the differential accessibility
400 between two alleles when one allele is more readily bound than the other. This is obtained from

401 our bulk ATAC-seq data which is available for all donors, thus highlighting the utility of a
402 combined bulk and single-cell approach. Moreover, the cell type-specificity of our scATAC-seq
403 data allows identification of the cell types in which these disease associations likely form.

404 Many studies have investigated the role of genes such as Phosphatidylinositol Binding
405 Clathrin Assembly Protein (*PICALM*)⁴⁰, Solute Carrier Family 24 Member 4 (*SLC24A4*)⁴¹,
406 Bridging Integrator 1 (*BIN1*)^{10,42}, and Membrane Spanning 4-Domains A6A (*MS4A6A*)⁴³ in AD
407 since their implication in the disease by GWAS. However, it remains unclear which
408 polymorphisms drive these associations. In the case of *PICALM*, our models predict a potential
409 functional variant (rs1237999) which resides within an oligodendrocyte-specific regulatory
410 element 35-kb upstream of *PICALM* and disrupts a putative FOS/AP1 factor binding site (Figure
411 4c-d). Moreover, rs1237999 shows striking allelic imbalance with the variant (effect) allele
412 showing diminished accessibility in bulk ATAC-seq data from heterozygotes across multiple brain
413 regions (Figure 4e). Lastly, rs1237999 shows 3D interaction with both *PICALM* and the *EED* gene,
414 a polycomb-group family member involved in maintaining a repressive transcriptional state. This
415 expands the potential functional role of this association to a novel gene and specifically points to
416 a role for oligodendrocytes which were not previously implicated in this phenotypic association⁴⁰.

417 Similarly, the *SLC24A4* locus harbors a small LD block with 46 SNPs that all reside within
418 an intron of *SLC24A4*. Previous work has implicated both *SLC24A4* and the nearby Ras And Rab
419 Interactor 3 (*RIN3*) gene in this association but the true mediator remains unclear^{44,45}. Our multi-
420 omic approach identifies a single SNP, rs10130373, which occurs within a microglia-specific peak,
421 disrupts an SPI1 motif, and communicates specifically with the promoter of the *RIN3* gene (Figure
422 4f-g). This is consistent with the role of *RIN3* in the early endocytic pathway which is crucial for
423 microglial function and of particular disease relevance in AD⁴⁶.

424 In the case of *BIN1*, our work and previous work¹⁰ predict SNP rs6733839 to disrupt a
425 MEF2 binding site in a microglia-specific enhancer located 28-kb upstream of the *BIN1* promoter
426 (Supplementary Fig. 8a). Our machine learning framework additionally implicates SNP
427 rs13025717 which we predict to disrupt a KLF4 binding motif in a microglia-specific putative
428 enhancer 21-kb upstream of *BIN1* (Supplementary Fig. 8b). Both of these SNPs have previously
429 been shown to have sequence-specific correlations with *BIN1* gene expression⁴⁷. Similarly, we
430 identified rs636317 in the *MS4A6A* locus which disrupts a microglia-specific CTCF binding motif
431 (Supplementary Fig. 8c-d). Cumulatively, these results annotate the most likely functional SNPs
432 mediating known disease associations in AD and PD (Supplementary Table 3). Importantly, these
433 predicted functional SNPs do not always affect the expected cell type nor target the closest gene,
434 further emphasizing the utility of our integrative multi-omic approach.

435 Nevertheless, the true promise in studying these noncoding polymorphisms is the
436 identification of novel genes affected by disease-associated variation. This is perhaps most
437 important in PD where identification of disease-associated genes is less mature. The *ITIH1* GWAS
438 locus occurs within a 600-kb LD block harboring 317 SNPs and no plausible gene association has
439 been made to date. We nominate rs181391313, a SNP occurring within a putative microglia-
440 specific intronic enhancer of the Stabilin 1 (*STAB1*) gene (Figure 5a). *STAB1* is a large

441 transmembrane receptor protein that functions in lymphocyte homing and endocytosis of ligands
442 such as low density lipoprotein, two functions that would be consistent with a role for microglia
443 in PD⁴⁸. This SNP is predicted to disrupt a KLF4 binding site, consistent with the role of KLF4 in
444 regulation of microglial gene expression⁴⁹ (Figure 5b). Similarly, the *KCNIP3* GWAS locus
445 resides in a 300-kb LD block harboring 94 SNPs. Our results identify two putative mediators of
446 this phenotypic association which lead to very different functional interpretations (Figure 5c).
447 First, rs7585473 occurs more than 250 kb upstream of the lead SNP and disrupts an
448 oligodendrocyte-specific SOX6 motif in a peak found to interact with the Myelin and Lymphocyte
449 (*MAL*) gene, a gene implicated in myelin biogenesis and function (Figure 5d). Alternatively, we
450 find rs3755519 in a neuronal-specific intronic peak within the *KCNIP3* gene with clear interaction
451 with the *KCNIP3* gene promoter. While this SNP does not show a robust machine learning
452 prediction, nor reside within a known motif, we do identify allelic imbalance supporting its
453 predicted functional alteration of transcription factor binding (Figure 5e). Together, these SNPs
454 provide competing interpretations of this locus, implicating oligodendrocyte- and neuron-specific
455 functions, and demonstrating the complexities of noncoding SNP interpretation.

456 Though many such anecdotes exist (Supplementary Table 3), we also noted a pattern
457 whereby many SNPs appear to disrupt binding sites related to the CCCTC-Binding Factor (*CTCF*)
458 protein. For example, SNP rs6781790 disrupts a predicted CTCFL binding site within the promoter
459 of the WD Repeat Domain 6 (*WDR6*) gene (Supplementary Fig. 9a-b). This SNP shows clear
460 allelic imbalance across a large number of bulk ATAC-seq samples (Supplementary Fig. 9c).
461 Similarly, SNP rs7599054 disrupts a putative CTCF binding site near the Transmembrane Protein
462 163 (*TMEM163*) gene (Supplementary Fig. 9d-e).

463 Taken together, this vertical integration of multi-omic data provides an unprecedented
464 resolution of the landscape of inherited noncoding variation in neurodegenerative disease.
465 Moreover, this framework and data can be applied to inform the molecular ontogeny of any brain-
466 related GWAS polymorphism, extending the applicability of this work to all neurological disease.
467

468 **Epigenomic dissection of the *MAPT* locus explains haplotype-specific changes in local gene 469 expression**

470 One of the most common PD-associated risk loci is the microtubule associated protein tau (*MAPT*)
471 gene locus. *MAPT* encodes tau proteins, a primarily neuronal set of isoforms whose pathological,
472 hyperphosphorylated aggregates form the neurofibrillary tangles of AD⁵⁰; however, despite the
473 long known genetic association, it remains unclear how the *MAPT* locus may play a role in PD.

474 The *MAPT* locus is present within a large 1.8-Mb LD block and manifests as two distinct
475 haplotypes, H1 and H2, which differ genetically in two primary ways: (i) more than 2000 SNPs
476 differ across the two haplotypes, and (ii) an approximately 1-Mb inversion that includes the *MAPT*
477 gene^{51,52} (Figure 6a). Previous reports have nominated multiple explanations for how these
478 alterations are associated with PD, including increased *MAPT* expression in the H1 haplotype^{53,54}
479 (Figure 6b), different ratios of splice isoforms⁵⁵⁻⁵⁷, and the use of alternative promoters⁵⁸. We
480 created a haplotype-specific map of chromatin accessibility and 3D chromatin interactions at the

481 *MAPT* locus (Figure 6c). Using data from heterozygote H1/H2 individuals, we split reads into H1
482 and H2 haplotypes based on the presence of one of the 2366 haplotype divergent SNP
483 (Supplementary Table 11; see methods). We tiled the region into non-overlapping 500-bp bins (to
484 avoid biases in peak calling) and performed a Wilcoxon rank sum test to identify regions that are
485 differentially accessible both between H1/H1 and H2/H2 homozygotes and between split reads
486 from H1/H2 heterozygotes (Supplementary Fig. 10a-b). This identified 28 bins including an H1-
487 specific putative enhancer 68 kb upstream of the *MAPT* promoter and the promoter of the KAT8
488 regulatory NSL complex subunit 1 (*KANSL1*) gene located 330 kb downstream of *MAPT* (Figure
489 6d (asterisks) and Supplementary Fig. 10c). Using our HiChIP data, we performed haplotype-
490 specific virtual 4C to determine if any of these changes in chromatin accessibility were
491 accompanied by changes in 3D chromatin interaction frequency. We identified H2-specific 3D
492 interactions between a putative domain boundary upstream of *MAPT* (labeled “A”) and the region
493 surrounding the *KANSL1* promoter (labeled “B”) spanning a distance of more than 600 kb inside
494 of the inversion breakpoints (Figure 6d). Additionally, the H1-specific putative enhancer upstream
495 of *MAPT* showed increased interaction with a second putative enhancer intronic to *MAPT* as well
496 as with the *MAPT* promoter (Figure 6d).

497 To better understand how these epigenetic changes impact local transcription, we used
498 RNA-sequencing data from the Genotype-Tissue Expression (GTEx) database to identify genes
499 that show significant haplotype-specific changes. In addition to the previously mentioned
500 haplotype-specific differences in *MAPT* expression (Figure 6b), we also identified significant
501 changes in the expression of genes near the largest changes in chromatin accessibility and 3D
502 interaction (points “A” and “B”; Figure 6e). These genes include a *KANSL1* antisense transcript
503 (*KANSL1-ASI*) and a pseudogene of the mitogen-activated protein kinase 8 interacting protein 1
504 (*MAPK8IP1P2*) (Supplementary Fig. 10d-e). These increases in gene expression could play a
505 functional role in pathologic changes mediated by the different *MAPT* haplotypes or, more likely,
506 could be a non-functional byproduct of the genomic inversion.

507 The above analyses help to understand how the genomic region inside of the *MAPT*
508 inversion breakpoints differs between the H1 and H2 haplotypes; however, the inversion also
509 changes the relative orientation of genes inside the breakpoints to enhancers and promoters outside
510 of the breakpoints. In this way, the inversion could alter the 3D architecture of the locus and thus
511 change which enhancers are able to communicate with the *MAPT* gene. In support of this
512 hypothesis, we find a long-distance putative enhancer located 650 kb upstream of the *MAPT* gene
513 that shows elevated interaction with the *MAPT* promoter specifically in the H1 haplotype (Figure
514 6f). We find support for this interaction both in HiChIP data from H1/H1 or H2/H2 homozygotes
515 and from H1/H2 heterozygotes where the reads have been split based on haplotype divergent SNPs
516 (Figure 6f). Indeed, we find multiple neuron-specific putative enhancers in this upstream region,
517 consistent with the known neuron-specific expression of *MAPT* (Supplementary Fig. 10f), and an
518 increase in overall 3D interaction between this upstream region and the region surrounding *MAPT*
519 inside of the inversion breakpoints (Supplementary Fig. 10g). In total, our epigenomic dissection
520 of the *MAPT* locus provides multiple plausible explanations for the haplotype-specific differences

521 in *MAPT* expression and nominates multiple other genes who may exert haplotype-specific effects
522 that are linked to differing PD phenotypes (Figure 6g).

523

524 DISCUSSION

525

526 Here, we provide a high-resolution epigenetic characterization of the role of inherited noncoding
527 variation in AD and PD. Our integrative multi-omic framework and machine learning classifier
528 predicted dozens of functional SNPs, nominating gene and cellular targets for each noncoding
529 GWAS locus. These predictions both inform well-studied disease-relevant genes, such as *BIN1* in
530 AD, and predict novel gene-disease associations, such as *STAB1* in PD. This greatly expands our
531 understanding of inherited variation in AD and PD and provides a roadmap for the epigenomic
532 dissection of noncoding variation in neurodegenerative and other complex genetic diseases.

533 Our work initially focused on two clinically similar but pathologically distinct groups. All
534 brain donors had been longitudinal participants in research cohorts, extensively evaluated within
535 two years of death, and scored as high performers by neuropsychological testing (average interval
536 between last evaluation and death was 362 days). We have shown previously that this cut off
537 minimizes interval conversion to cognitive impairment or dementia⁵⁹. One subset of these high
538 performers had no or low levels of AD or PD neuropathologic change, and are labeled clinico-
539 pathologic normal controls. Another subset of high performers showed neuropathologic changes
540 of AD sufficient to warrant suspicion of dementia; this not common occurrence has several
541 designations but is usually labeled resilient, meaning resilient to the clinical expression of
542 pathologically determined AD. There is intense interest in what underlies resilience to AD because
543 its mechanisms or adaptations may illuminate means to suppress disease expression and extend
544 healthspan. Interestingly, our bulk ATAC-seq data showed no statistically significant differences
545 in chromatin accessibility in any of the seven brain regions profiled for clinico-pathologic controls
546 vs. resilience to AD. This likely indicates that the differences between these two clinical groups is
547 minor, or potentially encoded in a rare cell type or a brain region not profiled in this work.

548 To inform inherited noncoding variation in neurodegenerative disease, we generated an
549 epigenomic resource that spans the cellular and regional diversity of the adult brain. We used bulk
550 ATAC-seq to profile seven distinct brain regions, identifying regional heterogeneity that is largely
551 based on changes in cell type composition. To mitigate the contribution of cellular diversity to our
552 analysis, we additionally performed scATAC-seq, profiling the chromatin accessibility of 70,631
553 individual cells. Cumulatively, this single-cell data identified 24 different cellular clusters which
554 map to 7 distinct broad cell types (excitatory neurons, inhibitory neurons, nigral neurons,
555 astrocytes, oligodendrocytes, OPCs, and microglia). Together, this resource captures the regional
556 and cellular gene regulatory machinery that governs phenotypic expression of noncoding variation,
557 thus allowing us to identify all polymorphisms that could putatively affect gene expression through
558 overlap with peaks of chromatin accessibility (Tier 3). To further refine these putative functional
559 variants, we identified the subset of polymorphisms that could be mapped to gene targets through
560 3D chromatin interactions or co-accessibility networks (Tier 2). Finally, we employed a machine

561 learning approach to predict the subset of polymorphisms that would be likely to perturb
562 transcription factor binding and validated these predictions with measurements of allelic imbalance
563 (Tier 1). In total we implicate approximately 5 times as many genes in the phenotypic association
564 of AD and PD and nominate functional noncoding variants for dozens of previously orphaned
565 GWAS loci.

566 Through our integrative analysis, we additionally provide a comprehensive epigenetic
567 characterization of the *MAPT* gene locus. The *MAPT* gene encodes tau isoforms, primarily
568 neuronal microtubule binding proteins that, under pathologic conditions, can adopt an abnormal
569 structure and extensive post translational modifications, a process called neurofibrillary
570 degeneration, which is a hallmark of AD and other neurodegenerative diseases, but not PD¹⁵.
571 Enigmatically, *MAPT* is a replicated risk locus for PD despite the absence of neurofibrillary
572 degeneration^{60,61}. The *MAPT* locus, found on chromosome 17, represents one of the largest LD
573 blocks in the human genome (1.8 Mb) and is present in two distinct haplotypes, H1 and H2, the
574 latter formed by an approximately 900 kb inversion of H1 that occurred about 3 million years ago
575 and is present mostly in Europeans⁵¹. Cumulatively, previous work supports *MAPT* haplotype-
576 specific impacts on transcript amount, transcript stability, and alternative splicing in several
577 neurodegenerative disorders^{54,56,57}. We highlight multiple epigenetic avenues through which the
578 *MAPT* gene is differentially regulated in the H1 and H2 haplotypes, thus explaining at least a
579 portion of the molecular underpinnings of the observed *MAPT* GWAS association in PD.

580 We developed a multi-omic framework that provides a robust and comprehensive
581 dissection of inherited variation in neurodegenerative disease. Moreover, the functional
582 predictions made through our machine learning classifier and integrative analytical approach
583 greatly expand our understanding of noncoding contributions to AD and PD. More broadly, this
584 work represents a systematic approach to understand inherited variation in disease and provides
585 an avenue towards the nomination of novel therapeutic targets that previously remained obscured
586 by the complexity of the regulatory machinery of the noncoding genome.
587

588 **DATA AVAILABILITY**

589 All data generated in this work is available through SRA (in progress).

590

591 **ACKNOWLEDGEMENTS**

592 This work was supported by NIH NS062684, AG057707 (to T.M.), HG007735 (to H.Y.C.),
593 HG009431 (to S.B.M./A.K.), and AG059918 (to M.R.C.). Sequencing data for this project was
594 generated on an Illumina HiSeq 4000 supported in part by NIH award S10OD018220. Additional
595 resources at the Stanford Center for Genomics and Personalized Medicine Sequencing Center were
596 supported by NIH S10OD025212. H.Y.C. is an Investigator of the Howard Hughes Medical
597 Institute.
598

599

600

601 **AUTHOR CONTRIBUTIONS**

602 M.R.C., H.Y.C., and T.J.M conceived of and designed the project. M.R.C. and T.J.M. compiled
603 the figures and wrote the manuscript with help and input from all authors. A.S. and M.R.C.
604 performed bulk ATAC-seq data processing and analysis. M.R.C. performed all HiChIP data
605 analysis with help from M.R.M and J.M.G. J.M.G., M.R.C., and A.S. performed all single-cell
606 ATAC-seq data processing and analysis with supervision from W.J.G., A.K., S.B.M. and H.Y.C.
607 M.J.G. performed GWAS locus curation, colocalization analysis, and GTEx analysis and L.F. and
608 B.L. performed all LD score regression analysis with supervision from S.B.M. S.K. and A.S.
609 performed all machine learning analysis with supervision from A.K. B.H.L., S.S., and M.R.C.
610 performed all ATAC-seq, scATAC-seq, and HiChIP data generation with help from S.T.B. and
611 M.R.M. K.S.M. curated the frozen tissue specimens used in this work.

612

613 **COMPETING FINANCIAL INTERESTS**

614 H.Y.C. is a co-founder of Accent Therapeutics, Boundless Bio, and an advisor to 10x Genomics,
615 Arsenal Biosciences, Spring Discovery.

616

617 **REFERENCES (MAIN TEXT)**

618

- 619 1. Kunkle, B. W. *et al.* Genetic meta-analysis of diagnosed Alzheimer's disease identifies
620 new risk loci and implicates A β , tau, immunity and lipid processing. *Nat. Genet.* **2019** *513*
621 **51**, 414 (2019).
- 622 2. Jansen, I. *et al.* Genetic meta-analysis identifies 10 novel loci and functional pathways for
623 Alzheimer's disease risk. *bioRxiv* 258533 (2018). doi:10.1101/258533
- 624 3. Lambert, J.-C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility
625 loci for Alzheimer's disease. *Nat. Genet.* **45**, 1452–1458 (2013).
- 626 4. Beecham, G. W. *et al.* Genome-Wide Association Meta-analysis of Neuropathologic
627 Features of Alzheimer's Disease and Related Dementias. *PLoS Genet.* **10**, (2014).
- 628 5. Pankratz, N. *et al.* Meta-analysis of Parkinson's Disease: Identification of a novel locus,
629 RIT2. *Ann. Neurol.* **71**, 370–384 (2012).
- 630 6. Chang, D. *et al.* A meta-analysis of genome-wide association studies identifies 17 new
631 Parkinson's disease risk loci. *Nat. Genet.* **49**, 1511–1516 (2017).
- 632 7. Nalls, M. A. *et al.* Identification of novel risk loci, causal insights, and heritable risk for
633 Parkinson's disease: a meta-analysis of genome-wide association studies. *Lancet Neurol.*
634 **18**, 1091–1102 (2019).
- 635 8. Escott-Price, V., Myers, A. J., Huentelman, M. & Hardy, J. Polygenic risk score analysis
636 of pathologically confirmed Alzheimer disease. *Ann. Neurol.* **82**, 311–314 (2017).
- 637 9. Gallagher, M. D. & Chen-Plotkin, A. S. The Post-GWAS Era: From Association to
638 Function. *Am. J. Hum. Genet.* **102**, 717–730 (2018).
- 639 10. Nott, A. *et al.* Brain cell type – specific enhancer – promoter interactome maps and
640 disease-risk association. *Science (80-.).* **1139**, 1134–1139 (2019).
- 641 11. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J.
642 Transposition of native chromatin for fast and sensitive epigenomic profiling of open
643 chromatin, DNA-binding proteins and nucleosome position. *Nat. Methods* **10**, 1213–1218
644 (2013).

- 645 12. Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human
646 immune cell development and intratumoral T cell exhaustion. *Nat. Biotechnol.* **37**, 925–
647 936 (2019).
- 648 13. Mumbach, M. R. *et al.* HiChIP: efficient and sensitive analysis of protein-directed genome
649 architecture. *Nat. Methods* **13**, 919–922 (2016).
- 650 14. Mumbach, M. R. *et al.* Enhancer connectome in primary human cells reveals target genes
651 of disease-associated DNA elements. *Nat. Genet.* **49**, 1602–1612 (2017).
- 652 15. Hyman, B. T. *et al.* National Institute on Aging-Alzheimer's Association guidelines for
653 the neuropathologic assessment of Alzheimer's disease. *Alzheimer's Dement.* **8**, 1–13
654 (2012).
- 655 16. Montine, T. J. *et al.* National institute on aging-Alzheimer's association guidelines for the
656 neuropathologic assessment of Alzheimer's disease: A practical approach. *Acta
657 Neuropathol.* **123**, 1–11 (2012).
- 658 17. Edwards, S. L., Beesley, J., French, J. D. & Dunning, M. Beyond GWASs: Illuminating
659 the dark road from association to function. *Am. J. Hum. Genet.* **93**, 779–797 (2013).
- 660 18. Granja, J. M. *et al.* Single-cell multiomic analysis identifies regulatory programs in
661 mixed-phenotype acute leukemia. *Nat. Biotechnol.* **37**, 1458–1465 (2019).
- 662 19. Pliner, H. A. *et al.* Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell
663 Chromatin Accessibility Data. *Mol. Cell* **71**, 858–871.e8 (2018).
- 664 20. McKeown, M. R. *et al.* Superenhancer analysis defines novel epigenomic subtypes of
665 non-APL AML, including an RAR α dependency targetable by SY-1425, a potent and
666 selective RAR α agonist. *Cancer Discov.* **7**, 1136–1153 (2017).
- 667 21. Stolt, C. C. *et al.* The Sox9 transcription factor determines glial fate choice in the
668 developing spinal cord. *Genes Dev.* **17**, 1677–1689 (2003).
- 669 22. Kuhlbrodt, K., Herbarth, B., Sock, E., Hermans-Borgmeyer, I. & Wegner, M. Sox10, a
670 novel transcriptional modulator in glial cells. *J. Neurosci.* **18**, 237–250 (1998).
- 671 23. Kondo, T. & Raff, M. Basic helix-loop-helix proteins and the timing of oligodendrocyte
672 differentiation. *Development* **127**, 2989–2998 (2000).
- 673 24. Nakatani, H. *et al.* Ascl1/Mash1 promotes brain oligodendrogenesis during myelination
674 and remyelination. *J. Neurosci.* **33**, 9752–9768 (2013).
- 675 25. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles.
676 *Nat. Methods* **12**, 1–10 (2015).
- 677 26. Lake, B. B. *et al.* Neuronal subtypes and diversity revealed by single-nucleus RNA
678 sequencing of the human brain. *Science (80-.).* **352**, 1586–1590 (2016).
- 679 27. van Bruggen, D., Aguirre, E. & Castelo-Branco, G. Single-cell transcriptomic analysis of
680 oligodendrocyte lineage cells. *Curr. Opin. Neurobiol.* **47**, 168–175 (2017).
- 681 28. Molofsky, A. V. *et al.* Astrocyte-encoded positional cues maintain sensorimotor circuit
682 integrity. *Nature* **509**, 189–194 (2014).
- 683 29. Zhao, T. *et al.* Genetic mapping of Foxb1-cell lineage shows migration from caudal
684 diencephalon to telencephalon and lateral hypothalamus. *Eur. J. Neurosci.* **28**, 1941–1955
685 (2008).
- 686 30. Bosse, A. *et al.* Identification of the vertebrate Iroquois homeobox gene family with
687 overlapping expression during early development of the nervous system. *Mech. Dev.* **69**,
688 169–181 (1997).
- 689 31. Hirata, T. *et al.* Zinc-finger genes Fez and Fez-like function in the establishment of
690 diencephalon subdivisions. *Development* **133**, 3993–4004 (2006).

- 691 32. Hemonnot, A. L., Hua, J., Ulmann, L. & Hirbec, H. Microglia in Alzheimer disease: Well-known targets and new opportunities. *Front. Cell. Infect. Microbiol.* **9**, 1–20 (2019).
- 692 33. Efthymiou, A. G. & Goate, A. M. Late onset Alzheimer's disease genetics implicates microglial pathways in disease risk. *Mol. Neurodegener.* **12**, 1–12 (2017).
- 693 34. Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* **523**, 486–490 (2015).
- 694 35. Deng, H. X. *et al.* Identification of TMEM230 mutations in familial Parkinson's disease. *Nat. Genet.* **48**, 733–739 (2016).
- 695 36. Ghandi, M. *et al.* GkmSVM: An R package for gapped-kmer SVM. *Bioinformatics* **32**, 2205–2207 (2016).
- 696 37. Shrikumar, A., Prakash, E. & Kundaje, A. GkmExplain: Fast and accurate interpretation of nonlinear gapped k-mer SVMs. *Bioinformatics* **35**, i173–i182 (2019).
- 697 38. Bromberg, Y. & Rost, B. Comprehensive in silico mutagenesis highlights functionally important residues in proteins. *Bioinformatics* **24**, 207–212 (2008).
- 698 39. Lee, D. *et al.* A method to predict the impact of regulatory variants from DNA sequence. *Nat. Genet.* **47**, 955–961 (2015).
- 699 40. Xu, W., Tan, L. & Yu, J. T. The Role of PICALM in Alzheimer's Disease. *Mol. Neurobiol.* **52**, 399–413 (2015).
- 700 41. Stage, E. *et al.* The effect of the top 20 Alzheimer disease risk genes on gray-matter density and FDG PET brain metabolism. *Alzheimer's Dement. Diagnosis, Assess. Dis. Monit.* **5**, 53–66 (2016).
- 701 42. Andrew, R. J. *et al.* Reduction of the expression of the late-onset Alzheimer's disease (AD) risk-factor BIN1 does not affect amyloid pathology in an AD mouse model. *J. Biol. Chem.* **294**, 4477–4487 (2019).
- 702 43. Ma, J., Yu, J. T. & Tan, L. MS4A Cluster in Alzheimer's Disease. *Mol. Neurobiol.* **51**, 1240–1248 (2015).
- 703 44. Rouka, E. *et al.* Differential recognition preferences of the three Src Homology 3 (SH3) domains from the adaptor CD2-associated Protein (CD2AP) and Direct Association with Ras and Rab Interactor 3 (RIN3). *J. Biol. Chem.* **290**, 25275–25292 (2015).
- 704 45. Larsson, M. *et al.* GWAS findings for human iris patterns: Associations with variants in genes that influence normal neuronal pattern development. *Am. J. Hum. Genet.* **89**, 334–343 (2011).
- 705 46. Kajihara, H. *et al.* RIN3: A novel Rab5 GEF interacting with amphiphysin II involved in the early endocytic pathway. *J. Cell Sci.* **116**, 4159–4168 (2003).
- 706 47. Novikova, G. *et al.* Integration of Alzheimer's disease genetics and myeloid genomics reveals novel disease risk mechanisms. *biorxiv* (2019). doi:10.1101/694281
- 707 48. Lecours, C. *et al.* Microglial implication in Parkinson's disease: Loss of beneficial physiological roles or gain of inflammatory functions? *Front. Cell. Neurosci.* **12**, 1–8 (2018).
- 708 49. Kaushik, D. K., Gupta, M., Das, S. & Basu, A. Krüppel-like factor 4, a novel transcription factor regulates microglial activation and subsequent neuroinflammation. *J. Neuroinflammation* **7**, 1–20 (2010).
- 709 50. Schellenberg, G. D. & Montine, T. J. The genetics and neuropathology of Alzheimer's disease. *Acta Neuropathol.* **124**, 305–323 (2012).
- 710 51. Stefansson, H. *et al.* A common inversion under selection in Europeans. *Nat. Genet.* **37**, 129–137 (2005).

- 737 52. Zody, M. C. *et al.* Evolutionary toggling of the MAPT 17q21.31 inversion region. *Nat. Genet.* **40**, 1076–1083 (2008).
- 738 53. Valenca, G. T. *et al.* The Role of MAPT Haplotype H2 and Isoform 1N/4R in
739 Parkinsonism of Older Adults. *PLoS One* (2016).
- 740 54. Allen, M. *et al.* Association of MAPT haplotypes with Alzheimer's disease risk and
741 MAPT brain gene expression levels. *Alzheimer's Res. Ther.* **6**, 1–14 (2014).
- 742 55. Pascale, E. *et al.* Genetic architecture of MAPT gene region in parkinson disease
743 subtypes. *Front. Cell. Neurosci.* **10**, 1–7 (2016).
- 744 56. Beevers, J. E. *et al.* MAPT Genetic Variation and Neuronal Maturity Alter Isoform
745 Expression Affecting Axonal Transport in iPSC-Derived Dopamine Neurons. *Stem Cell
746 Reports* **9**, 587–599 (2017).
- 747 57. Lai, M. C. *et al.* Haplotype-specific MAPT exon 3 expression regulated by common
748 intronic polymorphisms associated with Parkinsonian disorders. *Mol. Neurodegener.* **12**,
749 1–16 (2017).
- 750 58. Huin, V. *et al.* Alternative promoter usage generates novel shorter MAPT mRNA
751 transcripts in Alzheimer's disease and progressive supranuclear palsy brains. *Sci. Rep.* **7**,
752 1–10 (2017).
- 753 59. White, L. R. *et al.* Neuropathologic comorbidity and cognitive impairment in the Nun and
754 Honolulu-Asia Aging Studies. *Neurology* **86**, 1000–1008 (2016).
- 755 60. Simón-Sánchez, J. *et al.* Genome-wide association study reveals genetic risk underlying
756 Parkinson's disease. *Nat. Genet.* **41**, 1308–1312 (2009).
- 757 61. Nalls, M. A. *et al.* Large-scale meta-analysis of genome-wide association data identifies
758 six new risk loci for Parkinson's disease. *Nat. Genet.* **46**, 989–993 (2014).
- 759 62. Pankratz, N. *et al.* Genomewide association study for susceptibility genes contributing to
760 familial Parkinson disease. *Hum. Genet.* **124**, 593–605 (2009).
- 761 63. Corces, M. R. *et al.* An improved ATAC-seq protocol reduces background and enables
762 interrogation of frozen tissues. *Nat. Methods* **14**, 959–962 (2017).
- 763 64. Corces, M. R. *et al.* The chromatin accessibility landscape of primary human cancers.
764 *Science (80-.).* **362**, (2018).
- 765 65. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing
766 genomic features. *Bioinformatics* **26**, 841–842 (2010).
- 767 66. Heinz, S. *et al.* Simple Combinations of Lineage-Determining Transcription Factors Prime
768 cis-Regulatory Elements Required for Macrophage and B Cell Identities. *Mol. Cell* **38**,
769 576–589 (2010).
- 770 67. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-
771 wide association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 772 68. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies
773 disease-relevant tissues and cell types. *Nat. Genet.* **50**, 621–629 (2018).
- 774 69. Li, Z. *et al.* Genome-wide association analysis identifies 30 new susceptibility loci for
775 schizophrenia. *Nat. Genet.* **49**, 1576–1583 (2017).
- 776 70. Duncan, L. *et al.* Significant locus and metabolic genetic correlations revealed in genome-
777 wide association study of anorexia nervosa. *Am. J. Psychiatry* **174**, 850–858 (2017).
- 778 71. Demontis, D. *et al.* Discovery of the first genome-wide significant risk loci for attention
779 deficit/hyperactivity disorder. *Nat. Genet.* **51**, 63–75 (2019).
- 780 72. Otowa, T. *et al.* Meta-analysis of genome-wide association studies of anxiety disorders.
781 *Mol. Psychiatry* **21**, 1391–1399 (2016).
- 782

- 783 73. Okbay, A. *et al.* Genetic variants associated with subjective well-being, depressive
784 symptoms, and neuroticism identified through genome-wide analyses. *Nat. Genet.* **48**,
785 624–633 (2016).
- 786 74. Anney, R. J. L. *et al.* Genetic determinants of common epilepsies: A meta-analysis of
787 genome-wide association studies. *Lancet Neurol.* **13**, 893–903 (2014).
- 788 75. Zillikens, M. C. *et al.* Large meta-analysis of genome-wide association studies identifies
789 five loci for lean body mass. *Nat. Commun.* **8**, (2017).
- 790 76. Kemp, J. P. *et al.* Identification of 153 new loci associated with heel bone mineral density
791 and functional involvement of GPC6 in osteoporosis. *Nat. Genet.* **49**, 1468–1475 (2017).
- 792 77. Howson, J. M. M. *et al.* Fifteen new risk loci for coronary artery disease highlight arterial-
793 wall-specific mechanisms. *Nat. Genet.* **49**, 1113–1119 (2017).
- 794 78. Benner, C. *et al.* FINEMAP: Efficient variable selection using summary data from
795 genome-wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 796 79. Hormozdiari, F. *et al.* Colocalization of GWAS and eQTL Signals Detects Target Genes.
797 *Am. J. Hum. Genet.* **99**, 1245–1260 (2016).
- 798 80. Liu, B. *et al.* Genetic Regulatory Mechanisms of Smooth Muscle Cells Map to Coronary
799 Artery Disease Risk Loci. *Am. J. Hum. Genet.* **103**, 377–388 (2018).
- 800 81. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with
801 Harmony. *Nat. Methods* **16**, (2019).
- 802 82. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–1902.e21
803 (2019).
- 804 83. Cusanovich, D. A. *et al.* The cis-regulatory dynamics of embryonic development at single-
805 cell resolution. *Nature* **555**, 538–542 (2018).
- 806 84. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion
807 for RNA-seq data with DESeq2. *Genome Biol.* 1–21 (2014). doi:10.1186/s13059-014-
808 0550-8
- 809 85. Servant, N. *et al.* HiC-Pro: An optimized and flexible pipeline for Hi-C data processing.
810 *Genome Biol.* **16**, 1–11 (2015).
- 811 86. Bhattacharyya, S., Chandra, V., Vijayanand, P. & Ay, F. Identification of significant
812 chromatin contacts from HiChIP data by FitHiChIP. *Nat. Commun.* **10**, (2019).
- 813 87. Lee, D. LS-GKM: A new gkm-SVM for large-scale datasets. *Bioinformatics* **32**, 2196–
814 2198 (2016).
- 815 88. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**,
816 2825–2830 (2011).
- 817 89. Machiela, M. J. & Chanock, S. J. LDlink: A web-based application for exploring
818 population-specific haplotype structure and linking correlated alleles of possible
819 functional variants. *Bioinformatics* **31**, 3555–3557 (2015).
- 820 90. Krueger, F. & Andrews, S. R. SNPsplits: Allele-specific splitting of alignments between
821 genomes with known SNP genotypes [version 2; referees: 3 approved]. *F1000Research* **5**,
822 1–16 (2016).

823 **FIGURE LEGENDS**

824 **Figure 1 - ATAC-seq defines brain-regional epigenetic heterogeneity**

- 825
826 A. Schematic of the brain regions profiled in this study. Indicated colors are used
827 throughout.

- 829 B. Bar plot showing the number of reproducible peaks identified from samples in each brain
830 region. The “Merged” bar represents the final merged peak set used for all bulk ATAC-
831 seq analyses. Colors represent the type of genomic region overlapped by a given peak.
832 The numbers above each bar represent the total number of biological samples profiled for
833 each brain region.
- 834 C. t-SNE dimensionality reduction showing all samples profiled in this study, colored by the
835 region of the brain from which the data was generated. Each dot represents a single piece
836 of tissue with technical replicates merged where applicable.
- 837 D. Heatmap representation of binarized peaks from ATAC-seq data. Each row represents an
838 individual peak and each column represents an individual sample. Feature groups
839 containing more than 1000 peaks are randomly subsetted down to 1000 peaks for display
840 on the heatmap. Feature groups containing fewer than 50 peaks are not displayed.
841 Heatmap color represents the row-wise Z-score of normalized chromatin accessibility at
842 the peak region. Motif names and logos shown to the right of the plot represent motifs
843 enriched in the various peak sets.
- 844 E. Sequencing tracks of region-specific ATAC-seq peaks identified through feature
845 binarization. From left to right, *DRD2* (striatum-specific; chr11:113367951-113538919),
846 *IRX3* (substantia nigra-specific; chr16:54276577-54291319), and *KCNS1* (isocortex-
847 specific; chr20:45086706-45107665). Track heights are the same in each vertical panel.
- 848

849 **Figure 2 - Single-cell ATAC-seq identifies cell type-specific chromatin accessibility in the**
850 **adult brain**

- 851 A. Left; UMAP dimensionality reduction showing identified clusters of cells. Each dot
852 represents a single cell (N = 70,631). Right; Bar plot showing the number of cells per
853 cluster. Each cluster is labeled to the right of the bar plot and the predicted cell type
854 corresponding to each cluster is shown colorimetrically.
- 855 B. The same UMAP dimensionality reduction shown in Figure 2a but each cell is colored by
856 its gene activity score for the annotated lineage-defining gene. Grey represents a gene
857 activity score of 0 while purple represents the maximum gene activity score for the given
858 gene.
- 859 C. Cluster residence heatmap showing the percent of each cluster that is composed of cells
860 from each sample. Cell numbers were normalized across samples prior to calculating
861 cluster residence percentages.
- 862 D. Bar plot showing the overlap of bulk ATAC-seq and scATAC-seq peak calls. “Bulk”
863 represents the number of peaks from the bulk ATAC-seq merged peak set that are
864 overlapped by a peak called in our scATAC-seq merged peak set. “Single-cell”
865 represents the number of peaks from our scATAC-seq merged peak set that are
866 overlapped by a peak called in our bulk ATAC-seq merged peak set.
- 867 E. Heatmap representation of binarized peaks from scATAC-seq data. Each row represents
868 an individual pseudo-bulk replicate (3 per cell type) and each column represents an

- 869 individual peak. Feature groups containing fewer than 1000 peaks are not displayed.
870 Heatmap color represents the column-wise Z-score of normalized chromatin accessibility
871 at the peak region.
- 872 F. Motif enrichments of binarized peaks identified in Figure 2e. Due to redundancy in
873 motifs, TF drivers were predicted using average gene expression in GTEx brain samples
874 and accessibility at TF promoters in cell type-grouped scATAC-seq profiles. The final
875 list of TFs represents a trimmed set of all TFs with the most likely driving TF labeled
876 below. Color represents the p-value of the hypergeometric test for motif enrichment.
- 877 G. Footprinting analysis of the SPI1 (left) and JUND (right) transcription factors across the
878 7 major cell types. The motif logos are shown above and the Tn5 transposase insertion
879 biases are shown below.
- 880

881 **Figure 3 - HiChIP and scATAC-seq predict gene and cellular targets of disease-associated**
882 **polymorphisms**

- 883 A. LD score regression identifying the enrichment of GWAS SNPs from various brain- and
884 non-brain-related conditions in the peak regions of various cell types derived from
885 pseudo-bulk-based scATAC-seq data.
- 886 B. Heatmap representation of HiChIP interaction signal at 100-kb, 25-kb, and 5-kb
887 resolution at the *OLIG2* locus.
- 888 C. Characterization of GWAS loci in AD and PD according to the predicted effects of the
889 polymorphisms. For example, loci whose phenotypic association is likely mediated by
890 changes in coding regions are marked as “Likely coding”. Loci whose effect could be
891 mediated by either coding or noncoding mechanisms are marked as “Either coding or
892 noncoding” whereas loci with no polymorphisms overlapping a peak region or an exonic
893 region are marked as “Unknown”.
- 894 D. Histogram of the number of genes linked per GWAS locus. Each bar represents a bin of
895 length 1.
- 896 E. Venn diagram of (i) the number of genes linked through assessment of the nearest gene to
897 the lead SNP of each AD (top) and PD (bottom) GWAS locus and (ii) the number of
898 genes linked though HiChIP and scATAC-seq analyses of LD-expanded polymorphisms.
- 899

900 **Figure 4 - Machine learning predicts functional polymorphisms in AD and PD**

- 901 A. Schematic of the overall strategy for identification of putative functional SNPs and their
902 corresponding gene targets.
- 903 B. Schematic of the gkm-SVM machine learning approach used to predict which noncoding
904 SNPs alter transcription factor binding and chromatin accessibility.
- 905 C. Normalized scATAC-seq-derived pseudo-bulk tracks, HiChIP loop calls, co-accessibility
906 correlations, and machine learning predictions for LD-expanded SNPs in the *PICALM*
907 gene locus. For HiChIP, each line represents a loop connecting the points on each end.
908 Red lines contain one anchor overlapping the SNP of interest while grey lines do not.

- 909 D. GkmExplain importance scores for each base in the 50-bp region surrounding rs1237999
910 for the effect and non-effect alleles from the gkm-SVM model corresponding to
911 oligodendrocytes (Cluster 21). The predicted motif affected by the SNP is shown at the
912 bottom and the SNP of interest is highlighted in blue.
913 E. Dot plot showing allelic imbalance at rs1237999. The ATAC-seq counts for the
914 reference/non-effect (G) allele and variant/effect (A) allele are plotted. Each dot
915 represents an individual bulk ATAC-seq sample colored by the brain region from which
916 the sample was collected.
917 F. Sequencing tracks as shown in Figure 4c but for the *SLC24A4* locus.
918 G. GkmExplain importance scores for each base in the 50-bp region surrounding
919 rs10130373 for the effect and non-effect alleles from the gkm-SVM model corresponding
920 to microglia (Cluster 24). The predicted motif affected by the SNP is shown at the bottom
921 and the SNP of interest is highlighted in blue.
922

923 **Figure 5 - Vertical integration of multi-omic data and machine learning nominates novel
924 gene targets in AD and PD**

- 925 A. Normalized scATAC-seq-derived pseudo-bulk tracks, HiChIP loop calls, co-accessibility
926 correlations, and machine learning predictions for LD-expanded SNPs in the *ITIH1* gene
927 locus. For HiChIP, each line represents a loop connecting the points on each end. Red
928 lines contain one anchor overlapping the SNP of interest while grey lines do not.
929 B. GkmExplain importance scores for each base in the 50-bp region surrounding
930 rs181391313 for the effect and non-effect alleles from the gkm-SVM model
931 corresponding to microglia (Cluster 24). The predicted motif affected by the SNP is
932 shown at the bottom and the SNP of interest is highlighted in blue.
933 C. Sequencing tracks as shown in Figure 5a but for the *KCNIP3* locus.
934 D. GkmExplain importance scores for each base in the 50-bp region surrounding rs7585473
935 for the effect and non-effect alleles from the gkm-SVM model corresponding to
936 oligodendrocytes (Cluster 21). The predicted motif affected by the SNP is shown at the
937 bottom and the SNP of interest is highlighted in blue.
938 E. Dot plot showing allelic imbalance at rs3755519. The ATAC-seq counts for the
939 reference/non-effect (A) allele and variant/effect (T) allele are shown. Each dot
940 represents an individual bulk ATAC-seq sample colored by the brain region from which
941 the sample was collected.
942

943 **Figure 6 - Epigenetic deconvolution of *MAPT* locus explains haplotype-associated
944 transcriptional changes**

- 945 A. Schematic of the *MAPT* locus (chr17:44905000-46895000) showing all genes, the
946 predicted locations of the inversion breakpoints, and the 2366 haplotype-divergent SNPs
947 used for haplotype-specific analyses.

- 948 B. Gene expression of the *MAPT* gene shown as a box plot from GTEx cortex brain samples
949 subdivided based on *MAPT* haplotype. The lower and upper ends of the box represent the
950 25th and 75th percentiles. The whiskers represent 1.5 multiplied by the inter-quartile
951 range.
952 C. Schematic for the allelic analysis of the *MAPT* region. Data from homozygous H1 and
953 H2 individuals are directly compared. Data from heterozygous H1/H2 individuals are
954 first split based off of the presence of haplotype-divergent SNPs in the reads and then
955 compared.
956 D. HiChIP (top) and ATAC-seq (middle) sequencing tracks of the region representing the
957 *MAPT* locus inside of the predicted inversion breakpoints (chr17:45510000-46580000;
958 bottom). Each track represents the merge of all available H1 or H2 reads from all
959 heterozygotes. HiChIP and ATAC-seq tracks represent unnormalized data from
960 heterozygotes where reads were split based on haplotype. No normalization was
961 performed because each sample is internally controlled for allelic depth. HiChIP is shown
962 as a virtual 4C plot where the anchor is indicated by a dotted line and the signal
963 represents paired-end tag counts overlapping a 10-kb bin. Regions showing significant
964 haplotype bias in ATAC-seq are marked by an asterisk.
965 E. GTEx cortex gene expression of genes in the *MAPT* locus comparing H1 homozygotes to
966 H1/H2. Regions A and B are shown as in Figure 6d. *p < 0.05 after multiple hypothesis
967 correction.
968 F. HiChIP (top) and cell type-specific scATAC-seq (middle) sequencing tracks of the region
969 representing the *MAPT* locus outside of the predicted inversion breakpoints (bottom).
970 HiChIP tracks for bulk homozygote H1 or H2 samples (normalized based on reads-in-
971 loops) are shown at the top while haplotype-specific tracks from heterozygotes
972 (unnormalized) are shown below. In each HiChIP plot, the anchor represents the *MAPT*
973 promoter.
974 G. Schematic illustrating the predicted haplotype-specific change in long-distance
975 interaction between the *MAPT* promoter and the predicted distal enhancer identified in
976 Figure 6d. Regions marked A and B represent the same regions marked in Figure 6d-e.
977

978 SUPPLEMENTARY FIGURE LEGENDS

980 Supplementary Figure 1 - Analysis of bulk ATAC-seq data from adult brain identifies 981 brain-regional heterogeneity.

- 982 A. Principal component analysis of all samples. Each dot represents a single piece of tissue
983 with technical replicates merged where applicable. Color represents the brain region from
984 which the sample was isolated.
985 B. Dot plot showing the proportion of variance explained by each principal component.

- 986 C. Dot plot showing the significance of correlation between covariates and each of the top 5
987 principal components. Dot size represents the absolute value of the correlation while
988 color represents the principal component number.
989 D. Sample by sample Pearson correlation heatmap of all 140 samples profiled in this study.
990 Brain region, donor biological sex, and *APOE* genotype are indicated colorimetrically at
991 the top.
992 E. MA plots showing the change in normalized bulk ATAC-seq accessibility for each peak
993 in cognitively healthy control samples with low AD-associated pathology compared to
994 cognitively healthy control samples with high AD-associated pathology. Each dot
995 represents an individual peak from the merged bulk ATAC-seq peak set. Only peaks that
996 showed non-zero accessibility in at least one sample were tested for significance. From
997 left to right, samples from the caudate nucleus, hippocampus, parietal lobe, and superior
998 and middle temporal gyrus are shown.
999 F. MA plots showing the change in normalized bulk ATAC-seq accessibility comparing the
1000 parietal lobe (PARL) to all other brain regions. Each dot represents an individual peak
1001 from the merged bulk ATAC-seq peak set. Only peaks that showed non-zero accessibility
1002 in at least one sample were tested for significance.
1003

1004 **Supplementary Figure 2 - LD score regression of bulk ATAC-seq data identifies weak**
1005 **region-specific enrichment of AD and PD GWAS SNPs.**

- 1006 A. Bar plot of the enrichment of AD SNPs in peaks regions of bulk ATAC-seq data from
1007 various brain regions.
1008 B. Bar plot of the enrichment of PD SNPs in peak regions of bulk ATAC-seq data from
1009 various brain regions.
1010 C. Dot plots showing the TSS enrichment score and total number of fragments for each of
1011 the 10 samples profiled by scATAC-seq. Each dot represents an individual cell. Dot color
1012 represents density on the plot. Dotted lines represent the quality control cutoffs
1013 implemented.
1014 D. Heatmap of cell type-specific markers used to identify clusters. Color represents the row-
1015 wise Z-score of chromatin accessibility in the vicinity of each gene for each cluster.
1016

1017 **Supplementary Figure 3 - Region-centric scATAC-seq identifies cellular and regional**
1018 **heterogeneity in chromatin accessibility in adult brain**

- 1019 A. UMAP dimensionality reduction as shown in Figure 2a but colored by the sample from
1020 which each cell was generated.
1021 B. UMAP dimensionality reduction as shown in Figure 2a but colored by the brain region
1022 from which each cell was generated.
1023 C. UMAP dimensionality reduction as shown in Figure 2a but colored by the biological sex
1024 of the donor for each cell.

- 1025 D. UMAP dimensionality reduction as shown in Figure 2a but colored by the predicted cell
1026 type for each cell.
1027 E. Bar plot showing the number of cells identified in scATAC-seq from each of the
1028 annotated cell types.
1029 F. Bar plot showing the number of cells in scATAC-seq from each of the annotated
1030 donors/samples. Color represents the predicted cell type as shown in the legend next to
1031 Supplementary Fig. 3h.
1032 G. Bar plot showing the number of cells identified in scATAC-seq from each of the
1033 annotated cell types broken down by the brain region from which they originated. Color
1034 represents the predicted cell type as shown in the legend next to Supplementary Fig. 3h.
1035 H. Bar plot showing the percentage of each brain region composed by each cell type in
1036 scATAC-seq data.
1037 I. Bar plot showing the percentage of cells from each cell type that originated from each
1038 donor sample profiled by scATAC-seq. Color represents the biological sample from
1039 which the data was collected.

1040
1041 **Supplementary Figure 4 - Cell type-specific scATAC-seq data enables deconvolution of**
1042 **chromatin accessibility data from bulk regions in the adult brain.**

- 1043 A. Sequencing tracks of lineage-defining factors shown across all 24 scATAC-seq clusters.
1044 From left to right, *NEFL* (neurons; chr8:24933431-24966791), *AIF1* (aka *IBA1*,
1045 microglia; chr6:31607841-31617906), *MOG* (oligodendrocytes; chr6:29652183-
1046 29699713), *PDGFRA* (OPCs; chr4:54209541-54303643), and *GJB6* (astrocytes;
1047 chr13:20200243-20239571).
1048 B. Bar plot showing CIBERSORT deconvolution of bulk ATAC-seq data based on
1049 reference cell populations derived from scATAC-seq data. Clusters were subdivided into
1050 the 8 groups shown in the legend. These groups were used to preserve as much diversity
1051 as possible while merging clusters with little divergence (i.e. oligodendrocyte clusters
1052 #19-23). Bars represent the average of all bulk ATAC-seq samples profiled in the given
1053 brain regions.
1054 C. Bar plot showing CIBERSORT deconvolution of bulk ATAC-seq data based on clusters
1055 derived from scATAC-seq data. Color represents the cluster as shown in the legend of
1056 Supplementary Fig. 4g. Bars represent the average of all bulk ATAC-seq samples
1057 profiled in the given brain regions.
1058 D. Dot plot showing the performance of the CIBERSORT classifier by comparing the
1059 “ground truth” from scATAC-seq data and the CIBERSORT prediction on the bulk
1060 ATAC-seq data from the same tissue sample. Each dot represents a cell type (i.e. the
1061 merge of multiple clusters) from one of the 10 scATAC-seq samples profiled. Dots are
1062 colored by cell type according to the legend above the plot.
1063 E. Dot plot showing the performance of the CIBERSORT classifier by comparing the
1064 “ground truth” from scATAC-seq data and the CIBERSORT prediction on the bulk

1065 ATAC-seq data from the same tissue sample. Each dot represents a cluster from one of
1066 the 10 scATAC-seq samples profiled. Dots are colored by cluster according to the legend
1067 in Supplementary Fig. 4g.

- 1068 F. Bar plot showing CIBERSORT predictions across all bulk ATAC-seq data generated in
1069 this study. Samples are sorted and colored (bottom of plot) by the region from which they
1070 were profiled as indicated in the legend below Supplementary Fig. 4g. Bars are colored
1071 by the predicted cell type. Donor IDs are annotated below the plot.
1072 G. Bar plot showing CIBERSORT predictions across all bulk ATAC-seq data generated in
1073 this study. Samples are sorted and colored (bottom of plot) by the region from which they
1074 were profiled. Bars are colored by the predicted cluster. Donor IDs are annotated below
1075 the plot.

1076

1077 **Supplementary Figure 5 - scATAC-seq reveals epigenetic encoding of region-specific**
1078 **cellular gene regulatory programs**

- 1079 A. Pearson correlation heatmaps showing the correlation of cell types across brain regions.
1080 Cell type signals were generated by making at least 2 non-overlapping pseudo-bulk
1081 replicates of at least 150 cells. Cases where insufficient cells were present to make these
1082 pseudo-bulk replicates were excluded from analysis (ND) to avoid overinterpretation. All
1083 heatmaps use the same color scale.
1084 B. Volcano plot of peaks that show differential signal between astrocytes from the substantia
1085 nigra and astrocytes from the isocortex. Peaks below a log₂(fold change) threshold of 2
1086 were not considered. Peaks near genes that are predicted to be key lineage-defining genes
1087 are accented with larger colored dots.
1088 C. UMAP dimensionality reduction plots showing gene activity scores colorimetrically for
1089 the 4 lineage-defining genes identified in Supplementary Fig. 5b (*FOXG1*, *ZIC5*, *FOXB1*,
1090 *IRX1*).
1091 D. Sequencing tracks of the multiple genomic regions showing differential chromatin
1092 accessibility between astrocytes or OPCs in the isocortex and substantia nigra. From left
1093 to right: Isocortex-specific - *FOXG1* (chr14:28750000-28787000), and *ZIC2/ZIC5*
1094 (chr13:99937000-99999000); Substantia Nigra-specific:- *FOXB1* (chr15:59996000-
1095 60012000), *IRX1* (chr5:3589600-3607800), *IRX2* (chr5:2737000-2760000), *IRX3*
1096 (chr16:54277000-54292000), *IRX5* (chr16:54927000-54940000), and *PAX3*
1097 (chr2:222189500-222333500). Peaks called in scATAC-seq data are shown below each
1098 plot. Sequencing tracks were derived from merging of all single cells corresponding to
1099 the annotated cell types in the specified regions.
1100 E. Volcano plot of peaks that show differential signal between OPCs from the substantia
1101 nigra and OPCs from the isocortex. Peaks below a log₂(fold change) threshold of 2 were
1102 not considered. Peaks near genes that are predicted to be key lineage-defining genes are
1103 accented with larger colored dots.

- 1104 F. Same as Supplementary Fig. 5e but for oligodendrocytes in the substantia nigra and
1105 isocortex.
1106 G. Same as Supplementary Fig. 5e but of microglia in the substantia nigra and isocortex.
1107 H. Sequencing tracks of regions identified as differentially accessible in oligodendrocytes
1108 from the substantia nigra and isocortex. From left to right: Isocortex-specific - *SHC2*
1109 (chr19:409800-463200), and *INSM1* (chr20:20361000-20374000); Substantia nigra-
1110 specific - *RBFOX1* (chr16:5899200-7791000). Sequencing tracks were derived from
1111 merging of all single cells corresponding to the annotated cell types in the specified
1112 regions.
1113 I. Same as Supplementary Fig. 5e but for inhibitory neurons in the isocortex and striatum.
1114 J. Sequencing tracks of regions identified as differentially accessible in inhibitory neurons
1115 from the striatum and isocortex. From left to right: Isocortex-specific - *KCNJ6*
1116 (chr21:37583000-37955000), and *NCALD* (chr8:101673000-102141000); Striatum-
1117 specific - *DRD2* (chr11:113369000-113602000), and *FOXP1* (chr3:70922000-
1118 71622000). Sequencing tracks were derived from merging of all single cells
1119 corresponding to the annotated cell types in the specified regions.
1120

1121 **Supplementary Figure 6 - HiChIP implicates disease-relevant genes in AD and PD through
1122 linkage of noncoding GWAS SNPs to target genes**

- 1123 A. LD score regression identifying the enrichment of GWAS SNPs from various brain- and
1124 non-brain-related conditions in the peak regions of bulk ATAC-seq data from various
1125 hematopoietic cell types as indicated by color.
1126 B. Heatmap representation of HiChIP interaction signal at 100-kb, 25-kb, and 5-kb
1127 resolution at the SOX9 locus.
1128 C. Bar plots showing the number of valid interaction pairs identified in HiChIP data from all
1129 samples profiled in this study. Color represents the type of interaction identified.
1130 D. Bar plot showing the overlap of FitHiChIP loop calls from the 4 gross brain regions
1131 profiled. Color indicates whether the loop was identified in a single region (unique) or
1132 more than one region (shared).
1133 E. Bar plot showing the classification of FitHiChIP loop calls based on whether the loop call
1134 contained an ATAC-seq peak (bulk or single-cell) or TSS in one, both, or no anchor.
1135 F. Bar plots showing the number of Cicero-predicted co-accessibility-based peak links that
1136 are observed in HiChIP (left) or the number of HiChIP-based FitHiChIP loop calls that
1137 are predicted as peak links by Cicero.
1138 G. GO-term enrichments of genes linked to AD GWAS SNPs.
1139 H. GO-term enrichments of genes linked to PD GWAS SNPs.
1140

1141 **Supplementary Figure 7 - Machine learning and allelic imbalance predict functional
1142 noncoding SNPs in AD and PD**

- 1143 A. Flow chart of the analytical framework used to prioritize noncoding SNPs and predict
1144 functionality. The highest confidence SNPs (Tier 1) are supported by either machine
1145 learning predictions, allelic imbalance, or both. Moderate confidence SNPs (Tier 2) are
1146 supported by the presence of the SNP within a peak and a HiChIP loop or co-accessibility
1147 peak link that connects the SNP to a gene. Lower confidence SNPs (Tier 3) are only
1148 supported by the presence of the SNP in a peak.
- 1149 B. Box plot showing the area under the precision-recall curve for the gkm-SVM machine
1150 learning classifier. Performance for each cluster is shown with dots representing outliers.
1151 The lower and upper ends of the box represent the 25th and 75th percentiles. The
1152 whiskers represent 1.5 multiplied by the inter-quartile range.
- 1153 C. Box plot showing the area under the receiver-operating characteristics curve for the gkm-
1154 SVM machine learning classifier. Performance for each cluster is shown with dots
1155 representing outliers. The lower and upper ends of the box represent the 25th and 75th
1156 percentiles. The whiskers represent 1.5 multiplied by the inter-quartile range.
- 1157 D. GkmExplain importance scores shown across all 10 folds for each base across a 100-bp
1158 window surrounding rs636317 for the effect (left) and noneffect (right) bases.
- 1159 E. Dot plots showing comparison of the GkmExplain score, ISM score, and deltaSVM
1160 score. Each dot represents an individual SNP test in a given fold. Dot color represents the
1161 GWAS locus number. The only off-diagonal dots (circled) correspond to repetitive
1162 regions within the *MAPT* locus where the deltaSVM score appears to be particularly
1163 sensitive.
- 1164 F. Dot plot showing allelic imbalance across all bulk ATAC-seq data used in this study.
1165 ATAC-seq data was used to genotype individuals to identify heterozygotes. Allelic
1166 imbalance was defined as ratio of wildtype to variant reads that passes the binomial test
1167 with a p-value less than 0.05. Color indicates the average significance of the binomial test
1168 across all heterozygotes.
- 1169

1170 **Supplementary Figure 8 - Multi-omic characterization of well-studied AD-related GWAS**
1171 **loci pinpoints putative functional noncoding SNPs**

- 1172 A. Normalized scATAC-seq-derived pseudo-bulk tracks, HiChIP loop calls, co-accessibility
1173 correlations, and machine learning predictions for LD-expanded SNPs in the *BIN1* locus.
1174 For HiChIP, each line represents a loop connecting the points on each end. Red lines
1175 contain one anchor overlapping the SNP of interest while grey lines do not.
- 1176 B. GkmExplain importance scores for each base in the 50-bp region surrounding
1177 rs13025717 for the effect and non-effect alleles from the gkm-SVM model for microglia
1178 (Cluster 24). The predicted motif affected by the SNP is shown at the bottom and the
1179 SNP of interest is highlighted in blue.
- 1180 C. Sequencing tracks as shown in Supplementary Fig. 8a but for the *MS4A* gene locus.
- 1181 D. GkmExplain importance scores for each base in the 50-bp region surrounding rs636317
1182 for the effect and non-effect alleles from the gkm-SVM model for microglia (Cluster 24).

1183 The predicted motif affected by the SNP is shown at the bottom and the SNP of interest is
1184 highlighted in blue.
1185

1186 **Supplementary Figure 9 - Multi-omic characterization of noncoding SNPs identifies novel**
1187 **genes implicated in PD**

- 1188 A. Normalized scATAC-seq-derived pseudo-bulk tracks, HiChIP loop calls, co-accessibility
1189 correlations, and machine learning predictions for LD-expanded SNPs in the *IP6K2*
1190 locus. For HiChIP, each line represents a loop connecting the points on each end. Red
1191 lines contain one anchor overlapping the SNP of interest while grey lines do not.
- 1192 B. GkmExplain importance scores for each base in the 50-bp region surrounding rs6781790
1193 for the effect and non-effect alleles from the gkm-SVM model for astrocytes (Cluster 15).
1194 The predicted motif affected by the SNP is shown at the bottom and the SNP of interest is
1195 highlighted in blue.
- 1196 C. Dot plot showing allelic imbalance at rs6781790. The ATAC-seq counts for the
1197 reference/non-effect (C) allele and variant/effect (T) allele are plotted. Each dot
1198 represents an individual bulk ATAC-seq sample colored by the brain region from which
1199 the sample was collected.
- 1200 D. Sequencing tracks as shown in Supplementary Fig. 9a but for the *TMEM163* locus.
- 1201 E. GkmExplain importance scores for each base in the 50-bp region surrounding rs7599054
1202 for the effect and non-effect alleles from the gkm-SVM model for microglia (Cluster 24).
1203 The predicted motif affected by the SNP is shown at the bottom and the SNP of interest is
1204 highlighted in blue.

1205

1206 **Supplementary Figure 10 - Epigenomic dissection of the *MAPT* locus**

- 1207 A. Flowchart illustrating the analytical scheme used to identify bins with significant allelic
1208 imbalance across the H1 and H2 *MAPT* haplotypes.
- 1209 B. Heatmaps showing chromatin accessibility in 500-bp bins identified as having
1210 significantly different accessibility across *MAPT* haplotypes. Regions are shown for
1211 homozygous samples without allelic read splitting (left) and for heterozygous samples
1212 after allelic read splitting (right). Bin start coordinates are shown to the right.
- 1213 C. Box and whiskers plots for multiple regions which show differential chromatin
1214 accessibility across the H1 and H2 *MAPT* haplotypes. Each dot represents a single
1215 homozygous H1 or homozygous H2 sample. Heterozygotes are not shown. The lower and
1216 upper ends of the box represent the 25th and 75th percentiles. The whiskers represent 1.5
1217 multiplied by the inter-quartile range.
- 1218 D. Gene expression of the *KANSL1-AS1* gene shown as a box plot from GTEx cortex brain
1219 samples subdivided based on *MAPT* haplotype. The lower and upper ends of the box
1220 represent the 25th and 75th percentiles. The whiskers represent 1.5 multiplied by the
1221 inter-quartile range. ***p < 10⁻⁵.

- 1222 E. Gene expression of the *MAPK8IP1P2* gene shown as a box plot from GTEx cortex brain
1223 samples subdivided based on *MAPT* haplotype. The lower and upper ends of the box
1224 represent the 25th and 75th percentiles. The whiskers represent 1.5 multiplied by the
1225 inter-quartile range. ***p < 10⁻⁵.
- 1226 F. Sequencing tracks from pseudo-bulk data derived from predicted cell types in scATAC-
1227 seq data. This region represents a zoomed in view of the predicted distal enhancer region
1228 (chr17:45216500-45324000) that interacts with the *MAPT* promoter in the H1 haplotype.
1229 Putative neuron-specific enhancers are highlighted in blue.
- 1230 G. Box plots showing differential HiChIP interaction signal occurring between regions
1231 within the *MAPT* inversion and regions outside the inversion (“left” or “right”). The
1232 schematic at the top explains the analysis performed. The box plots show normalized
1233 HiChIP interaction counts for the H1 and H2 haplotypes for upstream/“left” interactions
1234 and downstream/“right” interactions.
- 1235

1236 SUPPLEMENTARY TABLES

1237

1238 **Supplementary Table 1** – Donor information and sequencing statistics for all samples profiled
1239 by bulk ATAC-seq, scATAC, and HiChIP.

1240

1241 **Supplementary Table 2** – Final merged peak set derived from all bulk ATAC-seq data.

1242

1243 **Supplementary Table 3** – All LD-expanded GWAS SNPs from AD and PD and their relevant
1244 metadata and characterizations.

1245

1246 **Supplementary Table 4** – Quality control information for all individual cells profiled by
1247 scATAC-seq and the cluster residence information for all clusters and samples.

1248

1249 **Supplementary Table 5** – Final merged peak set derived from all scATAC-seq data.

1250

1251 **Supplementary Table 6** – Results of feature binarization from scATAC-seq data showing cell
1252 type-specific peaks.

1253

1254 **Supplementary Table 7** – CIBERSORT signature matrices for the cell group-specific and
1255 cluster-specific classifiers.

1256

1257 **Supplementary Table 8** – Results of differential accessibility comparisons between the
1258 substantia nigra and isocortex for astrocytes, OPCs, oligodendrocytes, and microglia.

1259

1260 **Supplementary Table 9** – Results of all LD score regression analyses across all conditions and
1261 cell types.

1262

1263 **Supplementary Table 10** – All FitHiChIP loop calls overlapping a SNP on at least one anchor.

1264

1265 **Supplementary Table 11** – All SNPs that are divergent between the H1 and H2 haplotypes in
1266 the *MAPT* locus.

1267

1268 **METHODS**

1269

1270 **Code Availability**

1271 All custom code used in this work is available in the following GitHub repository:

1272 https://github.com/kundajelab/alzheimers_parkinsons.

1273

1274 **Publicly Available Data Used In This Work**

1275 All QTL analysis was performed using GTEx v8. Additionally, we downloaded full-genome

1276 summary statistics of GWAS associations for three Alzheimer's cohorts^{1–3} and three Parkinson's

1277 cohorts^{6,7,62}; however, it should be noted that these cohorts are not all mutually exclusive.

1278

1279 **Genome Annotations**

1280 All data is aligned and annotated to the hg38 reference genome.

1281

1282 **Sequencing**

1283 Bulk ATAC-seq, and HiChIP were sequenced using an Illumina HiSeq 4000 with paired-end 75-

1284 bp reads. Single-cell ATAC-seq was sequenced using an Illumina NovaSeq 6000 with an S4 flow

1285 cell with paired-end 99 bp reads.

1286

1287 **Sample acquisition and patient consent**

1288 Primary brain samples were acquired post-mortem with IRB-approved informed consent. Human

1289 donor sample sizes were chosen to provide sufficient confidence to validate methodological

1290 conclusions. Human brain samples were collected with an average post-mortem interval of 3.9

1291 hours (range 2.0 – 6.9 hours). Macrodissected brain regions were flash frozen in liquid nitrogen.

1292 Some samples were embedded in Optimal Cutting Temperature (OCT) compound. All samples

1293 were stored at -80°C until use.

1294

1295 **Isolation of nuclei from frozen tissue chunks**

1296 Nuclei were isolated from frozen tissue as described previously^{63,64}. This protocol is now available

1297 on protocols.io (dx.doi.org/10.17504/protocols.io.6t8herw). After isolation, nuclei were

1298 cryopreserved in BAM Banker (Wako Chemicals) and stored at -80°C for use in other assays such

1299 as scATAC-seq and HiChIP.

1300

1301 **Statistics**

1302 All statistical tests performed are included in the figure legends or methods where relevant.

1303

1304 **ATAC-seq Data Processing**

1305 The ENCODE DCC ATAC-seq pipeline (doi:10.5281/zenodo.211733) (V1.1.7) was used to
1306 process bulk ATAC-seq samples, starting from fastq files. The pipeline was executed with IDR
1307 enabled and the IDR threshold set to 0.05. The GRCh38 reference genome assembly was used,
1308 keeping only the primary chromosomes chr1 - chr22, chrX, chrY, chrM. The pipeline was executed
1309 with ATAQC enabled, using GENCODE version 29 TSS annotations. Biological replicates were
1310 analyzed individually, with the two technical replicates for each bio-rep provided as inputs to the
1311 “atac.bams” argument of the pipeline. Other arguments to the pipeline were kept at their defaults.
1312

1313 **ATAC-seq Peak Calling**

1314 Pipeline peak calls underwent several levels of filtering to identify credible peak sets. The IDR
1315 optimal peak set from the DCC pipeline for each biological replicate was determined. It was
1316 observed that although the IDR peaks for individual biological replicates were corrected for
1317 multiple testing, the high number of biological samples in the dataset served as another source of
1318 multiple testing error. To address this source of error, tagAlign files for all biological replicates
1319 for a given brain region/ condition were concatenated. The DCC pipeline (v1.1.7) was
1320 subsequently executed on the merged tagAlign files as single-replicate inputs. The pipeline
1321 generated pseudo-replicates from the input tagAlign files for each brain region/condition. Optimal
1322 IDR peaks were called from the pseudo-replicates. This set of IDR peaks was filtered to keep peaks
1323 supported by 30 percent or more of IDR peaks from the pipeline runs on individual biological
1324 replicates.

1325 Sample-by-peak count matrices were then generated from the resulting set of filtered peaks.
1326 Filtered peaks from the pooled tagAlign files were concatenated and truncated to within 200 base
1327 pairs of the summit (100 base pair flank kept upstream and downstream of the peak summit). These
1328 200 bp regions were merged with the bedtools⁶⁵ merge command to avoid merging peaks with low
1329 levels of overlap. The bedtools coverage -counts was used to compute the number of tagAlign
1330 reads that overlapped each peak region in the pseudo-replicates in the merged tagAlign dataset.
1331 This analysis yielded a total of n=186,559 peaks combined across the brain regions.
1332

1333 **Motif enrichment**

1334 Motif enrichment was performed using the hypergeometric test as described previously^{64,66}.
1335

1336 **Feature Binarization**

1337 Identification of “unique” peaks from ATAC-seq data was performed as described previously^{12,64}.
1338

1339 **Sequencing Tracks**

1340 Sequencing tracks were created using the WashU Epigenome Browser. All sequencing tracks of a
1341 given locus have the same y-axis. All tracks show data that has been normalized by “reads-in-

1342 peaks” (for ATAC-seq) or “reads-in-loops” for HiChIP to account for differences in signal-to-
1343 background ratios across multiple samples, unless otherwise stated. For all sequencing tracks,
1344 genes that are on the plus strand (i.e. 5’ to 3’ in the left to right direction) are shown in red and
1345 genes that are on the minus strand (i.e. 5’ to 3’ in the right to left direction) are shown in blue to
1346 enable identification of the TSS.

1347

1348 **LD score regression**

1349 We apply stratified LD score regression, a method for partitioning heritability from GWAS
1350 summary statistics, to sets of tissue or cell type specific ATAC-seq peaks to identify disease-
1351 relevant tissues and cell types across for Alzheimer's and Parkinson's diseases along with other
1352 brain-related GWAS traits. We used both bulk ATAC-seq and single cell ATAC-seq data. For
1353 bulk ATAC-seq we kept only peaks replicating in at least 30% of samples for each tissue part.
1354 ATAC-seq peaks were converted from hg38 to hg19 for analysis with GWAS data. We followed
1355 the LD score regression tutorial (<https://github.com/bulik/ldsc/wiki>) as used previously⁶⁷ for bulk
1356 data and as recently developed for single-cell specific analysis⁶⁸. We used brain related GWAS
1357 summary statistics such as Alzheimer's¹, Parkinson's⁶, Schizophrenia⁶⁹, Anorexia Nervosa⁷⁰,
1358 Attention Deficit Hyperactivity Disorder (ADHD)⁷¹, Anxiety⁷², Neuroticism⁷³ and Epilepsy⁷⁴. To
1359 serve as controls, we also used summary statistics for GWAS of traits not obviously linked to brain
1360 tissues such as Lean Body Mass⁷⁵, Bone Mineral Density⁷⁶ and Coronary Artery Disease⁷⁷. In
1361 particular, we looked at the regression coefficient p-value, indicative of the contribution of this
1362 annotation to trait heritability, conditional on the other annotations.

1363

1364 **Allelic imbalance from ATAC-seq data**

1365 Samples were first re-aligned to an N-masked version of the hg38 genome where all relevant SNP
1366 positions were changed to “N” to prevent mapping bias. Allelic depth at each desired position was
1367 obtained using samtools mpileup (v1.5) followed by varscan mpileup2snp (v2.4.3). Allele counts
1368 for the reference and variant alleles were extracted and compared using the binomial test to identify
1369 significant allelic imbalance.

1370

1371 **SNP selection for colocalization testing**

1372 A single test for colocalization of GWAS and eQTL association signals involves a locus, a GWAS,
1373 an eQTL tissue, and a gene expressed in that tissue. For each GWAS, we selected the set of all loci
1374 for which the lead GWAS variant had p-value < 1e-5. Using eQTLs from GTEx brain tissues in
1375 the GTEx v8 dataset, we then found all tissue-gene combinations for which the lead SNP at one
1376 of the GWAS loci had an eQTL SNP (association p-value < 1e-5) for that gene in that GTEx tissue.
1377 This resulted in a list of unique combinations of GWAS trait / genomic locus / eQTL tissue / eQTL
1378 gene, each to be tested individually for colocalization of GWAS and eQTL signals. The GWAS
1379 threshold of 1e-5 is less stringent than the threshold for genome-wide significance, but we favored
1380 sensitivity over specificity when selecting which SNPs to test, since colocalization with a strong

1381 eQTL signal may still suggest that a sub-threshold GWAS locus has an expression-mediated effect
1382 on disease.

1383

1384 **Colocalization analysis**

1385 For each colocalization test combination as defined above, we selected all 1000 Genomes Phase 3
1386 variants within a window of 500kb around the lead GWAS variant. We narrowed this list down to
1387 SNPs measured not only in the 1000 Genomes VCF, but also in the GWAS and eQTL summary
1388 statistics for the selected trait, tissue, and gene. We used a streamlined version of the FINEMAP
1389 tool⁷⁸ to compute posterior causal probabilities for each SNP at the locus in both the GWAS and
1390 eQTL studies, and then combined these probabilities as described in eCAVIAR⁷⁹ to compute a
1391 colocalization posterior probability (CLPP) score for this test locus. We considered a SNP weakly
1392 colocalized if its CLPP score exceeded 0.01 and strongly colocalized if its CLPP score exceeded
1393 0.05; although these seem like quite low probabilities, we have seen previously that loci exceeding
1394 this latter cutoff show strong likelihood of sharing causal variants⁸⁰.

1395

1396 **Selection of candidate SNPs for ATAC-seq overlap analysis, HiChIP interaction tests, and** 1397 **gkm-SVM model-based allelic effect scores**

1398 Our goal was to identify SNPs with a causal effect on any of the selected GWAS traits. To
1399 minimize the chances of excluding causal GWAS SNPs, we selected the set of all variants
1400 achieving a genome-wide significant p-value < 5e-8 for any GWAS trait. We then added in any
1401 lead SNPs from the colocalization analysis that achieved CLPP score of > 0.01, even those that
1402 did not pass the genome-wide significance value of p < 5e-8. We also included all trait-associated
1403 SNPs curated from two other Parkinson's studies^{6,7}. In these studies, full summary statistics were
1404 not publicly available for the entire genome because meta-analysis was applied only to the subset
1405 of SNPs reaching genome-wide significance in a previous Parkinson's GWAS. We then computed
1406 the full set of SNPs that had LD R^2 ≥ 0.8 with at least one of the SNPs in the set selected above.
1407 Together, these LD buddies plus the original set of trait-relevant SNPs comprised the set of SNPs
1408 tested in our subsequent functional analyses.

1409

1410 **Testing GWAS loci for overlap with ATAC-seq peaks**

1411 We tested all SNPs in the above set for overlap with ATAC-seq peaks from two different
1412 annotation formats. The first annotation consisted of bulk ATAC-seq peaks identified in one of 7
1413 brain regions. The second annotation consisted of cluster-specific peaks from single-cell ATAC-
1414 seq data. For each variant selected for functional analysis, we determined all cellular contexts in
1415 which an ATAC-seq peak contained this variant, as well as the nearest peak if no peak contained
1416 the variant.

1417

1418 **Single-cell ATAC-seq library generation**

1419 Cryopreserved nuclei were thawed on ice and 65,000 nuclei were transferred to a tube containing
1420 1 ml of RSB-T [10 mM Tris-HCl pH 7.5, 10 mM NaCl, 3 mM MgCl2, 0.1% Tween]. Nuclei were

1421 pelleted at 500 RCF for 5 minutes at 4°C in a fixed angle rotor. The supernatant was fully removed
1422 using two pipetting steps (p1000 to remove down to the last 100 ul, then p200 to remove all
1423 remaining supernatant). This pellet was then gently resuspended in 12 ul of 1x Nuclei Buffer (10x
1424 Genomics). To transpose, 5 ul of this nuclei suspension (containing 27,000 nuclei) was transferred
1425 to a tube containing 10 ul of transposition mix (10x Genomics). This reaction mixture was
1426 incubated at 37°C for 1 hour to transpose. The remainder of library generation was completed as
1427 described in the 10x Genomics Single Cell ATAC Regent Kits User Guide (v1 Chemistry).

1428

1429 **Single-cell ATAC-seq LSI clustering and visualization**

1430 To cluster our scATAC-seq data, we first identified a robust set of peak regions followed by
1431 iterative LSI clustering^{12,18}. Briefly, we created 1-kb windows tiled across the genome and
1432 determined whether each cell was accessible within each window (binary). Next, we identified the
1433 top 50,000 accessible windows across all samples (accounting for GC bias) and performed an LSI
1434 dimensionality reduction (TF-IDF transformation followed by Singular Value Decomposition
1435 SVD) on these windows followed by Harmony batch correction⁸¹. We then performed Seurat⁸²
1436 clustering (FindClusters v2.3) on the harmonized LSI dimensions at a resolution of 0.8, 0.4 and
1437 0.2, keeping the clustering for which the minimum cluster size was greater than 100 cells (0.2 if
1438 this condition is not met). For each cluster, we called peaks on the Tn5-corrected insertions (each
1439 end of the Tn5-corrected fragments) using the MACS2 callpeak command with parameters ‘--shift
1440 -75 --extsize 150 --nomodel --call-summits --nolambda --keep-dup all -q 0.05’. The peak summits
1441 were then extended by 250 bp on either side to a final width of 501 bp, filtered by the ENCODE
1442 hg38 blacklist ([https://www.encodeproject.org/ annotations/ENCSR636HFF/](https://www.encodeproject.org/annotations/ENCSR636HFF/)), and filtered to
1443 remove peaks that extend beyond the ends of chromosomes. We then created a non-overlapping
1444 set of extended summits across all of these peaks as described previously^{12,18}.

1445 We then counted the accessibility for each cell in these peak regions to create an
1446 accessibility matrix. We then adopted the iterative LSI clustering approach^{12,18} to unbiasedly
1447 identify clusters that are due to biological vs technical variation. Briefly, we computed the TF-IDF
1448 transformation as described by Cusanovich et. al.⁸³. To do this, we divided each index by the
1449 colSums of the matrix to compute the cell “term frequency”. Next, we multiplied these values by
1450 log(1 + ncol(matrix)/rowSums(matrix)), which represents the “inverse document frequency”. This
1451 yields a TF-IDF matrix that can be used as input to irlba’s SVD implementation in R. We then
1452 used Harmony to batch correct the LSI dimensions in R. Using the first 25 reduced dimensions as
1453 input into a Seurat object, crude clusters were identified using Seurat’s (v2.3) SNN graph
1454 clustering FindClusters function with a resolution of 0.2. We then calculated the cluster sums from
1455 the binarized accessibility matrix and then log-normalized using edgeR’s ‘cpm(matrix,
1456 log = TRUE, prior.count = 3)’ in R. Next, we identified the top 25,000 varying peaks across all
1457 clusters using ‘rowVars’ in R. This was done on the cluster log-normalized matrix rather than the
1458 sparse binary matrix because: (1) it reduced biases due to cluster cell sizes, and (2) it attenuated
1459 the mean-variability relationship by converting to log space with a scaled prior count. The 25,000
1460 variable peaks were then used to subset the sparse binarized accessibility matrix and recompute

1461 the TF-IDF transform. We used SVD on the TF-IDF matrix to generate a lower dimensional
1462 representation of the data by retaining the first 25 dimensions. We then used Harmony to batch
1463 correct the LSI dimensions in R. We then used these reduced dimensions as input into a Seurat
1464 object and crude clusters were identified using Seurat's (v.2.3) SNN graph clustering FindClusters
1465 function with a resolution of 0.6. This process was repeated a third time with a resolution of 1.0.
1466 Then, these same reduced dimensions were used as input to Seurat's 'RunUMAP' with default
1467 parameters and plotted in ggplot2 using R.

1468

1469 **Identification of clusters and cell types from scATAC-seq data**

1470 Different clusters and cell types were manually identified using promoter accessibility and gene
1471 activity scores for various lineage-defining genes. Microglia (Cluster 24) were identified based on
1472 accessibility near the *IBA1*, *CD14*, *CD11C*, *PTGS1*, and *PTGS2* genes. Astrocytes (Clusters 13-
1473 17) were identified based on accessibility near the *GFAP* and *FGFR3* genes. Excitatory neurons
1474 (Clusters 1, 3, and 4) were identified based on accessibility near the *SLC17A6* and *SLC17A7* genes.
1475 Inhibitory neurons (Cluster 2, 11, and 12) were identified based on accessibility near the *GAD2*
1476 and *SLC32A1* genes. Medium spiny neurons (most of Cluster 2) were identified based on
1477 accessibility near the *DARPP32* gene. Oligodendrocytes (Clusters 19-23) were identified based on
1478 accessibility near the *MAG* and *SOX10* genes. OPCs (Clusters 8-10) were identified based on
1479 accessibility near the *PDGFRA* gene. All neuronal subsets, for example nigral neurons (Cluster 5-
1480 6), were identified primarily as neurons based on accessibility near the *NEFL*, *RBFOX3*, *VGF*, and
1481 *GRIN1* genes and then subdivided based on the region of origin and the accessibility near other
1482 genes mentioned above.

1483

1484 **Single-cell ATAC-seq peak calling**

1485 For scATAC-seq peak calling from clusters or manually defined cell types, all single cells
1486 belonging to the given group were pooled together. These pooled fragment files were converted to
1487 the paired-end tagAlign format and processed with version 1.4.2 of the ENCODE DCC ATAC-
1488 seq pipeline. The conversion to tagAlign was performed as follows. For fragments on the positive
1489 strand, the read start coordinate was the fragment start coordinate, zero-indexed. The read end
1490 coordinate was the fragment start coordinate plus the read length (99 bp). For fragments on the
1491 negative strand, the read start coordinate was the fragment end coordinate, zero-indexed. The read
1492 start coordinate was the fragment end coordinate minus the read length (99 bp). Then, these
1493 tagAlign files were used as input to the DCC ATAC-seq pipeline. IDR optimal peak sets with an
1494 IDR threshold of 0.05 were determined for each cluster by the pipeline, using pseudo-bulk
1495 replicate tagAligns for the cluster. Other pipeline parameters were the same as for bulk ATAC-seq
1496 data (see above).

1497

1498 **Single-cell ATAC-seq gene activity scores**

1499 We calculated gene activity scores by summing the binarized accessibility, weighted by distance,
1500 in the 1-kb tiles within 100 kb. The distance weights were computed by determining the distance

1501 from the tile to the gene promoter start site and computing “ $\exp(-\text{abs}(\text{distance})/10000)$ ”. These
1502 were then scaled to 10,000 and log-normalized with a pseudo count of 1. For visualization
1503 purposes, the top and bottom 2.5% of scores were thresholded.
1504

1505 **Single-cell ATAC-seq pseudo-bulk replicate generation and differential accessibility 1506 comparisons**

1507 For differential comparisons of clusters or cell types, including Pearson correlation determination,
1508 non-overlapping pseudo-bulk replicates were generated from groups of cells. For each cell
1509 grouping (i.e a cluster or a cell type), a minimum of 300 cells was required in order to make at
1510 least two non-overlapping pseudo-bulk replicates of 150 cells each. A maximum of 3 pseudo-bulk
1511 replicates was made per group if the total number of cells per group was greater than 450 cells.
1512 Cells were randomly deposited into one of the pseudo-bulk replicates and all available cells were
1513 used. In this way, the non-overlapping pseudo-bulk replicates are agnostic to which donor the cell
1514 came from but aware of individual cells (i.e. all reads from a given cell are deposited into the same
1515 pseudo-bulk replicate). These pseudo-bulk replicates were then used for differential comparisons
1516 using DESeq2⁸⁴.
1517

1518 **CIBERSORT deconvolution**

1519 CIBERSORT²⁵ was used to deconvolve bulk ATAC-seq data using signature matrices generated
1520 from scATAC-seq data. Default parameters were used. For the cell type-specific classifier, pseudo-
1521 bulk replicates were generated for each of the 8 main cell types. For the cluster-specific classifier,
1522 pseudo-bulk replicates were generated for each of the 24 clusters.
1523

1524 **Transcription factor footprinting**

1525 Transcription factor footprinting was performed as described previously⁶⁴.
1526

1527 **HiChIP library generation**

1528 HiChIP library generation was performed as described previously¹³. One million cryopreserved
1529 nuclei were used per experiment. Enzyme MboI was used for restriction digest. Sonication was
1530 performed on a Covaris E220 instrument using the following settings: duty cycle 5, peak incident
1531 power 140, cycles per burst 200, time 4 minutes. All HiChIP was performed using H3K27ac as
1532 the target (Abcam ab4729).
1533

1534 **HiChIP data analysis**

1535 HiChIP paired-end sequencing data was processed using HiC-Pro⁸⁵ version 2.11.0 with a
1536 minimum mapping quality of 10. FitHiChIP⁸⁶ was used to identify “peak-to-all” interactions using
1537 peaks called from the one-dimensional HiChIP data. A lower distance threshold of 20 kb and an
1538 upper distance threshold of 2 Mb were used. Bias correction was performed using coverage-
1539 specific bias.
1540

1541 **HiChIP linkage of SNPs to genes**

1542 To link SNPs to genes, we identified FitHiChIP loops that contained a SNP in one anchor and a
1543 TSS in the other anchor. This was performed for all LD-expanded SNPs to identify the full
1544 complement of genes that could be putatively implicated in AD and PD.

1545

1546 **gkm-SVM machine learning classifier training and testing**

1547 For each of the 24 scATAC-seq clusters, we used a 10 fold cross-validation scheme to train
1548 weighted gapped k-mer Support Vector Machine (gkm-SVM) models to classify 1000 bp
1549 sequences into two classes - accessible (corresponding to sequences underlying peaks) and
1550 inaccessible (GC matched inaccessible genomic regions). The test sets for each of the 10 folds are
1551 as follows. Fold 0 consisted of chr 1. Fold 1 consisted of chr 2 and chr 19. Fold 2 consisted of chr
1552 3 and chr 20. Fold 3 consisted of chr 6, chr 13, and chr 22. Fold 4 consisted of chr 5, chr 16, and
1553 chr Y. Fold 5 consisted of chr 4, chr 15, and chr 21. Fold 6 consisted of chr 7, chr 14, and chr 18.
1554 Fold 7 consisted of chr 11, chr 17, and chr X. Fold 8 consisted of chr 9 and chr 12. Fold 9 consisted
1555 of chr 8 and chr 10.

1556 For each of the 24 scATAC-seq clusters, we merged the IDR peaks with identical genomic
1557 coordinates (peaks with multiple summits) while preserving the summit position and the MACS2
1558 p-value of the peak with the lowest p-value among the ones with the identical coordinates. Next,
1559 we ranked the peaks by the MACS2 p-value, expanded each peak by 500 bp on either side of the
1560 summit, to a total of 1000 bp, and eliminated those peaks with any 'N' bases in the 1000 bp. For
1561 each of 10 cross-validation folds, we kept up to 60,000 of the top peaks belonging to the training
1562 set and all of the peaks belonging to the much smaller test set, all of which comprised the positively
1563 labeled (accessible) examples for training.

1564 In order to generate the negative (inaccessible) examples for each of the cross-validation
1565 folds in each single-cell cluster, first, we used seqdataloader
1566 (<https://github.com/kundajelab/seqdataloader>) to generate all 1000 bp sequences obtained by tiling
1567 the hg38 genome 200 bp at a time, with a stride of 50 bp, keeping those 200 bp segments that have
1568 no IDR peak summits in that cluster, and then expanding those 200 bp segments by 400 bp on each
1569 side for a total of 1000 bp. Next, we calculated the GC content of the selected positive examples
1570 and all of the negative sequences. We matched each of the positive examples, both in the training
1571 set and the test set, with a negative sequence with the closest GC content, without replacement.

1572 For each of the 10 folds in each of the 24 clusters, we used the 1000-bp DNA sequences
1573 corresponding to the positive and GC-matched negative training examples as inputs to the
1574 gkmtrain function from the LS-GKM package⁸⁷ with the default options, producing a total of 240
1575 models; the default options for LS-GKM included the gapped *k*-mer + center weighted (wgkm)
1576 kernel (*t* = 4), a word length of 11 (*l* = 11), 7 informative columns (*k* = 7), 3 maximum mismatches
1577 to consider (*d* = 3), an initial value of the exponential decay function of 50 (*M* = 50), a half-life
1578 parameter of 50 (*H* = 50), a regularization parameter of 1.0 (*c* = 1.0), and a precision parameter of
1579 0.001 (*e* = 0.001). We used the resulting support vectors for each trained model to score the DNA
1580 sequences corresponding to the positive and GC-matched negative test set examples for each fold

1581 in each cluster by running gkmpredict, and used the scikit-learn python library⁸⁸ to calculate both
1582 auROC and auPRC accuracy metrics.

1583

1584 **gkm-SVM allelic scores of candidate SNPs**

1585 We intersected the coordinates of all LD-expanded candidate AD and PD GWAS and
1586 colocalization SNPs with those of the peaks for each single-cell ATAC-seq cluster to obtain the
1587 SNPs in each cluster that are in peaks. For each SNP in a peak in each of the clusters, we retrieved
1588 the 1000 bp DNA sequence around the SNP, with the SNP at its center, and created a sequence
1589 corresponding to the effect allele by replacing the 500th position of the sequence with the effect
1590 allele. Similarly, we created another sequence corresponding to the non-effect allele by replacing
1591 the 500th position of the sequence with the non-effect allele. Furthermore, we repeated the same
1592 procedure to also produce 50 bp sequences for each SNP with the effect allele and the non-effect
1593 allele by retrieving the 50 bp DNA sequence around each SNP and replacing the 25th position
1594 with the effect and the non-effect allele, respectively.

1595 For each SNP in a peak in each of the clusters, we computed **GkmExplain**³⁷ importance
1596 scores for each position in each of the 1000 bp effect and non-effect allele sequences using each
1597 of the 10 gkm-SVM³⁶ models for the respective cluster. GkmExplain is a method to infer the
1598 importance or predictive contribution of every base in an input sequence to its corresponding
1599 output prediction from a gkm-SVM model. Next, for each SNP in a given cluster, we computed
1600 the average score for each position across all 10 models (from the 10 folds) for that cluster for both
1601 the effect allele sequence and the non-effect allele sequence, producing a set of consensus
1602 importance scores for both the effect allele and the non-effect allele. Then, we subtracted the sum
1603 of these consensus importance scores corresponding to the central 50 bp of the non-effect allele
1604 sequence from that of the effect allele sequence to compute the GkmExplain score for each SNP
1605 in each cluster.

1606 To compute ***in silico* mutagenesis (ISM)** scores for each SNP in a peak in each of the
1607 clusters, we used each of the 10 fold gkm-SVM models from the respective cluster to compute
1608 model output prediction scores for the 50 bp effect and non-effect allele sequences by running
1609 gkmredict. Then, we subtracted the score of the non-effect allele sequence from that of the effect
1610 allele sequence to obtain the ISM score and computed the average ISM score for each SNP across
1611 all 10 folds in each cluster.

1612 To compute **deltaSVM** scores, we generated all possible non-redundant k-mers of size 11
1613 and scored each of them using each of the 240 models. Next, for each SNP in a peak in each of the
1614 clusters, we used each of the 10 sets of *k*-mer scores from the gkm-SVM models from the
1615 respective cluster to run deltaSVM³⁹ on the 50 bp effect and non-effect allele sequences. We
1616 computed the average of the resulting deltaSVM scores for each SNP across all 10 folds in each
1617 cluster.

1618

1619 **Statistical significance and high confidence sets of gkm-SVM based allelic scores for**
1620 **candidate SNPs**

1621 In order to obtain a statistical significance for each of the three gkm-SVM model based allelic SNP
1622 scores (GkmExplain, ISM and deltaSVM), we computed an empirical null distribution of scores.
1623 We expect most of the LD expanded candidate SNPs to be non functional. Hence, we simply use
1624 the distribution of the scores for all candidate SNPs as an empirical null distribution. For each type
1625 of score, in order to control for any arbitrary bias in the sign of the score, we included the negative
1626 value of each score to the list of scores to enforce symmetry. We found that the t-distribution was
1627 a good fit (based on KS test) to the empirical null distribution for all three scores. Hence, we used
1628 the fitted t-distributions (using SciPy python library <http://www.scipy.org/>) to each of the three
1629 sets of scores as the null distributions.

1630 To select SNPs with **statistically significant gkm-SVM allelic scores**, for each cluster,
1631 we selected those SNPs that fall outside the 95% confidence interval for all three null *t*-
1632 distributions fitted to the GkmExplain, ISM, and deltaSVM scores.

1633 Next, we developed a method to identify putative transcription factor binding sites around
1634 each gkm-SVM scored statistically significant candidate SNP, by identifying the subsequences
1635 around the SNP whose base-resolution importance scores are significantly above background. For
1636 each SNP, we defined the **active allele** as the allele for which the 50 bp sequence centered on the
1637 SNP has the higher gkmpredict output score (relative to the other allele) from the gkm-SVM
1638 model. We fitted a background null *t*-distribution to the consensus GkmExplain importance scores
1639 (averaged across models for all 10 folds) of all bases in the 200 bp sequence centered on the SNP
1640 and containing the active allele. We use this null distribution to identify bases around the SNP with
1641 high signal-to-noise ratio. Specifically, starting from the center of the positive allele's sequence,
1642 which is the location of the SNP, we continue advancing one pointer upstream and another
1643 downstream, each up to the position beyond which lie two consecutive bases that both have
1644 consensus importance scores that are within or lower than the 90% confidence interval for the
1645 distribution fitted to the consensus importance scores for that sequence. The subsequence between
1646 the terminal positions of the two pointers corresponds to one that underlies a series of bases with
1647 high GkmExplain importance scores that are significantly above scores of surrounding background
1648 sequence and potentially contains transcription factor binding sites and motifs that are relevant for
1649 the given cluster. We refer to these high-importance subsequences seqlets.

1650 Next, we defined two additional scores (prominence score and magnitude score) to further
1651 identify high confidence candidates from the gkm-SVM scored statistically significant candidate
1652 SNPs supported by seqlets that could potentially match identifiable transcription factor binding
1653 sites. We compute the sum of the non-negative consensus importance scores from the active
1654 allele's seqlet, which we refer to as the **active seqlet score**, and divide that score by the sum of the
1655 non-negative consensus importance scores from the entire central 200-bp region of the active
1656 allele's sequence; we refer to this ratio as the **active seqlet signal-to-noise ratio**. Similarly, we
1657 compute the **inactive seqlet score** as the sum of the non-negative consensus importance scores in
1658 the inactive allele's sequence from the same positions overlapping the active seqlet. We obtain a
1659 corresponding **inactive seqlet signal-to-noise ratio** by dividing the inactive seqlet score by the
1660 sum of the non-negative consensus importance scores from the entire central 200-bp region of the

1661 inactive allele's sequence. Then, for each SNP, we compute the **prominence score** by subtracting
1662 the non-effect allele's seqlet signal-to-noise ratio from the effect allele's seqlet signal-to-noise
1663 ratio. In addition, we also compute a **magnitude score** by subtracting the non-effect allele's seqlet
1664 score from the effect allele's seqlet score.

1665 To compute the statistical significance of the prominence and magnitude scores for
1666 candidate SNPs, for each cluster, we fit null *t*-distributions to the prominence scores and magnitude
1667 scores (using a KS test to test goodness of fit of the *t*-distribution to the empirical distribution of
1668 scores). For each type of score, in order to control for any arbitrary bias in the sign of the score,
1669 we include the negative value of each score to the list of scores to enforce symmetry before fitting
1670 the distribution.

1671 Finally, to prioritize SNPs that disrupt potential transcription factor binding sites, in each
1672 cluster, among the SNPs with statistically significant gkm-SVM allelic scores, we designate as
1673 high confidence SNPs those that have prominence scores outside the 95% confidence interval for
1674 the distribution fitted to the prominence scores. These are the SNPs that have an allele that
1675 completely destroys a prominent and high-scoring seqlet and, as a result, potentially disrupts an
1676 important transcription factor binding site. Next, among the confident SNPs that do not pass the
1677 high confidence threshold, we designated as medium confidence SNPs those that have either peak
1678 magnitude scores outside the 95% confidence interval or prominence scores outside the 80%
1679 confidence interval. The magnitude threshold is intended to capture those SNPs that have a
1680 significant deleterious effect on the seqlet score, even if those SNPs do not necessarily destroy the
1681 entire seqlet and even for cases where the seqlet around the SNP is not among the most prominent
1682 seqlets in the local 200 bp sequence window. In addition, the relaxed prominence threshold is
1683 intended to capture those SNPs that do not pass the stringent filter for the high confidence set, but
1684 nevertheless, demonstrate at least a partial deleterious effect on a moderately scoring seqlet around
1685 the SNP. Together, these two filters serve to increase the recall in the prioritization of the SNPs,
1686 allowing us to identify all promising SNPs that are worthy of in-depth evaluation, which can assess
1687 their potential regulatory effect through a case-by-case analysis. The remaining SNPs in the
1688 confident set, which fail to meet the threshold set for medium confidence, are designated as low
1689 confidence SNPs, as they include SNPs that significantly reduce the GkmExplain score, the ISM
1690 score, and the deltaSVM score, but do not have a clear impact on a seqlet around the SNP, making
1691 it unlikely for them to have a disruptive effect on a key transcription factor binding site.
1692

1693 **Identification of MAPT haplotypes**

1694 The MAPT haplotype block is part of one of the largest LD blocks in the human genome. To
1695 identify SNPs that belong exclusively to either the H1 or H2 haplotype, we used minor allele
1696 frequencies from dbSNP version 151. SNPs were required to be within the coordinates of the
1697 MAPT inversion breakpoints (hg38 chr17:45551578-46494237) and to have a minor allele
1698 frequency between 8.4% and 9%. While there are undoubtedly haplotype specific SNPs outside
1699 this frequency range, we chose this range to be as conservative as possible and to pick SNPs that
1700 showed minimal haplotype switching. Each SNP was verified to track with the predicted haplotype

1701 using LDLink⁸⁹. This resulted in 2366 SNPs that could be confidently called as haplotype
1702 divergent.
1703

1704 **MAPT locus differential expression analysis**

1705 A 900-kb block of variants in strong LD at the *MAPT* locus hampered the resolution of
1706 colocalization methods for identifying causal variants and/or genes at this locus. To probe this
1707 locus more deeply, we assembled a list of 2366 variants uniquely found in either the H1 or the H2
1708 haplotype of the *MAPT* locus (described above). For each of the 838 individuals genotyped in
1709 GTEx v8, we counted the number of variants in support of either haplotype. We designated
1710 individuals as homozygous if they possessed less than 1% of variants favoring the opposite
1711 haplotype and heterozygous if 45% to 55% of variants supported either haplotype. This determined
1712 the individual's haplotype in all but six cases, which were excluded from the remainder of the
1713 *MAPT* analysis. In total, we identified 539 individuals with the H1/H1 haplotype, 260 with H2/H1,
1714 and 33 with H2/H2. Our a priori gene of interest was *MAPT*, whose expression had
1715 previously been demonstrated to be higher in H1 than H2 haplotypes. At a nominal cutoff of $p <$
1716 0.05, we confirmed this expected direction of differential *MAPT* expression (higher in H1
1717 haplotypes) in multiple tissues, with the strongest contrasts in "Brain - Cortex".

1718 We then extended our analysis to include all genes expressed in any of the brain tissues
1719 from GTEx v8. We compared the log2-fold change of gene expression (TPM) between H1/H1 and
1720 H1/H2 individuals, given that these subgroups had the largest sample size. A change was
1721 considered statistically significant if a Wilcoxon rank-sum test between the two groups produced
1722 a p-value of $< 0.05 / (\text{total } \# \text{ genes}) / (\text{total } \# \text{ tissues})$. We also performed pairwise Wilcoxon rank-
1723 sum test comparisons for each gene in each brain tissue between all 3 pairings of haplotypes.
1724

1725 **MAPT haplotype-specific ATAC-seq and HiChIP analysis**

1726 For both ATAC-seq and HiChIP, reads from heterozygote donors were re-mapped to an N-masked
1727 genome (using bowtie2 or HiCPro, respectively) where all dbSNP v151 positions were masked to
1728 "N". After alignment, SNPs⁹⁰ was used to divide reads mapping to either the H1 or H2
1729 haplotypes based on the presence of one of the 2366 haplotype-divergent SNPs identified above.
1730 In this way, reads mapping to regions that lack a haplotype-divergent SNP could not be assigned
1731 in an allelic fashion to either the H1 or H2 haplotypes and were ignored. For track-based
1732 visualizations of haplotype-specific data, all available data from a given haplotype was merged
1733 agnostic to what brain region the data was derived from. To identify regions with haplotype-
1734 specific chromatin accessibility in the *MAPT* locus, the entire locus was tiled into non-overlapping
1735 500 bp bins and the number of Tn5 transposase insertions were counted for each haplotype in each
1736 bin for each sample. A Wilcoxon signed-rank test was used to determine if the difference between
1737 H1 and H2 for each bin was significant after multiple hypothesis correction (FDR < 0.01).

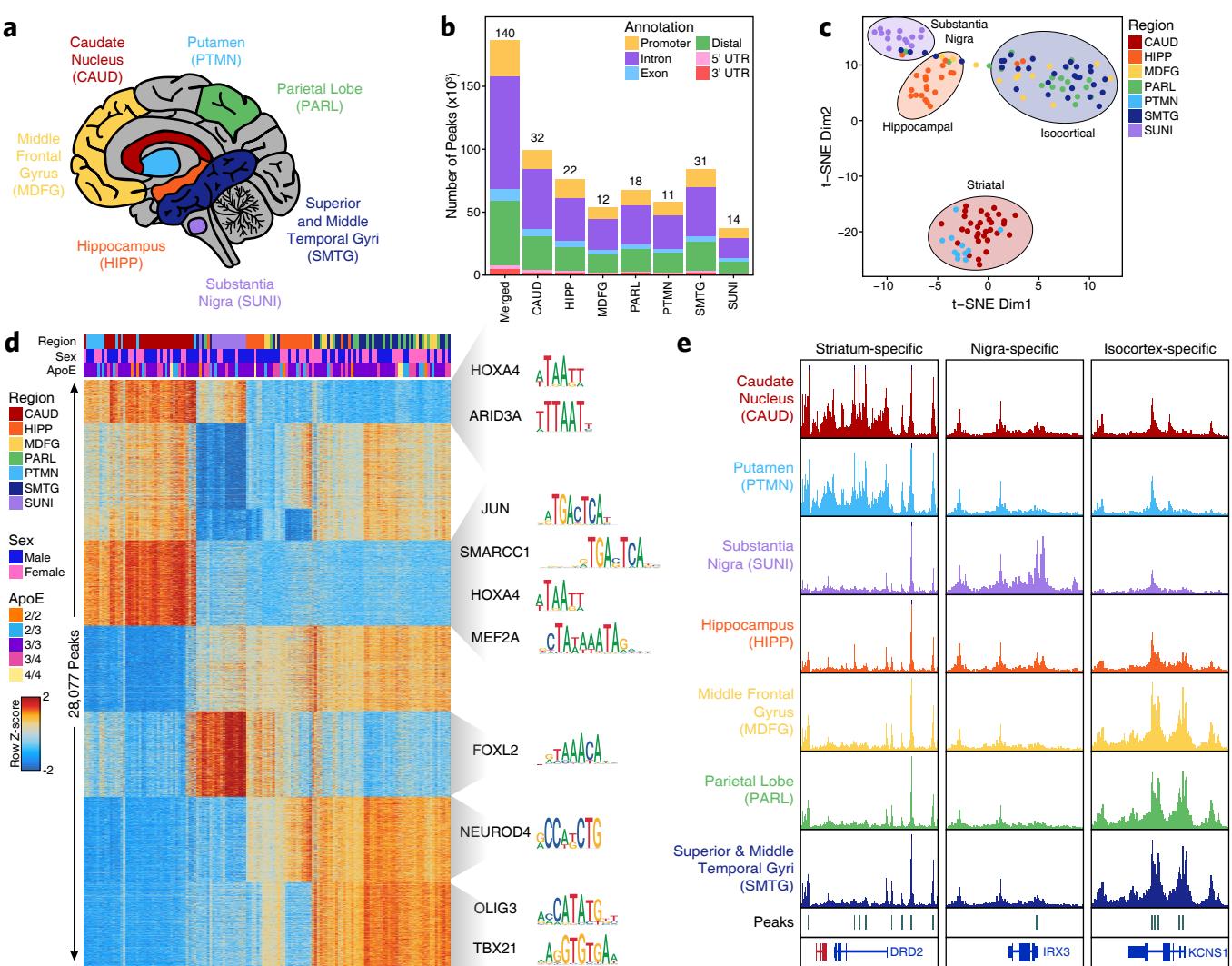


Figure 1

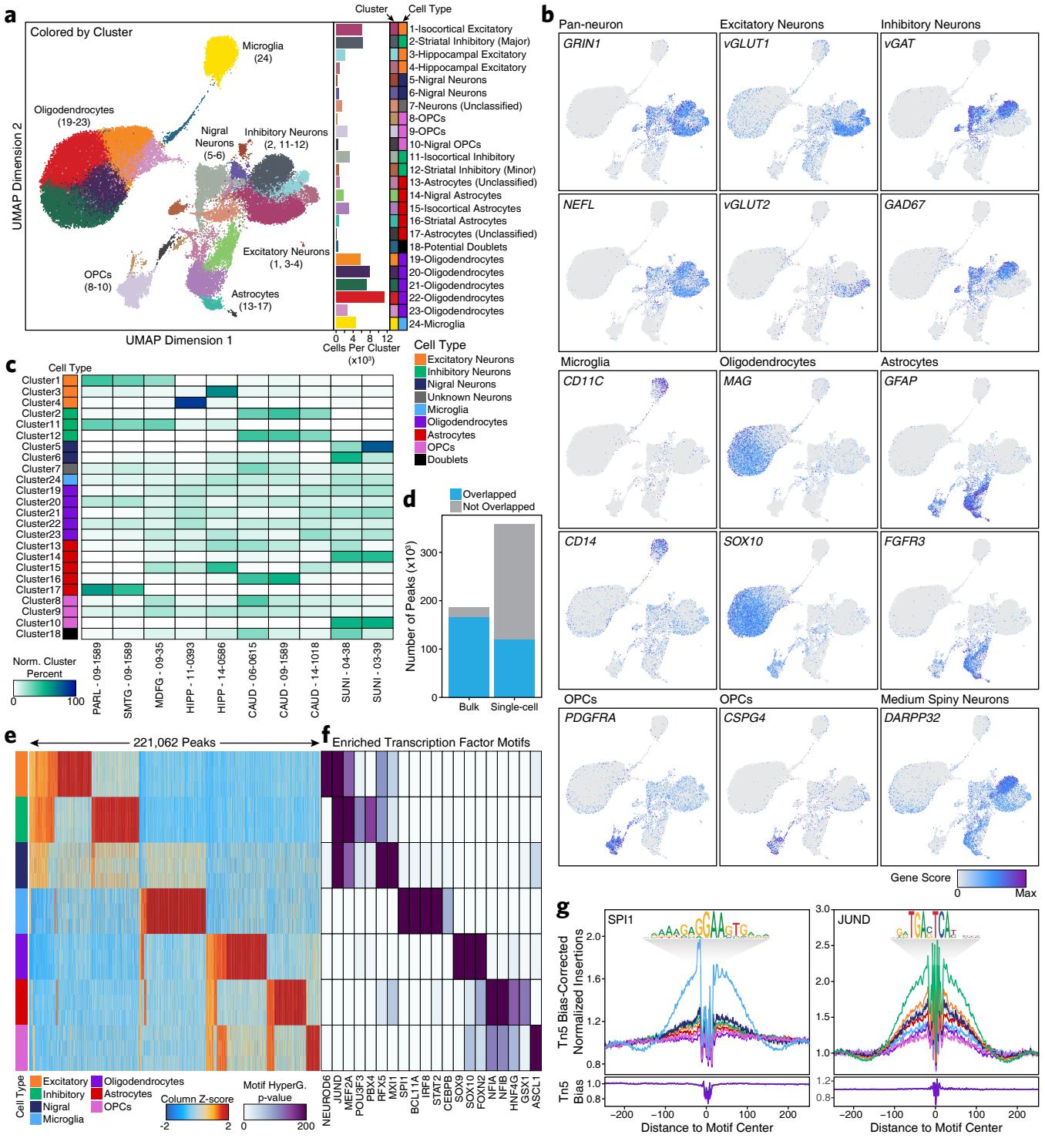


Figure 2

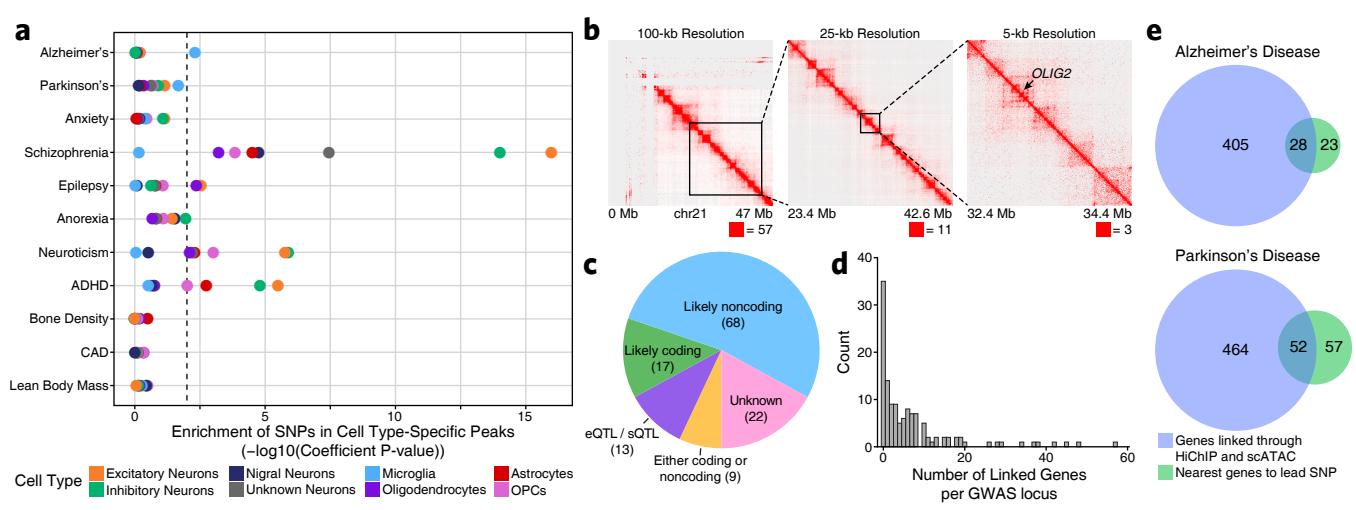


Figure 3

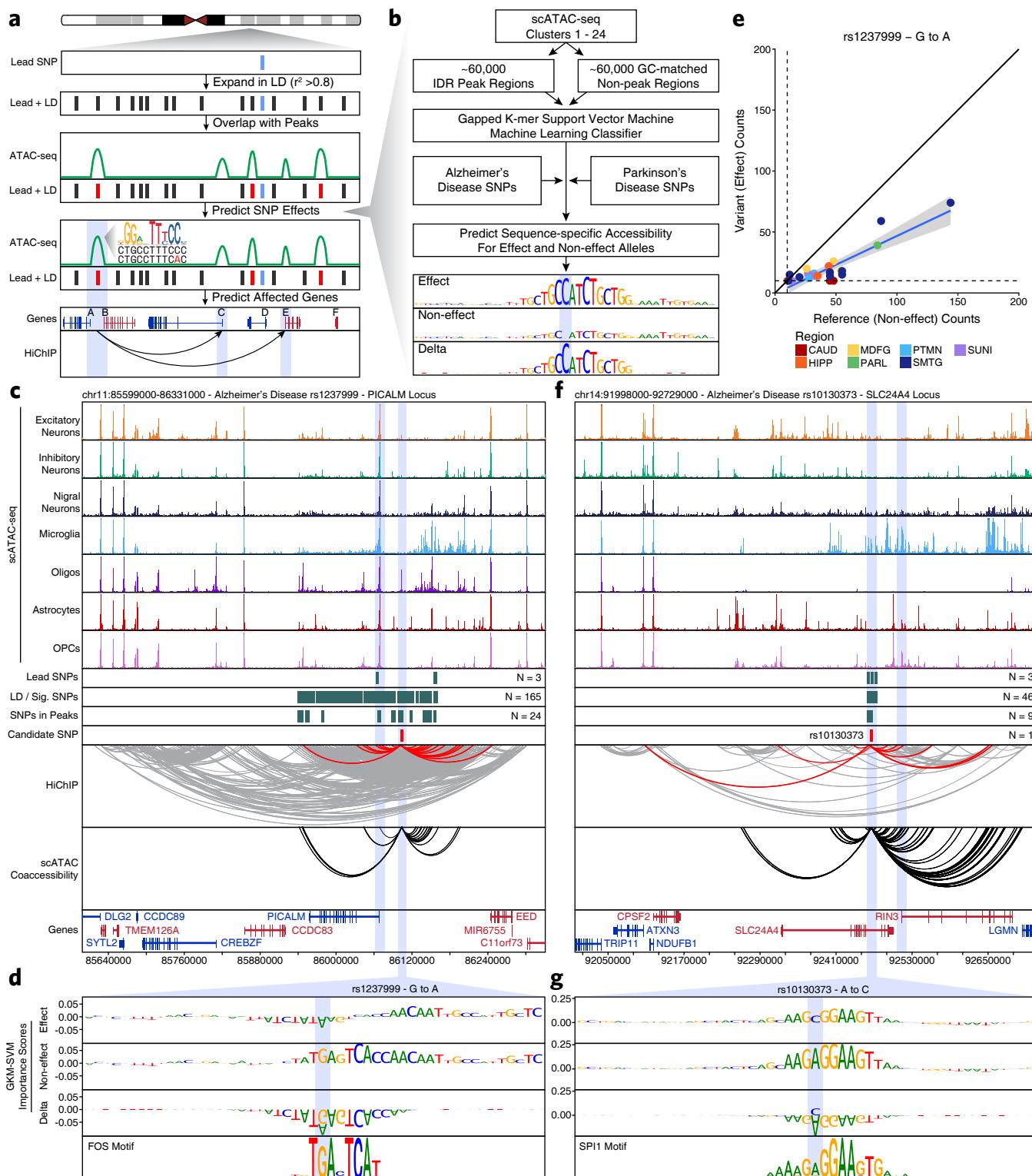


Figure 4

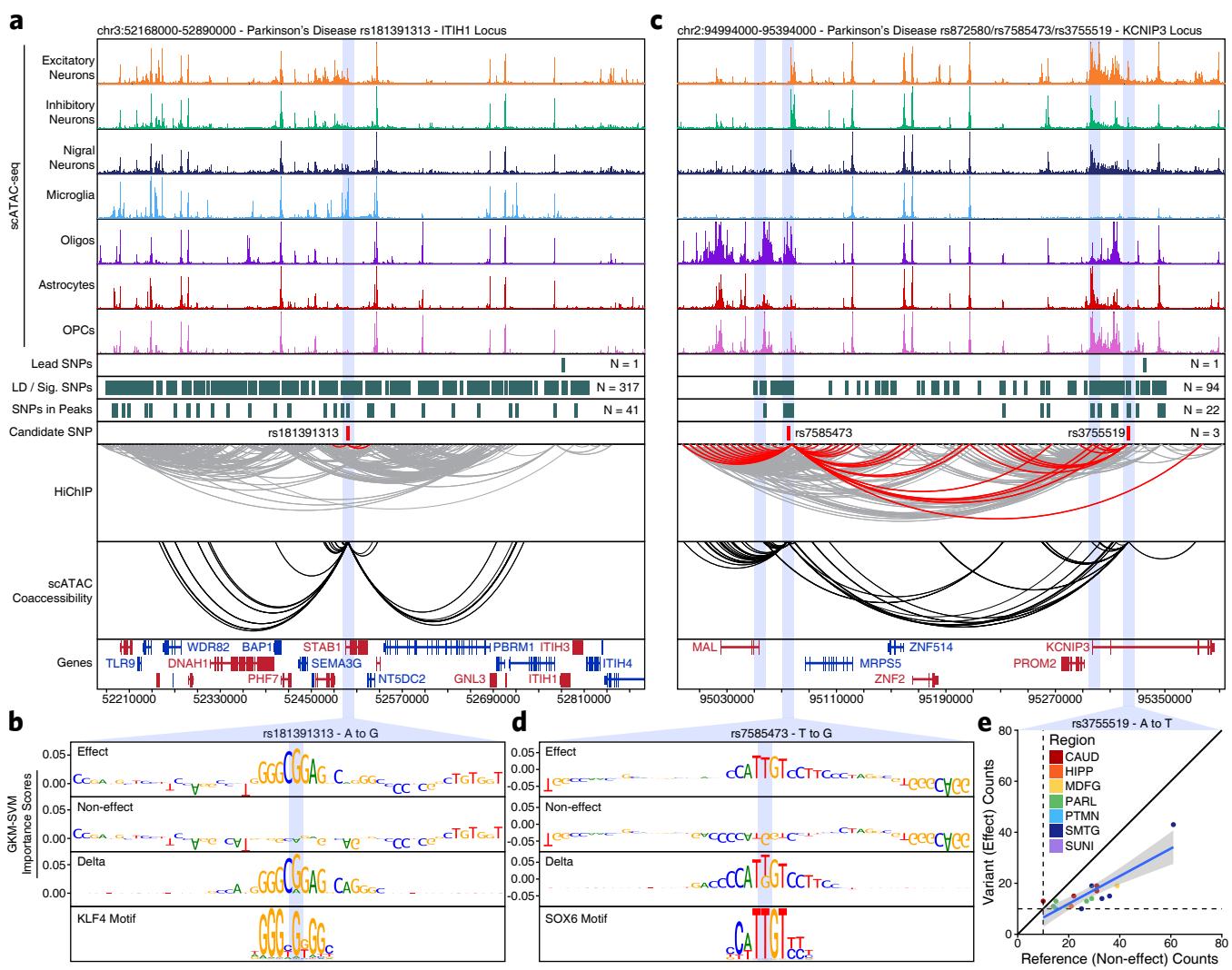


Figure 5

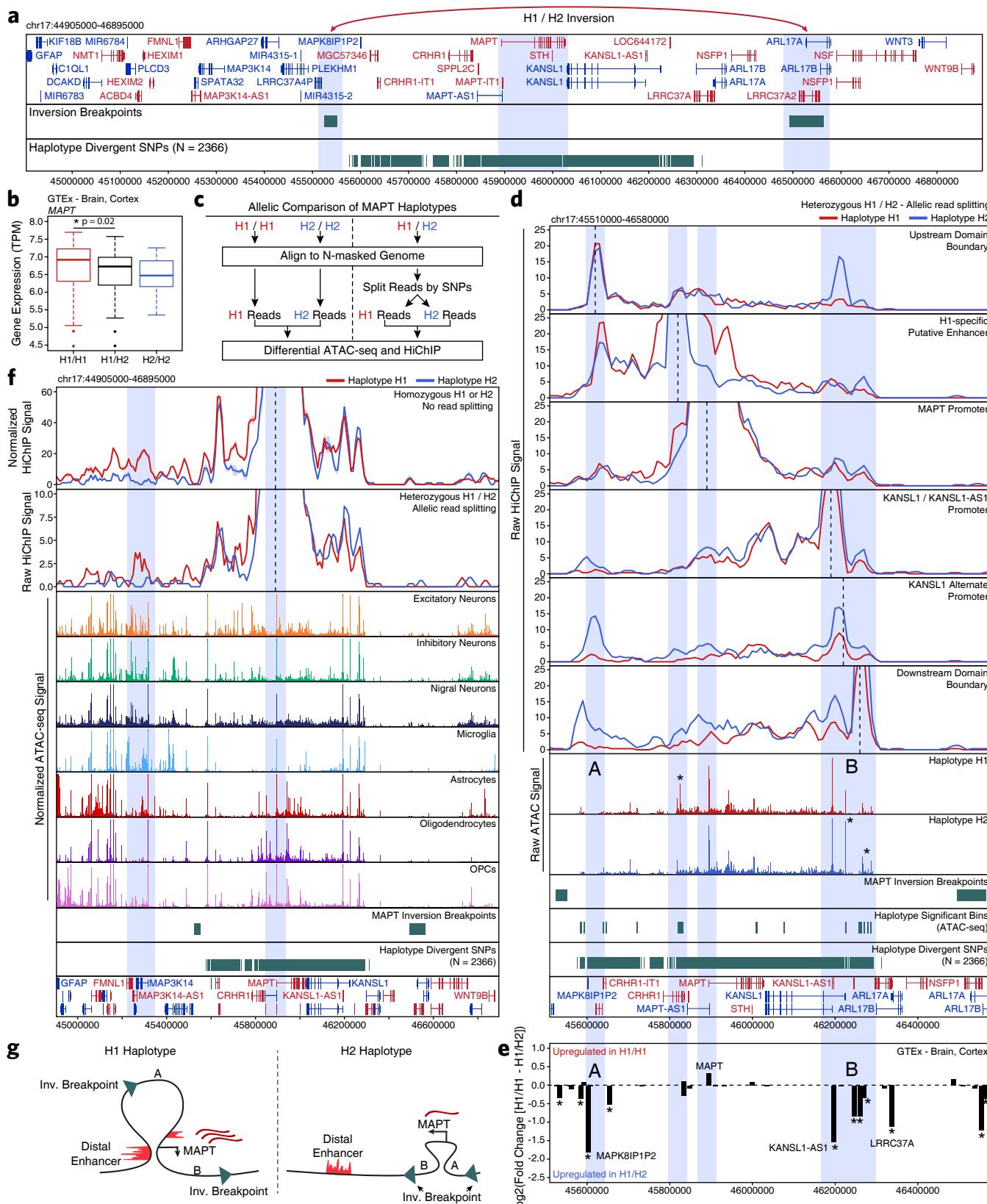
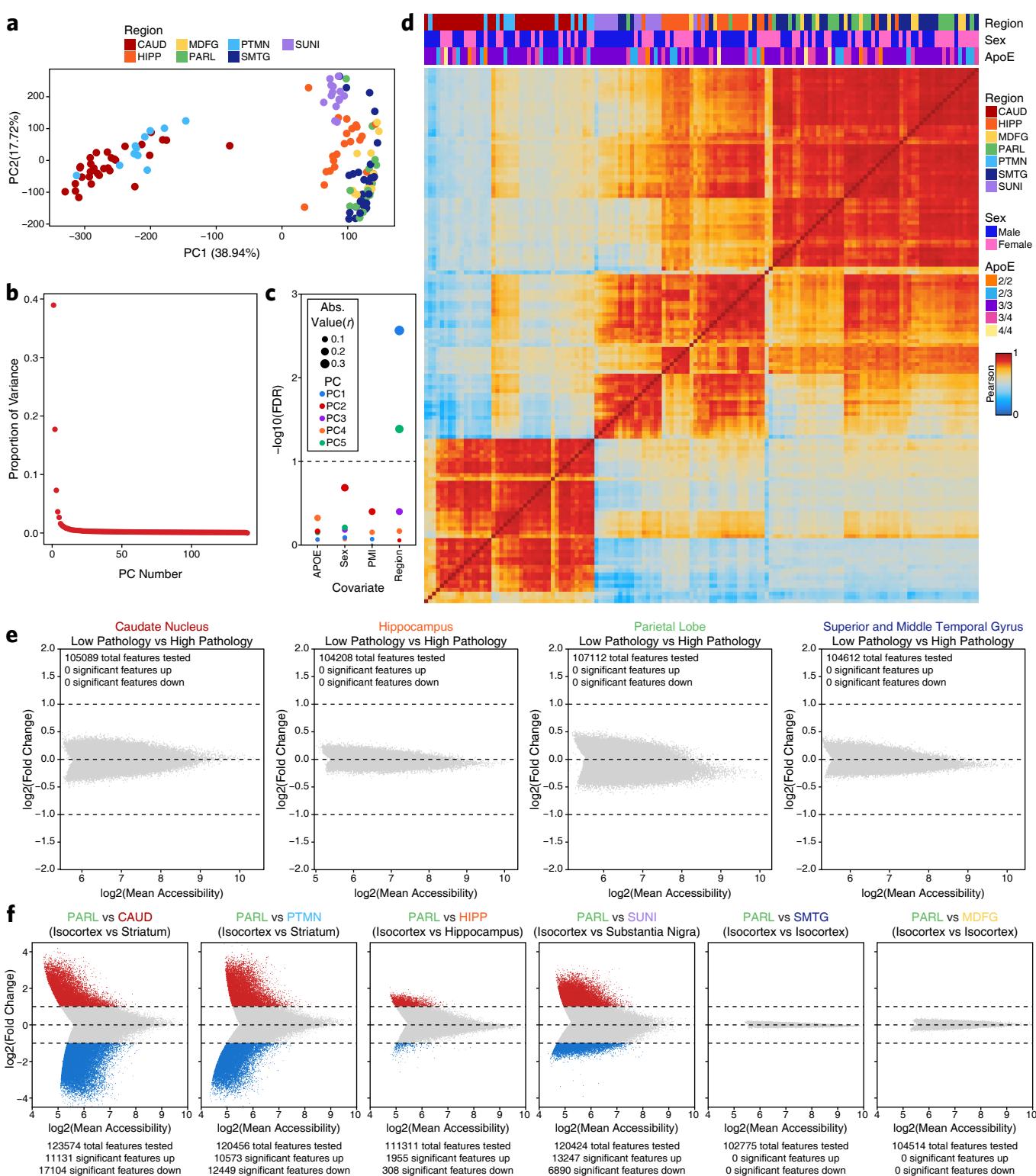
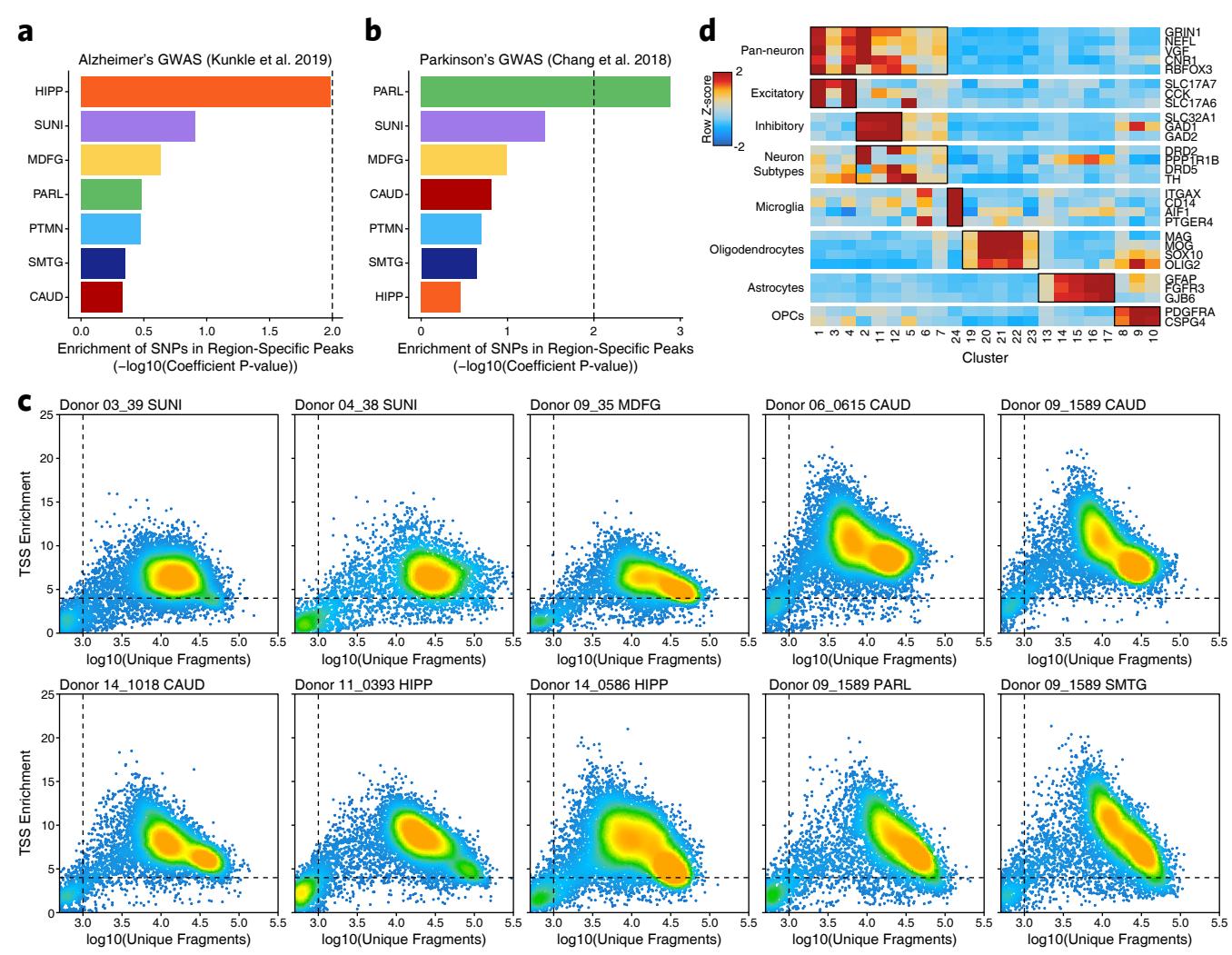


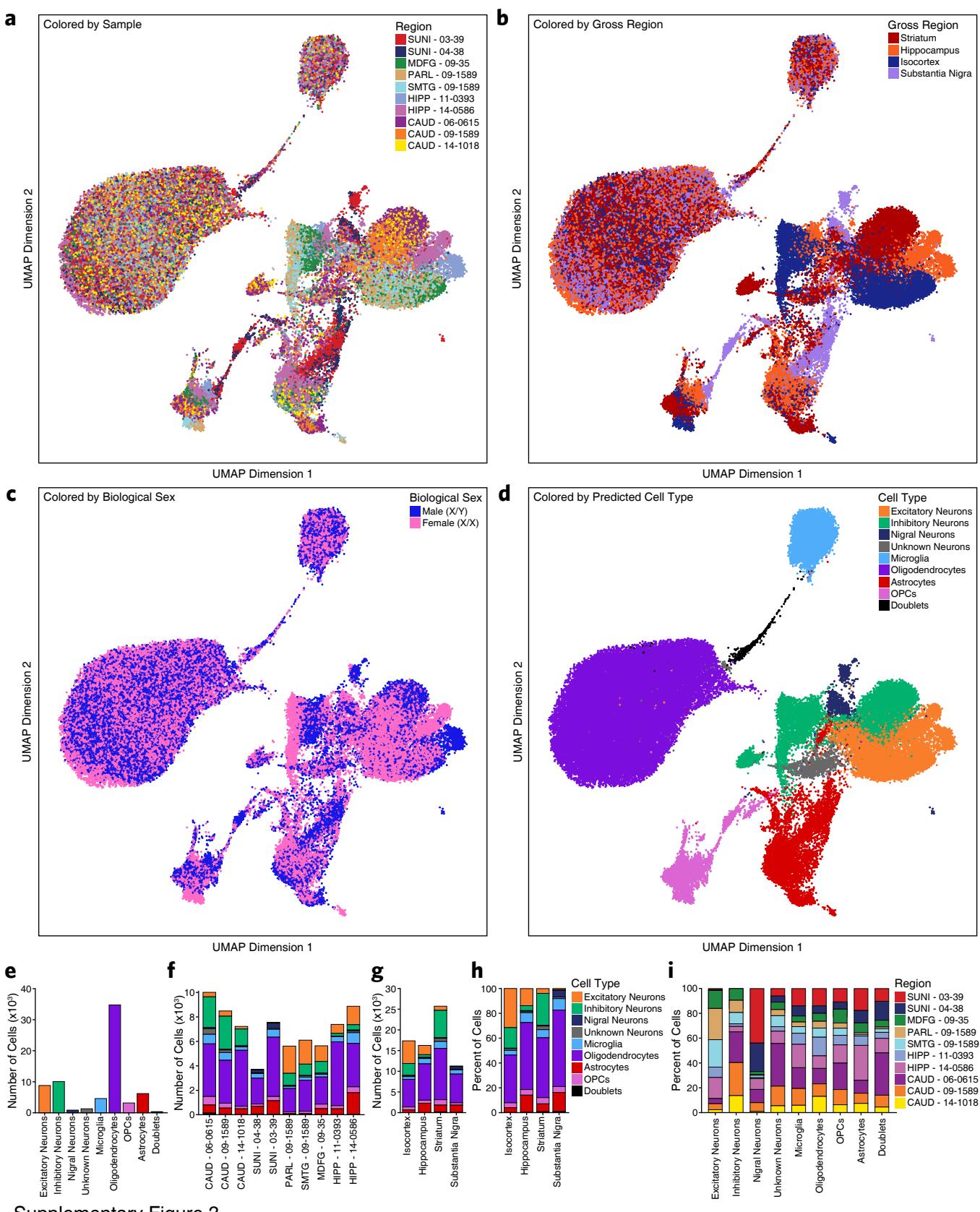
Figure 6



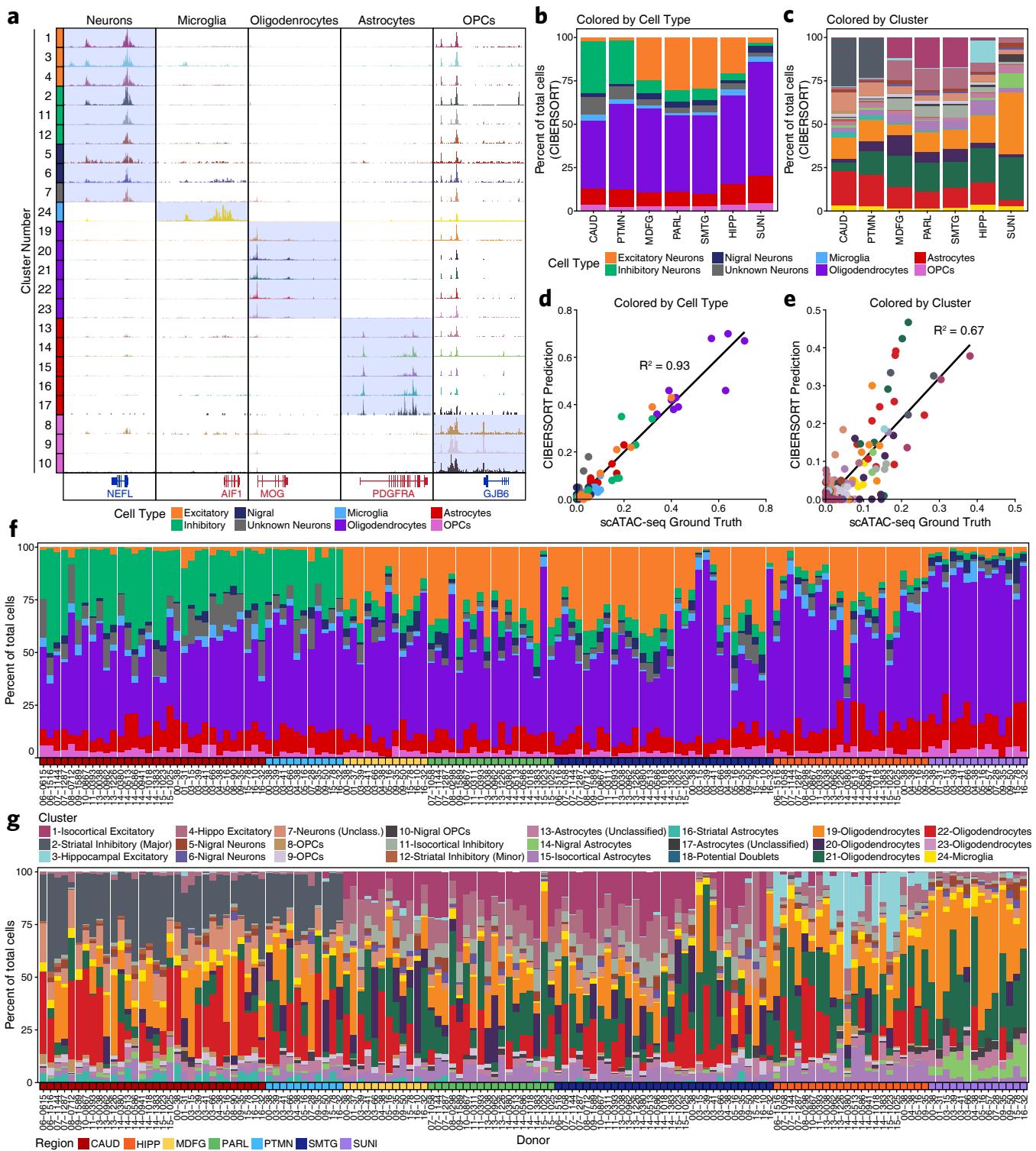
Supplementary Figure 1



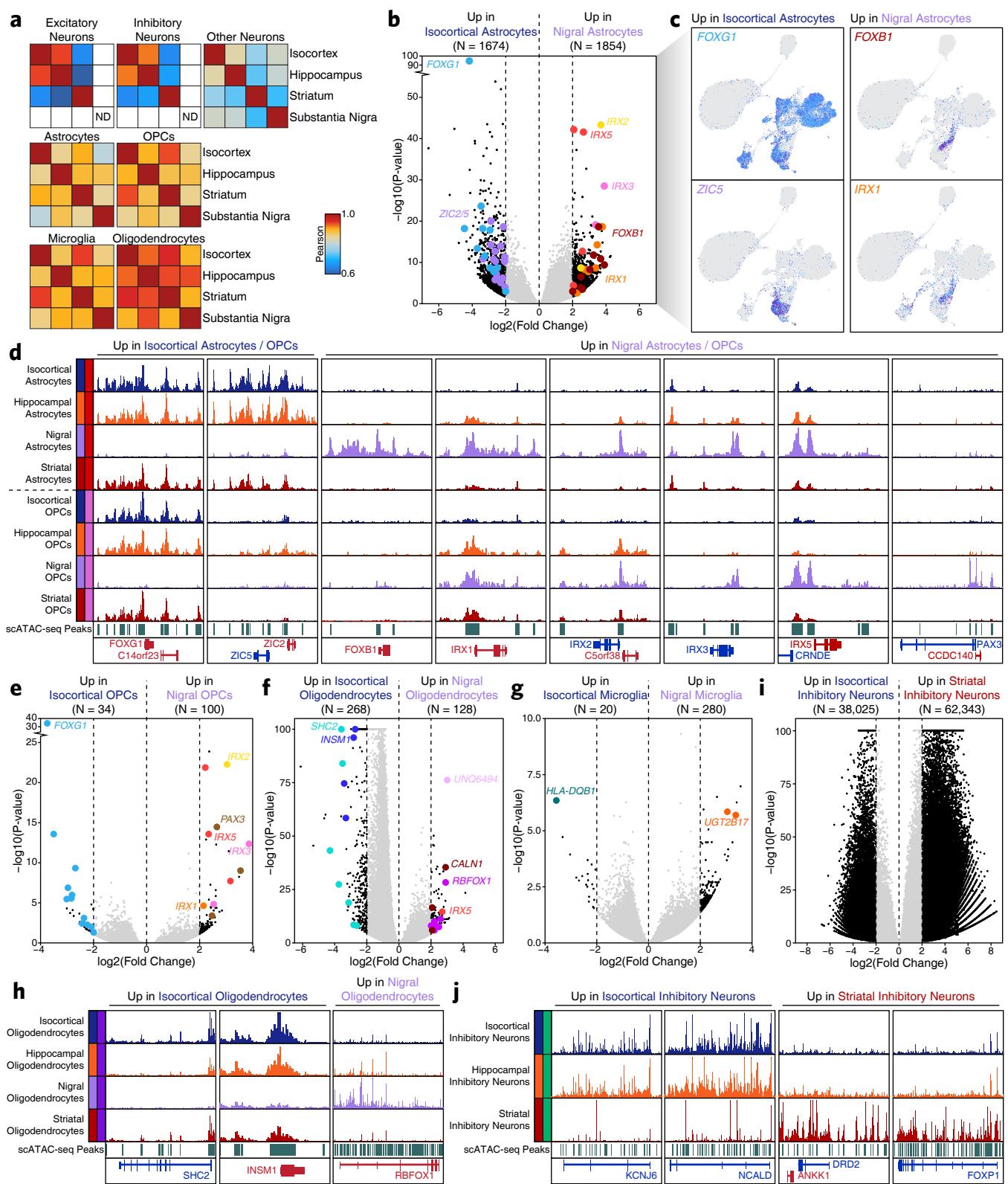
Supplementary Figure 2



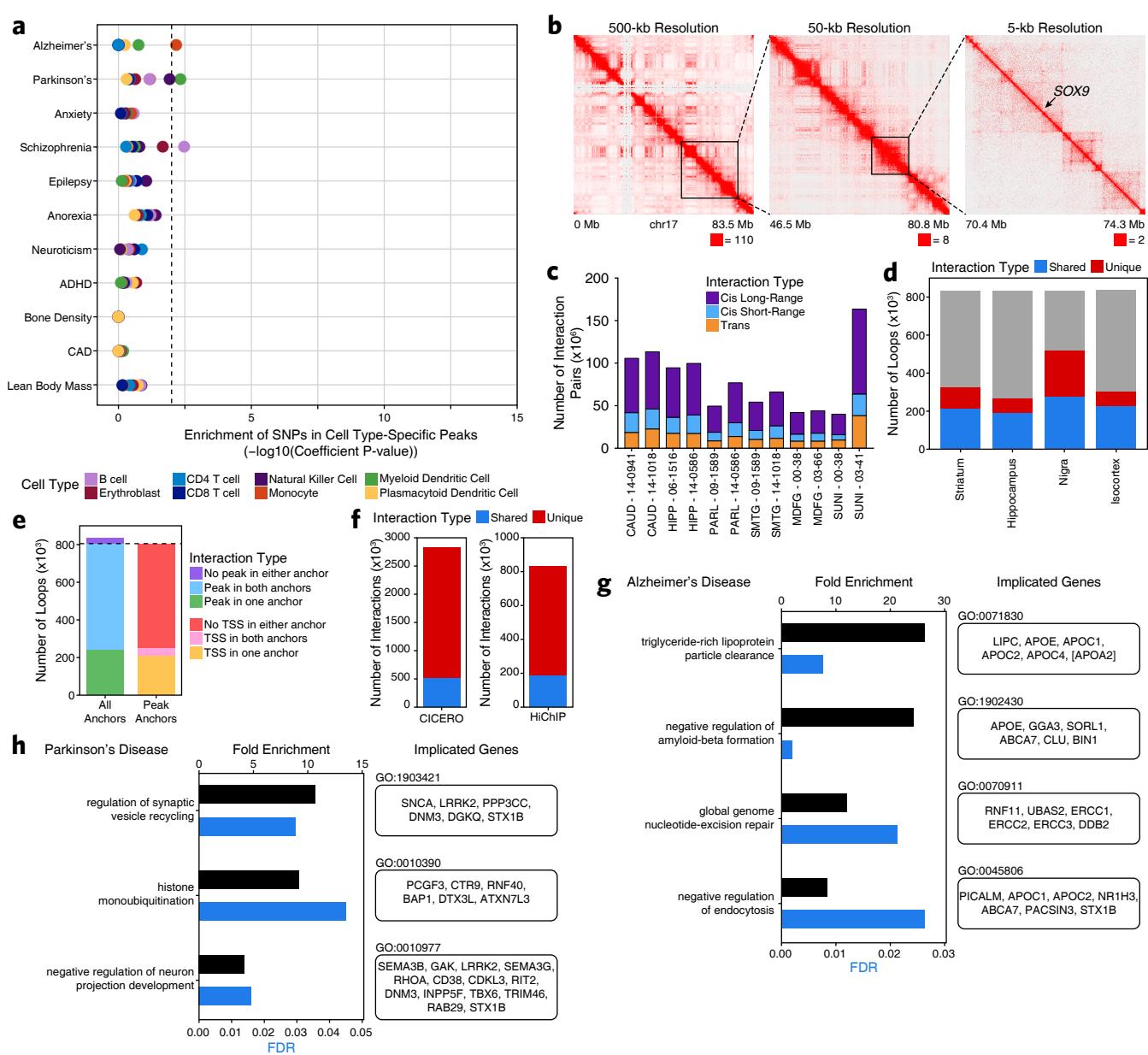
Supplementary Figure 3



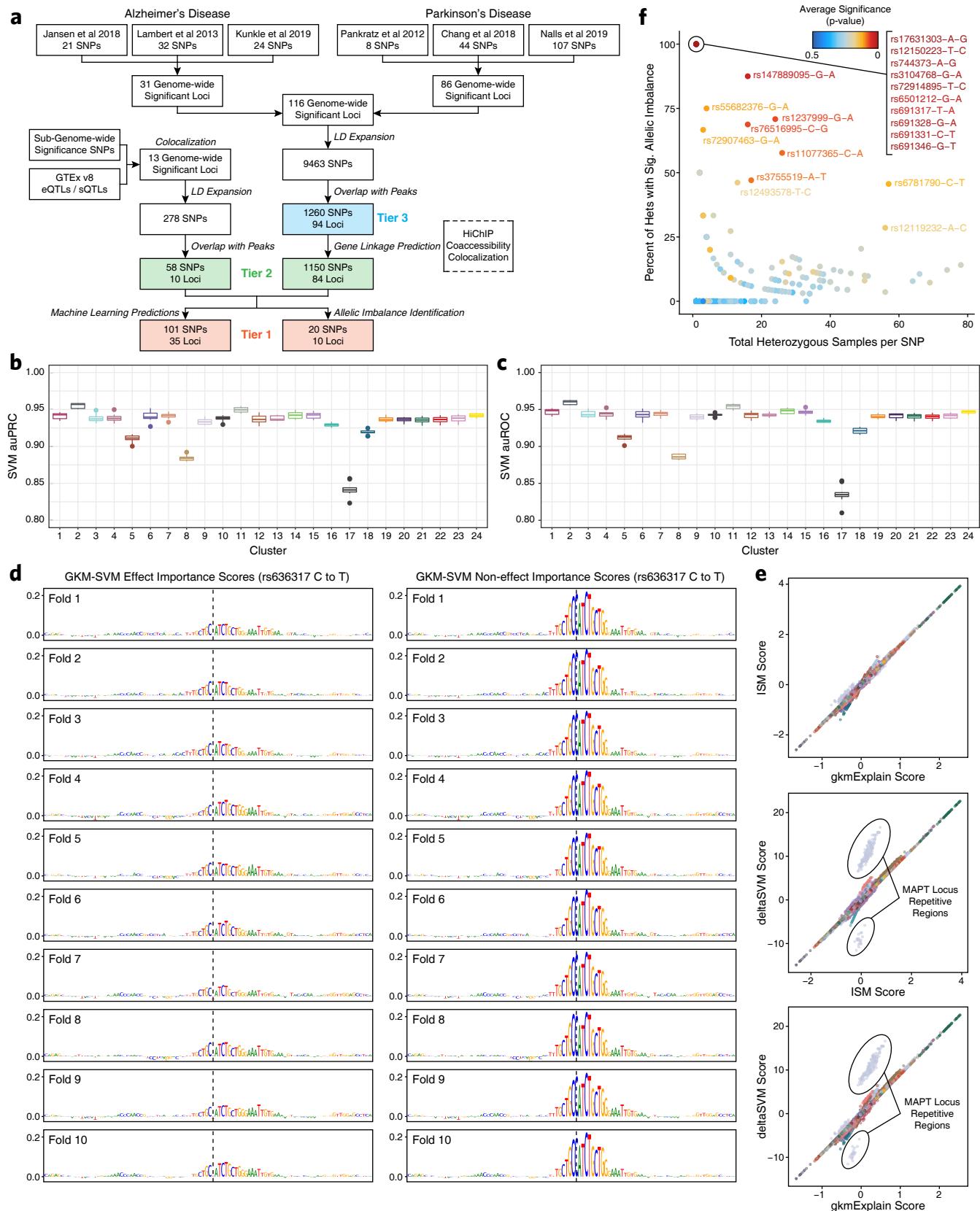
Supplementary Figure 4



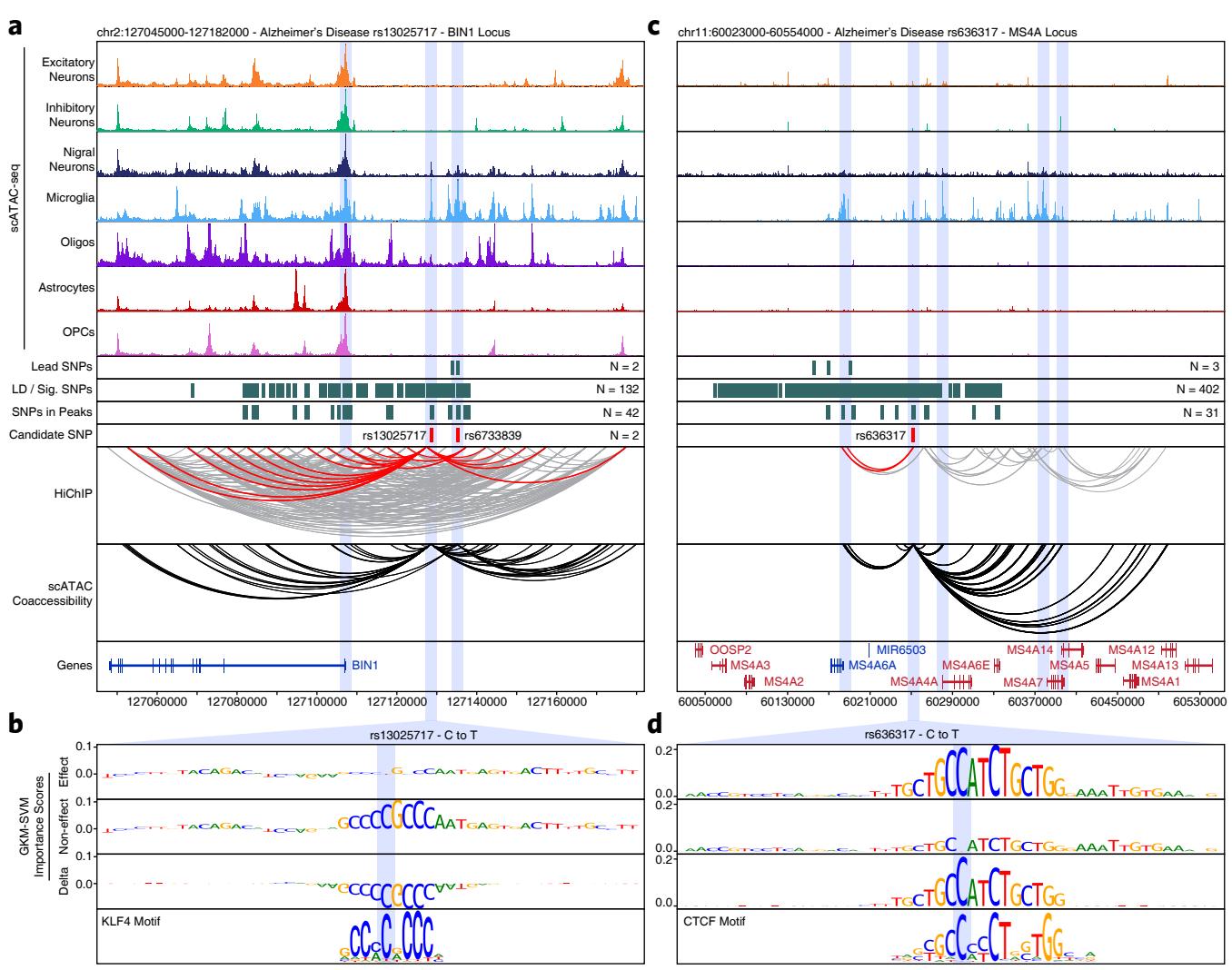
Supplementary Figure 5



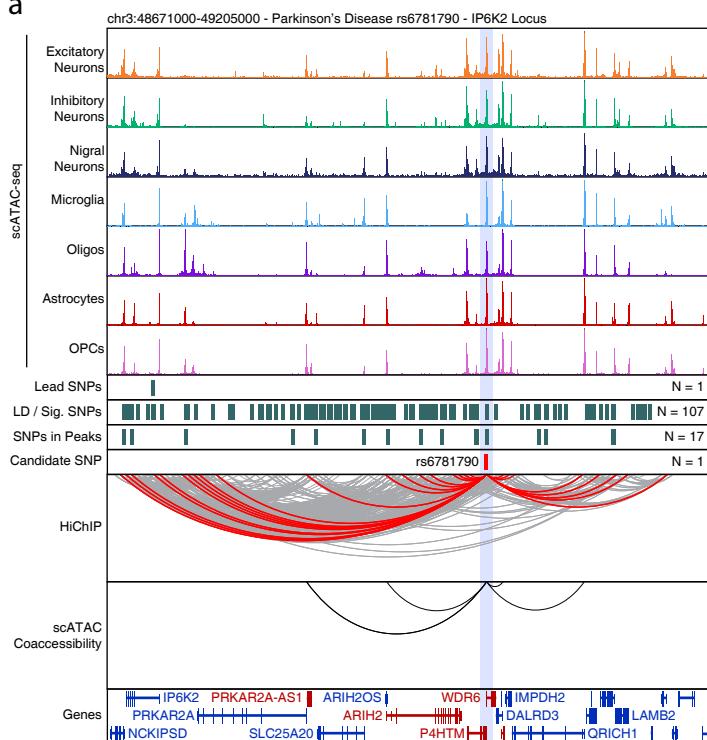
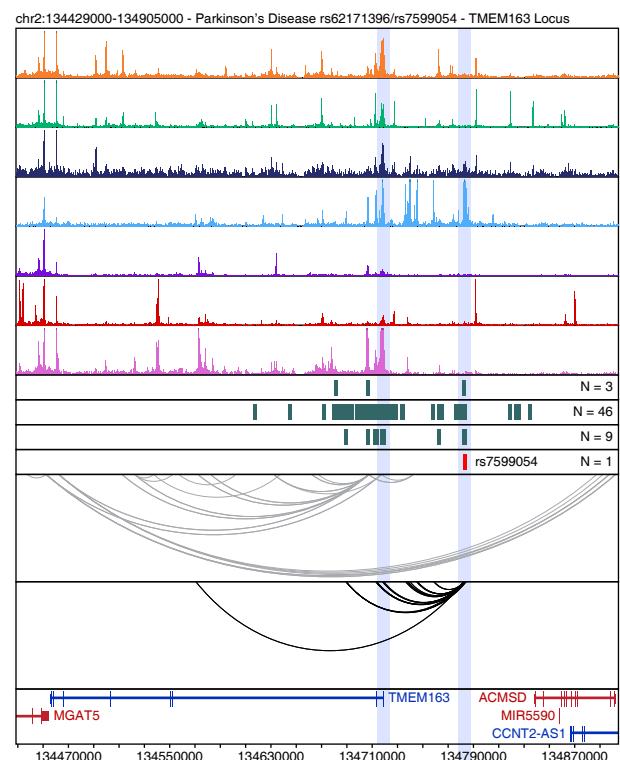
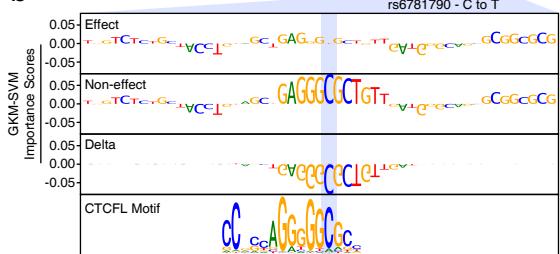
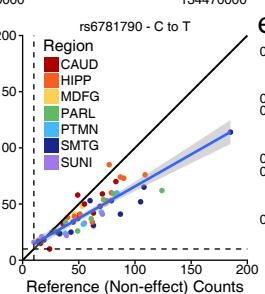
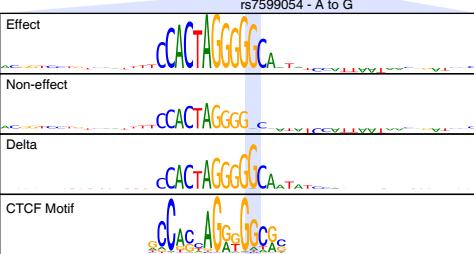
Supplementary Figure 6

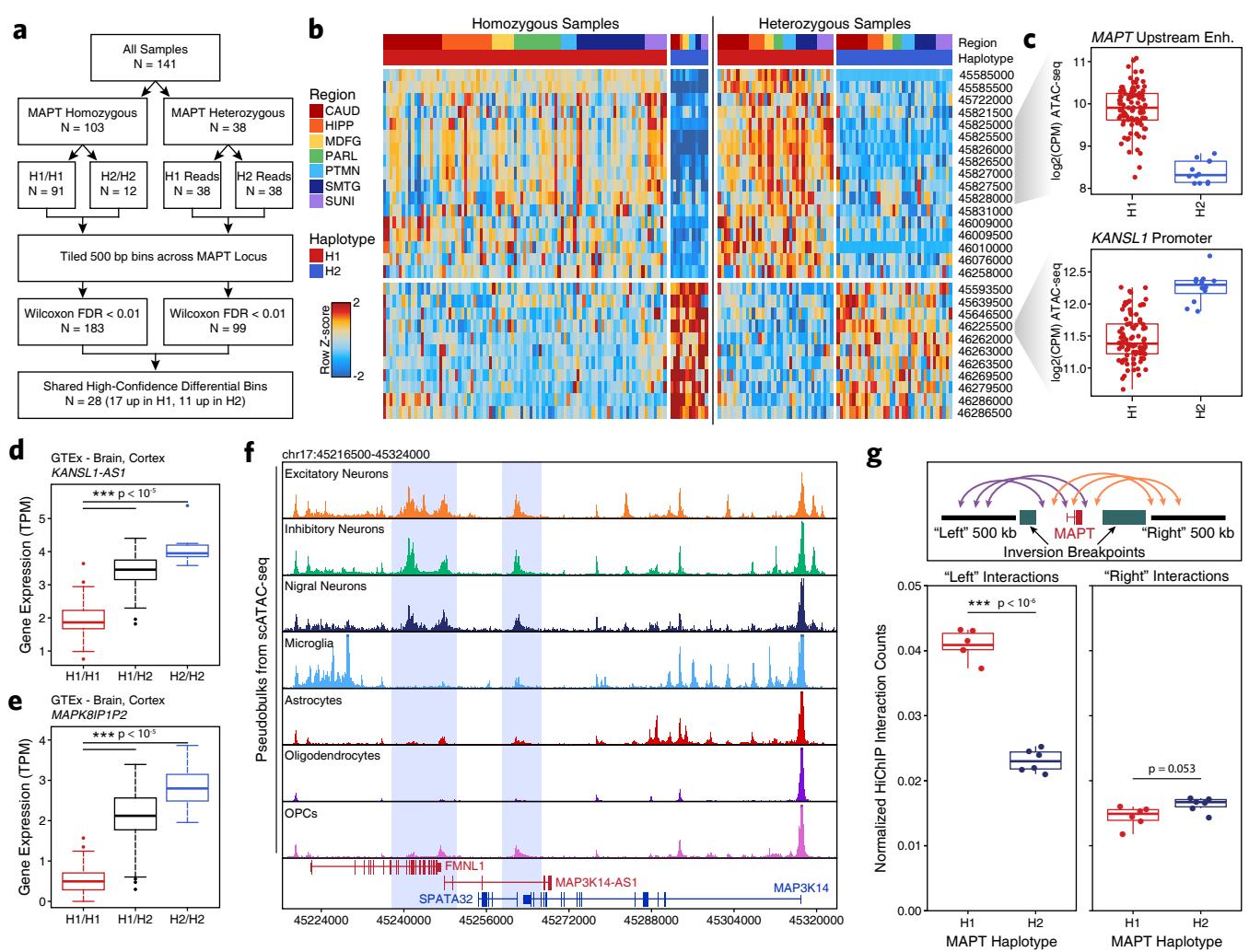


Supplementary Figure 7



Supplementary Figure 8

a**d****b****c****e****Supplementary Figure 9**



Supplementary Figure 10