

1 **Parkinson-associated *SNCA* enhancer variants revealed by open chromatin in mouse dopamine**
2 **neurons**

3 Sarah A. McClymont¹, Paul W. Hook¹, Alexandra I. Soto², Xylena Reed¹, William D. Law¹, Samuel J.
4 Kerans¹, Eric L. Waite¹, Nicole J. Briceno¹, Joey F. Thole¹, Michael G. Heckman³, Nancy N. Diehl³, Zbigniew
5 K. Wszolek⁴, Cedric D. Moore⁵, Heng Zhu⁵, Jennifer A. Akiyama⁶, Diane E. Dickel⁶, Axel Visel^{6,7,8}, Len A.
6 Pennacchio^{6,7,9}, Owen A. Ross^{2,10,11}, Michael A. Beer^{1,12}, Andrew S. McCallion^{1,13,14,*}

7
8 ¹ McKusick Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine,
9 Baltimore, Maryland, USA.

10 ² Department of Neuroscience, Mayo Clinic, Jacksonville, Florida, USA.

11 ³ Division of Biomedical Statistics and Informatics, Mayo Clinic, Jacksonville, Florida, USA.

12 ⁴ Department of Neurology, Mayo Clinic, Jacksonville, Florida, USA.

13 ⁵ Department of Pharmacology and Molecular Sciences, Johns Hopkins University School of Medicine,
14 Baltimore, Maryland, USA.

15 ⁶ Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory,
16 Berkeley, California, USA.

17 ⁷ US Department of Energy Joint Genome Institute, Walnut Creek, California, USA.

18 ⁸ School of Natural Sciences, University of California, Merced, Merced, California, USA.

19 ⁹ Comparative Biochemistry Program, University of California, Berkeley, California, USA.

20 ¹⁰ Mayo Graduate School, Neurobiology of Disease, Jacksonville, Florida, USA.

21 ¹¹ Department of Clinical Genomics, Jacksonville, Florida, USA.

22 ¹² Department of Biomedical Engineering, Johns Hopkins University School of Medicine, Baltimore,
23 Maryland, USA.

24 ¹³ Department of Comparative and Molecular Pathobiology, Johns Hopkins University School of
25 Medicine, Baltimore, Maryland, USA.

26 ¹⁴ Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA.

27 *, To whom correspondence should be addressed: andy@jhmi.edu

28

29

30 **ABSTRACT**

31 The progressive loss of midbrain (MB) dopaminergic (DA) neurons defines the motor features of
32 Parkinson disease (PD) and modulation of risk by common variation in PD has been well established
33 through GWAS. Anticipating that a fraction of PD-associated genetic variation mediates their effects
34 within this neuronal population, we acquired open chromatin signatures of purified embryonic mouse
35 MB DA neurons. Correlation with >2,300 putative enhancers assayed in mice reveals enrichment for MB
36 cis-regulatory elements (CRE), data reinforced by transgenic analyses of six additional sequences in
37 zebrafish and mice. One CRE, within intron 4 of the familial PD gene *SNCA*, directs reporter expression in
38 catecholaminergic neurons of transgenic mice and zebrafish. Sequencing of this CRE in 986 PD patients
39 and 992 controls reveals two common variants associated with elevated PD risk. To assess potential
40 mechanisms of action, we screened >20,000 DNA interacting proteins and identify a subset whose
41 binding is impacted by these enhancer variants. Additional genotyping across the *SNCA* locus identifies a
42 single PD-associated haplotype, containing the minor alleles of both of the aforementioned PD-risk
43 variants. Our work posits a model for how common variation at *SNCA* may modulate PD risk and
44 highlights the value of cell context-dependent guided searches for functional non-coding variation.

45 INTRODUCTION

46 Parkinson disease (PD) is a common progressive neurodegenerative disorder characterized by
47 preferential and extensive degeneration of dopaminergic (DA) neurons in the *substantia nigra*^{1,2}. This
48 loss of midbrain (MB) DA neurons disrupts the nigrostriatal pathway and results in the movement
49 phenotypes observed in PD. While this disorder affects approximately 1% of people over 70 years old
50 worldwide³, the mechanisms underlying genetic risk of sporadic PD in the population remains largely
51 unknown. Familial cases of PD with known pathogenic mutations are better understood but account for
52 $\leq 10\%$ of PD cases⁴.

53 The α -synuclein gene (*SNCA*) is commonly disrupted in familial PD through missense mutations
54 predicted to promote misfolding⁵⁻⁷ or genomic multiplications, resulting in an over-expression
55 paradigm⁸. The *SNCA* locus has also been shown by genome-wide association studies (GWAS) to harbour
56 common variation modulating risk of sporadic PD, with far stronger association signals observed in
57 comparison to all other nominated PD risk loci⁹. In the same way, common variation at over 40
58 additional loci have been implicated in PD¹⁰, but the genes modulated and causal variants responsible
59 for elevating risk remain largely undetermined.

60 That most GWAS-implicated variants are non-coding¹¹ is a major source of this uncertainty, obstructing
61 the identification of: 1) the causative variant at a locus; 2) the context in which the variation is acting
62 and; 3) the mechanism by which a variant asserts its effect on disease risk.

63 GWAS are inherently biologically agnostic and their exploitation of linkage disequilibrium (LD) structure
64 frequently results in many variants being implicated at any one locus, with no one variant prioritized
65 over those in LD. One method to prioritize non-coding variants is to examine the chromatin status at
66 that locus¹¹⁻¹³. Accessible chromatin is more likely to be functional and variants therein may impact that
67 activity, more so than those variants residing in inaccessible chromatin. Chromatin accessibility is
68 dynamic, often varying across cell types and developmental time. Understanding the cellular context in
69 which variation acts is therefore critical to begin the process of prioritizing variants and querying their
70 methods of action^{11,14-16}.

71 Exploiting the preferential vulnerability of MB DA neurons in PD, we have prioritized DA neurons as the
72 biological context in which a fraction of PD-associated variation likely acts. DA neurons in other brain
73 regions, such as the forebrain (FB), provide a related substrate that is less vulnerable to loss in PD. We
74 sought to use chromatin data from *ex vivo* populations of DA neurons to investigate the contributions of

75 non-coding variation to PD risk. To maximize the specificity of the biological context, we generated
76 chromatin signatures of purified mouse MB and FB DA neurons. We examined the resulting regulatory
77 regions for their ability to direct *in vivo* reporter expression and developed a regulatory sequence
78 vocabulary specific to DA neurons. In doing so, we identified a novel MB DA regulatory element that falls
79 within intron 4 of *SNCA* and demonstrate its ability to direct reporter expression in catecholaminergic
80 neurons of transgenic mice and zebrafish. Furthermore, this enhancer harbours two common variants
81 falling in a haplotype that we determine to be associated with PD risk. We demonstrate these enhancer
82 variants impact protein binding and we propose a model for how the variants and the haplotype at large
83 contribute to *SNCA* regulatory control. This work illustrates the power of cell context-dependent guided
84 searches for the identification of disease associated and functional non-coding variation.

85 RESULTS

86 ATAC-seq identifies open chromatin in MB and FB DA neurons

87 To identify regions of open chromatin in DA neurons, we performed ATAC-seq¹⁷ on ~50,000 fluorescent-
88 activated cell sorting (FACS)-isolated cells (per replicate) from microdissected regions of the MB and FB
89 of embryonic day 15.5 (E15.5) Tg(*Th*-EGFP)DJ76Gsat BAC transgenic mice¹⁸ (**Figure 1a**). This mouse line
90 expresses EGFP under the control of the tyrosine hydroxylase (*Th*) locus, labeling catecholaminergic
91 neurons (i.e.: DA, noradrenergic, and adrenergic neurons). To confirm capture of the corresponding
92 catecholaminergic neurons, we performed RT-qPCR on the isolated reporter-labelled cells, establishing
93 them to be enriched for DA neuronal markers relative to unlabelled populations from the same
94 dissected tissues (**Supplemental Figure 1**).

95 To evaluate the ATAC-seq libraries, we examined the called peaks and read pile-ups with the Integrative
96 Genomics Viewer (IGV)^{19,20} and quantified the correlation between brain regions and within replicates. A
97 representative browser trace at the *Th* locus in both MB and FB libraries is presented in **Figure 1b**.
98 Replicates are well correlated: MB library replicates have an average correlation of 0.72 (**Figure 1c**), and
99 FB replicates are more correlated at $r = 0.86$ (**Figure 1d**). Given the robust correlation between
100 replicates, we pooled all reads from the same brain region and called peaks on this unified set to
101 increase our power to detect regions of open chromatin. As a result, we identified 104,217 regions of
102 open chromatin in the MB DA neurons and 87,862 regions in the FB. MB and FB libraries are moderately
103 well correlated (average $r = 0.64$; **Supplemental Figure 2**), with approximately 60% of MB peaks also
104 represented in the FB libraries.

105 To assess these catalogues for characteristics of functionality, we examined the sequence constraint
106 underlying the called regions of open chromatin, excluding peaks that overlap promoters. Promoters are
107 typically accessible²¹ and thus, we aimed to reduce the inflation of sequence conservation due to highly
108 conserved promoter-overlapping ATAC-seq peaks. Despite removal of these highly conserved peaks, we
109 observed a high degree of sequence constraint underlying open chromatin peaks compared to
110 background (**Figure 1e**). The fact that elements in these libraries of putative cis-regulatory elements
111 (CREs) are constrained, highlights their likely functional significance.

112 To further examine the catalogues for biological relevance, we explored the gene ontology (GO) terms
113 of nearby genes. While CREs are not restricted to acting solely on the nearest gene, this restriction is
114 often used as a proxy in the absence of other information. To bolster our predictions, we also generated
115 bulk RNA-seq data on these same populations of sorted cells and used these data to examine the GO
116 terms of the nearest expressed gene (RPKM ≥ 1). While still imperfect, implementing this as a proxy for
117 function results in GO terms enriched for neuronal functions in both MB and FB catalogues (**Figures 1f,**
118 **g**). Thus, we establish these catalogues are enriched for putative CREs likely directing the expression of
119 genes with key roles in neuronal biology.

120 **Candidate regulatory regions are capable of directing expression *in vivo***

121 Although our candidate CRE catalogues appear to be enriched for functional elements on the basis of
122 sequence conservation and GO, both of these metrics are indirect surrogates for true measures of
123 function. To more directly measure the biological relevance of the catalogues and to identify enhancers,
124 we assessed the capability of the candidate CREs to direct expression *in vivo*.

125 We took advantage of the large repository of elements that have already been tested in *lacZ* reporter
126 assays *in vivo* and catalogued in the VISTA enhancer browser²² (accessed September 4, 2016). Overlap
127 between our catalogues and all 2,387 elements in the VISTA enhancer browser, which were scored for
128 their ability to direct *lacZ* reporter expression in E11.5 mice, was quantified (**Supplemental Table 1**). Of
129 the 1,264 elements in VISTA identified as enhancers, 786 were present in the MB catalogue and 719
130 were present in the FB catalogue (**Figure 2a**). Examining the overlap of the FB and MB catalogues with
131 enhancers demonstrated to direct expression in either non-neuronal or neuronal tissues, we observed
132 that 42-47% of enhancers reported to direct expression in non-neuronal tissues are present in the
133 catalogues. By contrast, 71-76% of enhancers that direct expression in one or more regions of the brain
134 overlapped the FB and MB catalogues (**Figure 2b**) confirming an abundance of brain enhancers in our

135 catalogues. Stratifying these confirmed neuronal enhancers on the basis of their expression patterns in
136 VISTA, we observed an abundance of MB-specific enhancers in our MB catalogue, and an abundance of
137 FB-specific enhancers in our FB catalogue, with 77% of MB- and FB-specific enhancers in VISTA captured
138 in our MB and FB catalogues, respectively (**Figure 2c**). Collectively, these data establish that our region-
139 specific catalogues capture region-specific, active CREs with high efficiency.

140 To extend our assessment of the biological activity of sequences within these catalogues, we focused on
141 an additional five candidate CREs not already tested in the VISTA browser and evaluated their ability to
142 act as enhancers in *lacZ* reporter mice and in transgenic zebrafish TdTomato reporter assays. All five
143 regions were represented by robust peaks in both the MB and FB catalogues (**Supplemental Figure 3**).
144 Two regions, one in the first intron of *Kcnq3* and the other downstream of *Foxg1*, were additionally
145 prioritized using H3K27Ac ChIP-seq from a variety of tissues from E11.5 and E15.5 embryonic mice,
146 seeking to limit our selection to candidate enhancers predicted to have neuronal-specific activity. The
147 remaining three candidate CREs were selected on their proximity to genes important in DA neuron
148 biology. We selected sequences at *Foxa2* and *Nr4a2*, as both are key transcription factors (TFs) in the
149 development and maintenance of DA neurons²³⁻²⁶. The final region, located in an intron of *Crhr1*, was
150 selected as this locus has been implicated in PD by GWAS⁹ and our group has recently prioritized this
151 gene as a candidate for PD risk²⁷. All selected sequences were lifted over to hg19, facilitating the
152 identification and assay of their corresponding human sequence intervals.

153 When tested in transgenic reporter mice at E11.5 (**Supplemental Figure 3**), two of the five regions
154 (those near *KCNQ3* and *FOXA2*) were validated as enhancers (**Figure 2c, g**). Recognizing that a disparity
155 exists between the developmental time at which we generated the catalogues (E15.5) and when the
156 data was assayed (E11.5), which may compromise validation rates, we also assayed each sequence
157 across multiple time points in zebrafish. All assayed regions except that at *KCNQ3* directed reporter
158 expression in mosaic transgenic zebrafish (**Figures 2d, e, f, g**). All five regions displayed enhancer activity
159 *in vivo* in neuronal tissues in one or both transgenic assays. Our transgenic animal experiments
160 corroborate the results of the retrospective VISTA enhancer browser intersection; implying that our
161 catalogues of candidate CREs are biologically active and enriched for sequences capable of driving
162 neural expression *in vivo*.

163

164 **Candidate CREs are enriched for TF motifs active in DA neurons**

165 To identify sequence modules (kmers) predicted to contribute regulatory activity of putative CREs in our
166 catalogues, we applied the machine learning algorithm, gkm-SVM²⁸. The resulting regulatory
167 vocabularies of kmers had high predictive power (auROC_{MB} = 0.915, auROC_{FB} = 0.927). We rank ordered
168 and collapsed related kmers to reveal motifs enriched in the putative CREs and their corresponding TFs
169 (**Figures 3a, e, i, m**). In the MB, the four most enriched motifs correspond to Rfx1, Foxa2, Ascl2, and
170 Nr4a2. Given the degeneracy of binding motifs within TF families, we consulted the bulk RNA-seq data
171 for each of the implicated TF families and examined the relative expression levels to prioritize which TFs
172 are most likely producing the observed motif enrichments (**Figures 3b, f, j, n**). While no member of the
173 Rfx family has been canonically associated with MB DA neurons, we anticipate Rfx3 and Rfx7, as the two
174 highest expressed Rfx genes, to likely be active in MB DA neurons and driving this motif enrichment
175 (**Figure 3b**). Foxa1, and more specifically, Foxa2 are both known to DA neuron biology^{23,29} and both are
176 highly expressed in the MB DA neurons (**Figure 3f**). Regarding enrichment for the Ascl family, Ascl1 is
177 known to be involved in DA neuron biogenesis³⁰ and is more highly expressed than any other TF in the
178 family (**Figure 3j**). Finally, Nr4a2 is both canonically associated with DA neurons and required for their
179 development²⁶; we observe it to be highly expressed in MB DA neurons (**Figure 3n**). Examining the
180 sequences underlying the CRE catalogues, we identified TF families known and unknown to DA neuron
181 biology and further refined the TF associations using expression data.

182 We also examined the qualities that differentiate MB CREs from FB CREs by examining the sequences
183 underlying MB-specific and FB-specific regions. We developed a vocabulary that discriminates MB and
184 FB regions with high predictive power (auROC = 0.926) and identified kmers enriched in MB-specific
185 peaks where the top corresponding TFs are Foxa1/2 and Nr4a2 (**Supplemental Figure 4**). We confirmed
186 this MB bias by again considering the bulk RNA-seq for these genes. As expected, these TFs are more
187 highly expressed in the MB where *Nr4a2* is present at 12-fold higher levels in the MB (135 RPKM in the
188 MB vs 11 RPKM in the FB) and *Foxa1/2* are not expressed in the FB, but are present in the MB (*Foxa1*: 28
189 RPKM, *Foxa2*: 7 RPKM). Not only do we identify Foxa1/2 and Nr4a2 as more active in MB DA neurons
190 than in the FB, we did so solely by comparing their role in the vocabulary of MB-specific candidate CREs
191 versus FB-specific CREs.

192 In a parallel strategy to identify TFs actively engaging the DNA in MB DA neurons, we performed TF
193 footprinting in a single deeply sequenced MB ATAC-seq library. Doing so, we confirm that two of the TFs
194 prioritized by gkm-SVM leave robust footprints. The motif corresponding to Rfx binding results in a

195 dearth of cuts directly over predicted binding sites (**Figure 3c**). The same can be seen to a lesser extent
196 for the motif corresponding to Foxa1/2 (**Figure 3g**). By contrast, motifs corresponding to Ascl1 or Nr4a2
197 fail to leave a robust mark on the chromatin availability (**Figures 3k, o**). These footprinting data
198 substantiate the claim that the Rfx family of TFs and Foxa1/2 are active in MB DA neuron CREs.

199 We confirmed that these sequences are indeed enriched in the catalogues by examining the pileup of
200 reads overlapping all genome-wide predicted motif binding sites for each motif identified by gkm-SVM.
201 We see an abundance of reads over predicted binding sites of all four motifs (**Figures 3d, h, l, p**), with
202 the strongest enrichment overlapping Rfx and Ascl1 motif sites (**Figures 3d, l**). Despite the less robust
203 footprint generated at the Ascl1, this TF clearly underlies a larger than expected proportion of CREs in
204 the MB catalogue. The integration of a support vector machine learning algorithm as applied to the
205 sequences underlying open chromatin regions with footprinting analysis in the same chromatin
206 substrate powerfully identifies TFs that are important for DA neuron biology and suggests the Rfx family
207 of TFs, Foxa1/2, Ascl1, and Nr4a2 are actively influencing gene expression in the MB DA neurons.

208 **A candidate CRE in intron 4 of SNCA is associated with PD risk**

209 Having established the biological robustness of the CRE catalogue, we moved to exploit these data to
210 investigate how non-coding variation therein may be contributing to PD risk; given α -synuclein's
211 established role in PD pathogenesis, we prioritized this locus for investigation. We first noted that *Snca*
212 expression differs significantly between the MB and FB DA neurons in our bulk RNA-seq (**Figure 4b**).
213 Examining the chromatin accessibility at the *Snca* locus, the MB and FB are largely the same with the
214 exception of one robust peak in intron 4 (mm9: chr6:60,742,503-60,744,726) that is present in the MB
215 and completely absent in the FB (**Figure 4a**). DNase hypersensitivity site (DHS) linkage^{21,31} suggests that
216 this putative CRE interacts with the *SNCA* promoter. Given the MB-specificity of this putative CRE and
217 indications that it interacts with the *SNCA* promoter, we anticipated this region to be a driving force
218 behind the MB-specific expression of *Snca*.

219 To test this hypothesis, we assayed whether the central portion of this putative CRE, when lifted over to
220 hg19 (chr4:90,721,063-90,722,122), is capable of directing appropriate reporter expression in transgenic
221 zebrafish and mouse reporter assays. Stable transgenesis of zebrafish indicates that this CRE directs
222 reporter expression at 72 hours post fertilization in the locus coeruleus, a key population of
223 catecholaminergic neurons preferentially degenerated in PD³², and along the catecholaminergic tract
224 through the hindbrain, which is largely composed of DA neurons³³ (**Figure 4c**). Additionally, we observe

225 reporter expression throughout the diencephalic catecholaminergic cluster with projections to the
226 subpallium, which is analogous to mammalian dopaminergic projections from the ventral midbrain to
227 the striatum³⁴. Reporter expression in these transgenic zebrafish is largely consistent with an enhancer
228 active in catecholaminergic populations.

229 To further evaluate this CRE in a mammalian system, we generated *lacZ* reporter mice and examined
230 reporter activity across developmental time. Whole mount E12.5 reporter mice indicate this enhancer
231 directs exquisitely restricted expression in Th+ populations, including the dorsal root ganglia, extending
232 into the sympathetic chain, and throughout the cranial nerves (particularly the trigeminal). Additional
233 diffuse staining is noted throughout the MB and FB (**Figure 4d**). Specifically examining the brains of *lacZ*
234 animals at E15.5, reporter expression is identified in the MB and hypothalamus, with strong expression
235 through the amygdala/piriform cortex and along the anterior portion of the sympathetic chain (**Figure**
236 **4e**); similar reporter patterns are seen at P7 (**Figure 4f**). At P30, we detect reporter activity in the
237 amygdala, hypothalamus, thalamus, periaqueductal grey area, brain stem, and importantly, in the
238 *substantia nigra* and ventral tegmental area (**Figure 4g**). By contrast, in aged *lacZ* reporter mice (574
239 days old, ~19 months), we only detect strong reporter expression in the brain stem and observe weak
240 reporter expression in the amygdala (**Figure 4h**). Collectively, the regions in which we detect reporter
241 activity reflect those compromised in PD; Lewy bodies (aggregates of α -synuclein) have been detected in
242 the locus coeruleus, sympathetic chain, amygdala, hypothalamus, ventral tegmental area,
243 periaqueductal grey area of PD patients³⁵⁻³⁹, and critically the preferential degradation of the *substantia*
244 *nigra* is the pathological hallmark of PD progression². This enhancer directs region-specific appropriate
245 expression throughout development in key locations concordant with SNCA activity in PD pathogenesis.

246 Following confirmation of this CRE's regulatory activity in brain regions associated with PD, we next
247 inspected this sequence for PD-associated variation. We sequenced across this interval in 986 PD
248 patients and 992 controls and identified 14 variants (**Supplemental Table 2**), 4 of which were common
249 and present in both cases and controls with a minor allele frequency greater than 5%. Of these, two
250 tightly linked variants ($r^2 = 0.934$; **Supplemental Table 3**), rs2737024 (OR = 1.25, 95% CI = 1.09-1.44, p-
251 value = 0.002) and rs2583959 (OR = 1.22, 95% CI = 1.06-1.40, p-value = 0.005), were significantly
252 associated with PD (**Table 1**). These data support a role for variation within the enhancer in conferring
253 PD risk.

254 To assess how these variants may impact enhancer function and thus PD risk, we assayed differential
255 protein binding at these variants for >20,000 proteins⁴⁰. In doing so, we identify five proteins whose

256 binding is robustly impacted by these implicated variants: NOVA1, APOBEC3C, PEG10, SNRPA, and
257 CHMP5 (**Figure 5a, b, c**). Of these, all are expressed at appreciable levels in both MB and FB DA neurons
258 (**Figure 5d**), excluding APOBEC3C (RPKM ≤ 1). Of the remaining four proteins, three (PEG10, SNRPA, and
259 CHMP5) demonstrate an increased binding affinity for the minor risk allele over the major allele; this
260 direction of effect is consistent with the over-expression paradigm by which *SNCA* confers PD risk⁸.
261 Interestingly, CHMP5 is the sole protein we identify whose binding affinity is impacted by variant
262 rs2583959, and our group has recently implicated one of its family members, CHMP7, in conferring PD
263 risk²⁷, perhaps indicating a role for this family of proteins in PD. Although no single protein stands out,
264 the increased affinity for the risk alleles of the identified enhancer variants by proteins expressed in DA
265 neurons is consistent with a potential mechanistic contribution to *SNCA* expression and therefore, PD
266 risk.

267 Finally, we set out to refine the haplotype structure and understand how this identified variation may be
268 interacting with other variants at this locus. A panel of common variants had previously been genotyped
269 across *SNCA* and PD-associated haplotypes were identified⁴¹. After genotyping our patients and controls
270 for a subset of this panel of variants in addition to all enhancer-associated variants identified by
271 sequencing (**Supplemental Table 4**), we identified a single haplotype that was significantly associated
272 with PD (p-value = 0.003), with a higher observed frequency in PD patients (28.3%) compared to controls
273 (23.4%; **Table 2**). This haplotype implicates some of the same variants as in Guella *et al.*⁴¹ (rs356220,
274 rs737029) but also implicates rs356225 and rs356168, and the two enhancer-associated variants.
275 Collectively, these data identify a catecholaminergic enhancer harbouring common variation that is part
276 of a larger haplotype associated with PD risk.

277 **DISCUSSION**

278 The identification and prioritization of biologically pertinent non-coding variation associated with
279 disease remains challenging. Recent studies by our and other groups have emphasized the importance
280 of cellular context in the identification of sequences harbouring biologically pertinent variation and the
281 genes they regulate. To this end, we used chromatin signatures from *ex vivo* isolated DA neurons to
282 reveal biologically active sequences that harbour non-coding variation contributing to PD risk. We
283 generated robust CRE catalogues for both MB and FB DA neurons, confirmed their capacity to act as
284 enhancers, identified motifs that confer their regulatory potential, and notably, identified two variants
285 located within a MB-specific enhancer that are associated with an increase in PD risk.

286 In contrast to strategies predicated solely on dissection of post-mortem tissues or on the differentiation
287 of cultured cells, we leveraged the use of transgenic reporter mice to specifically isolate Th-expressing
288 neurons from discrete neuroanatomical (FB and MB) domains. While our approach assays a more
289 refined population of DA neurons than would be achieved via gross dissection, recent single-cell RNA-
290 seq analyses of these same cells make clear that even within these highly restricted MB and FB
291 populations there exist two primary cellular phenotypes²⁷. The “homogenous” MB and FB populations
292 each are comprised of an immature neuroblast population and a more mature, domain specific, post-
293 mitotic population of DA neurons. As such, our CRE catalogues capture the chromatin accessibility from
294 both of these states. These catalogues are demonstrably biologically relevant for our purposes, but
295 future studies requiring even greater homogeneity may wish to consider single-cell ATAC-seq to refine
296 these domains further⁴².

297 In our *in silico* validation of the catalogues, we established them to be enriched for both sequence
298 constraint and biological relevance in a manner consistent with function and their FB/MB origin.
299 Furthermore, these sequences are frequently domain appropriate enhancers, with each catalogue
300 capturing a large fraction (77%) of previously validated MB and FB enhancers. Although an abundance of
301 regions are shown to direct neuronal expression compared to those annotated as negative or non-
302 neuronal, it is interesting to note that almost half of the sequences previously documented not to direct
303 expression *in vivo* are also represented in one or both of our catalogues.

304 Given the frequently dynamic nature of CRE activity, this overlap with negative regions likely results
305 from temporal differences in these assays. Our data indicates these regions are accessible at E15.5 but
306 the *lacZ* reporter assays were carried out at E11.5; regions that have been annotated as negative at
307 E11.5 may be active at later time points and, as such, appear in our catalogues. As we moved from these
308 unbiased functional comparisons to more highly selected ones, the potential impact of temporal
309 differences became more pronounced. In mouse transgenic reporter assays, two of five assayed
310 putative CREs direct detectable expression of *lacZ* in neuronal populations. Consistent with the
311 temporally dynamic nature of CREs, when these same regions are tested in zebrafish across multiple
312 developmental time points, we observe four of the five sequences to act as neuronal enhancers.

313 In examining the sequence composition underlying the ATAC-seq peaks, we illuminate powerful
314 vocabularies for both FB and MB DA neuron transcriptional regulatory control. Machine learning using
315 gkm-SVM prioritizes four transcription factor families (Rfx, Foxa1/2, Nr4a2, Ascl1/2) as those conveying
316 significant regulatory potential in the CRE catalogues. Of these, the Rfx family had not previously been

317 implicated in DA neuron biology. Although several of the Rfx family members have been annotated as
318 having expression in the cerebellum or fetal brain⁴³, a role specifically in MB DA neurons has not
319 previously been appreciated. By contrast, Nr4a2 is canonically associated with MB DA neurons^{25,26}, is
320 highly expressed in this population (139 RPKM), and was prioritized as a TF conferring regulatory
321 potential in these cells; however, TF footprinting fails to provide evidence supporting its activity. We
322 postulate that this lack of footprint may reflect the transient DNA binding dynamics of Nr4a2.
323 Transcription factors with short DNA residence times often fail to reveal footprints, and nuclear
324 receptors, such as Nr4a2, have markedly transient DNA interactions⁴⁴.

325 Taken collectively, these data establish a robust biological platform in which PD-associated variation can
326 be evaluated. To this end, an obvious candidate to interrogate was an apparent MB-specific open
327 chromatin domain within intron 4 of the known PD-associated gene, *SNCA*. We assayed the activity of
328 this putative CRE in zebrafish and across the life course of mice and found it to be active in key
329 catecholaminergic structures injured in PD (e.g.: the *substantia nigra* and locus coeruleus), from mid-
330 gestation until at least P30. Thereafter, the utilization of this enhancer in the brain is diminished and by
331 late life appears restricted to the brainstem and amygdala. By the time of clinical presentation, PD
332 patients have already lost a significant proportion ($\geq 30\%$) of their nigral DA neurons^{2,45}; the observed
333 biology of this CRE is consistent with a progressive pathogenic influence acting early in life, rendering
334 these populations preferentially vulnerable to loss over an extensive period of time.

335 Sequencing this interval in PD cases and controls revealed two common variants (rs2737024 and
336 rs2583959) therein, individually associated with an increased risk of PD. Testing these variants for their
337 effect on protein binding, we identify five proteins whose binding is affected, three of which, PEG10,
338 SNRPA, and CHMP5, display greater affinity for the risk allele. Furthermore, we identify a larger
339 haplotype containing these variants, also significantly associated with PD risk. While none of the other
340 SNPs in this haplotype overlap with CREs identified in the DA neuron catalogues, variant rs356168 has
341 significant functional evidence of its activity and contribution to PD risk⁴⁶. The same DHS correlation
342 analysis^{21,31} that suggests an interaction between the *SNCA* promoter and our identified CRE, also
343 suggests an interaction between the *SNCA* promoter and the rs356168 variant. Additionally, ChIA-PET
344 data^{31,47} indicates that sequence encompassing this variant may interact with our enhancer, suggesting a
345 potential co-operative mode of action; a paradigm recently proposed by Gupta and colleagues⁴⁸ at the
346 *EDN1* locus. We propose that the variants within the enhancer, independently or in concert with other

347 variation within the identified haplotype, may act throughout the lifespan to render key populations of
348 catecholaminergic neurons vulnerable, thus increasing PD risk in individuals harbouring this variation.

349 This work emphasizes the value of biologically informed, cell context-dependent guided searches for the
350 identification of disease associated and functional non-coding variation. Given the extent of non-coding
351 GWAS-identified variation, the need for strategies to prioritize variants for functional follow-up is
352 greater than ever. Here, we generate chromatin accessibility data from purified populations of DA
353 neurons to generate catalogues of putative CREs. We have demonstrated how these data can be used to
354 reveal non-coding variation contributing to PD risk; focusing on a single region of open chromatin at the
355 *SNCA* locus, we uncover PD-associated variation therein and propose a model through which this
356 sequence can contribute to normal DA neuronal biology and PD risk. There remains a plethora of
357 information still to be explored in these catalogues, either through further single locus investigations or
358 through massively parallel assays. For example, our MB DA neuron CRE catalogue overlaps SNPs at 10 of
359 26 (38%) PD-associated loci⁹, all of which can be investigated further for their mechanisms by which they
360 impact PD risk. Our work establishes a powerful paradigm, leveraging transgenic model systems to
361 systematically generate cell type specific chromatin accessibility data and reveal disease-associated
362 variation, in a manner that can be progressively guided by improved biological understanding.

363

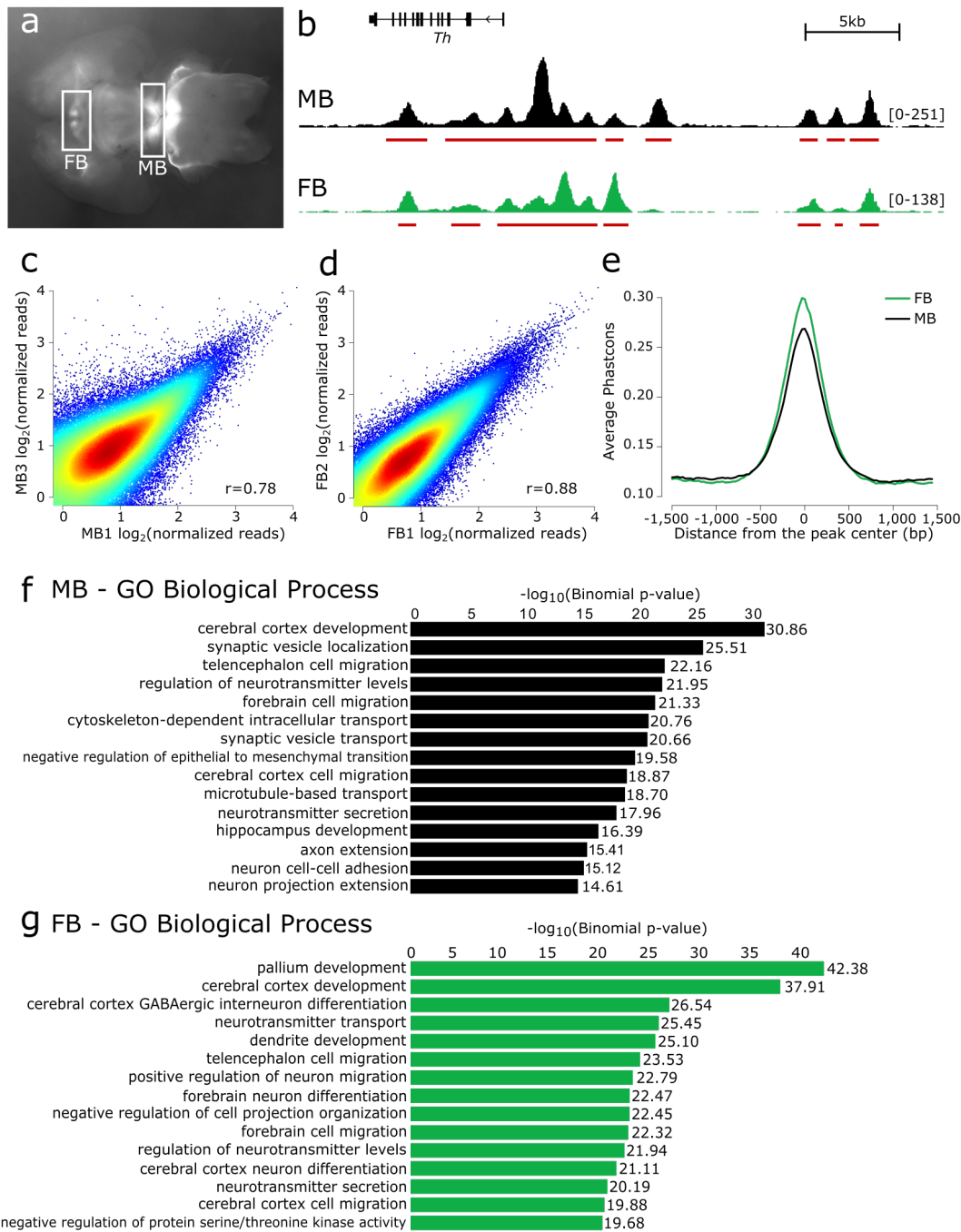


Figure 1 Preliminary validation of ATAC-seq catalogues generated from *ex vivo* DA neurons. **(a)** The midbrain (MB) and forebrain (FB) of E15.5 brains from Tg(Th-EGFP)DJ76Gsat mice are microdissected, dissociated, and isolated by FACS. **(b)** Read pile-up and called peaks for the MB and FB libraries at the *Th* locus. **(c, d)** Chromatin accessibility, genome-wide, is correlated between replicates. **(e)** The sequences underlying MB and FB peaks display a high degree of evolutionary sequence constraint as measured by PhastCon scores. **(f, g)** Gene ontology terms of the nearest expressed genes to all peaks in both the MB and FB reflect the neuronal origin and function of these catalogues.

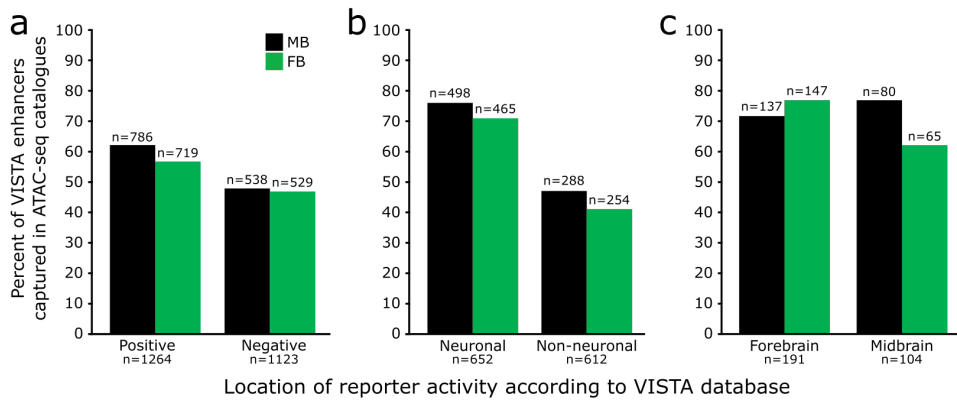
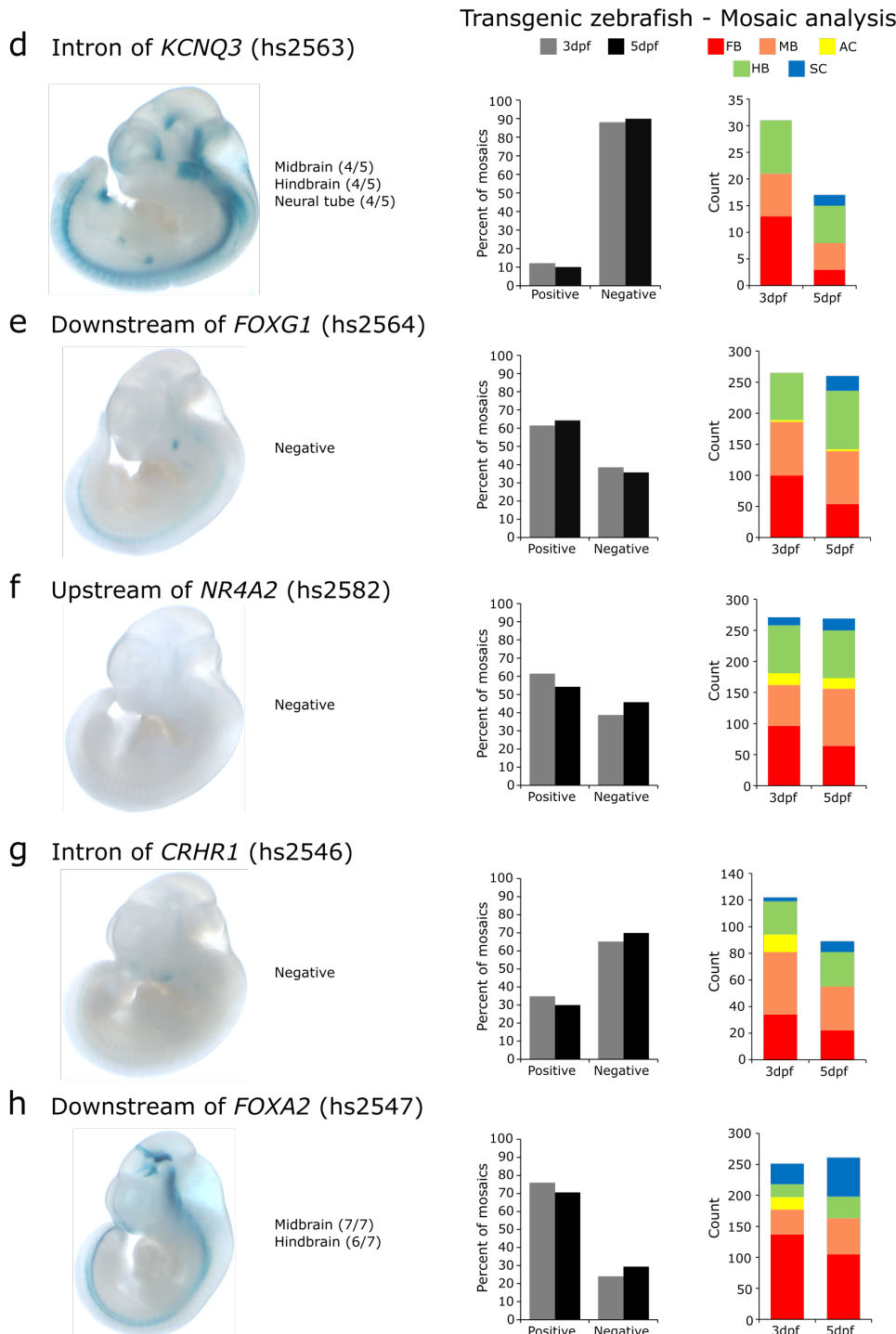


Figure 2 Validation of the putative CRE catalogues *in vivo*. **(a)** Of the elements annotated in VISTA as having enhancer activity, 62% and 56% of these are represented in the MB and FB catalogues, respectively **(b)** An abundance of open chromatin regions in the MB and FB catalogues overlap confirmed neuronal enhancers ($\geq 70\%$). **(c)** Stratifying neuronal enhancers, MB- and FB-specific enhancers are enriched in our MB and FB catalogues, respectively. **(d-h)** Testing five prioritized putative CREs *in vivo* identifies five neuronal enhancers. **(d)** A putative CRE in intron 1 of *KCNQ3* directs expression in the midbrain, hindbrain, and neural tube of E11.5 lacZ reporter mice. It fails to direct expression in a transgenic zebrafish assay at either 3 or 5 days post fertilization (dpf); reporter expression present in $\leq 25\%$ of mosaics. **(e, f, g)** Putative CREs downstream of *FOXP1*, upstream of *NR4A2*, and in an intron of *CRHR1* fail to direct expression in transgenic mice, however, they direct robust neuronal appropriate expression in transgenic zebrafish reporter assays (scored for expression in MB, FB, amacrine cells (AC), hindbrain (HB), spinal cord (SC)). **(h)** A putative CRE downstream of *FOXA2* directs neuronal expression in both transgenic mice and zebrafish assays. N mosaic zebrafish scored: ≥ 141 for 3dpf, ≥ 119 for 5dpf. All constructs have since been deposited in the VISTA database, under the hs numbers supplied.



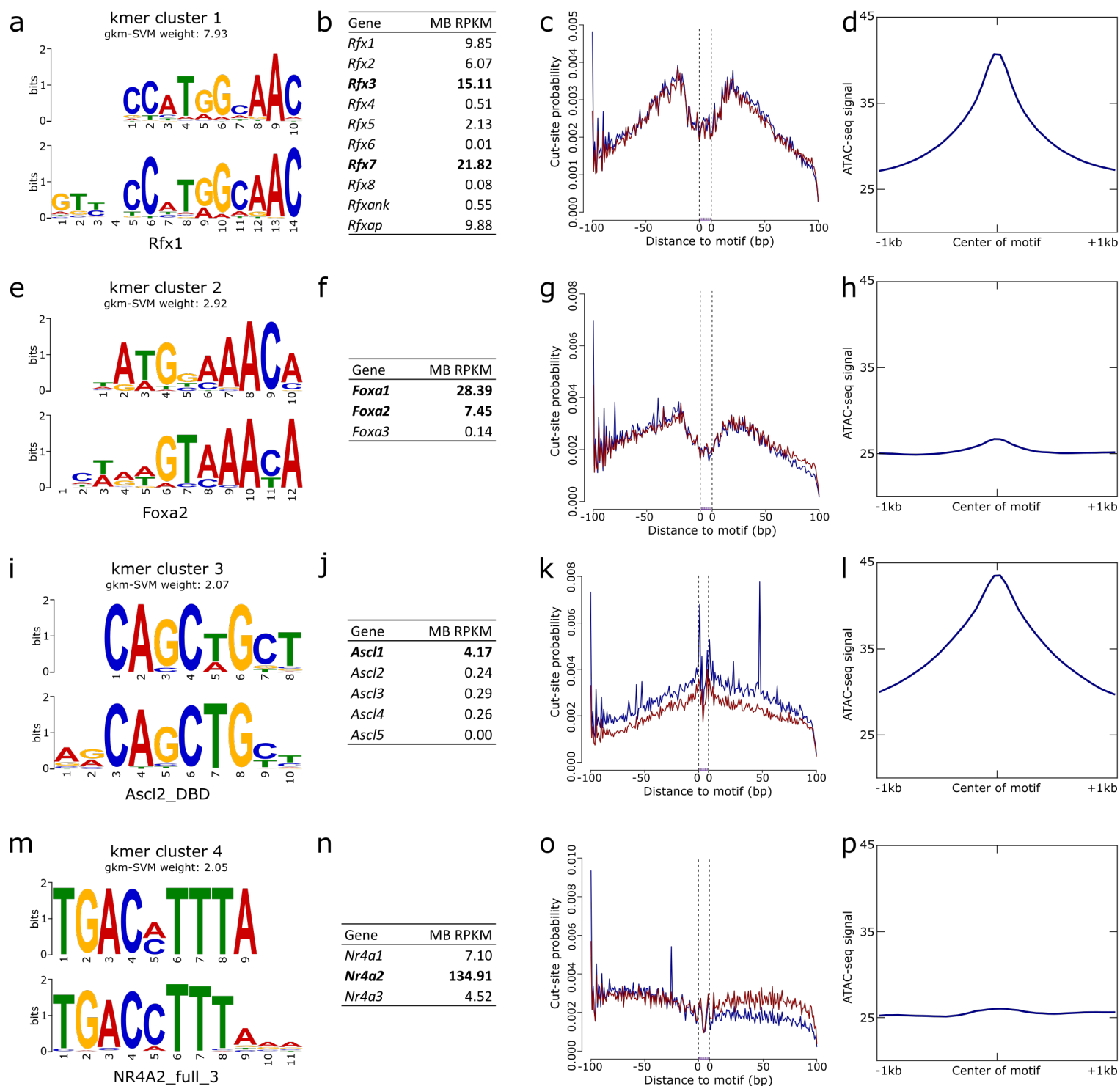


Figure 3 Identification of transcription factors (TFs) important to DA neurons. **(a)** The kmer predicted to have the greatest regulatory potential underlying MB ATAC-seq peaks corresponds to the Rfx family of TFs. **(b)** RNA-seq quantification in these same cells indicates this enrichment is likely due to Rfx3 or Rfx7 activity. Examining the ATAC-seq signal over predicted binding sites reveals a robust TF footprint **(c)** and a general enrichment of reads overlapping Rfx sites genome-wide **(d)**. **(e-h)** Similarly, a kmer corresponding to the TFs Foxa1/2 have similar evidence for their activity. **(i-j)** The third ranked motif likely corresponds to Ascl1, and while it fails to leave a robust TF footprint **(k)**, there is clear enrichment of ATAC-seq signal overlapping genome-wide predicted Ascl1 binding sites **(l)**. **(m-n)** Nr4a2, canonically associated with DA neuron biology, is identified as a highly expressed TF likely contributing to the regulatory potential of the putative CREs however, it fails to leave a TF footprint in the cut-site patterns around predicted motif sites **(o)** and is only mildly enriched for ATAC-seq reads over its predicted binding sites **(p)**.

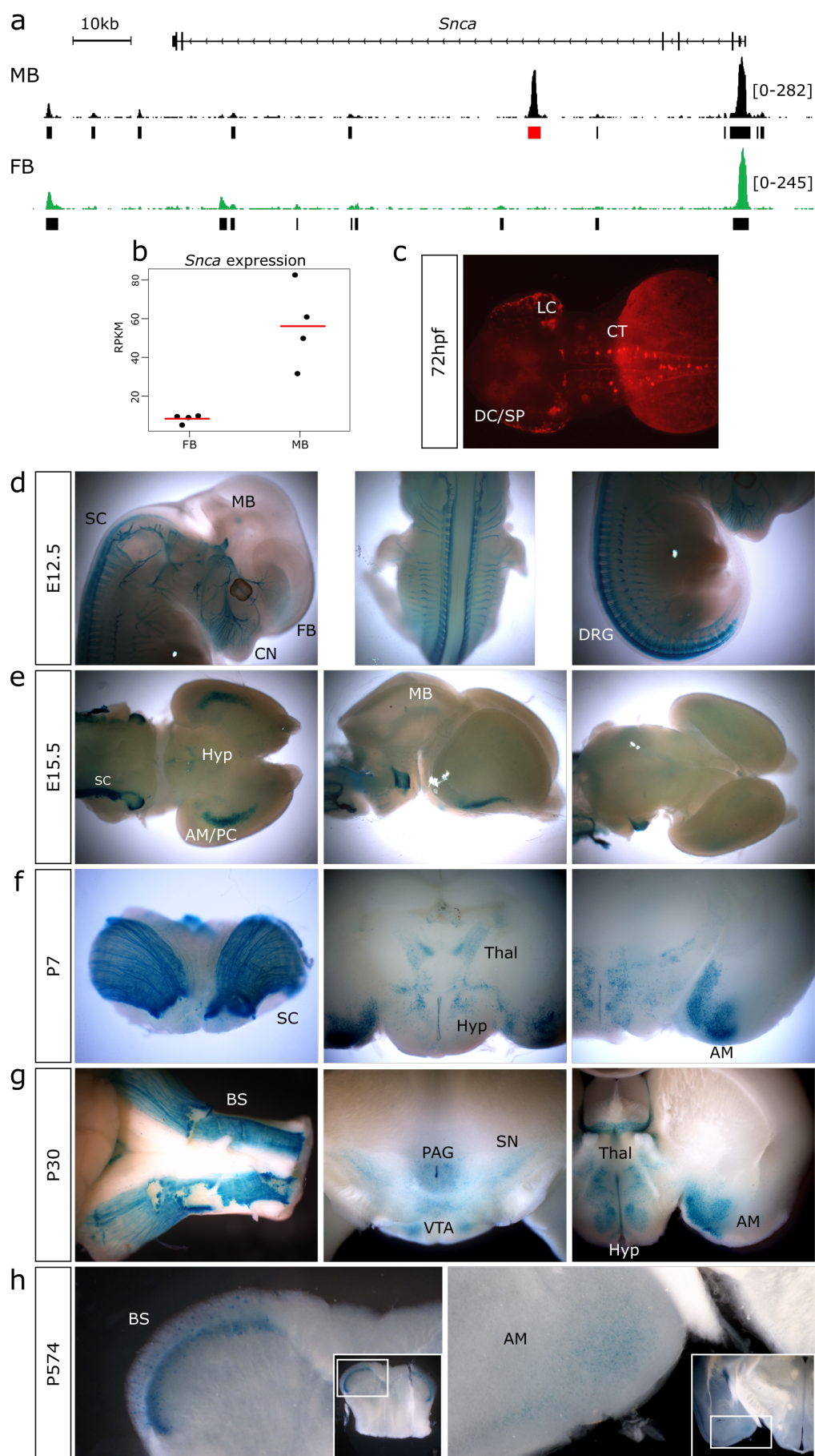


Figure 4 A MB-specific enhancer directs expression in catecholaminergic populations of neurons known to Parkinson disease biology. **(a)** IGV track indicating the location of the MB-specific region of open chromatin, located in intron 4 of *Snca*. **(b)** *Snca* is differentially expressed between the MB and FB DA neurons. Red bar is the mean expression of the four replicates (black dots). **(c)** At 72 hours post fertilization (hpf), stable transgenic zebrafish reporter assays indicate this putative CRE is capable of directing reporter expression in key catecholaminergic neuronal populations, including the locus coeruleus (LC), the catecholaminergic tract (CT) of the hindbrain, and the diencephalic cluster (DC) with projections to the subpallium (SP). **(d-g)** Further studies in *lacZ* reporter assays in embryonic (E) and post-natal (P) mice indicate dynamic enhancer usage across developmental time. **(d)** This enhancer directs expression throughout the MB, FB, dorsal root ganglia (DRG), sympathetic chain (SC), and cranial nerves (CN) of E12.5 mice. **(e)** By E15.5, reporter expression is observed in the amygdala and/or piriform cortex (AM/PC), sympathetic chain, MB, and hypothalamus (Hyp). **(f)** Patterns of reporter expression at P7 reflect those seen at E15.5. **(g)** Reporter activity is observed at P30 in the amygdala, hypothalamus and thalamus (Thal), brain stem (BS), *substantia nigra* (SN), ventral tegmental area (VTA), and the periaqueductal grey area (PAG). **(h)** In aged mice (P574), reporter expression is detected robustly in the brain stem and faintly in the amygdala.

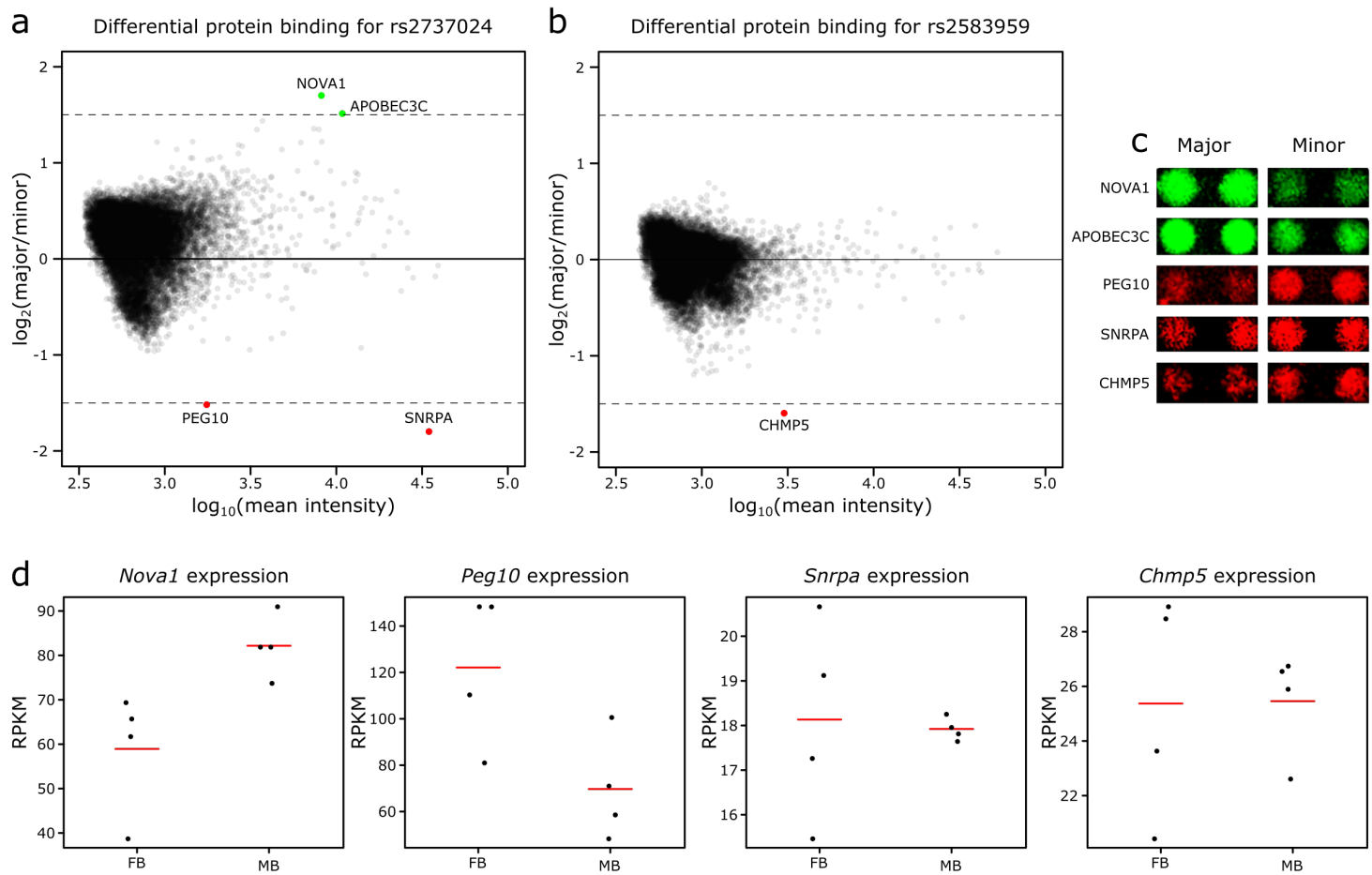


Figure 5 Identification of proteins whose binding is impacted by the implicated PD-risk SNPs. **(a, b)** MA plots for both rs2737024 and rs2583959 indicating the magnitude of the effect of the minor and major allele on binding. Cut-off for differential binding: $\log_2(\text{major}/\text{minor}) \geq 1.5$ or ≤ -1.5 . **(a)** NOVA1 and APOBEC3C (green circles) bind at rs2737024 with greater affinity for the major allele, while PEG10 and SNRPA (red circles) have a greater affinity for the minor allele. **(b)** CHMP5 (red circle) has a greater affinity for the minor allele of rs2583959. **(c)** Representative images of the protein binding for each of the differentially bound proteins. **(d)** Expression analysis in the MB and FB DA neurons for each of the differentially bound proteins indicate *Nova1*, *Peg10*, *Snrpa*, and *Chmp5* to be highly expressed in these populations, while none of the *Apoec* family member genes are expressed (RPKM ≤ 1 , data not shown). Red bar is the mean expression of the four replicates (black dots).

Table 1: Two tightly linked SNPs within the enhancer are significantly associated with PD risk

Variant	MA	MAF in PD cases (N=986)	MAF in controls (N=992)	Association with PD	
				OR (95% CI)	p-value
rs7684892	A	0.069	0.065	0.93 (0.72-1.20)	0.562
rs17016188	C	0.082	0.061	1.35 (1.04-1.75)	0.023
rs2583959	G	0.317	0.271	1.22 (1.06-1.40)	0.005 *
rs2737024	G	0.319	0.270	1.25 (1.09-1.44)	0.002 *

MA=minor allele; MAF=minor allele frequency; OR=odds ratio; CI=confidence interval
 Only variants with MAF > 0.05 were considered
 ORs, 95% CIs, and p-values result from additive logistic regression models adjusted for age at blood draw and sex
 p-values ≤ 0.0125 were considered as statistically significant after applying a Bonferroni correction for multiple testing (*)

Table 2: A single haplotype, containing the minor alleles of the implicated SNPs, is significantly associated with PD risk

		rs356220	rs356225	rs3857057	rs356168	rs10018362	rs2737029	rs62306323	rs2737024	rs2583959	rs17016188	rs7684892	rs7689942	Frequency in PD cases	Frequency in controls	p-value
Haplotypes spanning the <i>SNCA</i> locus	1	C	C	A	G	T	T	C	A	C	T	G	C	0.015	0.027	0.012
	2	C	G	A	A	T	T	T	A	C	T	G	C	0.092	0.110	0.029
	3	C	C	A	G	C	C	C	A	C	T	A	T	0.039	0.048	0.184
	4	T	C	A	G	T	T	C	A	C	T	G	C	0.037	0.040	0.480
	5	C	G	A	A	T	T	C	A	C	T	G	C	0.380	0.397	0.593
	6	T	C	A	G	T	T	T	A	C	T	G	C	0.010	0.011	0.698
	7	C	G	A	A	T	C	C	G	G	T	G	C	0.009	0.012	0.944
	8	C	C	A	G	T	C	C	G	G	T	G	C	0.016	0.015	0.768
	9	T	C	G	G	C	C	C	A	C	T	A	T	0.021	0.017	0.360
	10	T	C	G	G	T	C	C	A	C	C	G	C	0.014	0.009	0.189
	11	T	C	G	G	C	C	C	A	C	C	G	C	0.057	0.044	0.124
	12	T	C	A	G	T	C	C	G	G	T	G	C	0.283	0.234	0.003 *

Only haplotypes with frequency ≥ 0.01 were considered
 Black boxes indicate the minor allele in Europeans
 p-values result from score tests for association, performed under an additive model, adjusted for age at blood draw and sex
 p-values ≤ 0.0042 were considered as statistically significant after applying a Bonferroni correction for multiple testing (*)

366 **METHODS**

367 **Animal husbandry**

368 Tg(Th-EGFP)DJ76Gsat mice (Th-EGFP) were generated by the GENSAT Project and purchased through
369 the Mutant Mouse Resource and Research Centers Repository. Colony maintenance matings were
370 between hemizygous male Th-EGFP mice and female Swiss Webster (SW) mice, obtained from Charles
371 River Laboratories. This same mating scheme was used to establish timed matings, generating litters for
372 assay; day on which vaginal plug is observed, E0.5. Adult AB zebrafish lines were maintained in system
373 water according to standard methods⁴⁹. All work involving mice and zebrafish (husbandry, colony
374 maintenance, procedures, and euthanasia) were reviewed and pre-approved by the institutional care
375 and use committee.

376 **Neural dissociation and FACS**

377 Pregnant SW mice were euthanized at E15.5 and the embryos were removed and immediately placed in
378 chilled Eagle's Minimum Essential Media (EMEM) on ice. Embryos were decapitated and brains were
379 removed into Hank's Balanced Salt Solution without Mg²⁺ and Ca²⁺ (HBSS w/o) on ice. Under a
380 fluorescent microscope, EGFP+ brains were identified and microdissected to yield the desired forebrain
381 (FB) and midbrain (MB) regions desired. Microdissected regions were placed in fresh HBSS w/o on ice,
382 and pooled per litter for dissociation.

383 Pooled brain regions were dissociated using the Papain Dissociation System (Worthington Biochemical
384 Corporation). The tissue was dissociated in the papain solution for 30 minutes at 37°C, with gentle
385 trituration every 10 minutes using a sterile Pasteur pipette. Following dissociation, cells were passed
386 through a 40µm cell strainer into a 50mL conical, centrifuged for 5 minutes at 300g, resuspended in
387 albumin-inhibitor solution containing DNase, applied to a discontinuous density gradient, and
388 centrifuged for 6 minutes at 70g. The resulting cell pellet was resuspended in HBSS with Mg²⁺ and Ca²⁺
389 and submitted to FACS. Aliquots of 50,000 EGFP+ cells were sorted directly into 300µL HBSS with Mg²⁺
390 and Ca²⁺ with 10% FBS for ATAC-seq. Aliquots containing ≥50,000 EGFP+ cells were sorted into kit-
391 provided lysis buffer for RNA-seq. This procedure was repeated such that a single aliquot of cells from
392 each region per litter were submitted to either ATAC-seq or bulk RNA-seq three times over for each
393 region.

394 **ATAC-seq library preparation and quantification**

395 ATAC-seq library preparation generally follows the steps as set out in the original ATAC-seq paper¹⁷ with
396 minor modifications. Aliquots of 50,000 EGFP+ cells were centrifuged for 5 minutes at 4°C and 500g,
397 washed with 50µL of chilled PBS and centrifuged again for 5 minutes at 4°C and 500g. The cell pellet was
398 resuspended in lysis buffer, as set out in the protocol, and cells were left to lyse for 5 minutes at 4°C
399 before being centrifuged for 10 minutes at 4°C at 500g. The resulting nuclei pellet was transposed, as
400 written, using the transposase from the Nextera DNA Library Preparation Kit. Following transposition,
401 DNA was purified with the MinElute Reaction Clean-up Kit (Qiagen) and eluted in 10µL elution buffer.
402 Libraries were amplified according to the original ATAC-seq protocol¹⁷. The qPCR surveillance steps were
403 modified such that the additional number of cycles of amplification were calculated as ¼ maximum
404 intensity, so as to limit PCR duplication rates in the final libraries. Amplified libraries were purified with
405 Ampure XP beads (Beckman Coulter) following the Nextera DNA Library Prep Protocol Guide. Libraries
406 were quantified using the Qubit dsDNA High Sensitivity Assay (Invitrogen) in combination with the
407 Agilent 2100 Bioanalyzer using the High Sensitivity DNA Assay (Agilent).

408 **ATAC-seq sequencing, alignment, and peak calling**

409 Individual ATAC-seq libraries were sequenced on the Illumina MiSeq to a minimum depth of 20 million,
410 2x75bp reads per library. A single MB ATAC-seq library was sequenced on the Illumina HiSeq in Rapid
411 Run mode with 2x100bp reads, to a depth of ≥350 million paired-end reads.

412 Quality of sequencing was evaluated using FastQC (v0.11.2;
413 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Reads were aligned to mm9 using
414 Bowtie2⁵⁰ (v2.2.5), under --local mode. Reads aligning to the mitochondrial genome, unknown and
415 random chromosomes, and PCR duplicates were removed prior to peak calling (SAMtools⁵¹). Peaks were
416 called on individual libraries and on a concatenated file combining all MB and all FB libraries (“Joint”)
417 using MACS2⁵² (v2.1.1.20160309) “callpeak” with options: --nomodel --nolambda -B -f BAMPE --gsize
418 mm --keep-dup all. Peaks overlapping blacklisted regions called by ENCODE and in the original ATAC-seq
419 paper were removed^{17,47}.

420 **RNA-seq library preparation and quantification**

421 Total RNA was extracted using the Purelink RNA Micro Kit (Invitrogen). Following FACS isolation into kit-
422 provided lysis buffer, samples were homogenized and RNA extraction proceeded using manufacturer’s
423 recommendations. Total RNA integrity was determined using the RNA Pico Kit (Agilent). RNA samples

424 were sent to the Sidney Kimmel Comprehensive Cancer Center Next Generation Sequencing Core at
425 Johns Hopkins for library preparation, using the Ovation RNA-Seq System V2 (Nugen), and sequencing.

426 **RNA-seq sequencing, alignment, and transcript quantification**

427 Libraries were pooled and sequenced on Illumina's HiSeq 2500 in Rapid Run mode with 2x100bp reads
428 to an average depth of >90 million reads per library. Quality of sequencing was evaluated using FastQC.
429 FASTQ files were aligned to mm9 using HISAT2⁵³ (v2.0.1-beta) with --dta specified.

430 Aligned reads from individual samples were quantified against a reference transcriptome using the
431 Rsubread package⁵⁴⁻⁵⁶ (v1.22.3) function "featureCounts" with the following options: isPairedEnd =
432 TRUE, requireBothEndsMapped = TRUE, isGTFAnnotationFile = TRUE, useMetaFeature = TRUE. The
433 GENCODE vM9 GTF was downloaded⁵⁷ (date: March 30, 2016) and lifted over from the mm10 to the
434 mm9 genome using CrossMap (v0.2.2) with default parameters⁵⁸. This was used for quantification, in
435 which gene-level raw counts were converted to RPKM values and means for each region were
436 calculated.

437 **cDNA synthesis and RT-qPCR for DA neuron markers**

438 RNA was extracted using the RNeasy Mini Kit (Qiagen), after sorting 50,000 cells directly into Buffer RLT.
439 Aliquots of 50,000 non-fluorescing cells were also collected and processed in parallel. 100ng of each
440 RNA sample was submitted to first strand cDNA synthesis using the SuperScript III First-Strand Synthesis
441 System for RT-PCR (Invitrogen), following the Oligo(dT) method.

442 Primers (**Supplementary Table 5**) were designed using Primer-BLAST⁵⁹ under default parameters with
443 the requirement for exon-exon junction spanning specified. qPCR was performed using Power SYBR
444 Green Master Mix (Applied Biosystems). Reactions were run in triplicate, following default SYBR Green
445 Standard cycle specifications on the Viia7 Real-Time PCR System (Applied Biosystems). Relative
446 quantification followed the $2^{-\Delta\Delta CT}$ method, normalizing results to *Actb* in the EGFP- aliquot of cells for
447 each region, respectively.

448 **Correlation analysis between regions and within replicates**

449 Peaks from all six ATAC-seq libraries and the two "Joint" ATAC-seq libraries were concatenated together,
450 sorted on the basis of chromosomal location, merged into a unified peak set⁶⁰, and converted to
451 Simplified Annotation Format (SAF). Reads from each BAM file overlapping this unified peak set were
452 quantified with the Rsubread package "featureCounts" command, with the following options:

453 isPairedEnd = TRUE, requireBothEndsMapped = FALSE. Read counts were normalized for each library
454 using conditional quantile normalization⁶¹, accounting for library size, peak length, and peak GC content.
455 Pearson correlation co-efficients were calculated from this normalized count matrix and visualized using
456 corrplot⁶² and RColorBrewer⁶³ and LSD⁶⁴.

457 **Sequence constraint analysis**

458 Average phastCons⁶⁵ were calculated for the “Joint” peak file for both the MB and FB libraries using
459 Cistrome⁶⁶ (<http://cistrome.org/ap/>). Beforehand, peaks with overlap of exons or promoters (defined
460 here as +/-2000bp from the transcriptional start site) were removed. The exon and promoter BED files
461 were downloaded from the UCSC table browser⁶⁷ (Mouse genome; mm9 assembly; Genes and Gene
462 Predictions; RefSeq Genes track using the table refGene).

463 **Gene ontology of nearest expressed gene**

464 The Genomic Regions Enrichment of Annotations Tool⁶⁸ (GREAT; v3.0.0; <http://great.stanford.edu>)
465 predicted the GO term enrichment in the catalogues. Beforehand, peaks were processed to: a) remove
466 peaks overlapping commonly open regions; b) select the top 20,000 peaks and; c) overlap the nearest
467 expressed gene’s transcriptional start site (TSS).

468 First, regions that are commonly open were defined as those regions of the genome that are open in
469 >30% of ENCODE DNase hypersensitivity site (DHS) assays in mouse tissues. These ubiquitously open
470 regions were removed from the peak files. Next, to limit the number of regions submitted to GREAT
471 such that the binomial distribution for calculating fold enrichment values was still valid, peak files were
472 limited to the top 20,000 peaks on the basis of q-value.

473 Finally, in order to limit ourselves to the nearest expressed gene, we supplied a list of the TSSs of the
474 nearest expressed gene that are in the GREAT database. The list of genes and their TSSs used by GREAT
475 was downloaded from:

476 <http://bejerano.stanford.edu/help/download/attachments/2752609/mm9.great3.0.genes.txt>. Only
477 genes that are in this list with RPKM > 1 were considered as expressed. The nearest expressed gene to
478 each of the top 20,000 peaks was identified. Each peak is associated with its nearest expressed gene and
479 to ensure that GREAT only considered these nearest genes for analysis, we submitted these nearest
480 expressed gene’s TSSs as a proxy for each peak. These proxy peaks were submitted to GREAT using the
481 NCBI build 37 (mm9) assembly, under whole genome background regions, with the single nearest gene
482 as the association rule, including curated regulatory domains.

483 **Quantification of overlap between CRE catalogues and the VISTA Enhancer Browser**

484 All elements tested *in vivo* were downloaded from the VISTA Enhancer Browser
485 (<https://enhancer.lbl.gov>) on September 4, 2016. These regions were stratified into those annotated as
486 positive or negative. BED co-ordinates of these regions were extracted and intersected with the ATAC-
487 seq catalogues. Positive regions were further stratified into those with annotations for only forebrain,
488 only midbrain, only hindbrain, combinations of regions (“Multiple regions”), all three regions (“Whole
489 brain”), summing to the “Neuronal” category, or were annotated as positive but driving expression in
490 none of those three regions (“Non-neuronal”).

491 **Testing five putative CREs for *in vivo* reporter activity**

492 Prioritized regions were PCR amplified (**Supplementary Table 5**) from human gDNA and cloned into
493 either pENTR for mouse *lacZ* assays (Invitrogen) or pDONR221 for zebrafish assays (Invitrogen). Regions
494 were sequence validated and LR cloned (Invitrogen) into either an *hsp68-lacZ* vector or pXIG vector, with
495 a TdTomato cassette in place of GFP.

496 Generation of transgenic mice and E11.5 embryo staining was performed as previously described^{69–71}
497 using FVB strain mice. Embryos expressing the *lacZ* reporter gene were scored and annotated for their
498 expression patterns by multiple curators. For a construct to be considered positive, a minimum of three
499 embryos per construct were required to demonstrate reporter activity in the same tissue. Mouse
500 transient transgenic assays were approved by the Lawrence Berkeley National Laboratory Animal
501 Welfare and Research Committee.

502 Generation of transgenic zebrafish was performed as previously described⁷² in AB zebrafish. At 3dpf and
503 5dpf, reporter expression patterns were evaluated. For a construct to be considered as positive, $\geq 25\%$ of
504 mosaic embryos had to display reporter activity in one or more anatomical structures. Positive zebrafish
505 were quantified for reporter activity in five anatomical regions (forebrain, midbrain, hindbrain, amacrine
506 cells, spinal cord).

507 **Regulatory vocabulary development**

508 We applied the machine learning algorithm gkm-SVM to our MB and FB catalogues, under default
509 settings. We trained on the sequences underlying the summits ± 250 bp of non-ubiquitously open, top
510 10,000 peaks by signal intensity, versus five negative sets, matched for GC content, length, and repeat
511 content. Weights across all five tests were averaged for all 10-mers.

512 All 10-mers with weight ≥ 1.50 were clustered on sequence similarity using Starcode⁷³, using sphere
513 clustering with distance set to 3. clustalOmega⁷⁴ aligned the sequences within these clusters and
514 MEME⁷⁵, under default parameters, excepting -dna -maxw 12, generated position weight matrices of
515 these aligned clusters. Tomtom⁷⁶, querying the Jolma 2013, JASPAR Core 2014, and Uniprobe mouse
516 databases, identified the top transcription factors corresponding to these position weight matrices,
517 under default parameters excepting -no-ssc -min-overlap 5 -evaluate -thresh 10.0.

518 The same procedure was used to identify transcription factors specifically conveying regulatory potential
519 in the MB library relative to the FB library, except during gkm-SVM training, the positive set was
520 specified as the top 10,000 non-ubiquitously open MB summits and the negative set was specified to be
521 the top 10,000 non-ubiquitously open FB summits, both ± 250 bp.

522 **Transcription factor footprinting**

523 CENTIPEDE⁷⁷ was used to identify footprints. Sequences underlying the deeply sequenced MB library
524 peaks, less those ubiquitously open, were extracted. FIMO⁷⁸, with options --text --parse-genomic-coord,
525 identified all locations underlying ATAC-seq peaks of the motifs identified above. Additionally,
526 conservation data from 30-way vertebrate phastCons was considered in the CENTIPEDE calculations; for
527 each PWM site, those with mean conservation score greater than 0.9 were considered. Finally, the BAM
528 file read end co-ordinates were adjusted in response to the shift in co-ordinates due to the transposase
529 insertion⁷⁹. As such, following the original ATAC-seq method¹⁷, reads were adjusted +4bp on the positive
530 strand and -5bp on the negative strand.

531 **Genome-wide read pileup over predicted motif sites**

532 FIMO, as above, was used to identify all co-ordinates genome wide of the identified motifs. deepTools⁸⁰
533 “bamCoverage” tool was run under default conditions, to convert the deeply sequenced MB library BAM
534 to bigwig format. Following this, a matrix file was generated with “computeMatrix”, with options --
535 referencePoint center -b 1000 -a 1000 -bs 50 specified. Finally, “plotHeatmap” was used to generate
536 plots indicating ATAC-seq read pileup over predicted motif sites.

537 ***In vivo* validation of the MB-specific enhancer**

538 The MB-specific peak was PCR amplified (**Supplementary Table 5**) from human genomic DNA and TA
539 cloned into pCR8 (Invitrogen). Regions were sequence validated and LR cloned (Invitrogen) into either
540 an *hsp68-lacZ* vector or a modified pXIG vector, with a TdTomato cassette in place of GFP.

541 For zebrafish transgenesis, the modified pXIG vector was injected into 1-2 cell stage embryos as
542 previously described⁷² in AB zebrafish. TdTomato reporter expression was assayed at 72hpf and 5dpf;
543 mosaic embryos positive for TdTomato expression were selected and raised to adulthood and founders
544 were identified. Progeny of founders were screened at 72hpf for reporter activity.

545 For mouse transgenesis, the generated *hsp68-lacZ* vector was purified in a double CsCl gradient
546 (Lofstrand Labs Ltd) and stable mouse transgenesis was performed in C57BL/6 mice by Cyagen
547 Biosciences Inc. Multiple founder lines were generated. For lacZ staining, embryos were collected at
548 E12.5, and mouse brains were isolated at E15.5, P7, P30, and P574. Brains were roughly sectioned in
549 1mm sections at P7 and P30 and animals were perfused at P574 and fixed brains were sectioned
550 (200 μ m) with a vibratome. Specimens were subsequently fixed for 2 hours on ice in 1% formaldehyde,
551 0.2% glutaraldehyde, 0.02% Igepal CA-630 in PBS. Following fixation, tissues were permeabilized over
552 3x15 minute washes in 2mM MgCl₂ and 0.02% Igepal CA-630 in PBS at room temperature.
553 Embryos/tissues were incubated overnight at 37°C in staining solution, containing 320 μ g/mL X-Gal in
554 N,N-dimethyl formamide, 12mM K-ferricyanide, 12mM K-ferrocyanide, 0.002% Igepal CA-630, 4mM
555 MgCl₂ in PBS. Specimens were washed in 0.2% Igepal CA-630 in PBS over 2x30 minutes and finally stored
556 in 4% formaldehyde, 100mM sodium phosphate, and 10% methanol.

557 **Patient sequencing and genotyping at *SNCA***

558 A total of 986 PD patients and 992 controls who were seen at the Mayo Clinic in Jacksonville, FL were
559 sequenced across the putative enhancer and genotyped for 25 variants across the *SNCA* locus. The
560 variants chosen for genotyping were confirming those identified by sequencing of the enhancer as well
561 as assessing those identified in Guella *et al*⁴¹. For PD patients, median age at blood draw was 69 years
562 (Range: 28-97 years), median age at PD onset was 67 years (Range: 28-97 years), and 631 patients
563 (64.0%) were male. Median age at blood draw in controls was 67 years (Range: 18-92 years) and 415
564 subjects (41.8%) were male. Patients were diagnosed with PD using standard clinical criteria⁸¹. All
565 subjects are unrelated non-Hispanic Caucasians of European descent. The Mayo Clinic Institutional
566 Review Board approved the study and all subjects provided written informed consent.

567 Genomic DNA was extracted from whole blood using the Autogen FlexStar. Sanger sequencing of the
568 enhancer region was performed bidirectionally using the ABI 3730xl DNA analyzer (Applied Biosystems)
569 standard protocol. Sequence data was analyzed using SeqScape v2.5 (Applied Biosystems). Statistical
570 analyses were performed using both SAS and R⁸². Of the variants identified within the enhancer, only

571 those with minor allele frequency greater than 5% were evaluated for association with PD in single-
572 variant analysis. Associations between individual variants and PD were evaluated using logistic
573 regression models, adjusted for age at blood draw and sex, and where variants were considered, under
574 an additive model (i.e. effect of each additional minor allele). Odds ratios and 95% confidence intervals
575 were estimated and a Bonferroni correction for multiple testing, due to the four common variants that
576 were evaluated for association with PD, was utilized in single-variant analysis, after which p-values \leq
577 0.0125 were considered as statistically significant.

578 Genotyping the 25 SNPs across the SNCA locus was performed using the iPLEX Gold protocol on the
579 MassARRAY System and analysed with TYPER 4.0 software (Agena Bioscience). For the 25 SNPs
580 genotyped across the SNCA locus, all genotype call rates were $>95\%$ and there was no evidence for
581 departure from Hardy-Weinberg equilibrium (all χ^2 p-values > 0.05 after Bonferroni correction).
582 Haplotype frequencies in cases and controls was estimated using the haplo.stats package⁸³ function
583 “haplo.group”. Associations between haplotypes and risk of PD were evaluated using score tests of
584 association⁸⁴ using the “haplo.score” function. Tests were adjusted for age at blood draw and sex,
585 haplotypes occurring in less than 1% of subjects were excluded, and only individuals with no missing
586 genotype calls for any variants were included. A Bonferroni correction for multiple testing was applied,
587 after which p-values ≤ 0.0042 were considered as statistically significant, due to the 12 different
588 common haplotypes that were observed and tested for association with PD risk.

589 **Protein array testing differential binding**

590 HuProt v3.1 human proteome microarrays printed on the PATH surface containing $>20,000$ unique
591 proteins representing 16,152 genes (CDI laboratories)⁴⁰ were blocked with 25mM HEPES pH 8.0, 50mM
592 K Glutamate, 8mM MgCl₂, 3mM DTT, 10% glycerol, 0.1% Triton X-100, 3% BSA on an orbital shaker at
593 4°C for ≥ 3 hours. Allele specific protein-DNA binding interactions were identified through dye-swap
594 competition of major and minor alleles labeled with either Cy3 or Cy5. DNA fragments for rs2737024
595 and rs2583959 were synthesized with the SNP for each allele flanked by 15 nucleotides of the upstream
596 and downstream sequence and a common priming site at the 3' end (**Supplementary Table 5**).

597 The dsDNA fragments were created by separately annealing a primer containing a Cy3 or Cy5 label and
598 adding Klenow (NEB) with dNTP to fill-in the complementary strand for each allele⁸⁵. Cy3 labeled major
599 allele was mixed with Cy5 labeled minor allele (each at 40nM) in 1x hybridization buffer (10mM TrisCl
600 pH 8, 50mM KCl, 1mM MgCl₂, 1mM DTT, 5% glycerol, 10 μ M ZnCl₂, 3mg/mL BSA) and added to an array,

601 dyes were then swapped for each allele and the mixture was then added to a second array. DNA was
602 allowed to bind overnight at 4°C on an orbital shaker with protection from light. Chips were washed
603 once with cold 1xTBS (0.1% Triton X-100) for 5 minutes at 4°C, rinsed, and dried in the centrifuge. Cy5
604 and Cy3 images were taken separately on a Genepix 4000B scanner and, after alignment to the .gal file,
605 individual spot intensities were extracted using the Genepix Pro software.

606 Allele specific interactions were identified through dye swap analysis. The ratio of major/minor allele
607 binding was calculated using the duplicate spot average median foreground signal for each protein

608 according to the following equation:
$$\log_2 \sqrt{\frac{Cy3_{major} * Cy5_{major}}{Cy3_{minor} * Cy5_{minor}}}$$

609 Mean intensity was calculated by averaging the foreground signal for the Cy3 and Cy5 channels of the
610 major and minor alleles. MA plots were made for each allele using the calculated mean intensity and the
611 log ratio of the major/minor allele.

612

613 **DATA AVAILABILITY**

614 ATAC-sequencing and RNA-sequencing data will be available at the Gene Expression Omnibus (GEO)
615 under the accession number GSEXXXXXX.

616 **ACKNOWLEDGEMENTS**

617 This research undertaken at Johns Hopkins University School of Medicine was supported in part by
618 awards from NIH (NS62972 and MH106522 to ASM; GM111514 to HZ; HG007348 to MAB). The Mayo
619 Clinic collection was supported in part by a Morris K. Udall Center of Excellence in Parkinson's disease
620 Research (P50 NS072187), American Parkinson's Disease Association Center and The Mangurian
621 Foundation for Lewy body research. OAR is supported by NS078086 and NS10069 (NIH), W81XWH-17-1-
622 0249 (Department of Defense), The Michael J. Fox Foundation and The Little Family Foundation. ZKW is
623 supported by the Mayo Clinic Center for Regenerative Medicine, Mayo Clinic Center for Individualized
624 Medicine, Mayo Clinic Neuroscience Focused Research Team (Cecilia and Dan Carmichael Family
625 Foundation, and the James C. and Sarah K. Kennedy Fund for Neurodegenerative Disease Research at
626 Mayo Clinic in Florida), the gift from Carl Edward Bolch, Jr., and Susan Bass Bolch, The Sol Goldman
627 Charitable Trust, and Donald G. and Jodi P. Heeringa. Research conducted at the E.O. Lawrence Berkeley
628 National Laboratory was performed under U.S. Department of Energy contract DE-AC02-05CH11231,
629 University of California and was supported by HG003988 (NIH) to LAP.

630 **AUTHOR CONTRIBUTIONS**

631 SAM, ASM designed the study and wrote the paper. SAM, PWH, XR, WDL, SJK, and ELW performed
632 various experiments. Transgenic experiments were performed by SAM, NJB, and JFT (zebrafish) and by
633 SAM, JAA, DED, AV, and LP (mice). SAM and MAB performed the gkm-SVM analyses. PD patient
634 sequencing and analysis was performed by SAM, AIS, MGH, NND, ZKW, and OAR. SAM, CDM, and HZ
635 performed and analysed differential protein array binding assays. SAM implemented the computational
636 algorithms to process the raw data and conduct analyses thereof. SAM and ASM analyzed and
637 interpreted the resulting data. SAM contributed novel computational pipeline development. SAM and
638 ASM wrote the manuscript, with all other authors contributing. Correspondence to ASM
639 (andy@jhmi.edu).

640 **FINANCIAL INTERESTS STATEMENT**

641 The authors declare no competing financial interests.

642 **REFERENCES**

- 643 1. Ma, S. Y., Roytta, M., Rinne, J. O., Collan, Y. & Rinne, U. K. Correlation between
644 neuromorphometry in the substantia nigra and clinical features in Parkinson’s disease using
645 disector counts. *J. Neurol. Sci.* **151**, 83–87 (1997).
- 646 2. Fearnley, J. M. & Lees, A. J. Ageing and Parkinson’s disease: substantia nigra regional selectivity.
647 *Brain* **114 (Pt 5)**, 2283–2301 (1991).
- 648 3. Pringsheim, T., Jette, N., Frolkis, A. & Steeves, T. D. L. The prevalence of Parkinson’s disease: a
649 systematic review and meta-analysis. *Mov. Disord.* **29**, 1583–1590 (2014).
- 650 4. Thomas, B. & Beal, M. F. Parkinson’s disease. *Hum. Mol. Genet.* **16 Spec No**, R183-94 (2007).
- 651 5. Zarranz, J. J. *et al.* The new mutation, E46K, of alpha-synuclein causes Parkinson and Lewy body
652 dementia. *Ann. Neurol.* **55**, 164–173 (2004).
- 653 6. Kruger, R. *et al.* Ala30Pro mutation in the gene encoding alpha-synuclein in Parkinson’s disease.
654 *Nature Genetics* **18**, 106–108 (1998).
- 655 7. Polymeropoulos, M. H. *et al.* Mutation in the alpha-synuclein gene identified in families with
656 Parkinson’s disease. *Science* **276**, 2045–2047 (1997).
- 657 8. Singleton, A. B. *et al.* alpha-Synuclein locus triplication causes Parkinson’s disease. *Science* **302**,
658 841 (2003).
- 659 9. Nalls, M. A. *et al.* Large-scale meta-analysis of genome-wide association data identifies six new
660 risk loci for Parkinson’s disease. *Nat. Genet.* **46**, 989–993 (2014).
- 661 10. Chang, D. *et al.* A meta-analysis of genome-wide association studies identifies 17 new Parkinson’s
662 disease risk loci. *Nat. Genet.* **49**, 1511–1516 (2017).
- 663 11. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in
664 regulatory DNA. *Science* **337**, 1190–1195 (2012).
- 665 12. Schaub, M. A., Boyle, A. P., Kundaje, A., Batzoglou, S. & Snyder, M. Linking disease associations
666 with regulatory information in the human genome. *Genome Res.* **22**, 1748–1759 (2012).
- 667 13. Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature*
668 **473**, 43–49 (2011).

- 669 14. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An Expanded View of Complex Traits: From Polygenic to
670 Omnigenic. *Cell* **169**, 1177–1186 (2017).
- 671 15. Lee, D. *et al.* A method to predict the impact of regulatory variants from DNA sequence. *Nat.*
672 *Genet.* **47**, 955–961 (2015).
- 673 16. Praetorius, C. *et al.* A polymorphism in IRF4 affects human pigmentation through a tyrosinase-
674 dependent MITF/TFAP2A pathway. *Cell* **155**, 1022–1033 (2013).
- 675 17. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition of native
676 chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins
677 and nucleosome position. *Nat. Methods* **10**, 1213–8 (2013).
- 678 18. Heintz, N. Gene expression nervous system atlas (GENSAT). *Nat. Neurosci.* **7**, 483 (2004).
- 679 19. Robinson, J. T. *et al.* Integrative genomics viewer. *Nature Biotechnology* **29**, 24–26 (2011).
- 680 20. Thorvaldsdottir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV): high-
681 performance genomics data visualization and exploration. *Brief. Bioinform.* **14**, 178–192 (2013).
- 682 21. Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome. *Nature* **489**, 75–
683 82 (2012).
- 684 22. Visel, A., Minovitsky, S., Dubchak, I. & Pennacchio, L. A. VISTA Enhancer Browser--a database of
685 tissue-specific human enhancers. *Nucleic Acids Res.* **35**, D88–92 (2007).
- 686 23. Stott, S. R. W. *et al.* Foxa1 and foxa2 are required for the maintenance of dopaminergic
687 properties in ventral midbrain neurons at late embryonic stages. *J. Neurosci.* **33**, 8022–8034
688 (2013).
- 689 24. Arenas, E. Foxa2: the rise and fall of dopamine neurons. *Cell Stem Cell* **2**, 110–112 (2008).
- 690 25. Prakash, N. & Wurst, W. Development of dopaminergic neurons in the mammalian brain. *Cell.*
691 *Mol. Life Sci.* **63**, 187–206 (2006).
- 692 26. Smits, S. M., Ponnio, T., Conneely, O. M., Burbach, J. P. H. & Smidt, M. P. Involvement of Nurr1 in
693 specifying the neurotransmitter identity of ventral midbrain dopaminergic neurons. *Eur. J.*
694 *Neurosci.* **18**, 1731–1738 (2003).
- 695 27. Hook, P. W. *et al.* Single-Cell RNA-Seq of Mouse Dopaminergic Neurons Informs Candidate Gene

- 696 Selection for Sporadic Parkinson Disease. *Am. J. Hum. Genet.* **102**, 427–446 (2018).
- 697 28. Ghandi, M., Lee, D., Mohammad-Noori, M. & Beer, M. A. Enhanced regulatory sequence
698 prediction using gapped k-mer features. *PLoS Comput. Biol.* **10**, e1003711 (2014).
- 699 29. Kittappa, R., Chang, W. W., Awatramani, R. B. & McKay, R. D. G. The *foxa2* gene controls the birth
700 and spontaneous degeneration of dopamine neurons in old age. *PLoS Biol.* **5**, e325 (2007).
- 701 30. Caiazzo, M. *et al.* Direct generation of functional dopaminergic neurons from mouse and human
702 fibroblasts. *Nature* **476**, 224–227 (2011).
- 703 31. Wang, Y. *et al.* The 3D Genome Browser: a web-based browser for visualizing 3D genome
704 organization and long-range chromatin interactions. *bioRxiv* (2017). at
705 <<http://biorxiv.org/content/early/2017/02/27/112268.abstract>>
- 706 32. Zarow, C., Lyness, S. A., Mortimer, J. A. & Chui, H. C. Neuronal loss is greater in the locus
707 coeruleus than nucleus basalis and substantia nigra in Alzheimer and Parkinson diseases. *Arch.*
708 *Neurol.* **60**, 337–341 (2003).
- 709 33. Kasthuber, E., Kratochwil, C. F., Ryu, S., Schweitzer, J. & Driever, W. Genetic dissection of
710 dopaminergic and noradrenergic contributions to catecholaminergic tracts in early larval
711 zebrafish. *J. Comp. Neurol.* **518**, 439–458 (2010).
- 712 34. Rink, E. & Wullimann, M. F. The teleostean (zebrafish) dopaminergic system ascending to the
713 subpallium (striatum) is located in the basal diencephalon (posterior tuberculum). *Brain Res.* **889**,
714 316–330 (2001).
- 715 35. Seidel, K. *et al.* The brainstem pathologies of Parkinson’s disease and dementia with Lewy bodies.
716 *Brain Pathol.* **25**, 121–135 (2015).
- 717 36. Wakabayashi, K., Mori, F., Tanji, K., Orimo, S. & Takahashi, H. Involvement of the peripheral
718 nervous system in synucleinopathies, tauopathies and other neurodegenerative proteinopathies
719 of the brain. *Acta Neuropathol.* **120**, 1–12 (2010).
- 720 37. Wakabayashi, K. & Takahashi, H. Neuropathology of autonomic nervous system in Parkinson’s
721 disease. *Eur. Neurol.* **38 Suppl 2**, 2–7 (1997).
- 722 38. Braak, H. *et al.* Amygdala pathology in Parkinson’s disease. *Acta Neuropathol.* **88**, 493–500
723 (1994).

- 724 39. Langston, J. W. & Forno, L. S. The hypothalamus in Parkinson disease. *Ann. Neurol.* **3**, 129–133
725 (1978).
- 726 40. Jeong, J. S. *et al.* Rapid identification of monospecific monoclonal antibodies using a human
727 proteome microarray. *Mol. Cell. Proteomics* **11**, O111.016253 (2012).
- 728 41. Guella, I. *et al.* alpha-synuclein genetic variability: A biomarker for dementia in Parkinson disease.
729 *Ann. Neurol.* **79**, 991–999 (2016).
- 730 42. Preissl, S. *et al.* Single-nucleus analysis of accessible chromatin in developing mouse forebrain
731 reveals cell-type-specific transcriptional regulation. *Nat. Neurosci.* **21**, 432–439 (2018).
- 732 43. Sugiaman-Trapman, D. *et al.* Characterization of the human RFX transcription factor family by
733 regulatory and target gene analysis. *BMC Genomics* **19**, 181 (2018).
- 734 44. Sung, M.-H., Guertin, M. J., Baek, S. & Hager, G. L. DNase footprint signatures are dictated by
735 factor dynamics and DNA sequence. *Mol. Cell* **56**, 275–285 (2014).
- 736 45. Greffard, S. *et al.* Motor score of the Unified Parkinson Disease Rating Scale as a good predictor
737 of Lewy body-associated neuronal loss in the substantia nigra. *Arch. Neurol.* **63**, 584–588 (2006).
- 738 46. Soldner, F. *et al.* Parkinson-associated risk variant in distal enhancer of alpha-synuclein
739 modulates target gene expression. *Nature* **533**, 95–99 (2016).
- 740 47. The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human
741 genome. *Nature* **489**, 57–74 (2012).
- 742 48. Gupta, R. M. *et al.* A Genetic Variant Associated with Five Vascular Diseases Is a Distal Regulator
743 of Endothelin-1 Gene Expression. *Cell* **170**, 522–533.e15 (2017).
- 744 49. Westerfeld, M. *The Zebrafish Book. A Guide for the Laboratory Use of Zebrafish (Danio rerio).*
745 (Univ. Oregon Press, Eugene, 2007).
- 746 50. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–
747 359 (2012).
- 748 51. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079
749 (2009).
- 750 52. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).

- 751 53. Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory
752 requirements. *Nat. Methods* **12**, 357–360 (2015).
- 753 54. Liao, Y., Smyth, G. K. & Shi, W. The Subread aligner: fast, accurate and scalable read mapping by
754 seed-and-vote. *Nucleic Acids Res.* **41**, e108 (2013).
- 755 55. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor. *Nat.*
756 *Methods* **12**, 115–121 (2015).
- 757 56. Gentleman, R. C. *et al.* Bioconductor: open software development for computational biology and
758 bioinformatics. *Genome Biol.* **5**, R80 (2004).
- 759 57. Mudge, J. M. & Harrow, J. Creating reference gene annotation for the mouse C57BL6/J genome
760 assembly. *Mamm. Genome* **26**, 366–378 (2015).
- 761 58. Zhao, H. *et al.* CrossMap: a versatile tool for coordinate conversion between genome assemblies.
762 *Bioinformatics* **30**, 1006–1007 (2014).
- 763 59. Ye, J. *et al.* Primer-BLAST: a tool to design target-specific primers for polymerase chain reaction.
764 *BMC Bioinformatics* **13**, 134 (2012).
- 765 60. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
766 *Bioinformatics* **26**, 841–842 (2010).
- 767 61. Hansen, K. D., Irizarry, R. A. & Wu, Z. Removing technical variability in RNA-seq data using
768 conditional quantile normalization. *Biostatistics* **13**, 204–216 (2012).
- 769 62. Wei, T. & Simko, V. corrplot: Visualization of a Correlation Matrix. (2016). at <[https://cran.r-](https://cran.r-project.org/package=corrplot)
770 [project.org/package=corrplot](https://cran.r-project.org/package=corrplot)>
- 771 63. Neuwirth, E. RColorBrewer: ColorBrewer Palettes. (2014). at <[https://cran.r-](https://cran.r-project.org/package=RColorBrewer)
772 [project.org/package=RColorBrewer](https://cran.r-project.org/package=RColorBrewer)>
- 773 64. Schwalb, B., Tresch, A., Torkler, P., Duemcke, S. & Demel, C. LSD: Lots of Superior Depictions.
774 (2015). at <<https://cran.r-project.org/package=LSD>>
- 775 65. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
776 genomes. *Genome Res.* **15**, 1034–1050 (2005).
- 777 66. Liu, T. *et al.* Cistrome: an integrative platform for transcriptional regulation studies. *Genome Biol.*

- 778 **12**, R83 (2011).
- 779 67. Karolchik, D. *et al.* The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.* **32**, D493-6
780 (2004).
- 781 68. McLean, C. Y. *et al.* GREAT improves functional interpretation of cis-regulatory regions. *Nat.*
782 *Biotechnol.* **28**, 495–501 (2010).
- 783 69. Pennacchio, L. A. *et al.* In vivo enhancer analysis of human conserved non-coding sequences.
784 *Nature* **444**, 499–502 (2006).
- 785 70. Poulin, F. *et al.* In vivo characterization of a vertebrate ultraconserved enhancer. *Genomics* **85**,
786 774–781 (2005).
- 787 71. Kothary, R. *et al.* Inducible expression of an hsp68-lacZ hybrid gene in transgenic mice.
788 *Development* **105**, 707–714 (1989).
- 789 72. Fisher, S. *et al.* Evaluating the biological relevance of putative enhancers using Tol2 transposon-
790 mediated transgenesis in zebrafish. *Nat. Protoc.* **1**, 1297–1305 (2006).
- 791 73. Zorita, E., Cusco, P. & Fillion, G. J. Starcode: sequence clustering based on all-pairs search.
792 *Bioinformatics* **31**, 1913–1919 (2015).
- 793 74. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence alignments
794 using Clustal Omega. *Mol. Syst. Biol.* **7**, 539 (2011).
- 795 75. Bailey, T. L. & Elkan, C. Fitting a mixture model by expectation maximization to discover motifs in
796 biopolymers. *Proceedings. Int. Conf. Intell. Syst. Mol. Biol.* **2**, 28–36 (1994).
- 797 76. Gupta, S., Stamatoyannopoulos, J. A., Bailey, T. L. & Noble, W. S. Quantifying similarity between
798 motifs. *Genome Biol.* **8**, R24 (2007).
- 799 77. Pique-Regi, R. *et al.* Accurate inference of transcription factor binding from DNA sequence and
800 chromatin accessibility data. *Genome Res.* **21**, 447–455 (2011).
- 801 78. Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif.
802 *Bioinformatics* **27**, 1017–1018 (2011).
- 803 79. Adey, A. *et al.* Rapid, low-input, low-bias construction of shotgun fragment libraries by high-
804 density in vitro transposition. *Genome Biol.* **11**, R119 (2010).

- 805 80. Ramirez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data analysis.
806 *Nucleic Acids Res.* **44**, W160-5 (2016).
- 807 81. Hughes, A. J., Daniel, S. E., Kilford, L. & Lees, A. J. Accuracy of clinical diagnosis of idiopathic
808 Parkinson's disease: a clinico-pathological study of 100 cases. *J. Neurol. Neurosurg. Psychiatry* **55**,
809 181–184 (1992).
- 810 82. R Core Team. R: A language and environment for statistical computing. (2017). at
811 <<https://www.r-project.org/>>
- 812 83. Sinnwell, J. P. & Schaid, D. J. haplo.stats: Statistical Analysis of Haplotypes with Traits and
813 Covariates when Linkage Phase is Ambiguous. (2016). at <[https://cran.r-](https://cran.r-project.org/package=haplo.stats)
814 [project.org/package=haplo.stats](https://cran.r-project.org/package=haplo.stats) >
- 815 84. Schaid, D. J., Rowland, C. M., Tines, D. E., Jacobson, R. M. & Poland, G. A. Score tests for
816 association between traits and haplotypes when linkage phase is ambiguous. *Am. J. Hum. Genet.*
817 **70**, 425–434 (2002).
- 818 85. Hu, S. *et al.* DNA methylation presents distinct binding sites for human transcription factors. *Elife*
819 **2**, e00726 (2013).
- 820