

# 1 *De novo* human brain enhancers created by single 2 nucleotide mutations

3  
4 Shan Li<sup>1,2</sup>, Sridhar Hannenhalli<sup>2,\*</sup>, Ivan Ovcharenko<sup>1,\*</sup>

5 <sup>1</sup> Computational Biology Branch, National Center for Biotechnology Information, National  
6 Library of Medicine, National Institutes of Health, Bethesda, MD, 20892, USA.

7 <sup>2</sup> Cancer Data Science Laboratory, Center for Cancer Research, National Cancer Institute,  
8 National Institutes of Health, Bethesda, MD, 20892, USA.

9 \*Correspondence: [ovcharen@nih.gov](mailto:ovcharen@nih.gov) and [sridhar.hannenhalli@nih.gov](mailto:sridhar.hannenhalli@nih.gov).

10  
11

## 12 **Abstract**

13 Advanced human cognition is attributed to increased neocortex size and complexity, but the  
14 underlying gene regulatory mechanisms are unknown. Using deep learning model of embryonic  
15 neocortical enhancers, and human and macaque embryonic neocortex H3K27ac data, we  
16 identified ~4000 enhancers gained *de novo* in the human, largely attributable to single-  
17 nucleotide essential mutations. The genes near *de novo* gained enhancers exhibit increased  
18 expression in human embryonic neocortex relative to macaque, are involved in critical neural  
19 developmental processes, and are expressed specifically in the progenitor cells and  
20 interneurons. The gained enhancers, especially the essential mutations, are associated with  
21 central nervous system disorders/traits. Integrative computational analyses suggest that the  
22 essential mutations establish enhancer activities through affecting binding of key transcription  
23 factors of embryonic neocortex. Overall, our results suggest that non-coding mutations may  
24 have led to *de novo* enhancer gains in the embryonic human neocortex, that orchestrate the  
25 expression of genes involved in critical developmental processes associated with human  
26 cognition.

27

## 28 **Introduction**

29 The neocortex is a mammalian innovation enabling complex cognitive and motor tasks  
30 (Geschwind and Rakic 2013; Emera et al. 2016). The substantial expansion and functional  
31 elaboration of the neocortex provides an essential basis for the advanced cognitive abilities of  
32 humans (Geschwind and Rakic 2013), which includes an increase in the proliferative capacity of  
33 the progenitor cells (Dehay et al. 2015; Namba and Huttner 2017; Sousa et al. 2017), an  
34 increase in the duration of their proliferative, neurogenic and gliogenic phases (Lewitus et al.  
35 2014; Otani et al. 2016), an increase in the number and diversity of progenitors, modification of  
36 neuronal migration, and establishment of new connections among functional areas (Geschwind  
37 and Rakic 2013).

38

39 Critical events in corticogenesis, including specification of cortical areas and differentiation of  
40 cortical layers require precise spatiotemporal orchestration of gene expression (Rakic et al.  
41 2009). Modifications in gene regulation are thus hypothesized to be a major source of  
42 evolutionary innovation during cortical development (Rakic 2009; Rakic et al. 2009; Geschwind  
43 and Rakic 2013). Among these are gain and loss of enhancers, repurposing of existing  
44 enhancers, rewiring of enhancer-gene interaction networks, and modifications of crosstalk  
45 between enhancers operating within the same cis-regulatory landscape (Long et al. 2016).  
46 However, several fundamental questions remain open: to what extent the evolutionary gain  
47 and loss of enhancers has contributed to human-specific features of corticogenesis?  
48 Specifically, how often enhancer gain is associated with an increased expression of the target  
49 gene involved in human corticogenesis? To what extent the emergence of human-specific  
50 enhancers could be explained by a single or a few single-nucleotide mutations? How often do  
51 such mutations establish an enhancer from neutral DNA through creation of binding sites of  
52 activators as opposed to the disruption of binding sites of repressors? What are the  
53 transcription factors (TFs) mediating critical enhancer gains and losses and what gene  
54 regulatory networks are induced by such mutations? A previous study identified Human Gained  
55 Enhancers (termed HGEs) (Reilly et al. 2015) that exhibit increased regulatory activity in human  
56 relative to macaque and mouse. In contrast, our focus is *de novo* gained enhancers in human  
57 that presumably originate from neutral non-coding sequence via minimum number of single-  
58 nucleotide substitutions along the human lineage. Besides the availability of enhancer activity  
59 profiles in the developing brain of humans and macaques (Reilly et al. 2015), a quantitative  
60 model that can accurately estimate enhancer activity from DNA sequence, with single-  
61 nucleotide sensitivity, is critical to answering the questions above.

62  
63 In this study, we developed a deep learning model (DLM) able to learn the sequence encryption  
64 of human and primate embryonic neocortex enhancers, enabling us to quantify the functional  
65 effect of single nucleotide mutations on enhancer activity. Leveraging the DLM and the recently  
66 available enhancer activity profiles in developing neocortex in humans and macaques (Reilly et  
67 al. 2015), we identified single-nucleotide mutations that potentially drive human-specific  
68 regulatory innovations. We observed that a single-nucleotide mutation is often sufficient to  
69 give rise to an enhancer, leading to increased expression of the proximal target gene. As a  
70 group, *de novo* gained enhancers induce genes that are critical to cognitive function and are  
71 expressed preferentially in the progenitor and interneuron cells of the developing neocortex.  
72 *De novo* gained enhancers and their target genes induce and mediate a potential core  
73 regulatory network in the developing human neocortex, with POU3F2 occupying a central  
74 position. Essential single-nucleotide mutations resulting in *de novo* enhancer gain exhibit  
75 relaxed negative, or potentially adaptive, selection. Interestingly, the essential mutations and  
76 *de novo* gained enhancers are enriched for cognitive traits; in particular, the *de novo* gained  
77 enhancers associated with regulation of key TFs are enriched for *de novo* mutations in patients  
78 with the autism spectrum disorder (ASD). Compared to HGEs, although *de novo* gained  
79 enhancers have relatively weaker enhancer activity, they are more likely to turn on gene  
80 expression in human and regulate genes associated with brain development. Integrating a DLM  
81 with epigenomic data allowed us not only to identify *de novo* gained human-specific enhancers

82 that might underlie advanced cognition, but also gauge the impact of single-nucleotide  
83 mutations in this process.

84

85 Overall, our results, based on the H3K27ac profiles in developing human and macaque brain,  
86 and a novel sequence-specific deep learning model of embryonic neocortical enhancers,  
87 suggests a wide-spread *de novo* gain in enhancers, largely driven by single nucleotide  
88 mutations, in the progenitors and interneurons of the developing human neocortex, that  
89 together induce a core regulatory network that associated with human cognitive abilities as  
90 well as cognitive disorders.

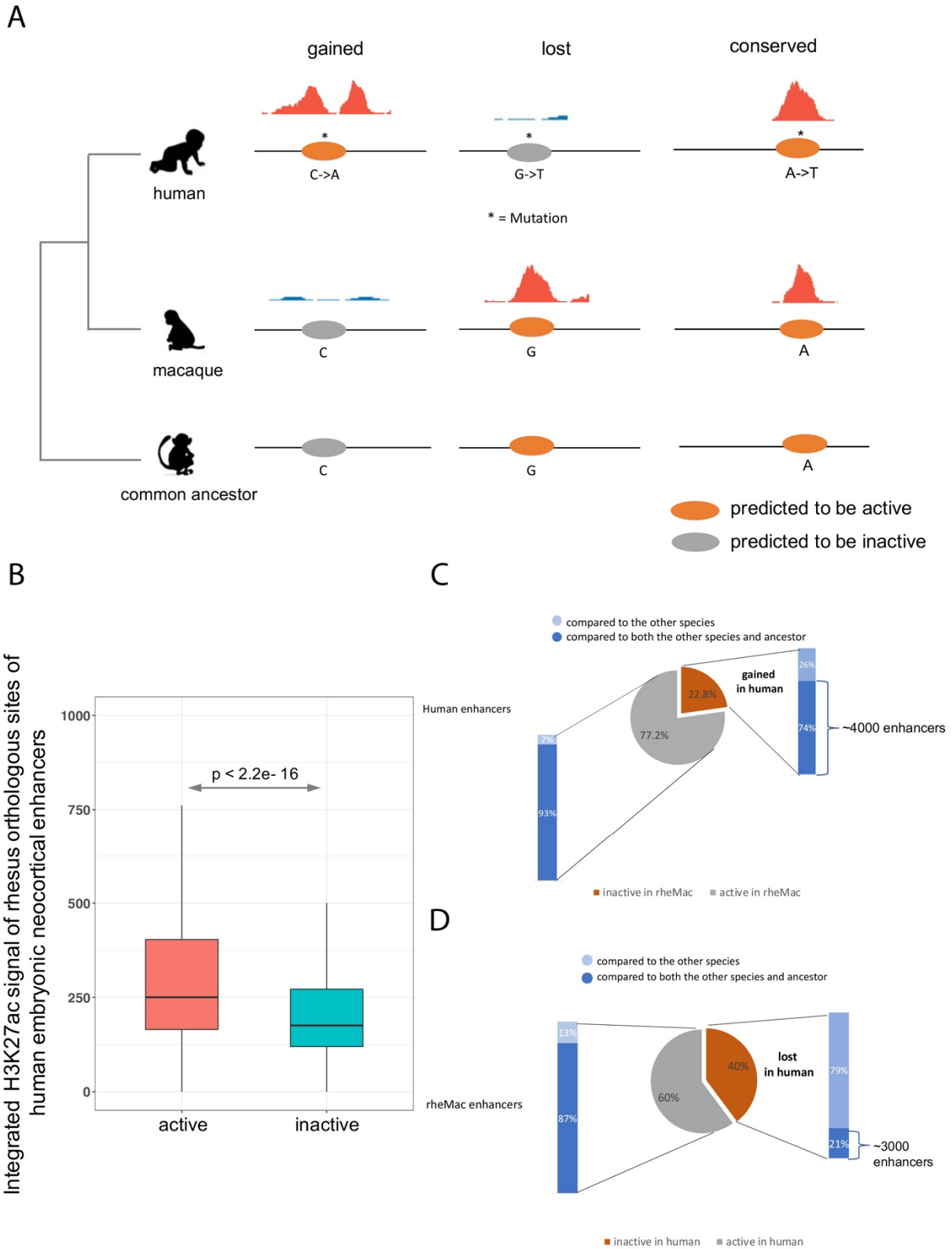
91

## 92 **Results**

93

### 94 **Identifying *de novo* enhancer gain and essential human mutations - Overview**

95 To assess functional impact of single nucleotide mutations on enhancer activity, we leveraged  
96 the H3K27ac ChIP-seq data during human and macaque corticogenesis as a proxy for active  
97 enhancers (Reilly et al. 2015) and built a DLM to learn the regulatory code encrypted in the  
98 enhancer sequences (Figure S1A-C, Methods). Next, by integrating the predicted enhancer  
99 activities in human, macaque, and the human-macaque common ancestor inferred from  
100 multiple sequence alignment (Paten et al. 2008) based on a probabilistic model (Holmes and  
101 Bruno 2001; Holmes 2003; Bradley and Holmes 2007) with the observed enhancer activities in  
102 human and macaque, we identified human-specific *de novo* gains and losses of enhancers  
103 (Figure 1A, Methods). We then prioritized the single-nucleotide human-macaque mutations in  
104 the *de novo* gained and lost enhancers based on the difference of the DLM scores between the  
105 macaque sequence and the intermediate sequence with one or more introduced human  
106 allele(s). For an enhancer with multiple mutations, which was either gained or lost in the  
107 human genome, we first introduced each human-specific allele to its matching macaque  
108 sequence and estimated its impact on enhancer activity using the difference in the DLM score  
109 attributed to the human allele. By iteratively increasing the number of introduced human-  
110 specific alleles and scoring the modified sequence, we evaluated the impact of combinations of  
111 mutations and determined the minimal number of mutations needed for an enhancer to be  
112 gained or lost in the human lineage.



113  
114  
115  
116  
117

**Figure 1. *De novo* gained and lost enhancers.** A) Identification of *de novo*-gained, lost, and conserved enhancers. If a human enhancer scored highly by the DLM and scored low both in macaque and in the common ancestor, and was not detected by H3K27ac in macaque, it was considered to be gained in humans. If a macaque enhancer having high DLM score, scored high in



118 common ancestor, scored low in human and was undetectable by H3K27ac in human it was considered a loss in human. The  
119 enhancers that are detected by H3K27ac in both human and macaque and scored highly in all three genomes were called  
120 conserved enhancers. B) Comparison of embryonic macaque neocortex integrated H3K27ac signal intensities (within the 1kb  
121 enhancers) between the predicted active and inactive macaque orthologs of human embryonic neocortex enhancers. C) The  
122 fraction of *de novo* gained human embryonic neocortex enhancers by comparing human to both macaque and their common  
123 ancestor. Specifically, 74% of human enhancers that are inactive in macaque are active in the common ancestor and 93% of  
124 human enhancers that are active in macaque are active in the common ancestor. D) The fraction of lost human embryonic  
125 neocortex enhancers by comparing human to both rhesus macaque and their common ancestor. Specifically, 21% of macaque  
126 enhancers that are inactive in human are active in the common ancestor and 87% of macaque enhancers that are active in  
127 human are active in the common ancestor. Light blue refers to relative to the other species, dark blue refers to relative to both  
128 the other species and common ancestor.

129

### 130 **An accurate DLM of embryonic neocortex enhancers in human and macaque**

131 The human embryonic neocortex H3K27ac ChIP-seq peaks were obtained from the four  
132 temporal/spatial groups: the whole cortex at 7 post conception weeks (p.c.w.) (CS16) and 8.5  
133 p.c.w. (CS23) and primitive frontal and occipital tissues from 12 p.c.w. (F2F and F2O) (Reilly et  
134 al. 2015). We trained a DLM separately for each set of enhancers (Methods). The DLM was able  
135 to discriminate human embryonic neocortex enhancers from accessible regions devoid of non-  
136 fetal-brain-enhancer with high accuracy: the area under the receiver operating characteristic  
137 curve (auROC) ranges from 0.9 to 0.94 (Figure S1B), and the area under the precision-recall  
138 curve (auPRC, expectation = 0.091) ranges from 0.56 to 0.63 for the four datasets (Figure S1C).  
139 The consistently high accuracy of all models showed the ability of DLMs in capturing sequence  
140 signatures of brain enhancers similarly to previous modeling of enhancers in other cells and  
141 tissues (Supplementary Results 1), and prompted us to conjecture that the four groups of  
142 enhancers tend to share either genomic locations or sequence characteristics. To assess their  
143 sequence similarity, we trained the DLM on one set and predicted those from all other sets. We  
144 observed both high auROCs and auPRCs (Figure S1D), strongly suggesting a shared sequence  
145 characteristics across the four enhancer sets. However, the genomic overlap between any two  
146 groups of enhancers is relatively low (20-40%) (Figure S1E), indicating that the four sets of  
147 enhancers overlap only partially but share sequence characteristics.

148

149 We proceeded to investigate the *de novo*-gain and loss of enhancers by comparing human 8.5  
150 p.c.w (CS23) sample and macaque sample at approximately matching time point (e55) (Reilly et  
151 al. 2015), as the DLM trained on CS23 has not only high auROC (0.92) but also the highest  
152 precision at a low false positive rate (FPR = 0.1) (Figure S1BC). To ascertain that the DLM trained  
153 on CS23 can accurately predict the enhancer activity in macaque, we scored the macaque  
154 orthologs of CS23 enhancers and compared the e55 H3K27ac signal intensities of the macaque  
155 orthologs predicted to be active with those predicted to be inactive (Methods). The predicted  
156 active regions indeed have significantly stronger H3K27ac signal (Figure 1B), suggesting that the  
157 DLM learned from human embryonic neocortical enhancers can accurately gauge the enhancer  
158 activity in macaque from its genomic sequence.

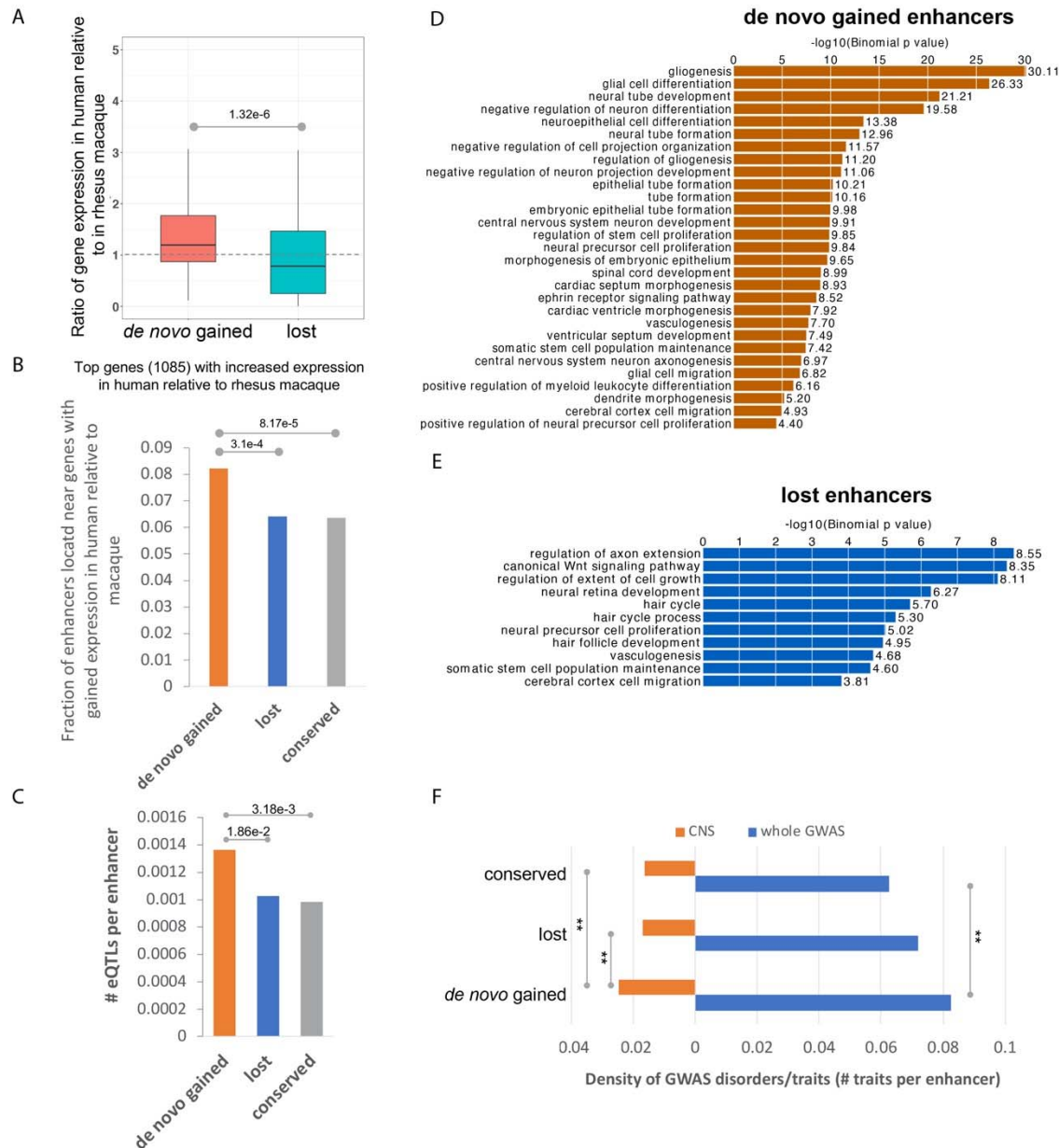
159

160 We next identified the enhancers *de novo*-gained, lost, or conserved in human relative to both  
161 macaque and human-macaque common ancestor based on the H3K27ac profile and DLM  
162 scores (Methods, Figure 1A). In total, we identified 4,066 *de novo* gained (Figure 1C), 2,925 lost,  
163 and 23,119 conserved neocortical enhancers (Figure 1D). Although the majority of the

164 developmental neocortical enhancers remained active since the divergence of human and  
165 macaque from their common ancestor, there are certain groups of enhancers that are gained or  
166 lost in the human lineage, prompting us to conjecture that these gain and loss events may  
167 correlate with the human-specific features of corticogenesis, which we investigate next.  
168

### 169 ***De novo* gained enhancers are associated with critical cortical developmental functions**

170  
171 Next, to investigate whether *de novo* enhancer gains are accompanied by an increase in the  
172 expression of their putative target genes, we compared the human-to-macaque ratios of gene  
173 expression near gained enhancers versus those near lost enhancers and observed that the  
174 genes near gained enhancers show a human-specific increase in expression while a reverse  
175 trend is exhibited by genes near lost enhancers (Figure 2A); this trend holds when we rely on  
176 Hi-C contact data to map an enhancer to its target genes (Figure S2). Consistently, gained  
177 enhancers are enriched near the genes with top 5% highest expression relative to macaque  
178 (Figure 2B). Notably, the fetal brain eQTLs (O'Brien et al. 2018) are significantly enriched in *de*  
179 *nov* gained enhancers compared to lost and conserved enhancers (Figure 2C and Figure S3).  
180 These results together support a causal link between enhancer gain and an increase in the  
181 expression of their target genes. Furthermore, the *de novo* gained enhancers are primarily  
182 associated with gliogenesis, neural tube development, and neural precursor cell proliferation,  
183 among other central nervous system (CNS) related developmental processes (Figure 2D, Figure  
184 S4A, and Table S1). In contrast, lost enhancers are associated with only a small number of CNS  
185 related essential biological processes, including regulation of axon extension, neural retina  
186 development, neural precursor cell proliferation, and cerebral cortex cell migration (Figure 2E,  
187 Figure S4B, and Table S2). Lost enhancers are enriched for far fewer processes than the *de novo*  
188 gained enhancers (Figure 2DE); at a stringent enrichment p-value threshold of  $10^{-9}$ , lost  
189 enhancers are not enriched for any process while gained enhancers are enriched for 17  
190 functions (Figure 2DE). As expected, conserved enhancers, which constitute the majority (72%)  
191 of all enhancers considered, are enriched for a large range of CNS developmental processes  
192 (Figure S5A and Table S3). Finally, we found that CNS related GWAS traits (Table S4-6) are  
193 enriched among *de novo* gained enhancers compared to conserved and lost enhancers (Figure  
194 2F), suggesting an essential role of *de novo* gained enhancers in establishing cognitive traits.



195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205

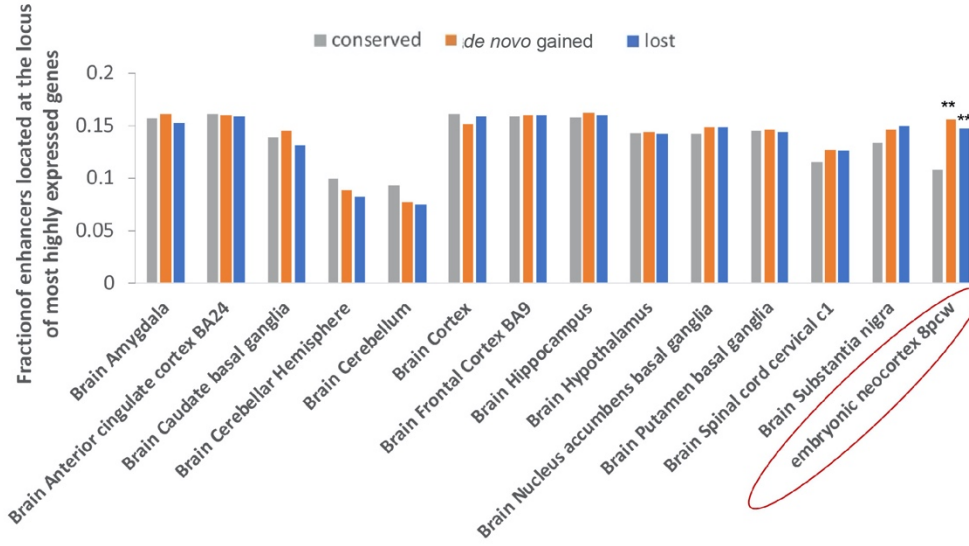
**Figure 2. *De novo* gained enhancers are associated with essential biological pathways.** (A) The expression level of genes near the *de novo* gained enhancers is increased. (B) Gained enhancers are enriched near the genes that are mostly highly expressed in humans as compared to rhesus macaque. (C) Average number of eQTLs per enhancer. (D) Biological processes that are associated with gained enhancers based on whole-genome region enrichment analysis performed using the GREAT tool (McLean et al. 2010). (E) Biological processes that are associated with lost enhancers based on GREAT whole-genome region enrichment. (F) The CNS related GWAS traits are enriched in the gained enhancers compared to both lost and conserved enhancers.

206 We further observed that, relative to conserved enhancers, *de novo* gained and lost enhancers  
207 are significantly enriched near genes that are specifically expressed in the embryonic neocortex

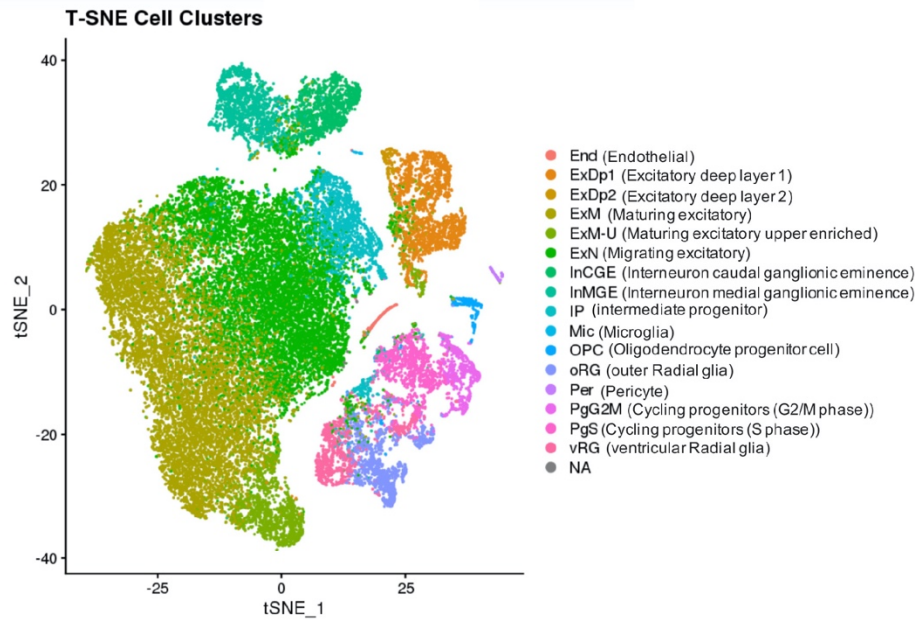
208 (8 pcw), but not adult brain (Figure 3A, Methods), implicating them specifically in brain  
209 development. To fine map gained and lost enhancer activities to specific cell types of the  
210 developing human brain, we leveraged the single-cell transcriptomic data of developing human  
211 neocortex during mid-gestation (Polioudakis et al. 2019). Among the 16 transcriptionally  
212 distinct cell types/states (Figure 3B), *de novo* gained enhancers are primarily enriched near the  
213 genes specifically expressed in progenitor cells including radial glia (oRG, vRG), cycling  
214 progenitors in G2/M phase (PgG2M) and S phase (PgS), intermediate progenitors (IP), as well as  
215 interneurons (InCGE and InMGE), which connect different brain regions and are involved in  
216 cell/axon migration (Figure 3C and Figure S6). Although lost enhancers are enriched near genes  
217 specifically expressed in excitatory neurons (excitatory deep layers ExDp1 and ExDp2, maturing  
218 excitatory neurons ExM, ExM-u and migrating excitatory neurons ExN), *de novo* gained  
219 enhancers also exhibited a comparable level of enrichment in the same loci, thus arguing for  
220 compensatory impact on either the target gene expression or the phenotypic change to a large  
221 extent. Thus, the unique enrichment of *de novo* gained enhancers in the progenitor cells and  
222 interneurons might have contributed to the expansion of cortical surface and to an increased  
223 complexity of connections in the human cerebral neocortex, both of which together underpin  
224 the advanced cognition in humans. As such, in the following, we focus specifically on the *de*  
225 *novo* gained enhancers and investigate their emergence and functional consequences.

226  
227  
228  
229

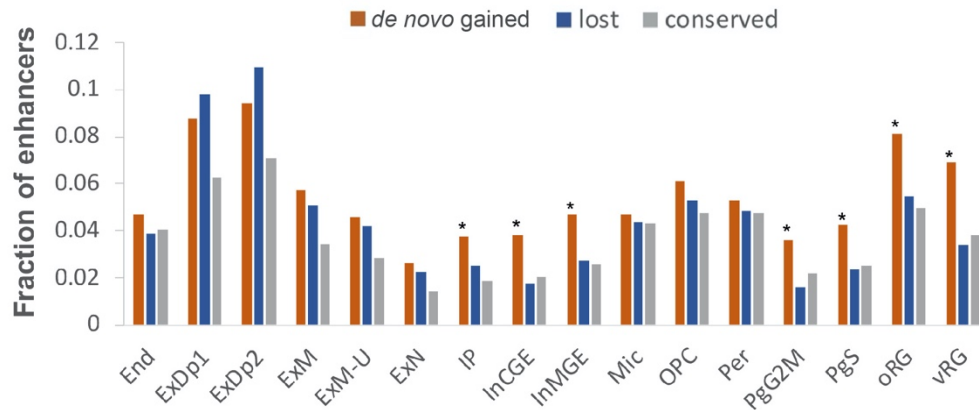
A



B



C



230  
231

232 **Figure 3.** The *de novo* gained enhancers are enriched in the progenitor cells and interneurons. A) The *de novo* gained enhancers  
233 are significantly enriched in the most highly expressed genes of embryonic human neocortex but no other adult brain regions.  
234 \*\* indicates Fisher's exact test P-value < 1e-3. B) Scatterplot visualization of cells after principal-component analysis and t-  
235 distributed stochastic neighbor embedding (tSNE), colored by Seurat clustering and annotated by major cell types. C) Fraction  
236 of enhancers near genes that are most highly expressed in all the cell clusters.

237

238

### 239 **A single essential mutation is often sufficient to create a human neocortical enhancer**

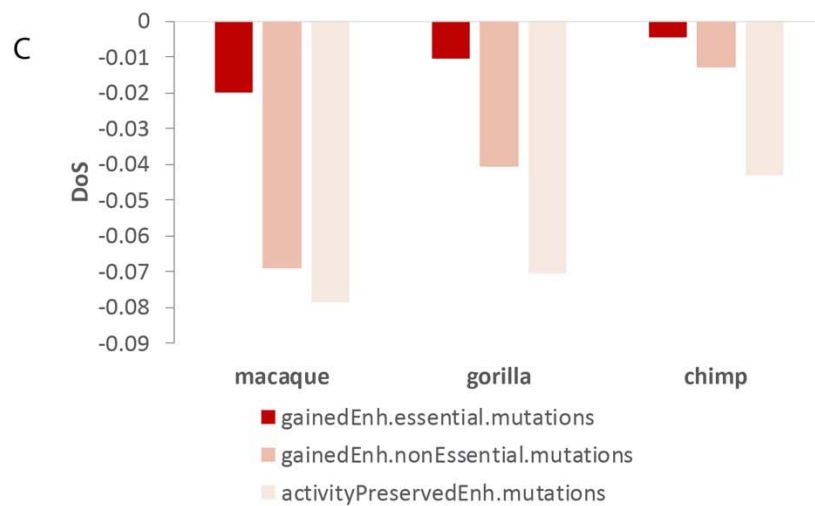
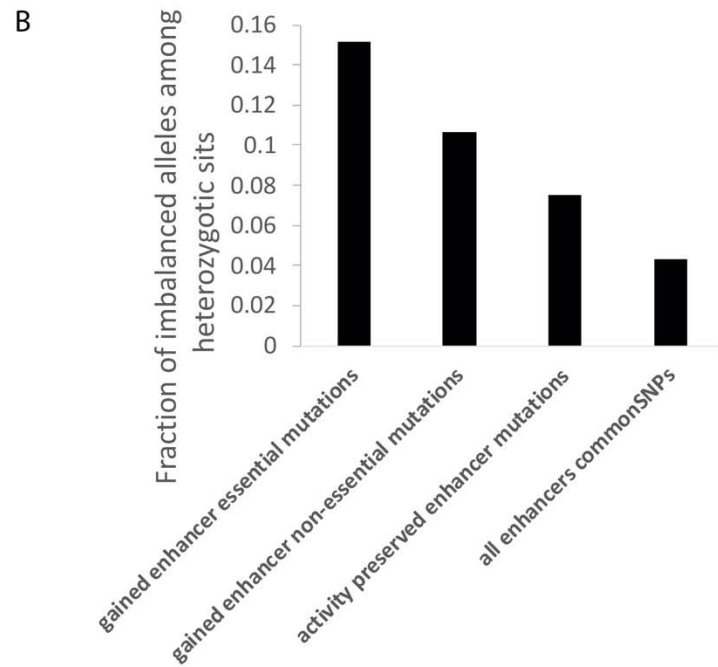
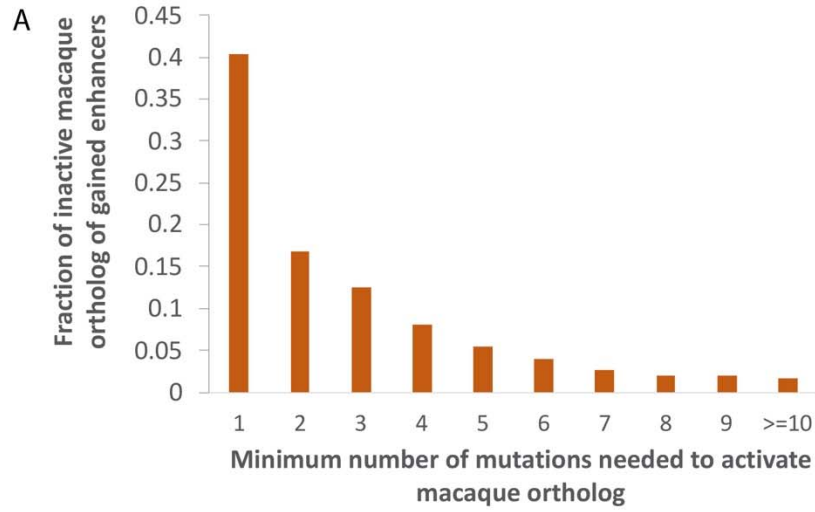
240

241 To investigate the extent to which the enhancer gains could be explained by single-nucleotide  
242 mutations and to identify the minimal number of mutations needed to activate a neutral DNA  
243 sequence, we first compared the number of human-macaque mutations in *de novo* gained and  
244 conserved enhancers. The number of human-macaque mutations in *de novo* gained and  
245 conserved enhancers are comparable -- ~50 in a 1 kb enhancer (Figure S7). Recall that our DLM  
246 is trained to distinguish fetal brain enhancers from accessible non-fetal-brain-enhancer regions  
247 and not necessarily to assess the effect of single nucleotide changes. Therefore, we first  
248 performed a series of analyses to ensure that the DLM score (i) tracks enhancer activity and (ii)  
249 can accurately predict allele specific effects on H3K27ac signals (Supplementary Results 2). To  
250 identify critical mutations, we applied our DLM to prioritize human-macaque mutations in *de*  
251 *nov*o gained enhancers based on the mutations' impact on enhancer activity by iteratively  
252 introducing them into the potentially inactive macaque sequence orthologous to human CS23  
253 enhancers. We were thus able to assess the minimal number of mutations capable of activating  
254 an enhancer (Methods). Even though only ~1.8% of all mutations in *de novo* gained enhancers  
255 are independently able to activate an enhancer (we call these essential mutations), ~40% of the  
256 *de novo* gained enhancers contain at least one essential mutation (Figure 4A). As expected, the  
257 smaller the minimal number of mutations needed to create an enhancer, the larger is their  
258 individual impact as per the DLM (Figure S8). To validate the impact of essential mutations on  
259 enhancer activity, we assessed their allelic imbalance of H3K27ac reads at the heterozygous  
260 sites. We hypothesized that the human reference allele at essential positions should exhibit  
261 larger H3K27ac read coverage than the macaque reference allele (Methods). Indeed, compared  
262 to three other groups of mutations/SNPs as controls, essential mutation positions are  
263 significantly associated with imbalance of H3K72ac reads coverage with the human reference  
264 allele (Figure 4B). This result strongly supports a causal link between the essential mutations  
265 and enhancer gain.

266

267 We next examined the evolutionary constraints on essential mutations by applying the  
268 direction of selection (*DoS*) (Stoletzki and Eyre-Walker 2011) test, which is a refinement of  
269 McDonald-Kreitman (MK) test (Stoletzki and Eyre-Walker 2011), to measure the direction and  
270 degree of departure from neutral selection (Methods). *DoS* test is applied to a pair of species  
271 and a positive and negative *DoS* indicate positive and negative selection respectively. We  
272 estimated the *DoS* values for three sets of mutations -- essential mutations, non-essential  
273 mutations in *de novo* gained enhancers, and mutations within activity preserved enhancers  
274 (Methods) -- comparing human with macaque, gorilla, and chimp. As shown in Figure 4C,  
275 compared to other mutation classes, essential mutations have the highest *DoS* values,  
276 consistent with a relaxed negative selection, or potentially a subset of sites being under positive

277 selection, both of which manifest as accelerated evolutionary rate (Cai and Petrov 2010; Hunt  
278 et al. 2011; Calderoni et al. 2016; Persi et al. 2016; Liu and Robinson-Rechavi 2018).  
279  
280





282 **Figure 4.** Essential mutations show larger impact on enhancer activity. A) Fraction of *de novo* gained enhancers  
283 that could be activated by specific number of mutations. B) Fraction of mutation/SNP sites that are in allelic  
284 imbalance. C) DoS score of the mutated sites, using macaque, gorilla, and chimp as comparison species.  
285

## 286 **Essential mutations are associated with cognition and neurodevelopmental disorders**

287  
288 Given our observation that the essential mutations are causally linked to enhancer activity in  
289 the embryonic neocortex, we assessed whether the essential mutations are preferentially  
290 associated with CNS-related GWAS traits (Methods). Indeed, we observed a ~2-fold enrichment  
291 of CNS related traits at the essential mutation positions as compared to non-essential mutation  
292 sites (Figure 5A, Table S7-8). Specifically, 7 out of 28 GWAS traits overlapping essential  
293 mutations are CNS related, and more importantly, 6 of those are associated with cognition  
294 (Table S7). We further investigated three such cases where the nearest genes are protein-  
295 coding genes with available expression data at approximate developmental stages (Zhu et al.  
296 2018).

297  
298 One essential mutation site coinciding with the common SNP rs9574096 is tightly linked to the  
299 tag SNP (rs9574095; correlation = 0.93) associated with the trait “Mathematical ability”. Both  
300 variants are located in the intronic region of the gene neurobeachin (NBEA), which is an autism-  
301 linked gene that fine-tunes signals at neuronal junctions (Nuytens et al. 2013). Mice missing one  
302 copy of NBEA show autism-like behavior (Nuytens et al. 2013). We found that NBEA exhibits a  
303 significantly higher embryonic neocortex expression in human compared to macaque at a  
304 similar early developmental stage (Zhu et al. 2018) (Figure 5B). Interestingly, the macaque allele  
305 A appears to be bound by another autism risk transcription factor, RFX3 (Harris et al. 2021),  
306 whereas the human allele T does not (Methods), suggesting a loss of RFX3 binding resulting in  
307 an increased enhancer activity and NBEA gene expression. Consistently, RFX3 expression is  
308 negatively correlated with that of NBEA in the embryonic neocortex across human and  
309 macaque individuals (Spearman  $\rho = -0.26$ ). In addition, NBEA is specifically expressed in sub-  
310 brain regions including excitatory neurons (ExDp1, ExDp2, ExM, ExM-U) and interneurons  
311 (InMGE) (Polioudakis et al. 2019), suggesting a link between these sub-brain regions and  
312 autism.

313  
314 Other two essential mutation positions coincide with two common SNPs rs747759 and  
315 rs1535043, both of which are in perfect LD with each other. Notably, rs747759 is the tag SNP of  
316 the GWAS trait “Neuroticism”. The nearest gene of the two SNPs is CD40, which again displays a  
317 much higher expression in humans as compared to macaque (Figure 5C). CD40 is a major  
318 regulator of dendrite growth and elaboration in the developing brain (Carriba and Davies 2017)  
319 and contributes to synaptic degeneration in Alzheimer’s disease (AD) (Ye et al. 2019), which  
320 may have developmental origins (Arendt et al. 2017). The human allele T at the tag SNP  
321 rs747759 either causes a potential binding site gain of NFYA or a potential binding site loss of  
322 NHLH1 (Table S9). NFYA is an AD associated gene (Leslie et al. 2014; Nazarian et al. 2018;  
323 Nazarian et al. 2019). On the other hand, NHLH1 is known to play important roles in neuronal  
324 and glial differentiation and maturation (Dennis et al. 2019). However, the chance for NHLH1 to

325 be a repressor of CD40 is dampened by their strong positive correlation of gene expression  
326 across human and macaque individuals (Spearman  $\rho = 0.58$ ). By contrast, NFYA expression is  
327 positively correlated with CD40 expression (Spearman  $\rho = 0.29$ ). At rs1535043, the human  
328 allele T is associated with the gain of an EHF binding site. However, its links with CNS traits are  
329 unclear.

330  
331 Together, these results suggest a link between essential mutations in *de novo* gained enhancers  
332 and cognition-related traits as well as neurodevelopmental disorders in humans.

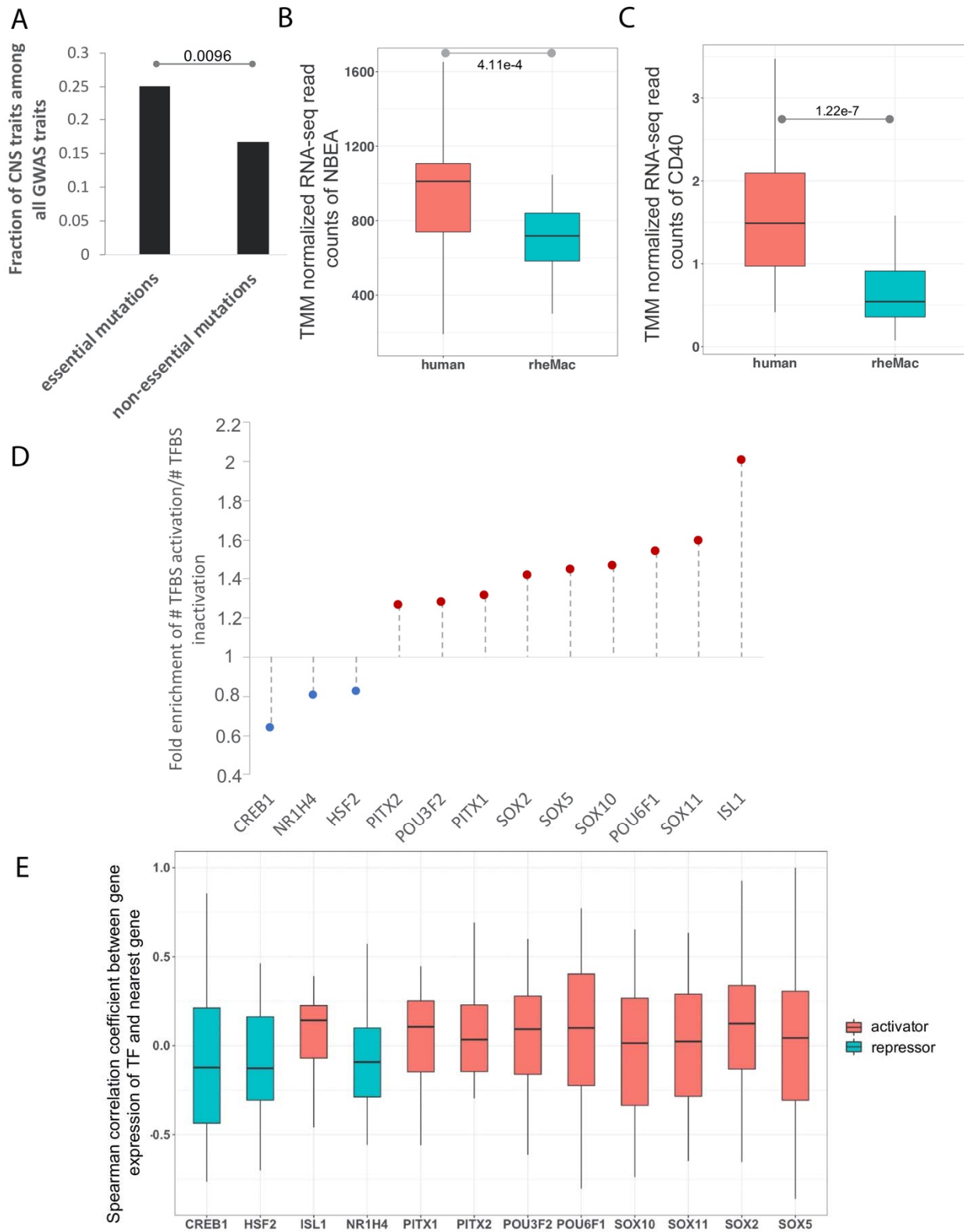
333

### 334 **Essential mutations tend to create binding sites of activating transcription factors**

335

336 Next, we investigated the relative prevalence and importance of binding site gain versus loss in  
337 the *de novo* gained enhancers. Toward this, we focused on the TFs whose binding sites are  
338 enriched in the *de novo* gained enhancers compared to the conserved ones (using both human  
339 and macaque sequences to avoid allelic bias) (Table S10) and quantified the global tendency of  
340 essential mutations to lead to binding site gain versus loss (Methods). Overall, we observed that  
341 9 TFs including POU3F2, PITX2, PITX1, SOX2, SOX5, SOX10, POU6F1, SOX11, and ISL1 tend to  
342 gain binding sites mediated by essential mutations in human (Figure 5D), suggesting an  
343 activator role of these TFs. Conversely, three TFs, CREB1, HSF2 and NR1H4, are more likely to  
344 lose their binding sites (Figure 5D), suggesting potentially repressive roles. Moreover, the  
345 overall positive or negative correlation of gene expression between these putative cognate TFs  
346 of the essential mutations and their nearest genes further validates their activator or repressor  
347 roles, respectively (Figure 5E). In short, the *de novo* gained enhancers are more likely to be  
348 activated by the creation of binding sites of activators due to the essential mutations.

349



350

351 **Figure 5.** Essential mutations are associated with cognition related traits and tend to create binding sites of  
 352 activators. A) Fraction of GWAS traits at the mutation sites which are CNS related. B) Comparison of TMM

353 normalized expression of NBEA between embryonic human and rhesus macaque individuals. P-values are based on  
354 the Wilcoxon test. C) Comparison of TMM normalized expression of CD40 between embryonic human and rhesus  
355 macaque individuals. P-values are based on the Wilcoxon test. D) Enrichment of ratio of binding site gain to loss  
356 caused by essential mutations overlapping enriched TFBSs as compared to those caused by common SNPs. E)  
357 Spearman correlation coefficient of expression between the cognate TF of essential mutation and its nearest gene.

358

359

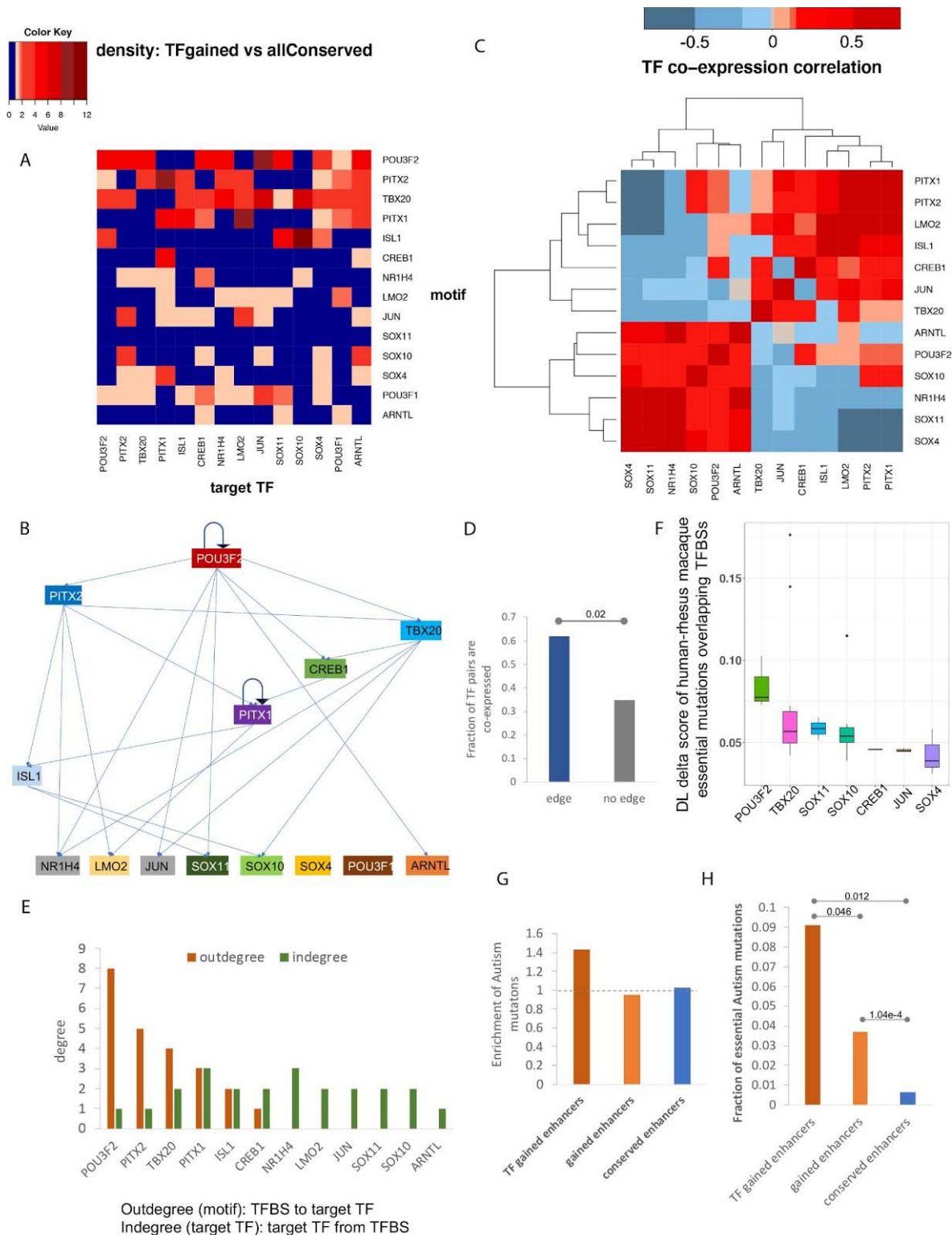
### 360 ***De novo* gained enhancers induce a potential human-specific TF regulatory network**

361 Transcriptional programs driving cell state are governed by a core set of TFs (also called master  
362 regulators), that auto- and cross-regulate each other to maintain a robust cell state. The  
363 ensemble of core TFs and their regulatory loops constitutes core transcriptional regulatory  
364 circuitry (Hnisz et al. 2013; Hnisz et al. 2015; Saint-André et al. 2016). Interestingly, the genes  
365 near *de novo* gained enhancers are enriched for transcriptional regulators (Figure S5B). We  
366 hypothesized that the TFs regulated by the *de novo* gained enhancers form a core regulatory  
367 network in the human embryonic neocortex. Toward this, first, we identified 24 TF genes (Table  
368 S11) near *de novo* gained enhancers and performed a motif scan for each of the 14 TFs having a  
369 known binding motif among all enhancers near the 24 TF genes (Methods). We found that the  
370 majority of the 14 TF motifs are enriched in the *de novo* gained enhancers near TF genes  
371 compared to the conserved enhancers in the same loci (Figure S9), suggesting a core regulatory  
372 network formed by these TFs. Next, we established a putative regulatory relationship for each  
373 TF pair based on the enrichment of the density of one TF's motif in the *de novo* gained  
374 enhancer near another TF, including autoregulation, using conserved enhancers associated with  
375 the 24 TFs as the background (Figure 6AB). The inferred links are supported by our observation  
376 that linked TF pairs tend to have correlated expressions, as compared to those which are not  
377 (Figure 6CD). Based on the number of TFs each TF regulates, POU3F2 is likely to be the master  
378 regulator, with PITX2, TBX20, and PITX1 playing critical roles (Figure 6E). Moreover, we found  
379 the essential mutations that create a binding site for the TFs at higher hierarchical levels have a  
380 larger impact on the enhancer activity according to the DLM (Figure 6F). Interestingly, the *de*  
381 *nov* non-coding mutations in Autism patients (Zhou et al. 2019) are specifically enriched in the  
382 set of *de novo* gained enhancers associated with TF activity (Figure 6G). Remarkably, the *de*  
383 *nov* Autism mutations within this subset of *de novo* gained enhancers are more likely to be  
384 essential, which alone can deactivate an enhancer, as compared to those other *de novo* gained  
385 and conserved enhancers (Fig 6H). Together, these results suggest that essential mutations and  
386 the resulting enhancer gains may have helped create a core transcriptional regulatory network,  
387 with POU3F2 in a central position, to mediate a novel gene expression program in the  
388 developing human neocortex, associated with cognitive traits.

389

390

391



392  
393  
394  
395  
396  
397

**Figure 6.** A hierarchical regulatory network of TFs induced by *de novo* gained enhancers. A) Density of TFBSs of the 14 TFs in the locus of the 14 TF genes. B) The inferred hierarchical structure of the 14 TFs. C) Spearman correlation coefficient of the 14 TF genes across the embryonic human and macaque individuals. D) Comparison of fraction of TF pairs that are co-expressed (Spearman correlation coefficient > 0.3) between the pairs with links and those without links. P-value is calculated using Fisher's exact test. E) Out-degree and in-degree of each TFs. F) DL delta score of essential mutations overlapping TFBSs. G) Enrichment of Autism mutations. H) Fraction of essential Autism mutations.

398 Distribution of DLM delta score caused by the essential mutations overlapping the 14 TFs. G) Fraction of Autism *de*  
399 *nov*o mutations located within each set of enhancers normalized by the fraction of common SNPs falling into the  
400 same set of enhancers. H) Fraction of Autism *de novo* mutations within each set of enhancers, which are essential.  
401

402

## 403 Discussion

404

405 Higher cognition in humans is attributed to substantial expansion of the cortical surface and  
406 increased complexity of cortical connections during early development. Such phenotypic  
407 changes are likely to be mediated, in significant part, by changes in transcriptional regulation  
408 during brain development (Geschwind and Rakic 2013). Recent availability of genome  
409 sequencing and epigenomic data in the developing brain of humans and a close relative –  
410 rhesus macaque – has opened the possibility to probe key regulatory changes underlying the  
411 cognitive innovations in humans.

412

413 Here, we focused on one critical component of transcriptional regulation, namely, cis-  
414 regulatory enhancers. Our results suggest that single-nucleotide mutation in the human  
415 lineage, by creating binding sites for key TFs, may have induced novel enhancers which,  
416 mediated by a core regulatory network, involving POU3F2, PITX2, TBX20, and PITX1, underlie an  
417 increased expression in the developing neocortex of key genes involved in gliogenesis, neural  
418 tube development, and neuron differentiation. Further, analysis of scRNA-seq data from the  
419 developing human brain shows that the *de novo* gained enhancers are likely to be active  
420 specifically in the progenitor cells and interneurons, which notably, are thought to underlie the  
421 expansion of the cortical surface and connectivity in the human neocortex, respectively. Given  
422 that corticogenesis in human differ from other species mainly with respect to an increased  
423 duration of neurogenesis, increases in the number and diversity of progenitors, introduction of  
424 new connections among functional areas, and modification of neuronal migration (Schwartz et  
425 al. 1991; Rakic 2009), our results are highly suggestive of a mechanistic link between enhancer  
426 gains and higher cognition in humans. We also find that the *de novo* mutations in autistic  
427 individuals are especially enriched in the *de novo* gained enhancers associated with  
428 transcription activator activities, suggesting a shared basis between human cognition and  
429 autism.

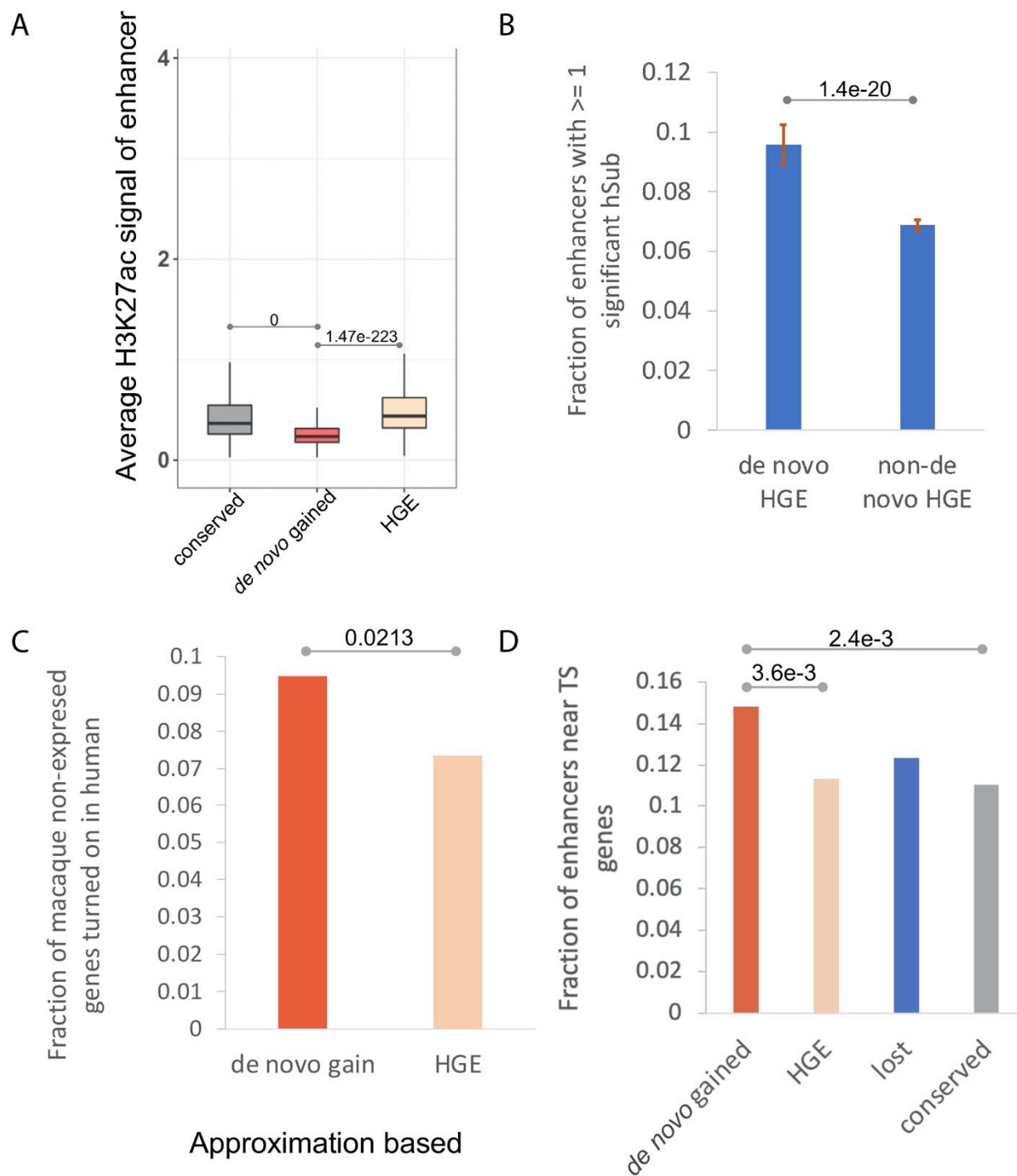
430

431 Our *de novo* gained enhancers differ significantly from previously identified human gained  
432 enhancers (HGEs) (Reilly et al. 2015), both conceptually as well as in terms of various functional  
433 properties. Reilly et al. defined HGEs based on a comparative analysis of enhancer-associated  
434 epigenetic marks (H3K27ac and H3K4me2) in human with rhesus macaque and mouse. In other  
435 words, HGEs are enhancers with increased activity in human compared to both macaque and  
436 mouse. In sharp contrast, our “*de novo*” gained enhancers originate from presumably “neutral”  
437 non-coding sequence (i.e., without a detectable enhancer activity) in either macaque or the  
438 common ancestor of humans and macaques. In fact, HGEs are largely a subset of what we  
439 consider conserved enhancers in our study (85% CS23 HGEs overlap our conserved enhancers)  
440 and not *de novo* gained enhancers (only 11.9% overlap *de novo* gained enhancers). Notably, *de*

441 *de novo* gained enhancers exhibit weaker H3K27ac signals compared to the HGEs and conserved  
442 enhancers (Figure 7A), as they are largely activated by single-nucleotide mutations that  
443 potentially create binding sites of essential TFs in the developing brain (Figure 5DE). Therefore,  
444 it is not surprising that the *de novo* gained enhancers are more vulnerable to human  
445 substitutions that significantly alter enhancer activity (termed hSubs) according to Massively  
446 Parallel Reporter Assay (MPRA) targeted HGEs (Figure 7B, Methods) (Uebbing et al. 2021).  
447 Importantly, the *de novo* gained enhancers are more likely to turn on the expression of a gene  
448 in human compared to the HGEs (Figure 7C and Figure S10A), as the macaque counterpart of  
449 the human *de novo* gained enhancers are inactive in embryonic neocortex, whereas the  
450 macaque counterpart of HGEs is also an active enhancer, albeit relatively weaker. Therefore, all  
451 else being equal, the macaque orthologs of human genes associated with *de novo* gained  
452 enhancers are more likely to be silent.

453  
454 In addition, the brain morphology related functions were reported to be associated with HGEs  
455 by the earlier study (specific functions in neuronal proliferation, migration, and cortical-map  
456 organization) (Reilly et al. 2015), which differs notably from our findings, which implicate  
457 human *de novo* gained enhancers specifically in human neocortex development. Furthermore,  
458 GO enrichment analysis based on either nearby genes (Figure S11) or genes linked via Hi-C  
459 contacts (Table S1 and Table S12) consistently shows that *de novo* gained enhancers are more  
460 likely to be associated with more tissue-specific functions of the developing human brain  
461 compared to HGEs (Supplementary Results 3). Indeed, *de novo* gained enhancers are more  
462 likely to reside near (Figure 7D) or at 3D contact positions (Figure S10B) with the most tissue-  
463 specific genes in embryonic neocortex (Table S13). As mentioned above, *de novo* gained  
464 enhancers exhibit weaker H3K27ac signals compared to the HGEs. Previous studies have  
465 implicated weaker enhancers to be specifically critical during development (Farley et al. 2015),  
466 further suggesting a link between *de novo* gained enhancers and brain development.

467  
468  
469  
470  
471



472  
473  
474  
475  
476  
477  
478  
479  
480  
481

**Figure 7.** *de novo* gained enhancers versus HGEs. (A) *de novo* gained enhancers exhibit weaker enhancer signal. (B) Fraction of enhancers with  $\geq 1$  significant hSubs. Bar plot shows the median and standard deviation of fraction of enhancers with at least one significant hSub by 90% Bootstrapping for 50 times. P-value is based on t-test. *De novo* HGEs refers to HGEs that overlap *de novo* gained enhancers, and non-*de novo* HGEs refers to CS23 HGEs that do not overlap *de novo* gained enhancers. (C) Fraction of enhancers located near genes whose RPKM  $< 1$  in macaque and  $> 1$  in human. (D) *de novo* gained enhancers are more likely to locate near tissue-specific (TS) genes.



## 482 **Methods**

483

### 484 **Data Availability**

485 We downloaded the gene expression data in the prenatal neocortex of human and macaque  
486 (Zhu et al. 2018). For human, we chose the time-points at 8 p.c.w and 12 p.c.w; for macaque,  
487 we selected the approximately matching time-points at E60 and E82 (Table S14). The data is  
488 shared by the authors at <http://evolution.psychencode.org/#>. The single-cell transcriptomic  
489 data of developing human neocortex during mid-gestation (Polioudakis et al. 2019) is shared by  
490 the authors at <http://solo.bmap.ucla.edu/shiny/webapp/>. In addition to assigning enhancers to  
491 their nearby genes, we also used the Hi-C loops in the developing brain of human (Won et al.  
492 2016) and macaque (Luo et al. 2021) to link enhancer to their gene targets. The CS23 HGEs  
493 were obtained from the study (Reilly et al. 2015). All the potential fetal brain enhancers (the  
494 merged ATAC-seq peaks from the germinal zone and cortical plate of the human developmental  
495 brain) were obtained from the study (de la Torre-Ubieta et al. 2018). The fetal brain eQTLs were  
496 obtained from the study (O'Brien et al. 2018).

497

### 498 **Embryonic neocortex enhancers in human and rhesus macaque**

499 The H3K27ac peaks of both species were obtained from the previous study (Reilly et al. 2015).  
500 The enhancers were defined as H3K27ac peaks extended to 1 kb from its original center.  
501 Integrating wider sequence context is critical because sequence surrounding the variant  
502 position determines the regulatory properties of the variant, as in vivo TF binding depends  
503 upon sequence beyond traditionally defined motifs (Deplancke et al. 2016; Inukai et al. 2017).  
504 Enhancers overlapping promoters (including all alternative promoters) and promoters (intervals  
505 [-1000 bp, 1000 bp] surrounding the transcription start site) were removed from the enhancer  
506 set. Overall, we identified 32,201 human enhancers, and 43,997 macaque enhancers.

507

### 508 **A deep convolutional neural network model for enhancer prediction**

509 We built a deep convolutional neural network to predict tissue-specific enhancer activity  
510 directly from the enhancer DNA sequence. The DLM comprises 5 convolution layers with 320,  
511 320, 240, 240, and 480 kernels, respectively (Table S15). Higher-level convolution layers receive  
512 input from larger genomic ranges and are able to represent more complex patterns than the  
513 lower layers. The convolutional layers are followed by a fully connected layer with 180 neurons,  
514 integrating the information from the full length of 1,000 bp sequence. In total, the DLM has  
515 3,631,401 trainable parameters. We used the Python library Keras version 2.4.0  
516 (<https://github.com/keras-team/keras>) to implement our model.

517

518 The model was trained for each of the four temporal-spatial groups of enhancers (CS16, CS23,  
519 F2F, and F2O). The positive sets contain the human embryonic enhancers of each group. The  
520 Dnase I-hypersensitive sites (DHSs) profiles of non-CNS-related and non-embryonic tissues  
521 from Roadmap Epigenomics projects (Kundaje et al. 2015), which do not overlap the positive  
522 sets, were collected as the negative training set of the DL model. The reason we used DHS sites  
523 not overlapping embryonic neocortex H3K27ac peaks as negative control regions is that we aim  
524 to identify tissue-specific enhancers of embryonic neocortex, and DHS is a good representation  
525 of active chromatin. The fact that DHS in general overlaps H3K27ac makes it a stringent control,

526 and in fact, our choice of DHS as the control is analogous to DeepSEA, which utilizes the  
527 genomic regions not overlapping the positive set and with at least one TF binding as the  
528 negative set, which broadly overlap with DHS regions.

529  
530 Training and testing sets were split by chromosomes. Chromosome 8 and 9 were excluded from  
531 training to test prediction performances. Chromosome 6 was used as the validation set, and  
532 the rest of the autosomes were used for training. Each training sample consists of a 1,000-bp  
533 sequence (and their reverse complement) from the human GRCh37 (hg19) reference genome.  
534 Larger DL score of the genomic sequence corresponds to a higher propensity to be an active  
535 enhancer. The genomic sequence with DLM score  $\geq 0.197$  (FPR  $\leq 0.1$ ) are predicted to be  
536 active enhancers. We used the difference of the DLM score induced by a human-macaque  
537 single-nucleotide mutation to estimate its impact on enhancer activity.

538  
539 Given a human (hg19) or macaque (rheMac2) enhancer, we used liftOver (Hinrichs et al. 2006)  
540 to identify their orthologs. Only the reciprocal counterparts with their lengths difference no  
541 more than 50 bp were considered to be ortholog pairs. For a human sequence with  $n$  mutations  
542 relative to its macaque ortholog, to score the impact of combinations of  $m$  ( $m < n$ ) mutations on  
543 enhancer activity, all possible combinations of  $m$  ( $n$  choose  $m$ ) human alleles at the human-  
544 macaque mutation sites were introduced to the macaque orthologs if the total number of  
545 combinations ( $n$  choose  $m$ ) is no more than 10,000, otherwise, we randomly sample 10,000  
546 combinations of  $m$  human alleles from the human-macaque mutation sites and introduce them  
547 to the macaque ortholog. The change of DL score caused by the set of introduced human  
548 mutations were used to estimate their impact on enhancer activity.

549  
550 **Gain and loss of enhancers**  
551 Briefly, if a human enhancer having a high DLM score scored low both in macaque and in the  
552 common ancestor, and was not detected by H3K27ac in macaque, it was considered to be a *de*  
553 *novo* gain in humans (Figure 1A). Likewise, if a macaque enhancer having high DL score scored  
554 high in common ancestor, scored low in human and was undetectable by H3K27ac in human, it  
555 was considered a loss in human (Figure 1A). The enhancers that are detected by H3K27ac in  
556 both human and macaque, and scored highly in all three genomes were called conserved  
557 enhancers (Figure 1A).

558  
559 **Normalization of gene expression data**  
560 We applied 'tmm' built-in normalization method of edgeR to normalize human and macaque  
561 embryonic neocortex gene expression and to remove differences across species and batch  
562 effects. To identify the most tissue-specific genes of human embryonic neocortex, the  
563 expression data of human individuals were averaged and quantile normalized together with the  
564 gene expression profile downloaded from GTEx. The top 2000 genes with the highest ratios of  
565 the human embryonic expression to the mean of the GTEx expression were identified as the  
566 most specifically highly expressed genes in human embryonic neocortex (Table S13).

567  
568 ***De novo* single-nucleotide substitutions in autism spectrum disorder (ASD)**

569 We obtained 127,141 *de novo* single-nucleotide mutations in ASD from a previous study (Zhou  
570 et al. 2019), which were identified from Simons Simplex Collection of whole-genome  
571 sequencing data for 1790 families that were available via the Simons Foundation Autism  
572 Research Initiative (SFARI).

573

#### 574 **Functional enrichment analysis using GREAT and DAVID tools**

575 To probe the potential functional roles of gained and lost enhancers we first tested for  
576 functional enrichment among genes near the enhancer loci using the online Genomic Regions  
577 Enrichment of Annotations Tool (GREAT) version 3.0.0 (McLean et al. 2010) using single nearest  
578 gene association rule with more strict settings than default. Specifically, the GO terms will be  
579 considered as enriched if it has at least 10 gene hits with FDR threshold set as 0.01. Two  
580 background options were used when using GREAT. Figure 2DE, Figure S5 and Figure S11 are  
581 based on enrichment against whole genome region. Next, we performed GO enrichment  
582 analysis using all potential fetal brain enhancers (the merged ATAC-seq peaks from the  
583 germinal zone and cortical plate of the human developmental brain) (de la Torre-Ubieta et al.  
584 2018) as the background and obtained consistent observations (Figure S4AB and Figure S12AB).  
585 The exception is the conserved enhancers, which are not enriched for CNS related biological  
586 processes (Figure S12A). The tissue-specific signal of conserved enhancers is dampened, as  
587 expected, by using the fetal brain enhancers as the background, as the conserved enhancers  
588 constitute the majority of the fetal brain enhancers.

589

590 We also applied DAVID (Huang da et al. 2009b; Huang da et al. 2009a) to do functional  
591 enrichment of the genes with Hi-C loops to different sets of enhancers.

592

#### 593 **Enrichment analysis of GWAS traits and eQTLs**

594 The NHGRI-EBI GWAS Catalog (Buniello et al. 2019) was downloaded. To study the enrichment  
595 of a set of SNPs coinciding with CNS related GWAS traits, the tag SNPs were first expanded by  
596 linkage disequilibrium (LD) ( $r^2 > 0.8$ , maximum distance of 500 kb) using Plink  
597 (<http://pngu.mgh.harvard.edu/purcell/urcell/plink/>; (Purcell et al. 2007)) with the following  
598 parameters:

599 `'--r2 --ld-window-kb 500 --ld-window-r2 0.8'`

600 We overlapped the LD-expanded GWAS traits with the human-macaque mutation sites of the  
601 gained enhancers where the human alternative alleles are the same as the macaque reference  
602 alleles. The CNS-related GWAS traits are listed in Table S4-6. We then use the fraction of CNS-  
603 related traits among the total GWAS traits overlapping the essential mutations, as compared to  
604 that of the non-essential mutations to estimate the enrichment of CNS-related traits in the  
605 essential mutation positions (Figure 5A).

606 As for the overall enrichment of the CNS-related GWAS traits in the three sets of enhancers  
607 (Figure 2F), we used the density (average number of LD-expanded GWAS traits per enhancer) to  
608 estimate the enrichment. As the density of common SNPs in the three sets of enhancers  
609 (average number of SNPs per enhancer) is comparable (gained: 4.1, lost: 4.05, conserved: 4.6)  
610 and would not change the trend of the enrichment upon normalization, we did not normalize  
611 the GWAS density by SNP density.

612

613 As for the enrichment of the fetal brain eQTLs (O'Brien et al. 2018) in the three sets of  
614 enhancers, we first compared the density of eQTLs (average number of eQTLs per enhancer) in  
615 the three sets of enhancers (Figure 2C). Next, we normalized the fraction of eQTLs fallen within  
616 a set of enhancers by the fraction of common SNPs fallen within that set of enhancers (Figure  
617 S3).

618

619

### 620 **Identification of potential TFBSs in the *de novo* gained enhancers**

621 To identify potential binding sites, we used FIMO (Bailey et al. 2009) to scan the profiles of  
622 binding sites for vertebrate TF motifs in Jaspar (Mathelier et al. 2014), CIS-BP (Weirauch et al.  
623 2014), SwissRegulon (Pachkov et al. 2007), HOCOMOCO (Kulakovskiy et al. 2016), and  
624 UniPROBE (Hume et al. 2015) databases, along the enhancer sequences. We identified motif-  
625 specific thresholds to limit the false discovery rate to no more than five false positives in 10<sup>2</sup>kb  
626 of sequence, by scanning each motif on random genomic sequences using FIMO (Bailey et al.  
627 2009). Enrichment of a motif in *de novo* gained (foreground) relative to conserved (background)  
628 enhancers were ascertained using Fisher's exact test. The occurrence of a particular TFBS in the  
629 set of *de novo*-gained/conserved sequences was normalized by the total number of *de novo*-  
630 gained/conserved regions.

631

632 However, when identifying TFs whose motifs are enriched in *de novo* gained enhancers relative  
633 to conserved enhancers, we included both the human and the macaque ortholog sequences, to  
634 avoid allelic bias in our following analysis of activation/repression of enhancers by single-  
635 nucleotide mutations. Next, we assessed whether a mutation (in a *de novo* gained enhancer)  
636 creates a binding site of a potential activator or disrupts binding of a potential repressor, we  
637 estimated, for each enriched TF, the ratio of binding site gain to loss caused by essential  
638 mutations within *de novo* gained enhancers relative to the same ratio caused by common SNPs.  
639 If the gain/loss (loss/gain, respectively) ratio caused by essential mutations was greater than  
640 1.2-fold that for common SNPs, the TF was considered activator (repressor, respectively).

641

### 642 **Identification of allelic imbalance in H3K27ac data**

643 We used BWA (Li and Durbin 2010) to map two replicates of CS23 H3K27ac data (Reilly et al.  
644 2015) to hg19 human reference sequence. At the mutation/SNP sites, the H3K27ac reads were  
645 extracted using BaalChIP (de Santiago et al. 2017). Allelic counts over heterozygous sites of the  
646 two replicates were merged, and variants that had at least 6 reads were further processed for  
647 allele specific enhancer activity analysis with Binomial test. We use the heterozygous sites  
648 within the activity preserved enhancers (the ratio between human and macaque H3K27ac  
649 signal is no more than 1.2) as the background. For a heterozygous site, if the ratio of reads  
650 number of the human allele to that of the macaque allele is over 1.3 and the Binomial p-value  
651  $\leq 1e-3$ , the position is considered to have allelic imbalance.

652

### 653 **Single-cell clustering and visualization**

654 Clustering was performed using Seurat (v2.3.4) (Stuart et al. 2019). Read depth normalized  
655 expression values were mean centered and variance scaled for each gene, and the effects of

656 number of UMI (sequencing depth), donor, and library preparation batch were removed using a  
657 linear model with Seurat ('ScaleData' function). Highly variable genes were then identified and  
658 used for the subsequent analysis (Seurat 'MeanVarPlot' function). Briefly, average expression  
659 and dispersion are calculated for each gene, genes are placed into bins, and then a z-score for  
660 dispersion within each bin is determined. Principal component analysis (PCA) was then used to  
661 reduce dimensionality of the dataset to the top 13 PCs (Seurat 'RunPCA' function). Clustering  
662 was then performed using graph-based clustering implemented by Seurat ('FindClusters'  
663 function). Cell clusters with fewer than 30 cells were omitted from further analysis. Clusters  
664 were annotated using the Seurat function 'group.by'.  
665 For visualization, t-distributed stochastic neighbor embedding (tSNE) coordinates were  
666 calculated in PCA space, independent of the clustering, using Seurat ('RunTSNE' function). tSNE  
667 plots were then colored by the cluster assignments derived above, gene expression values, or  
668 other features of interest. Gene expression values are mean centered and variance scaled  
669 unless otherwise noted.

670

### 671 **Direction of selection test**

672 The *DoS* test was designed to measure the direction and extent of departure from neutral  
673 selection based on the difference between the proportion of substitution and polymorphism in  
674 the selective sites. *DoS* is positive when there is evidence of adaptive evolution, is zero if there  
675 is only neutral evolution, and is negative when there are slightly deleterious mutations  
676 segregating (Stoletzki and Eyre-Walker 2011). Here, we used the mutated four-fold degenerate  
677 sites as the background to measure the selection on the mutations within *de novo* gained  
678 enhancers (formula 1). Note that all sites in our three mutational site classes are, by design,  
679 mutated in human relative to macaque. Therefore, to avoid ascertainment bias, we uniformly  
680 applied the same criteria of human-macaque mutation to select a subset of all fourfold  
681 degenerate sites.

682

683 Let, *n* represent the 'non-synonymous' sites, i.e. the essential or non-essential mutations within  
684 the *de novo* gained enhancers. *S* represents the 'synonymous' sites, i.e. the mutated four-fold  
685 degenerate sites. *D* means 'diverged' sites, i.e. mutations (or substitutions) that are fixed in the  
686 human populations, and *P* means 'polymorphic' sites, i.e. both the ancestor allele and the  
687 mutations are preserved in the human populations (Table 1).

688

$$689 \text{DoS} = D_n / (D_n + D_s) - P_n / (P_n + P_s) \quad (1)$$

690

691

692 Table 1. Contingency table of number of fixed mutations and polymorphic mutations at the  
693 foreground and background sites.

694

	Fixed	Polymorphic
Mutated four-fold degenerate sites	Ds	Ps
Mutated sites within gained enhancers	Dn	Pn

695

696 Ds: the number of fixed mutations at mutated four-fold degenerate sites  
697 Dn: the number of fixed mutations within *de novo* gained enhancers  
698 Ps: the number of polymorphic mutations at mutated four-fold degenerate sites  
699 Pn: the number of polymorphic mutations within *de novo* gained enhancers  
700

## 701 **Comparing *de novo* gained enhancers and HGEs using MPRA data**

702 Overlapping the significant human substitutions (relative to chimp, termed hSubs) from a MPRA  
703 targeted HGEs (Uebbing et al. 2021) with the *de novo* gained enhancers, we found that 141 *de*  
704 *nov*o gained enhancers overlapping HGEs (dubbed *de novo* HGEs) were tested by this assay. In  
705 total, 14 of the 141 (10%) *de novo* HGEs harbor at least one hSub. For the 1,019 CS23 HGEs that  
706 do not overlap *de novo* gained enhancers (dubbed non-*de novo* HGEs), 74 (7%) HGEs have at  
707 least one hSubs (Figure 7B). We applied 90% bootstrapping 50 times to estimate the statistical  
708 significance of the difference between the two fractions (Figure 7B).

709

710

## 711 **Acknowledgement**

712

713 This work utilized the computational resources of the NIH HPC Biowulf cluster and was  
714 supported by the Intramural Research Program of the National Cancer Institute, Center for  
715 Cancer Research, and National Library of Medicine, NIH. We would like to thank Di Huang,  
716 Vishaka Gopalan, and Arashdeep Singh for their feedback. We would also to thank James  
717 Noonan and Jian Zhou for nicely providing supplementary materials of their works.

718

## 719 **Competing interests**

720 The authors have no competing interests.

721

## 722 **References**

723

724 Arendt T, Stieler J, Ueberham U. 2017. Is sporadic Alzheimer's disease a developmental  
725 disorder? *J Neurochem* **143**: 396-408.

726 Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS. 2009.  
727 MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res* **37**: W202-208.

728 Bradley RK, Holmes I. 2007. Transducers: an emerging probabilistic framework for modeling  
729 indels on trees. *Bioinformatics* **23**: 3258-3262.

730 Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, McMahon A,  
731 Morales J, Mountjoy E, Sollis E et al. 2019. The NHGRI-EBI GWAS Catalog of published  
732 genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic*  
733 *Acids Res* **47**: D1005-d1012.

734 Cai JJ, Petrov DA. 2010. Relaxed purifying selection and possibly high rate of adaptation in  
735 primate lineage-specific genes. *Genome Biol Evol* **2**: 393-409.

- 736 Calderoni L, Rota-Stabelli O, Frigato E, Panziera A, Kirchner S, Foulkes NS, Kruckenhauser L,  
737 Bertolucci C, Fuselli S. 2016. Relaxed selective constraints drove functional modifications  
738 in peripheral photoreception of the cavefish *P. andruzzii* and provide insight into the  
739 time of cave colonization. *Heredity (Edinb)* **117**: 383-392.
- 740 Carriba P, Davies AM. 2017. CD40 is a major regulator of dendrite growth from developing  
741 excitatory and inhibitory neurons. *Elife* **6**.
- 742 de la Torre-Ubieta L, Stein JL, Won H, Opland CK, Liang D, Lu D, Geschwind DH. 2018. The  
743 Dynamic Landscape of Open Chromatin during Human Cortical Neurogenesis. *Cell* **172**:  
744 289-304.e218.
- 745 de Santiago I, Liu W, Yuan K, O'Reilly M, Chilamakuri CS, Ponder BA, Meyer KB, Markowitz F.  
746 2017. BaalChIP: Bayesian analysis of allele-specific transcription factor binding in cancer  
747 genomes. *Genome Biol* **18**: 39.
- 748 Dehay C, Kennedy H, Kosik KS. 2015. The outer subventricular zone and primate-specific cortical  
749 complexification. *Neuron* **85**: 683-694.
- 750 Dennis DJ, Han S, Schuurmans C. 2019. bHLH transcription factors in neural development,  
751 disease, and reprogramming. *Brain Res* **1705**: 48-65.
- 752 Deplancke B, Alpern D, Gardeux V. 2016. The Genetics of Transcription Factor DNA Binding  
753 Variation. *Cell* **166**: 538-554.
- 754 Emera D, Yin J, Reilly SK, Gockley J, Noonan JP. 2016. Origin and evolution of developmental  
755 enhancers in the mammalian neocortex. *Proc Natl Acad Sci U S A* **113**: E2617-2626.
- 756 Farley EK, Olson KM, Zhang W, Brandt AJ, Rokhsar DS, Levine MS. 2015. Suboptimization of  
757 developmental enhancers. *Science* **350**: 325-328.
- 758 Fulco CP, Nasser J, Jones TR, Munson G, Bergman DT, Subramanian V, Grossman SR, Anyoha R,  
759 Doughty BR, Patwardhan TA et al. 2019. Activity-by-contact model of enhancer-  
760 promoter regulation from thousands of CRISPR perturbations. *Nat Genet* **51**: 1664-1669.
- 761 Geschwind DH, Rakic P. 2013. Cortical evolution: judge the brain by its cover. *Neuron* **80**: 633-  
762 647.
- 763 Harris HK, Nakayama T, Lai J, Zhao B, Argyrou N, Gubbels CS, Soucy A, Genetti CA, Suslovitch V,  
764 Rodan LH et al. 2021. Disruption of RFX family transcription factors causes autism,  
765 attention-deficit/hyperactivity disorder, intellectual disability, and dysregulated  
766 behavior. *Genet Med* doi:10.1038/s41436-021-01114-z.
- 767 Hinrichs AS, Karolchik D, Baertsch R, Barber GP, Bejerano G, Clawson H, Diekhans M, Furey TS,  
768 Harte RA, Hsu F et al. 2006. The UCSC Genome Browser Database: update 2006. *Nucleic  
769 Acids Res* **34**: D590-598.
- 770 Hnisz D, Abraham BJ, Lee TI, Lau A, Saint-Andre V, Sigova AA, Hoke HA, Young RA. 2013. Super-  
771 enhancers in the control of cell identity and disease. *Cell* **155**: 934-947.
- 772 Hnisz D, Schuijers J, Lin CY, Weintraub AS, Abraham BJ, Lee TI, Bradner JE, Young RA. 2015.  
773 Convergence of developmental and oncogenic signaling pathways at transcriptional  
774 super-enhancers. *Mol Cell* **58**: 362-370.
- 775 Holmes I. 2003. Using guide trees to construct multiple-sequence evolutionary HMMs.  
776 *Bioinformatics* **19 Suppl 1**: i147-157.
- 777 Holmes I, Bruno WJ. 2001. Evolutionary HMMs: a Bayesian approach to multiple alignment.  
778 *Bioinformatics* **17**: 803-820.

- 779 Huang da W, Sherman BT, Lempicki RA. 2009a. Bioinformatics enrichment tools: paths toward  
780 the comprehensive functional analysis of large gene lists. *Nucleic Acids Res* **37**: 1-13.
- 781 Huang da W, Sherman BT, Lempicki RA. 2009b. Systematic and integrative analysis of large gene  
782 lists using DAVID bioinformatics resources. *Nat Protoc* **4**: 44-57.
- 783 Hume MA, Barrera LA, Gisselbrecht SS, Bulyk ML. 2015. UniPROBE, update 2015: new tools and  
784 content for the online database of protein-binding microarray data on protein-DNA  
785 interactions. *Nucleic Acids Res* **43**: D117-122.
- 786 Hunt BG, Ometto L, Wurm Y, Shoemaker D, Yi SV, Keller L, Goodisman MA. 2011. Relaxed  
787 selection is a precursor to the evolution of phenotypic plasticity. *Proc Natl Acad Sci U S A*  
788 **108**: 15936-15941.
- 789 Inukai S, Kock KH, Bulyk ML. 2017. Transcription factor-DNA binding: beyond binding site  
790 motifs. *Curr Opin Genet Dev* **43**: 110-119.
- 791 Kulakovskiy IV, Vorontsov IE, Yevshin IS, Soboleva AV, Kasianov AS, Ashoor H, Ba-Alawi W, Bajic  
792 VB, Medvedeva YA, Kolpakov FA et al. 2016. HOCOMOCO: expansion and enhancement  
793 of the collection of transcription factor binding sites models. *Nucleic Acids Res* **44**: D116-  
794 125.
- 795 Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, Kheradpour P, Zhang Z,  
796 Wang J, Ziller MJ et al. 2015. Integrative analysis of 111 reference human epigenomes.  
797 *Nature* **518**: 317-330.
- 798 Leslie R, O'Donnell CJ, Johnson AD. 2014. GRASP: analysis of genotype-phenotype results from  
799 1390 genome-wide association studies and corresponding open access database.  
800 *Bioinformatics* **30**: i185-194.
- 801 Lewitus E, Kelava I, Kalinka AT, Tomancak P, Huttner WB. 2014. An adaptive threshold in  
802 mammalian neocortical evolution. *PLoS Biol* **12**: e1002000.
- 803 Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows-Wheeler transform.  
804 *Bioinformatics* **26**: 589-595.
- 805 Liu J, Robinson-Rechavi M. 2018. Adaptive Evolution of Animal Proteins over Development:  
806 Support for the Darwin Selection Opportunity Hypothesis of Evo-Devo. *Mol Biol Evol* **35**:  
807 2862-2872.
- 808 Long HK, Prescott SL, Wysocka J. 2016. Ever-Changing Landscapes: Transcriptional Enhancers in  
809 Development and Evolution. *Cell* **167**: 1170-1187.
- 810 Luo X, Liu Y, Dang D, Hu T, Hou Y, Meng X, Zhang F, Li T, Wang C, Li M et al. 2021. 3D Genome of  
811 macaque fetal brain reveals evolutionary innovations during primate corticogenesis. *Cell*  
812 **184**: 723-740.e721.
- 813 Mathelier A, Zhao X, Zhang AW, Parcy F, Worsley-Hunt R, Arenillas DJ, Buchman S, Chen CY,  
814 Chou A, Ionescu H et al. 2014. JASPAR 2014: an extensively expanded and updated  
815 open-access database of transcription factor binding profiles. *Nucleic Acids Res* **42**:  
816 D142-147.
- 817 McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM, Bejerano G. 2010.  
818 GREAT improves functional interpretation of cis-regulatory regions. *Nat Biotechnol* **28**:  
819 495-501.
- 820 Namba T, Huttner WB. 2017. Neural progenitor cells and their role in the development and  
821 evolutionary expansion of the neocortex. *Wiley Interdiscip Rev Dev Biol* **6**.



- 822 Nazarian A, Arbeev KG, Yashkin AP, Kulminski AM. 2019. Genetic heterogeneity of Alzheimer's  
823 disease in subjects with and without hypertension. *Geroscience* **41**: 137-154.
- 824 Nazarian A, Yashin AI, Kulminski AM. 2018. Methylation-wide association analysis reveals AIM2,  
825 DGUOK, GNAI3, and ST14 genes as potential contributors to the Alzheimer's disease  
826 pathogenesis. *bioRxiv* **322503**.
- 827 Nuytens K, Gantois I, Stijnen P, Iscru E, Laeremans A, Serneels L, Van Eylem L, Liebhaber SA,  
828 Devriendt K, Balschun D et al. 2013. Haploinsufficiency of the autism candidate gene  
829 Neurobeachin induces autism-like behaviors and affects cellular and molecular  
830 processes of synaptic plasticity in mice. *Neurobiol Dis* **51**: 144-151.
- 831 O'Brien HE, Hannon E, Hill MJ, Toste CC, Robertson MJ, Morgan JE, McLaughlin G, Lewis CM,  
832 Schalkwyk LC, Hall LS et al. 2018. Expression quantitative trait loci in the developing  
833 human brain and their enrichment in neuropsychiatric disorders. *Genome Biol* **19**: 194.
- 834 Otani T, Marchetto MC, Gage FH, Simons BD, Livesey FJ. 2016. 2D and 3D Stem Cell Models of  
835 Primate Cortical Development Identify Species-Specific Differences in Progenitor  
836 Behavior Contributing to Brain Size. *Cell Stem Cell* **18**: 467-480.
- 837 Pachkov M, Erb I, Molina N, van Nimwegen E. 2007. SwissRegulon: a database of genome-wide  
838 annotations of regulatory sites. *Nucleic Acids Res* **35**: D127-131.
- 839 Paten B, Herrero J, Fitzgerald S, Beal K, Flicek P, Holmes I, Birney E. 2008. Genome-wide  
840 nucleotide-level mammalian ancestor reconstruction. *Genome Res* **18**: 1829-1843.
- 841 Persi E, Wolf YI, Koonin EV. 2016. Positive and strongly relaxed purifying selection drive the  
842 evolution of repeats in proteins. *Nat Commun* **7**: 13570.
- 843 Polioudakis D, de la Torre-Ubieta L, Langerman J, Elkins AG, Shi X, Stein JL, Vuong CK,  
844 Nichterwitz S, Gevorgian M, Opland CK et al. 2019. A Single-Cell Transcriptomic Atlas of  
845 Human Neocortical Development during Mid-gestation. *Neuron* **103**: 785-801.e788.
- 846 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar P, de Bakker  
847 PI, Daly MJ et al. 2007. PLINK: a tool set for whole-genome association and population-  
848 based linkage analyses. *Am J Hum Genet* **81**: 559-575.
- 849 Rakic P. 2009. Evolution of the neocortex: a perspective from developmental biology. *Nat Rev*  
850 *Neurosci* **10**: 724-735.
- 851 Rakic P, Ayoub AE, Breunig JJ, Dominguez MH. 2009. Decision by division: making cortical maps.  
852 *Trends Neurosci* **32**: 291-301.
- 853 Reilly SK, Yin J, Ayoub AE, Emera D, Leng J, Cotney J, Sarro R, Rakic P, Noonan JP. 2015.  
854 Evolutionary genomics. Evolutionary changes in promoter and enhancer activity during  
855 human corticogenesis. *Science* **347**: 1155-1159.
- 856 Saint-André V, Federation AJ, Lin CY, Abraham BJ, Reddy J, Lee TI, Bradner JE, Young RA. 2016.  
857 Models of human core transcriptional regulatory circuitries. *Genome Res* **26**: 385-396.
- 858 Schwartz ML, Rakic P, Goldman-Rakic PS. 1991. Early phenotype expression of cortical neurons:  
859 evidence that a subclass of migrating neurons have callosal axons. *Proc Natl Acad Sci U S*  
860 *A* **88**: 1354-1358.
- 861 Sousa AMM, Meyer KA, Santpere G, Gulden FO, Sestan N. 2017. Evolution of the Human  
862 Nervous System Function, Structure, and Development. *Cell* **170**: 226-247.
- 863 Stoletzki N, Eyre-Walker A. 2011. Estimation of the neutrality index. *Mol Biol Evol* **28**: 63-70.

864 Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, Hao Y, Stoeckius M,  
865 Smibert P, Satija R. 2019. Comprehensive Integration of Single-Cell Data. *Cell* **177**: 1888-  
866 1902.e1821.

867 Uebbing S, Gockley J, Reilly SK, Kocher AA, Geller E, Gandotra N, Scharfe C, Cotney J, Noonan JP.  
868 2021. Massively parallel discovery of human-specific substitutions that alter enhancer  
869 activity. *Proc Natl Acad Sci U S A* **118**.

870 Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, Najafabadi HS,  
871 Lambert SA, Mann I, Cook K et al. 2014. Determination and inference of eukaryotic  
872 transcription factor sequence specificity. *Cell* **158**: 1431-1443.

873 Won H, de la Torre-Ubieta L, Stein JL, Parikshak NN, Huang J, Opland CK, Gandal MJ, Sutton GJ,  
874 Hormozdiari F, Lu D et al. 2016. Chromosome conformation elucidates regulatory  
875 relationships in developing human brain. *Nature* **538**: 523-527.

876 Ye X, Zhou W, Zhang J. 2019. Association of CSF CD40 levels and synaptic degeneration across  
877 the Alzheimer's disease spectrum. *Neurosci Lett* **694**: 41-45.

878 Zhou J, Park CY, Theesfeld CL, Wong AK, Yuan Y, Scheckel C, Fak JJ, Funk J, Yao K, Tajima Y et al.  
879 2019. Whole-genome deep-learning analysis identifies contribution of noncoding  
880 mutations to autism risk. *Nat Genet* **51**: 973-980.

881 Zhou J, Troyanskaya OG. 2015. Predicting effects of noncoding variants with deep learning-  
882 based sequence model. *Nat Methods* **12**: 931-934.

883 Zhu Y, Sousa AMM, Gao T, Skarica M, Li M, Santpere G, Esteller-Cucala P, Juan D, Ferrández-  
884 Peral L, Gulden FO et al. 2018. Spatiotemporal transcriptomic divergence across human  
885 and macaque brain development. *Science* **362**.

886

887

888

889

## 890 **Supplementary Materials**

891

## 892 **Supplementary Results**

893

### 894 **1. Performance and further validation of DLM of embryonic neocortex enhancers**

895

896 Here we provide benchmarking and comparison of our enhancer model with DeepSEA (Zhou and  
897 Troyanskaya 2015).

898

899 To directly compare our model performance with DeepSEA, we applied our model to the training and  
900 testing H3K27ac data sets used by DeepSEA. Our model achieved a very similar (although slightly higher)  
901 accuracy (both auROC and auPRC) compared to DeepSEA across multiple datasets (Figure S13BC).

902

903 We have shown that the human embryonic neocortex DLM can accurately estimate the enhancer  
904 activity (independently) in macaque from its genomic sequence (Fig 2A). To further validate our model,  
905 we applied the model trained on the human embryonic neocortex (CS23) enhancers (H3K27ac peaks)  
906 and tested it on the mouse embryonic neocortex enhancers (H3K27ac peaks) (Reilly *et al* 2015,  
907 PMID:25745175), using random genomic regions (due to a lack of available multi-tissue DHS profile) that  
908 do not overlap H3K27ac peaks as the negative testing set. Even for this more distant species, the model  
909 achieves an auROC of 0.9 at e11 (Figure S13D).

910

### 911 **2. DLM can accurately predict allele specific effects on histone marks H3K27ac**

912

913 Our DLM is trained to distinguish enhancer region from non-enhancer regions in a specific context.  
914 However, its application to identify *de novo* enhancer gains driven by single nucleotide mutations  
915 requires the DLM score to be sensitive to single nucleotide changes. We performed additional analyses  
916 to ensure that DLM score indeed (i) represent the enhancer activity level, and (ii) is sensitive to single  
917 nucleotide changes.

918

919 First, we computed the direct correlation between the predicted DLM score (DL score) of the enhancers  
920 and the log of their average H3K27ac signal intensity. We observed a significant positive correlation  
921 between the two (correlation = 0.4, empirical p-value = 3.18e-6).

922

923 Next, DeepSEA was shown to work well in identifying variants at loci that affect histone signals (hQTLs of  
924 H3K27ac or H3K4me3) (Zhou and Troyanskaya 2015). As our approach is very similar to DeepSEA (just a  
925 different neural net architecture), and we aim to identify variants that create enhancers, we trained our  
926 model on H3K27ac peaks in a lymphoblastoid cell line, GM12878, and applied it to predict the same set  
927 of hQTLs of H3K27ac in lymphoblastoid cell lines (McVicker, G. et al. Science 342, 747-749 (2013)) as did  
928 DeepSEA. Our model shows similar accuracy as DeepSEA (Figure S14).

929

930 To further show the ability of our DLM to accurately predict chromatin features from sequence with  
931 single-nucleotide sensitivity, we applied our CS23 model to evaluate the 2,578 allelically imbalanced  
932 SNPs within the CS23 H3K27ac peaks, which were identified using the R-package BaalChIP (de Santiago  
933 et al. 2017). Our model makes similarly accurate predictions on this set of SNPs as well (Figure S15).

934

### 935 **3. Using Hi-C loops to link enhancers to their potential target genes**

936

937 In the main result sections, we opted to use proximity as the criterion to identify the enhancer-  
938 associated gene for several reasons. First, the available human Hi-C contacts (Won et al. 2016) are very  
939 sparse: only 23% of human embryonic neocortex enhancers are covered. The 3D contacts in macaque  
940 (Luo et al. 2021) are even sparser, where 8,399 and 15,048 loops were identified in the germinal zone  
941 and cortical plate, covering only 2.68% of total macaque enhancers. Second, in the study of 'Activity-by-  
942 Contact model' (Fulco et al. 2019), based on a small number of experiments, the authors concluded that  
943 it is rare for an enhancer to skip the nearest gene (Fulco et al. 2019). Finally, for the enhancers included  
944 in Hi-C loops, around 60% of *de novo* gained enhancers contact their nearest genes, and more than 50%  
945 of both lost and conserved enhancers are in contact with their nearest genes (Figure S16), suggesting  
946 that our findings based on the nearest genes are robust.

947  
948 Nevertheless, we examined the results when the enhancers were mapped to their putative targets  
949 based on Hi-C loops. The findings based on the Hi-C loops are consistent with the ones based on the  
950 proximity rule. For example, the *de novo* gained enhancers tend to associate with an increase in the  
951 expression of their target gene, whereas the lost enhancers show the reverse trend (Figure 2A and  
952 Figure S2). Enhancers are more likely to regulate the tissue-specific genes of embryonic neocortex either  
953 based on proximation rule (Figure 7B) or Hi-C contacts (Figure S10B). In addition, using either gene  
954 proximation rule (Figure 7C) or Hi-C contact (Figure S10A), we observed that *de novo* gained enhancers  
955 are more likely to turn on gene expression compared to HGEs.

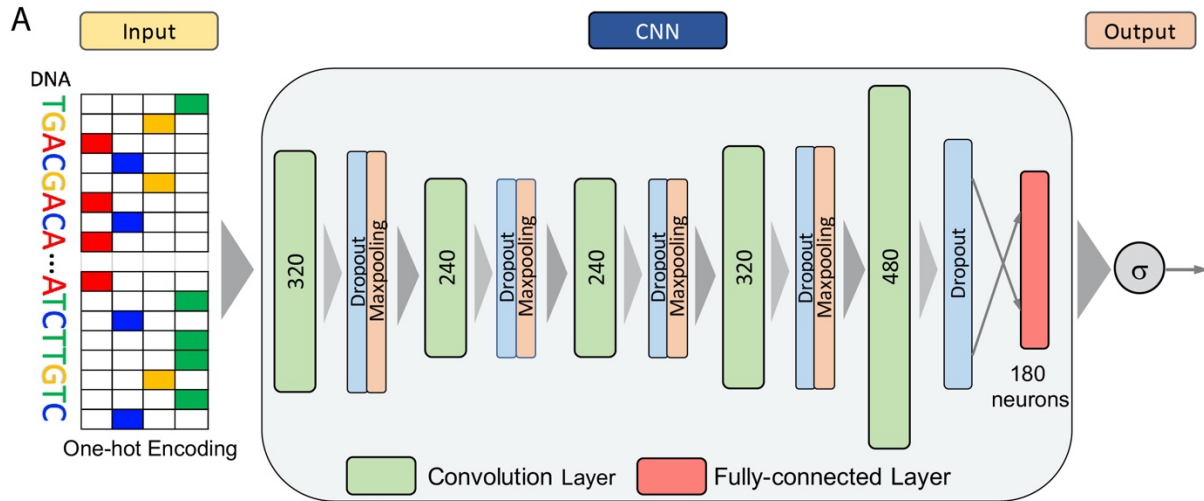
956  
957  
958  
959  
960  
961

962 **Supplementary Figures**

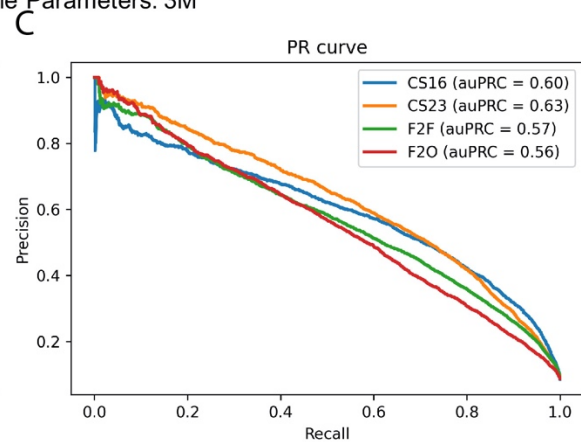
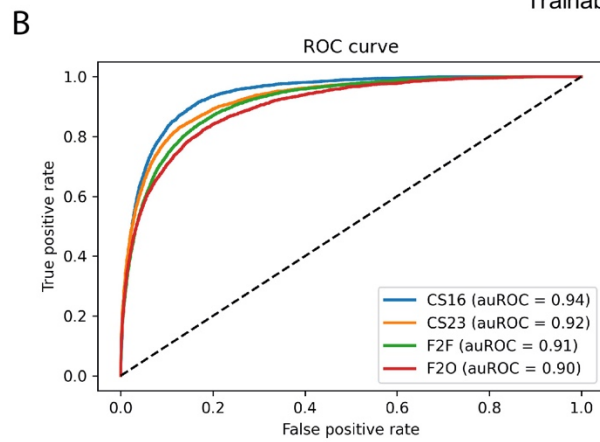
963

964 Figure S1

965



Trainable Parameters: 3M



**D**

		Predict				
		CS16	CS23	F2F	F2O	
Model	auROC	CS16	0.942	0.896	0.890	0.890
	CS23	0.912	0.924	0.900	0.886	
	F2F	0.922	0.907	0.914	0.887	
	F2O	0.897	0.902	0.894	0.904	

**E**

	CS16	CS23	F2F	F2O
CS16	1	0.322	0.278	0.236
CS23	0.322	1	0.448	0.415
F2F	0.278	0.448	1	0.490
F2O	0.236	0.415	0.490	1

$$\text{Similarity}(A, B) = \frac{A \cap B}{\min(A, B)}$$

966

967

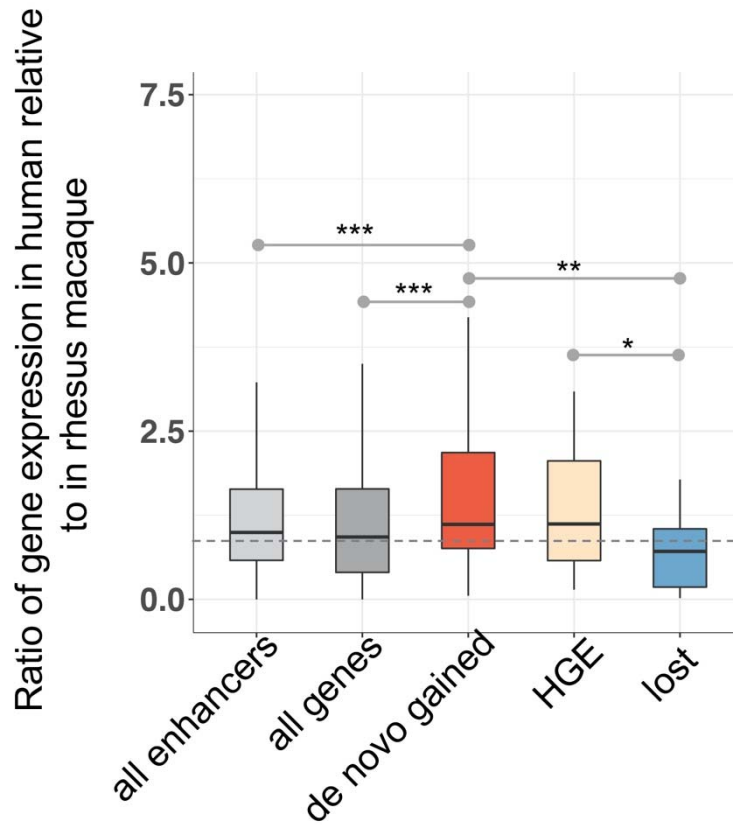
968

969

970

971 **Figure S1. Deep learning model of human embryonic neocortex enhancers used to score**  
972 **enhancer activity.** A) Structure of the deep convolutional model. The number within each  
973 convolutional layer indicates the number of kernels. B) ROC curve of the model. C) PR curve of  
974 the model. D) Model performance across four stages. E) Similarity between enhancer sets  
975 across stages.  
976  
977  
978

979 Figure S2



990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019

Figure S2. The expression level of genes with Hi-C loops to the *de novo* gained enhancers is increased, and so is the previously published enhancers that increase activity in human (HGEs, PMID: 25745175). By contrast, the genes in contact with the lost enhancers show the reverse trend. “all enhancers” refer to the genes link to all enhancers. \*Wilcoxon p-value  $\leq 0.01$ . \*\* Wilcoxon p-value  $\leq 1e-3$ . \*\*\* Wilcoxon p-value  $\leq 1e-5$ .

1020 Figure S3

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

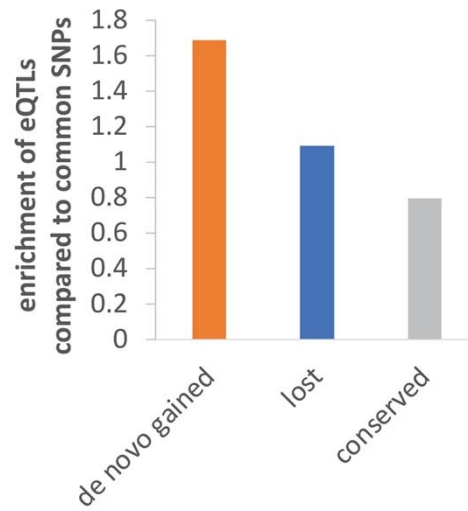
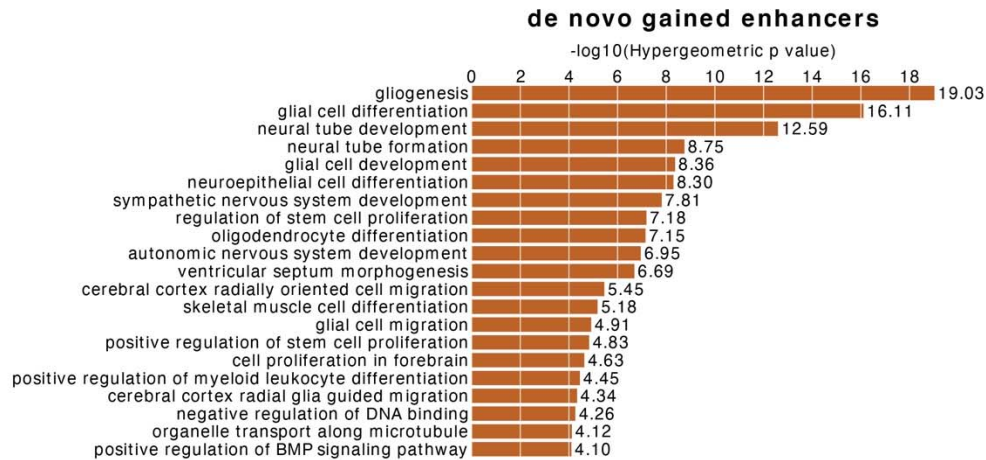


Figure S3. Enrichment of eQTLs compared to common SNPs in the three sets of enhancers. Specifically, the enrichment = fraction of eQTLs in enhancers/fraction of SNPs in enhancers.

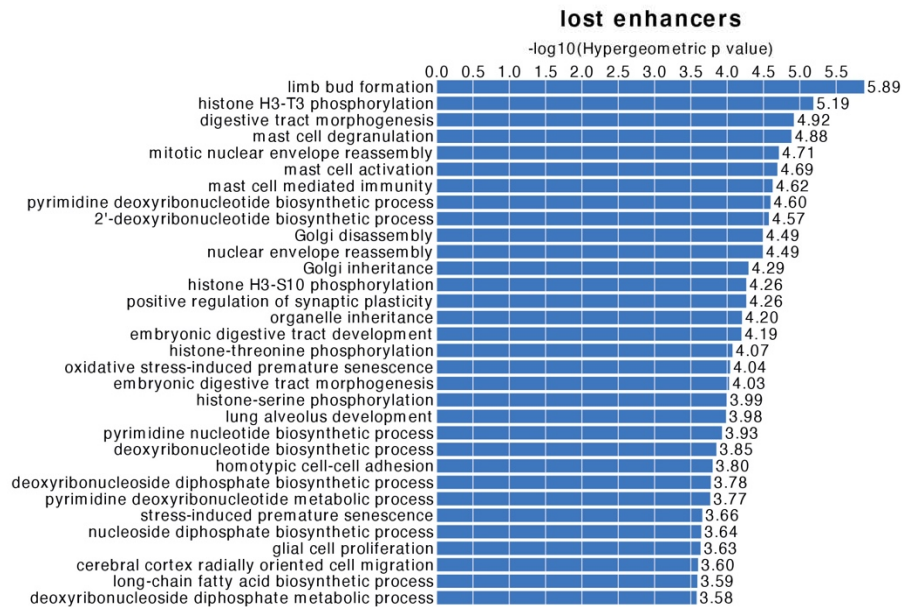


1046 Figure S4  
1047

A



B

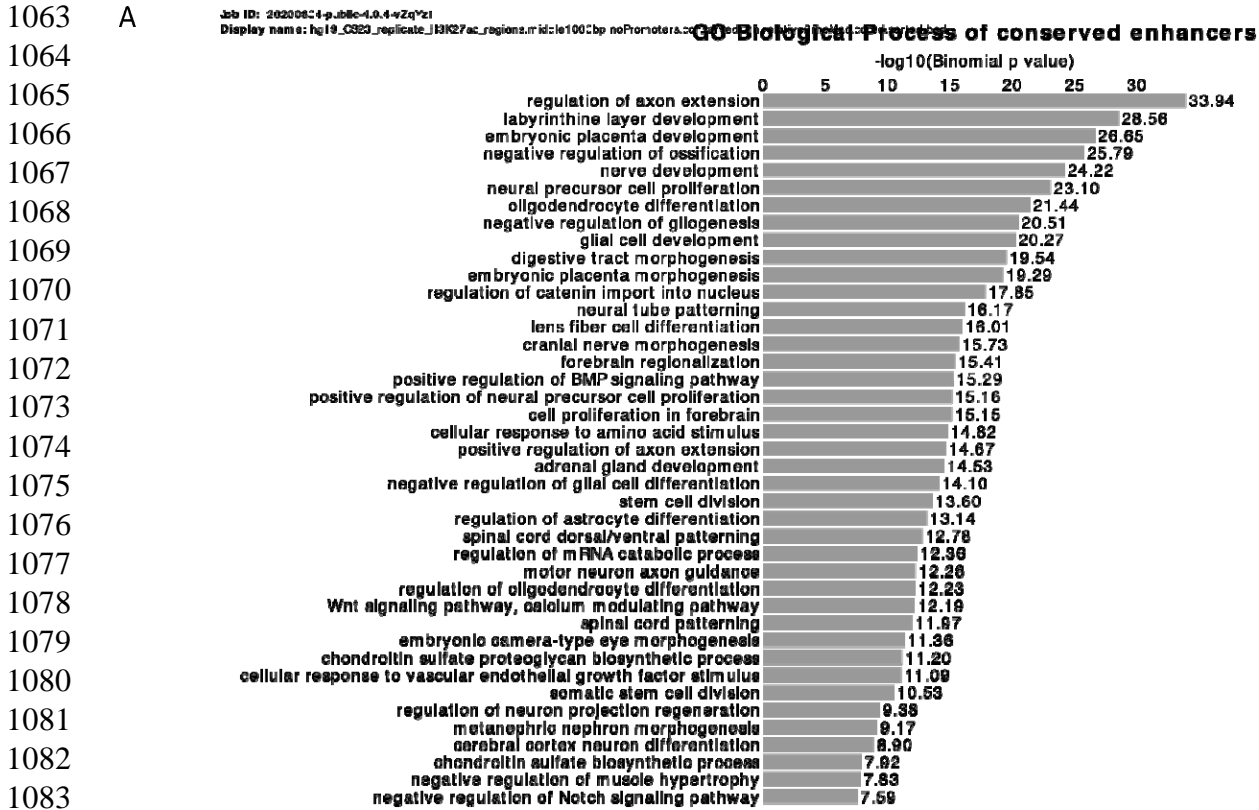


1048  
1049

1050 Figure S4. *De novo* gained enhancers are associated with essential CNS-related biological  
1051 processes, using all fetal brain enhancers (de la Torre-Ubieta et al. 2018) as the background. (A)  
1052 GO terms of *de novo* gained enhancers. (B) GO terms of lost enhancers. We apply GREAT with  
1053 the single nearest gene association rule to do functional enrichment of genes near enhancers.  
1054 The GO terms will be considered as enriched if it has at least 10 gene hits with FDR threshold  
1055 set as 0.01.

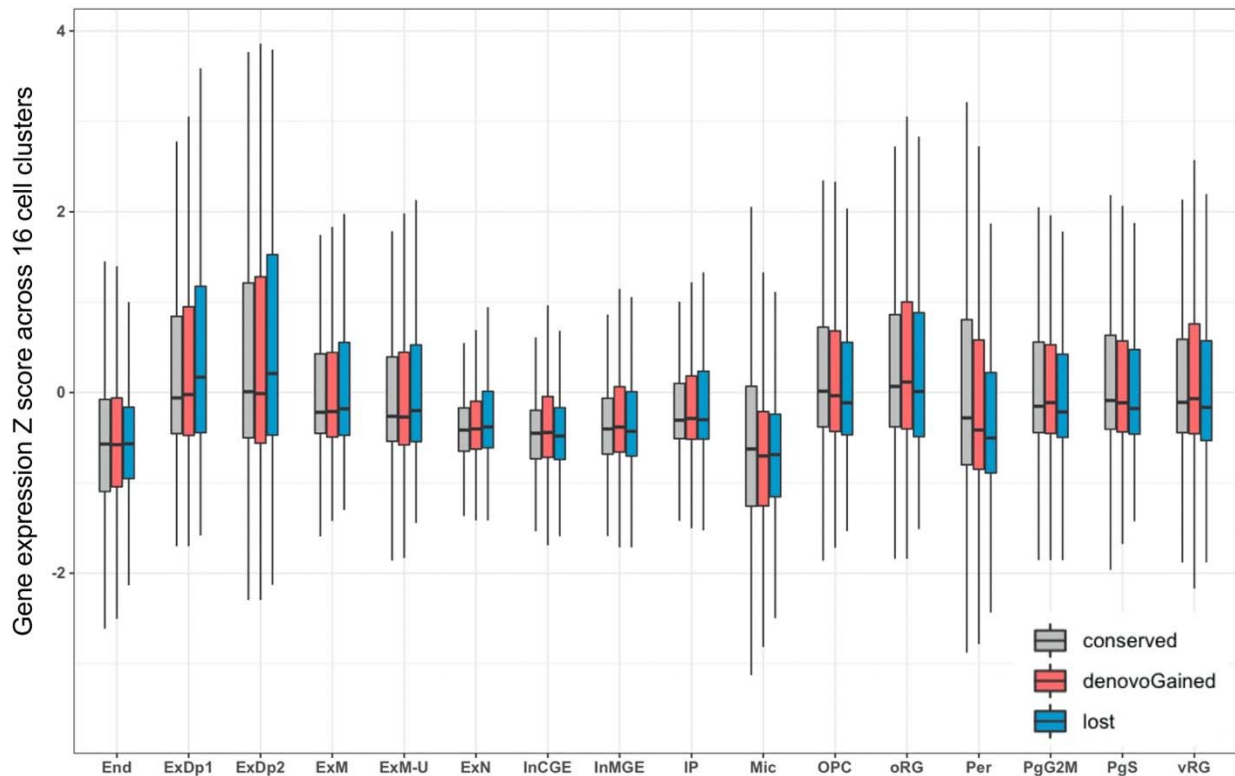
1056  
1057  
1058

1059  
1060  
1061 Figure S5  
1062



1101 Figure S5. A) Enriched GO Biological Processes terms of conserved enhancers. B) Enriched GO  
1102 Molecular Function terms of the three sets of enhancers.

1103 Figure S6  
1104



1105  
1106

1107 Figure S6. Z scores of expression of genes nearby the three sets of enhancers across 16 cell  
1108 clusters. The lack of statistical significance may partly be due to the high variability/noise in  
1109 single cell gene expression data, and also because only a subset of the genes near *de novo*  
1110 gained enhancers are likely to drive cluster-specific expression as revealed in our fractional  
1111 analysis (Figure 4C) but obscured in our analysis of z-scores for all genes.

1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127

1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171

Figure S7

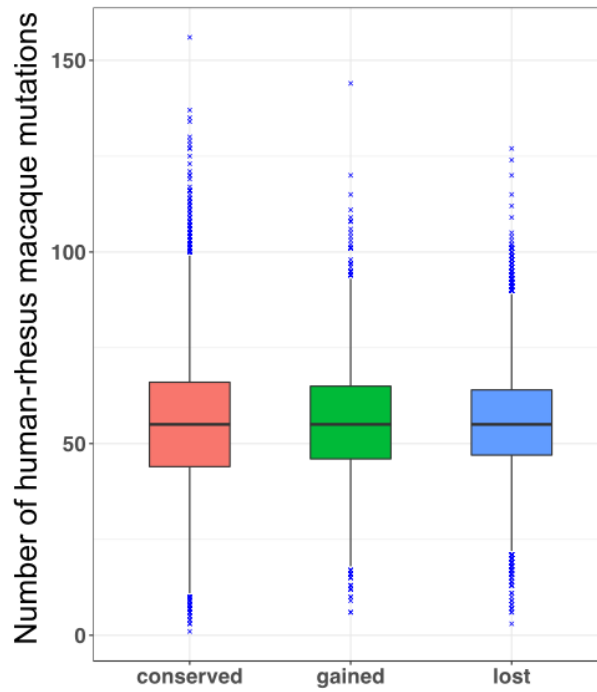
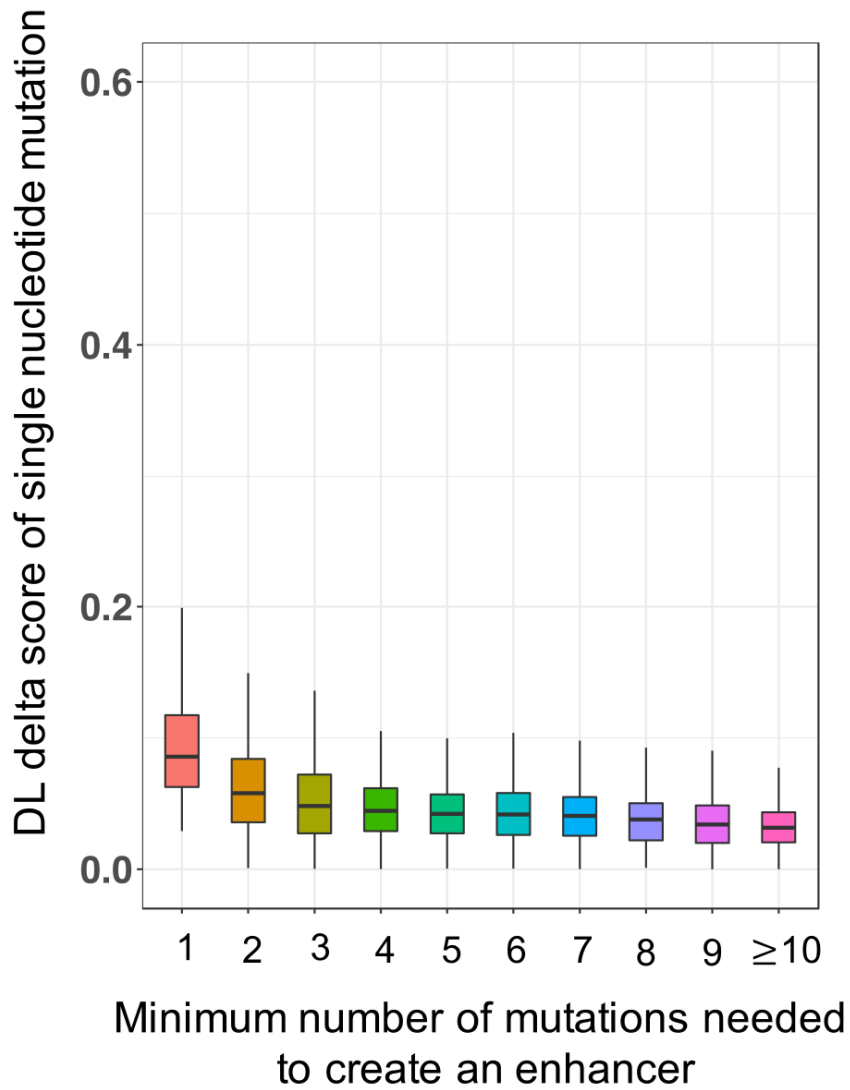


Figure S7. Number of human-macaque mutations within enhancers.

1172  
1173  
1174  
1175  
1176

Figure S8

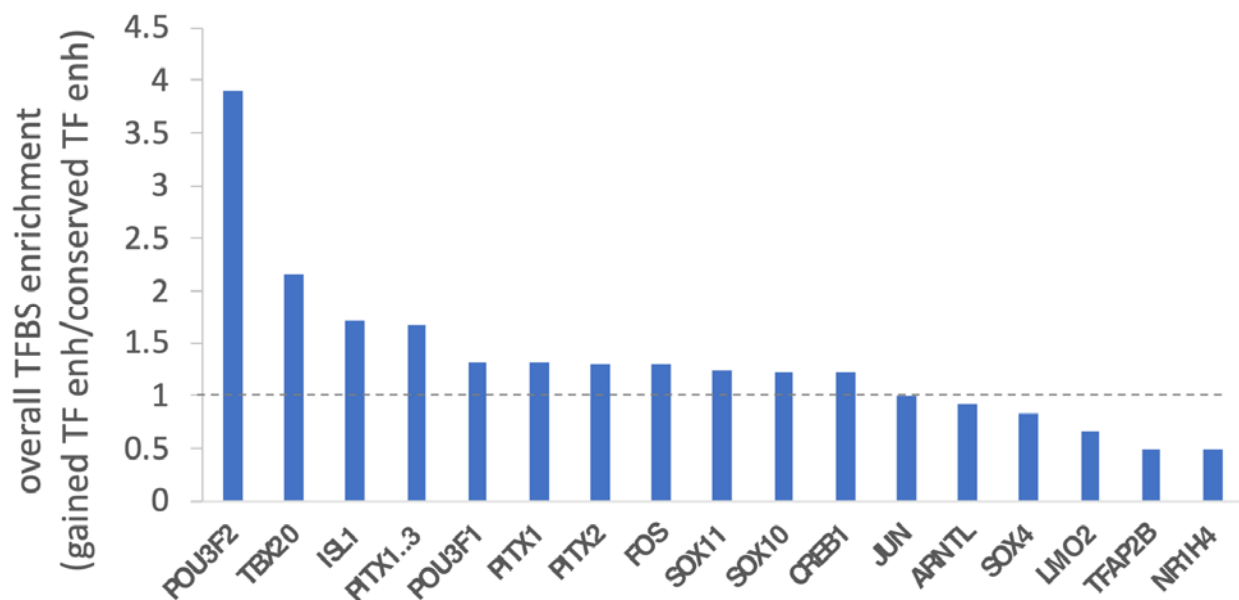


1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Figure S8. Distribution of delta score of the single nucleotide mutations that are minimally needed to create an enhancer.

1188  
1189  
1190  
1191  
1192  
1193  
1194

Figure S9



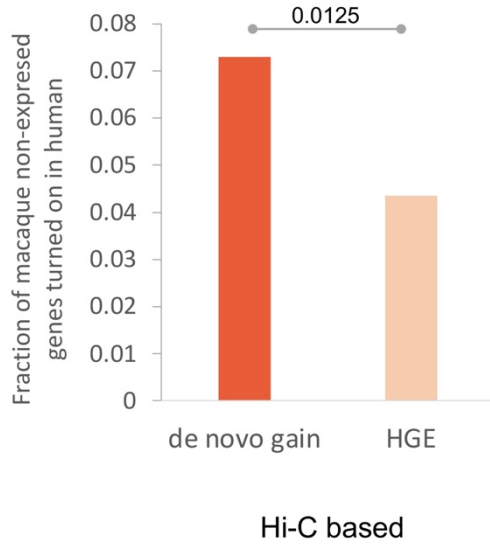
1195  
1196  
1197  
1198  
1199  
1200  
1201

Figure S9. TFBS enrichment of gained enhancers associated with TFs, as compared to the conserved enhancers associated with TFs.

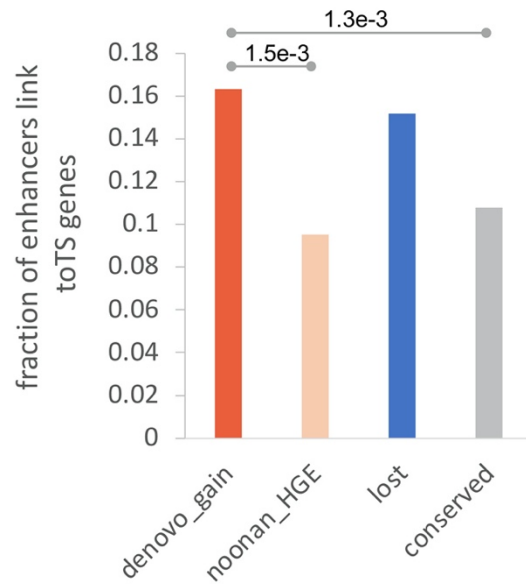
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225

Figure S10

A



B

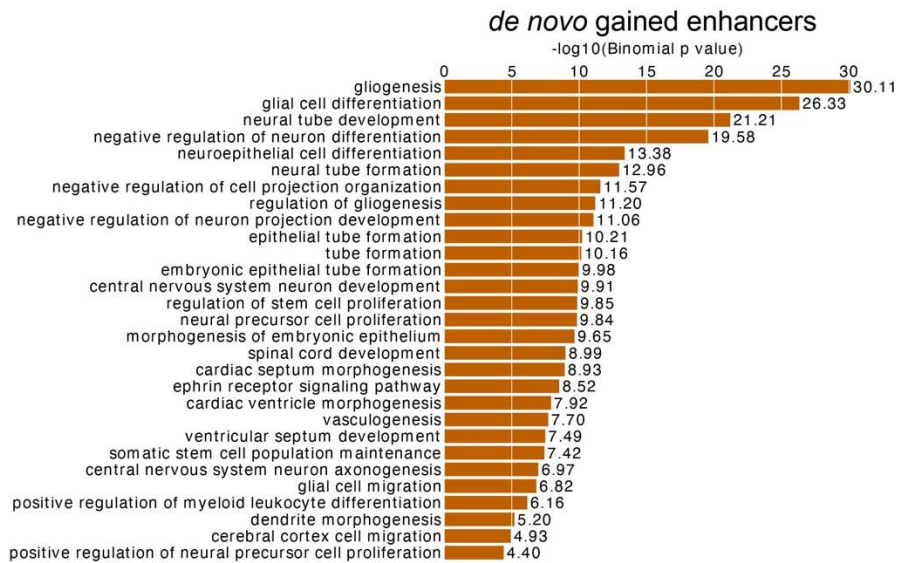


1226 Figure S10. (A) Fraction of enhancers in contact with genes whose RPKM < 1 in macaque and >  
1227 1 in human. (B) Fraction of enhancers in 3D contact with the most tissue-specific genes.

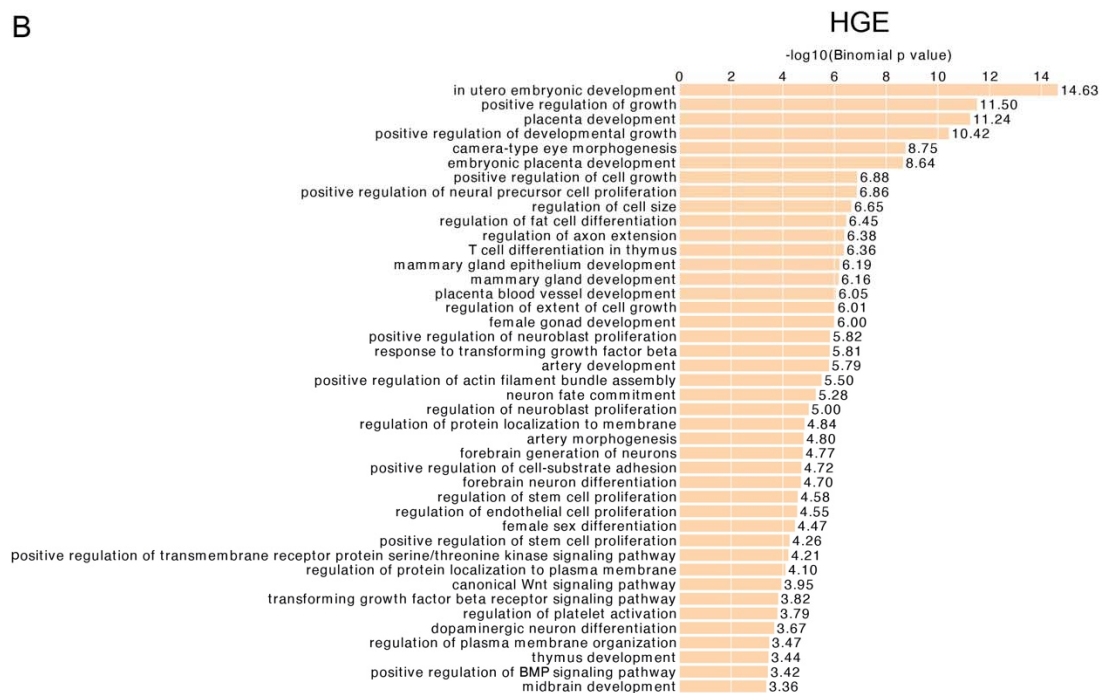
1228  
1229  
1230  
1231  
1232  
1233  
1234

1235 Figure S11  
1236

A



B

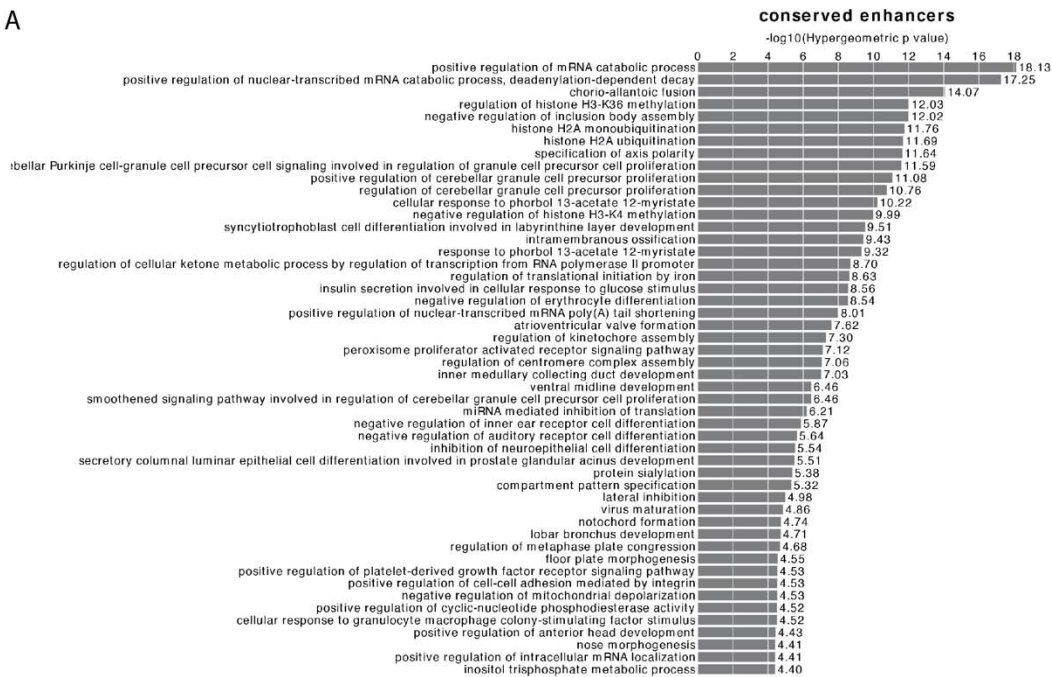


1237 Figure S11. Enriched GO biological processes of *de novo* gained enhancers (A), and HGEs (B)  
1238 using whole genome as the background.  
1239  
1240

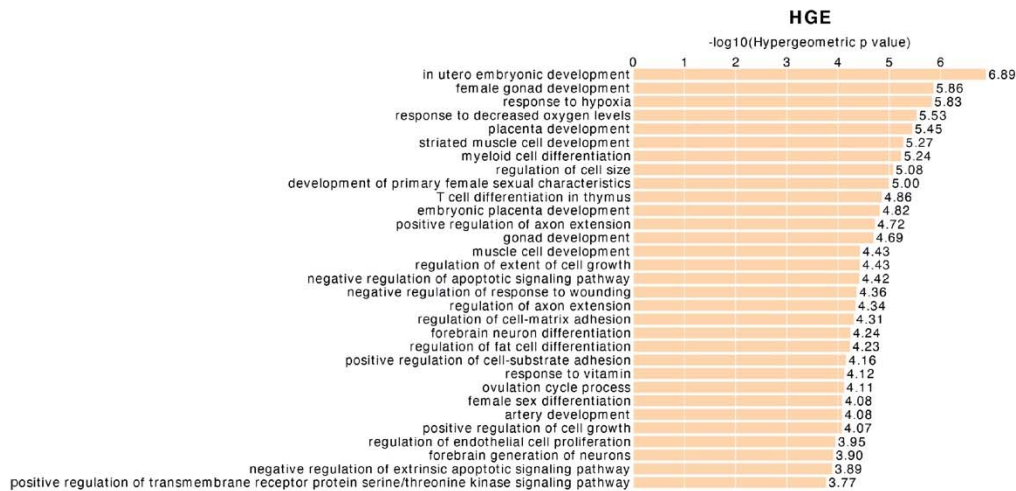


1241

A



B



1242

1243

1244

1245 Figure S12. Enriched biological processes of a set of enhancers, using all fetal brain enhancers

1246 (de la Torre-Ubieta et al. 2018) as the background. (A) GO terms of conserved enhancers. (B)

1247 GO terms of HGEs. We apply GREAT with the single nearest gene association rule to do

1248 functional enrichment of genes near enhancers. The GO terms will be considered as enriched if

1249 it has at least 10 gene hits with FDR threshold set as 0.01.

1250

1251

1252

1253

1254 Figure S13

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

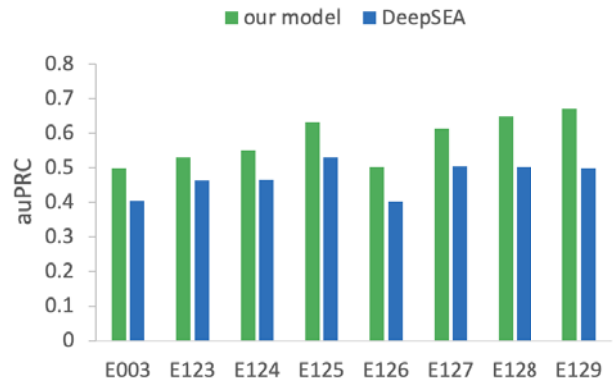
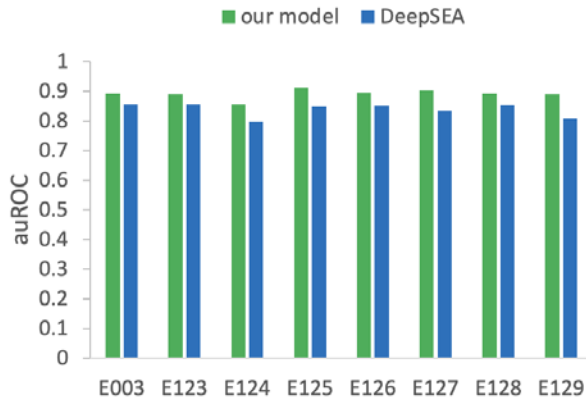
1300

1301

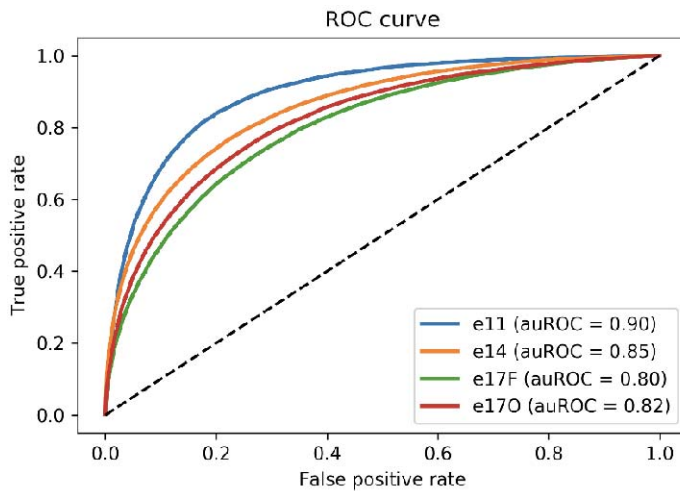
1302

A.

B.



C.



Performance of the DLM. (A) auROC and (B) auPRC of our model in predicting H3K27ac in 8 tissues which are tested by DeepSEA. (C). ROC curve of CS23 model tested on mouse embryonic neocortex enhancers corresponding to different stages of development (e11, e14, e17F, e17O). The E numbers on the x-axis are the tissue IDs defined by the Roadmap Epigenomic Project. E003: H1 Cell Line, E123: K562 Leukemia Cell Line, E124: Monocytes-CD14+ RO01746 Cell Line, E125: NH-A Astrocytes Cell Line, E126: NHDF-Ad Adult Dermal Fibroblast Primary Cells, E127: NHEK-Epidermal Keratinocyte Primary Cells, E128: NHLF Lung Fibroblast Primary Cells, E129: Osteoblast Primary Cells.

1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330

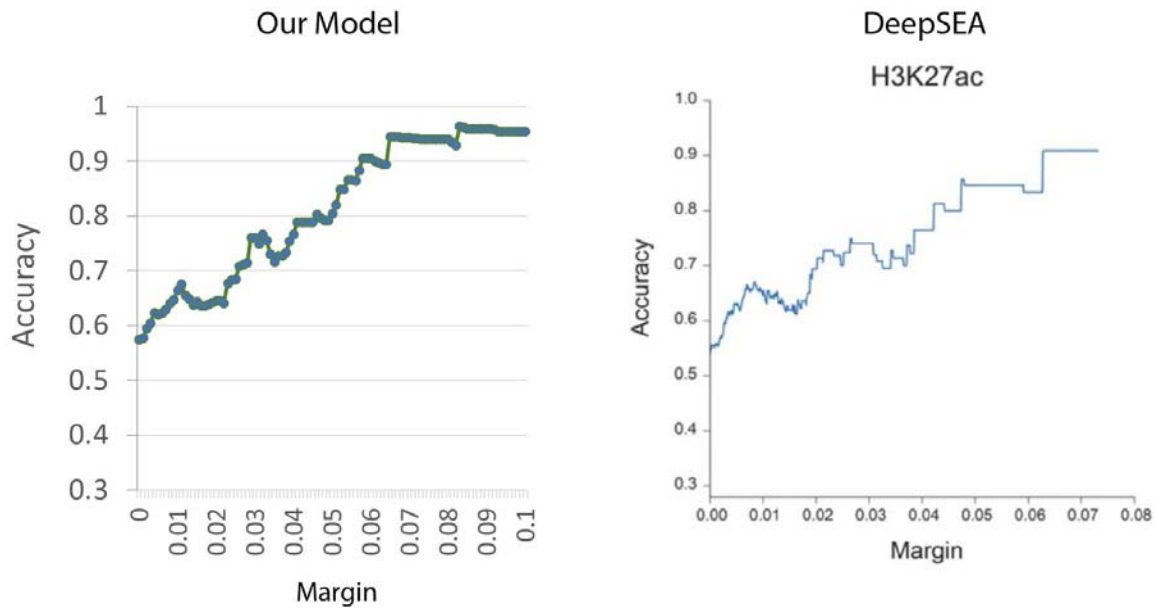


Figure S14. Deep learning histone mark classifiers provided accurate prediction of allele specific effects on histone marks H3K27ac (the allele with stronger histone mark signals). The predictions were evaluated with histone mark QTLs identified with FDR < 0.1 in Yoruba lymphoblastoid cell lines (McVicker, G. et al. Science 342, 747-749 (2013)). Margin shown on the x axis is the threshold of predicted probability differences between the two alleles for classifying high-confidence predictions. Performance is measured by accuracy (y-axis) of predicting the allele with higher read counts based on DLM score difference above certain threshold (x-axis).

1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356

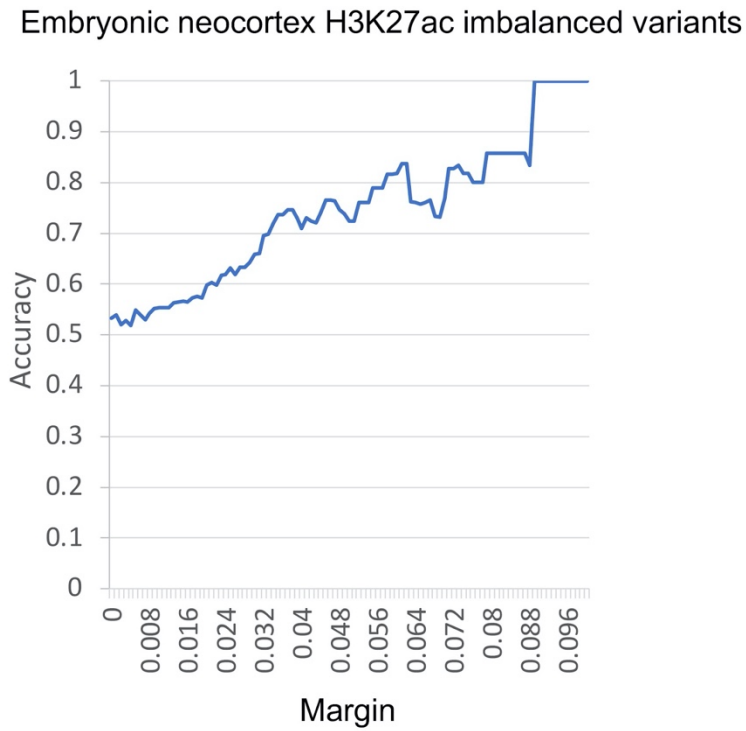
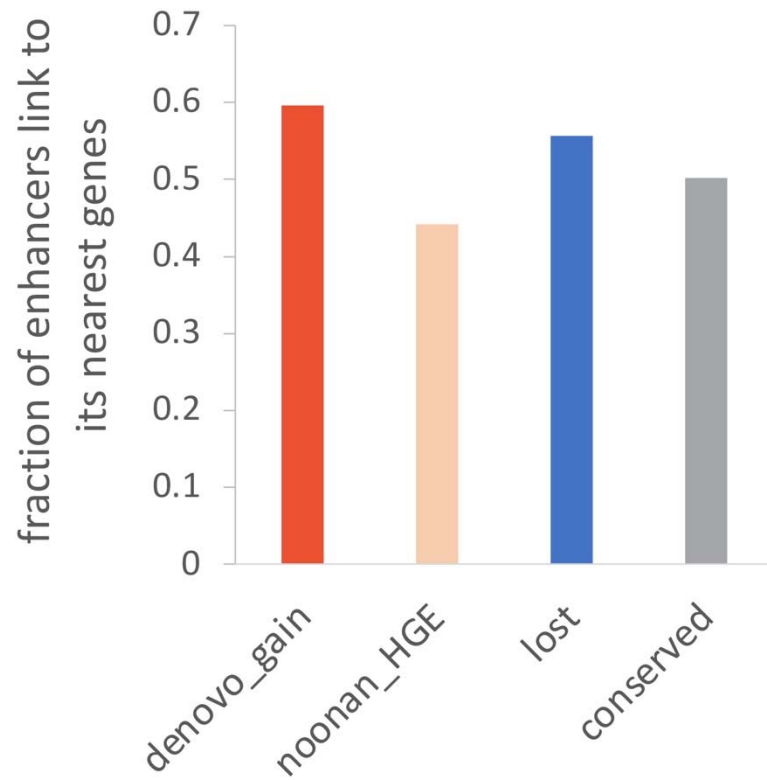


Figure S15. The DLM of CS23 H3K27ac accurately predict allelic imbalanced heterozygous variants within CS23 H3K27ac peaks.

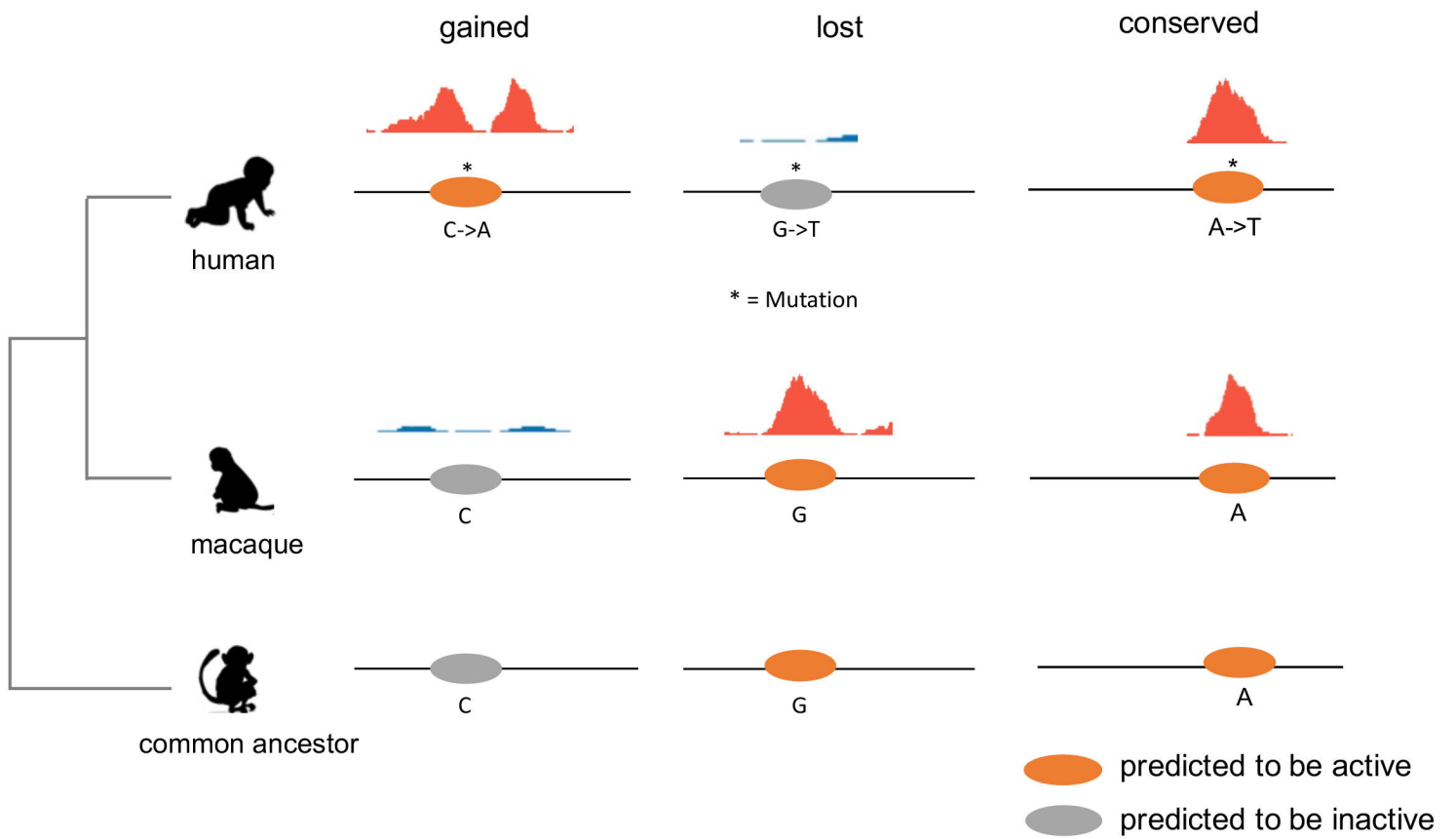
1357



1358  
1359

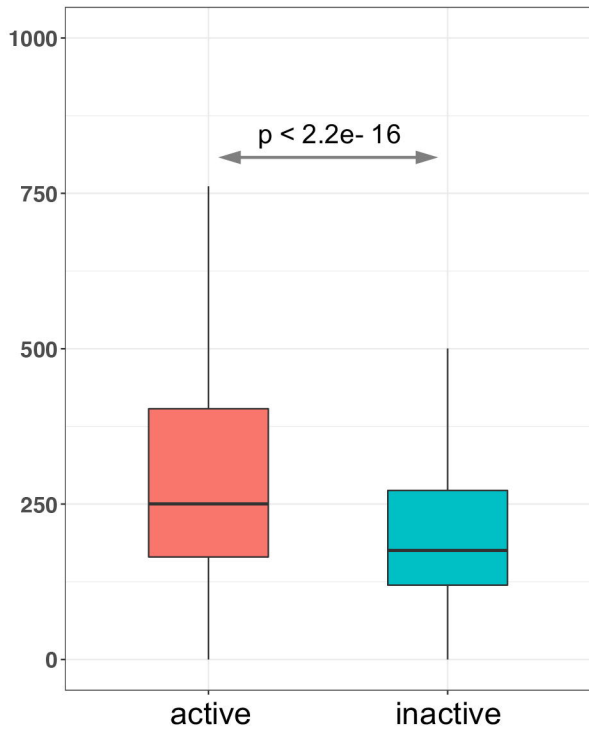
Figure S16. Fractions of enhancers that contact their nearest gene.

**A**

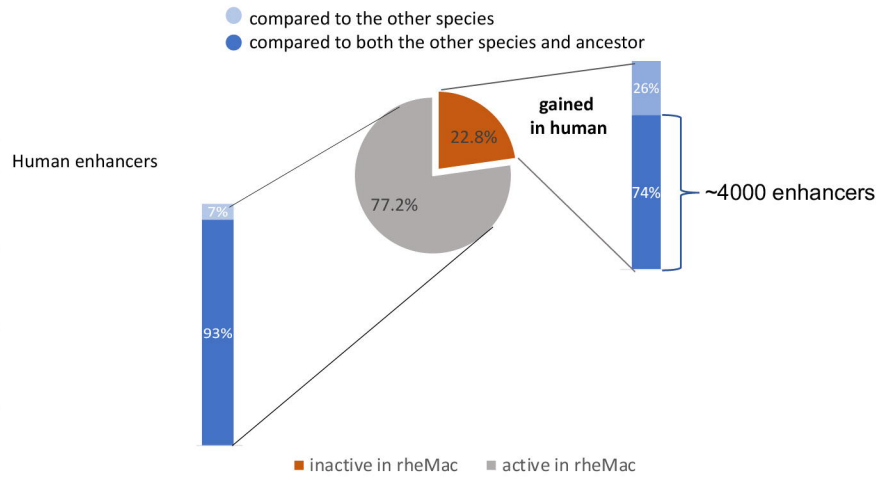


**B**

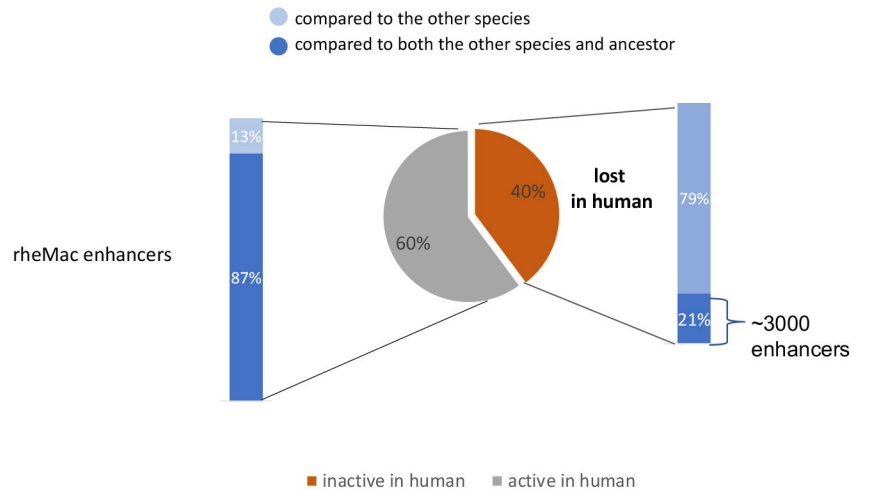
Integrated H3K27ac signal of rhesus orthologous sites of human embryonic neocortical enhancers



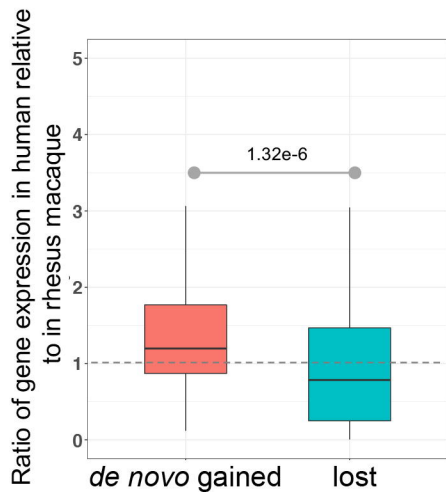
**C**



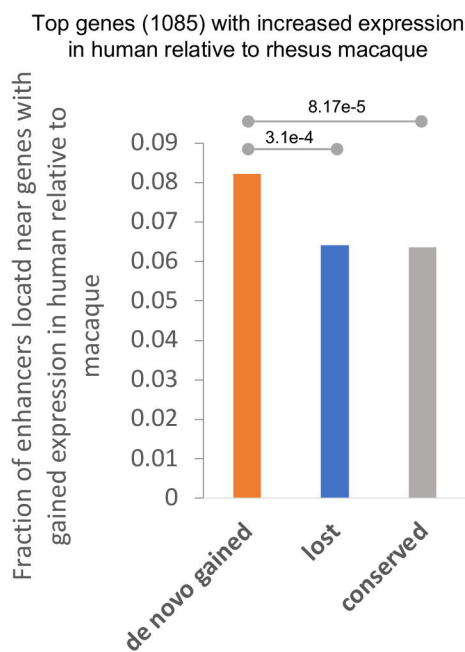
**D**



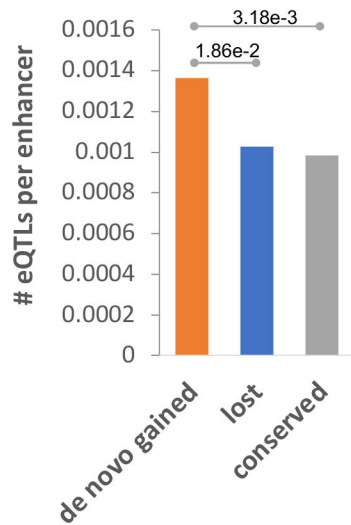
A



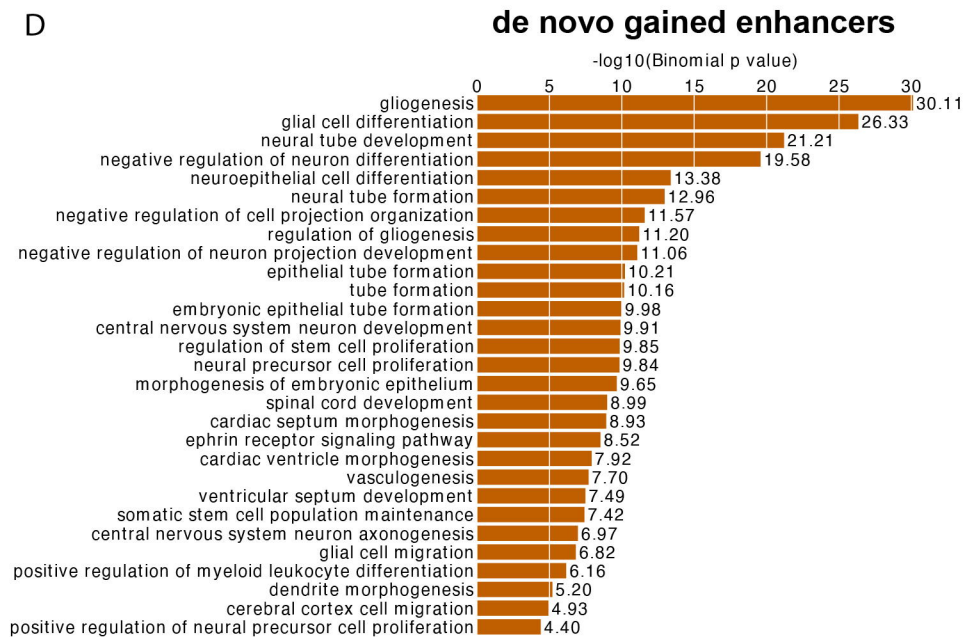
B



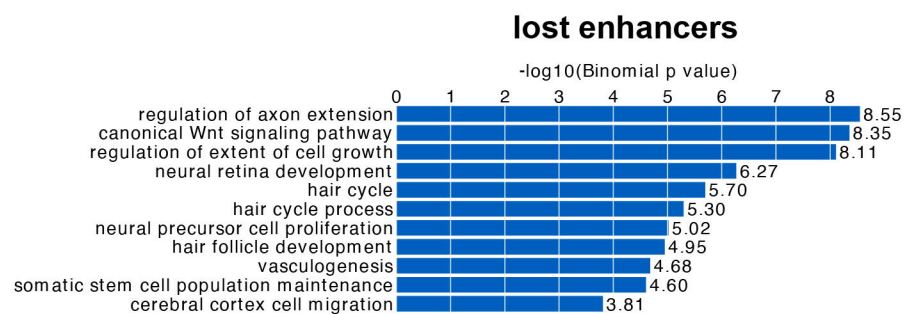
C



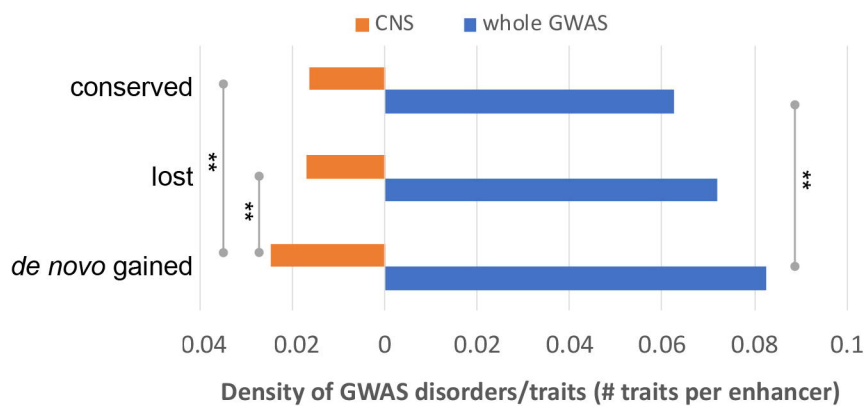
D



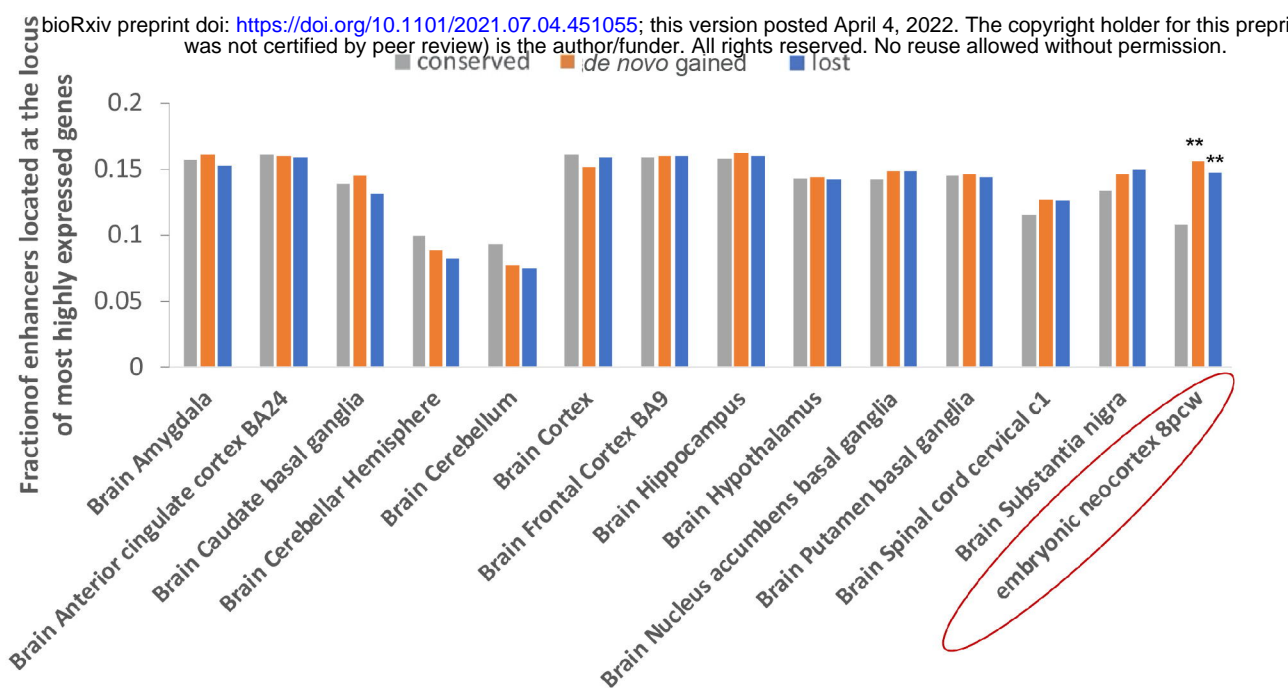
E



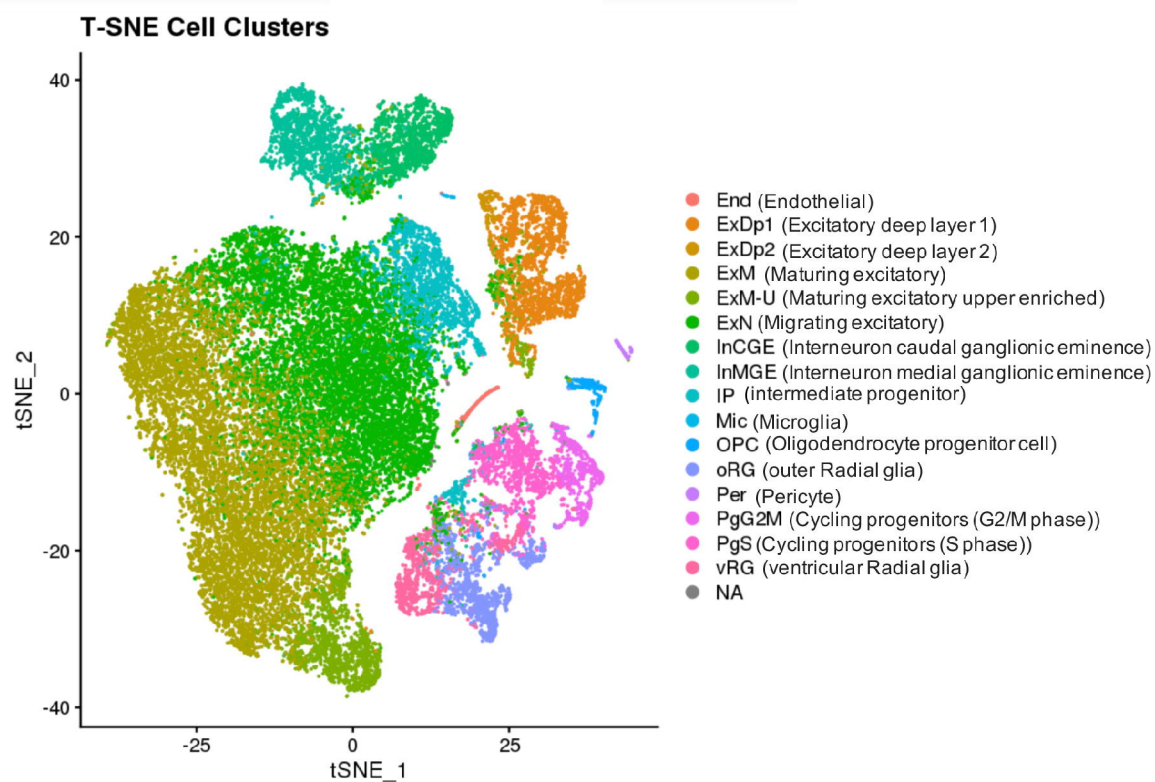
F



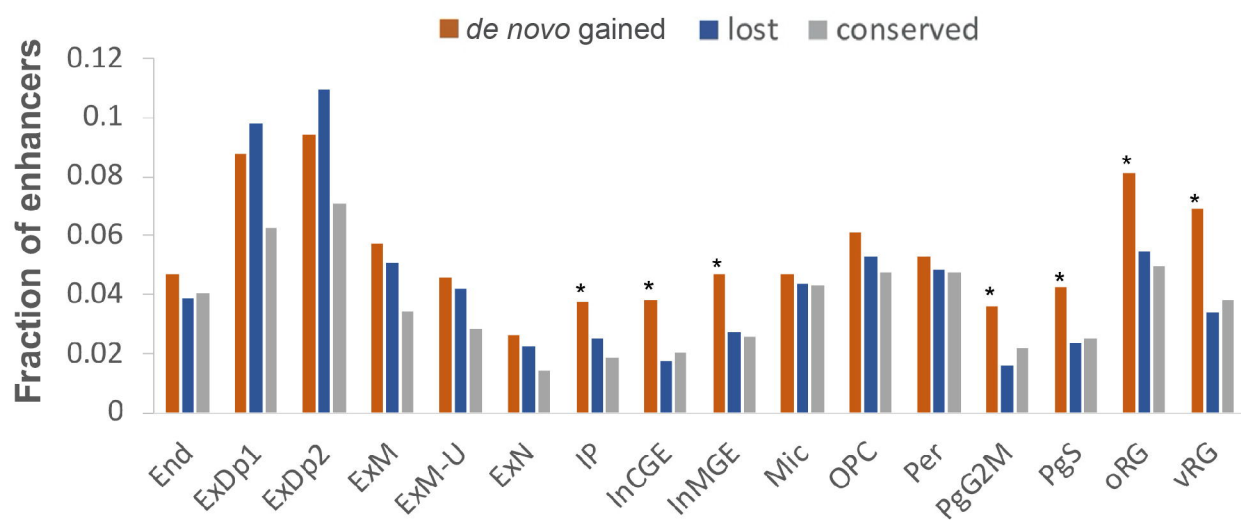




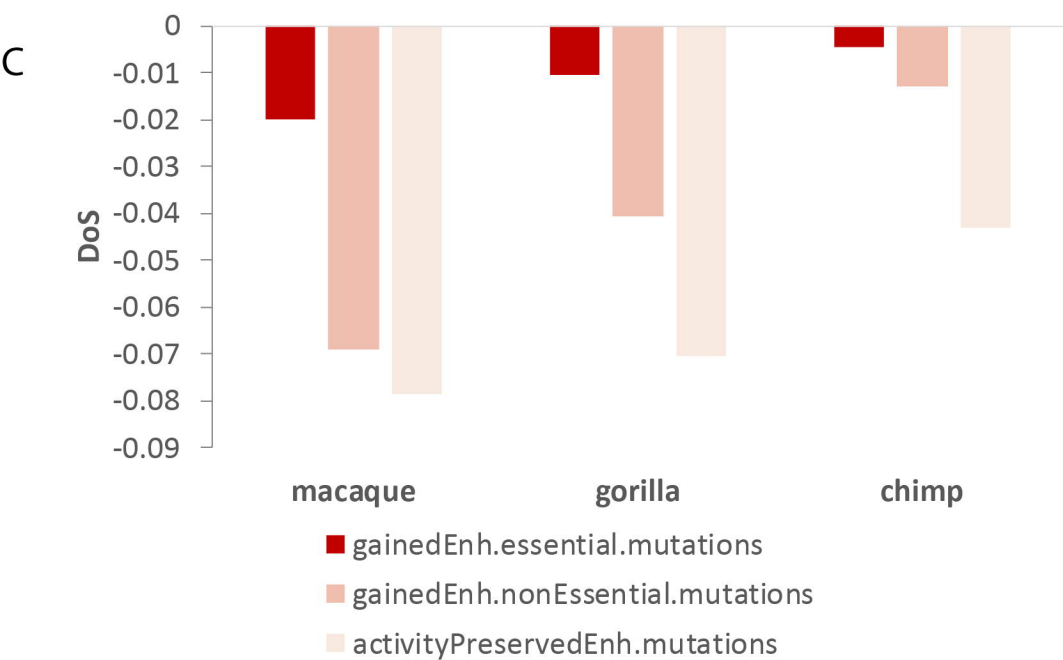
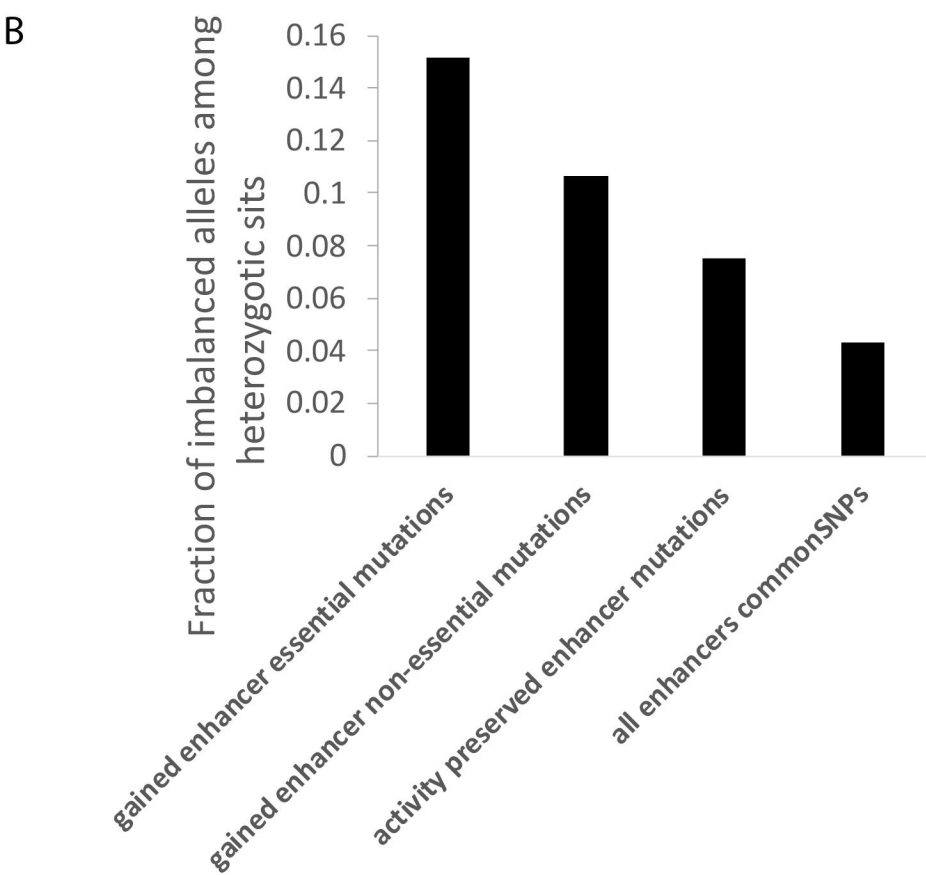
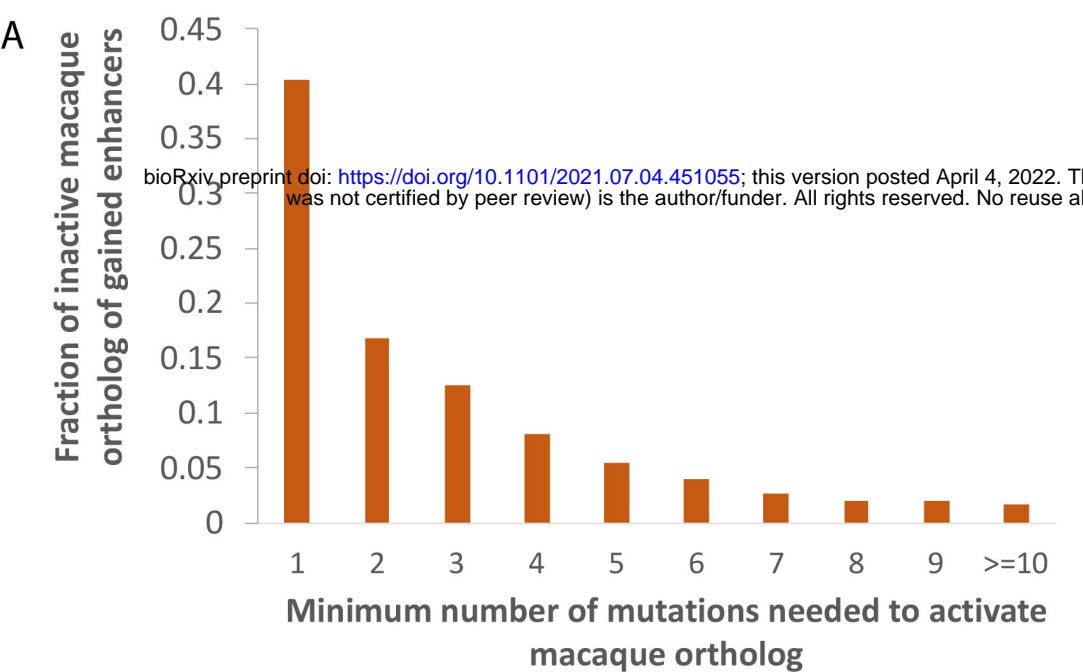
B

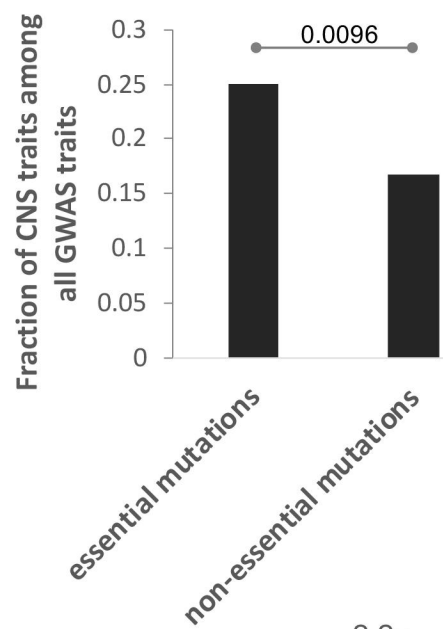
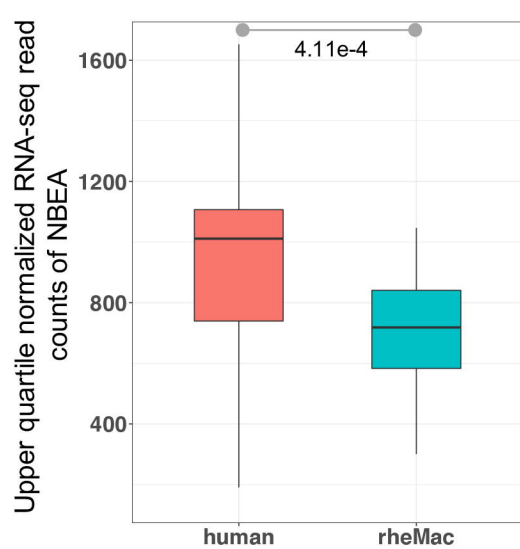
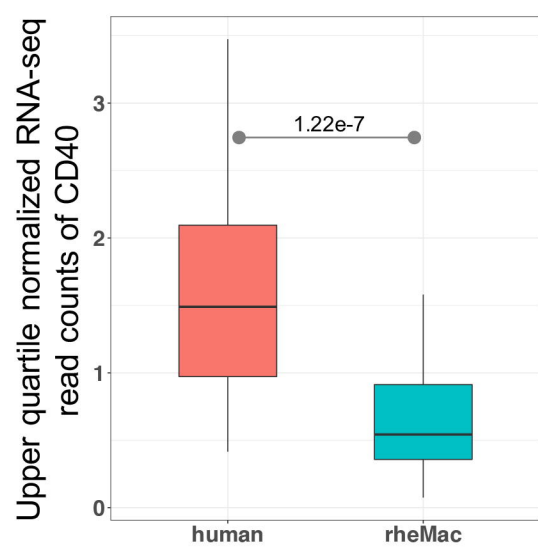
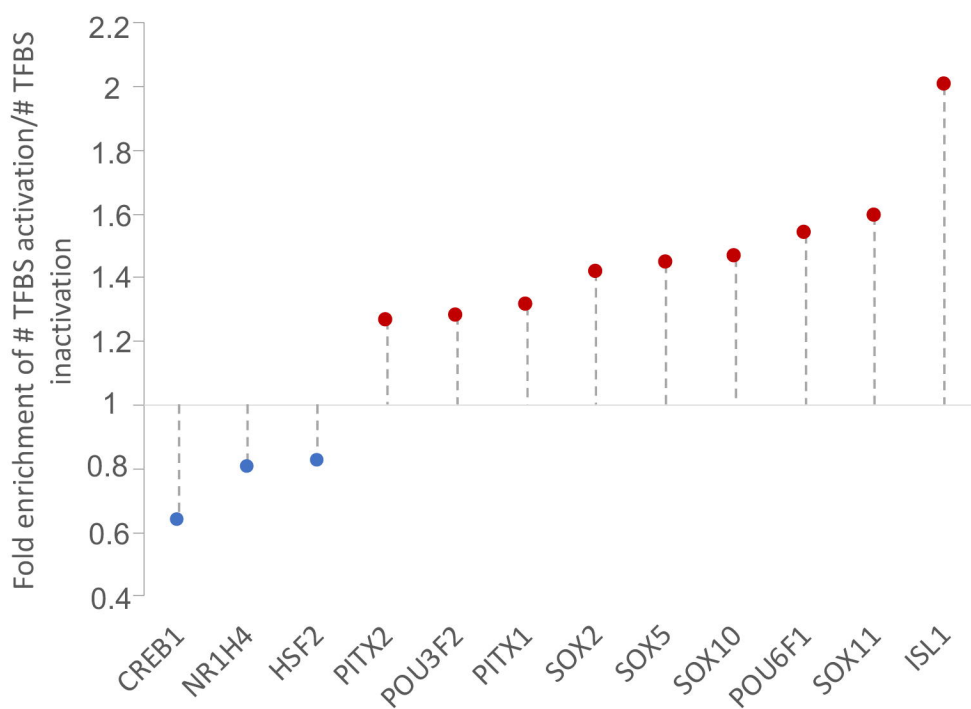
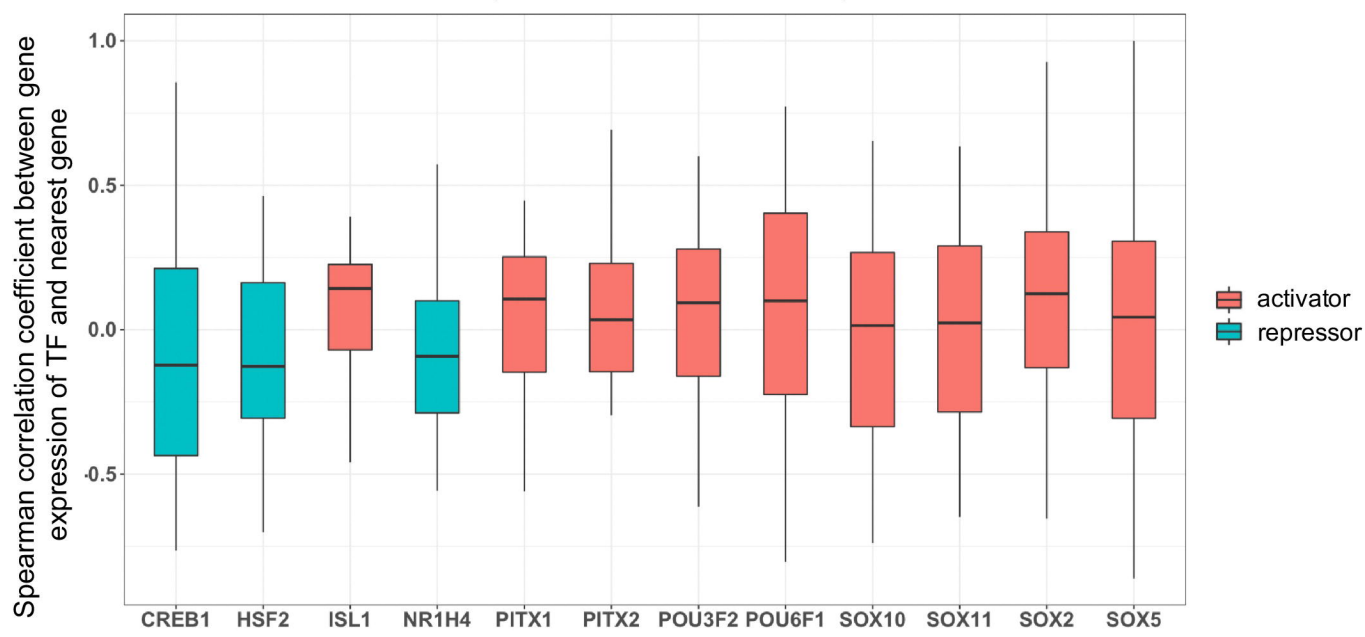


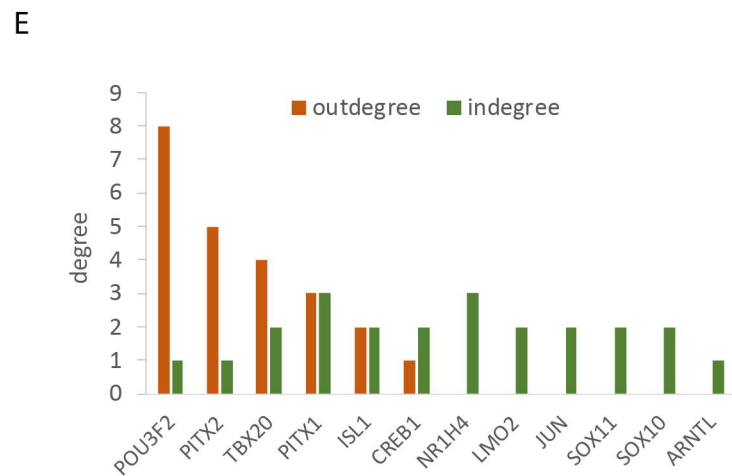
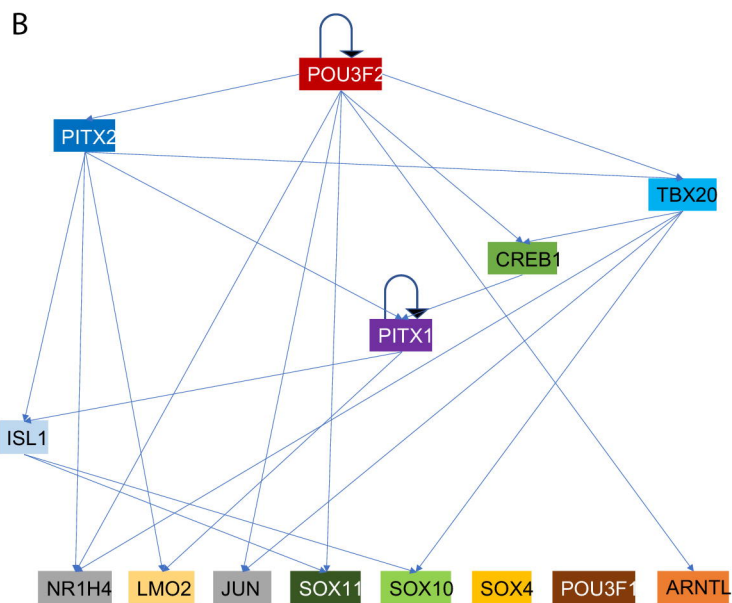
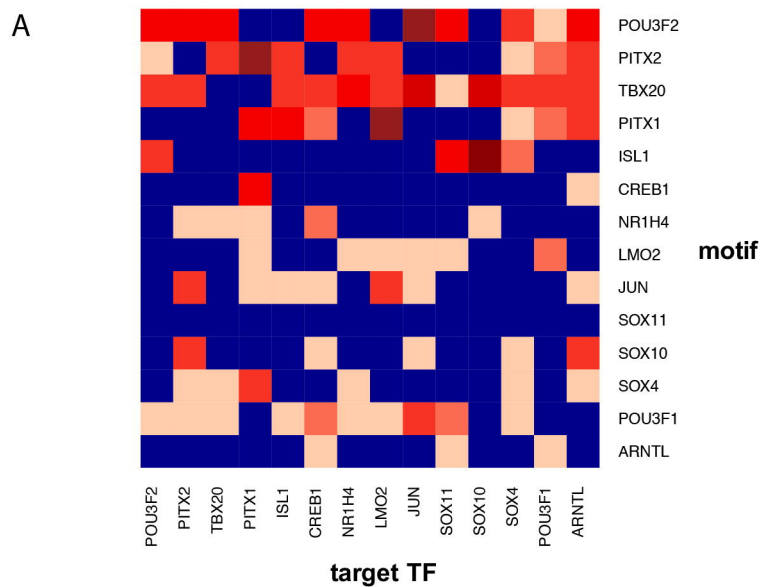
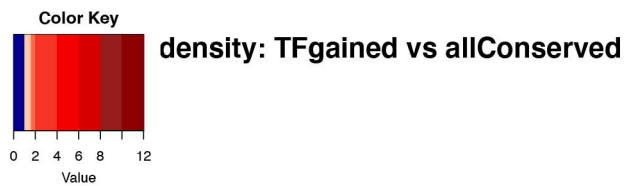
C



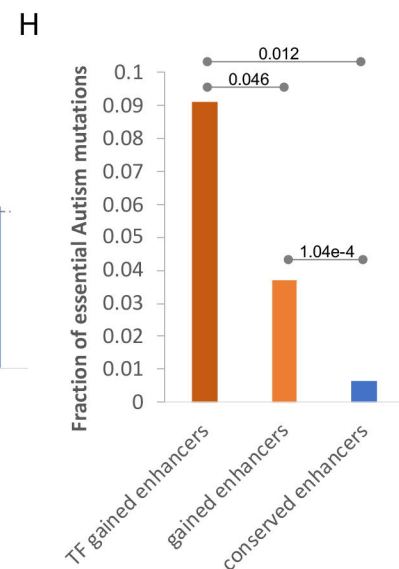
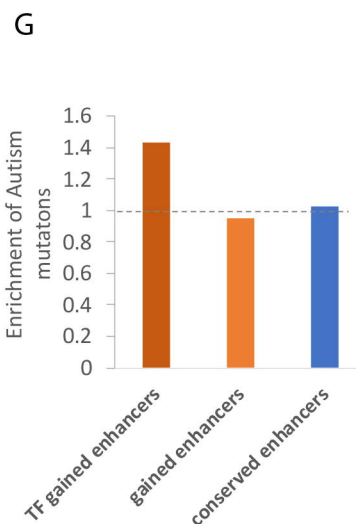
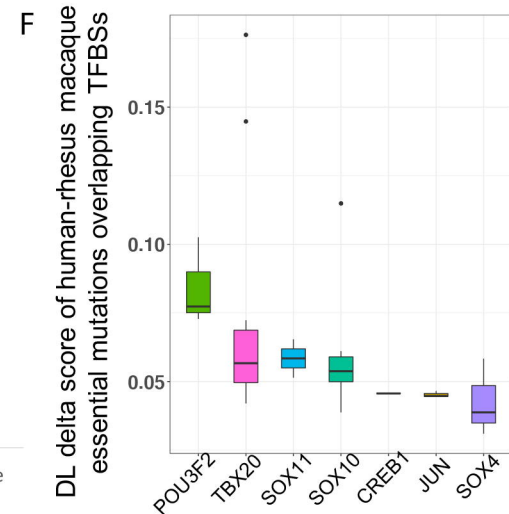
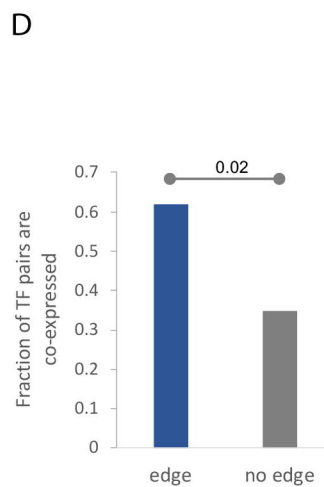
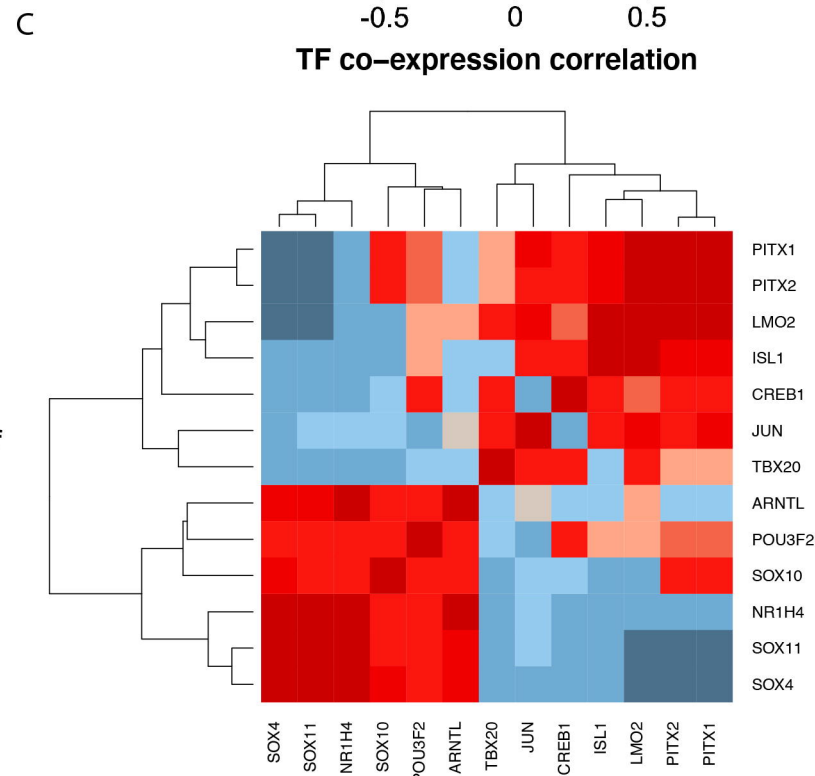




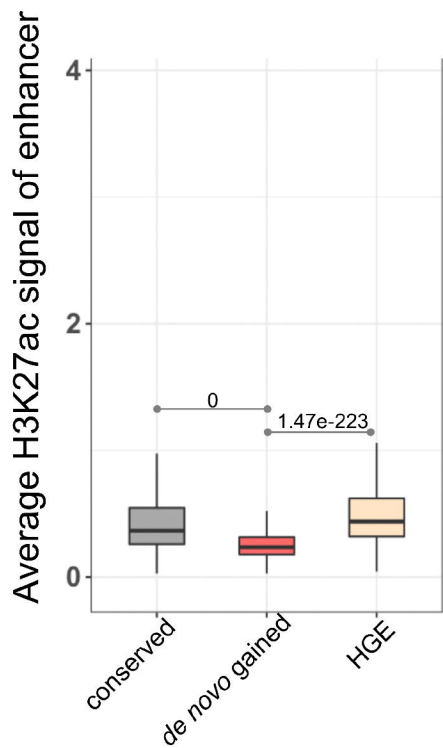
**A****B****C****D****E**



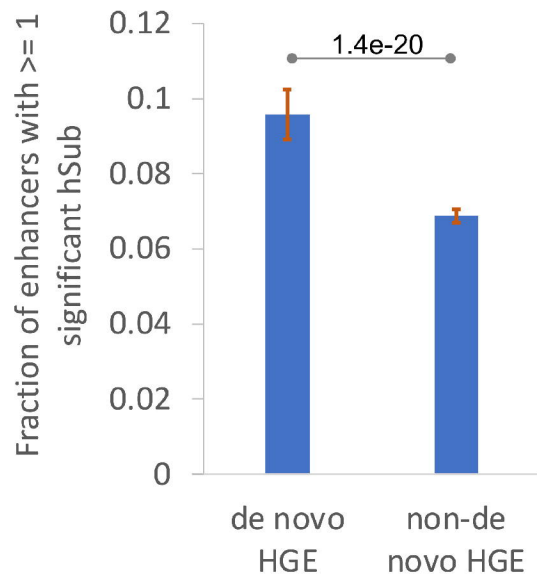
Outdegree (motif): TFBS to target TF  
 Indegree (target TF): target TF from TFBS



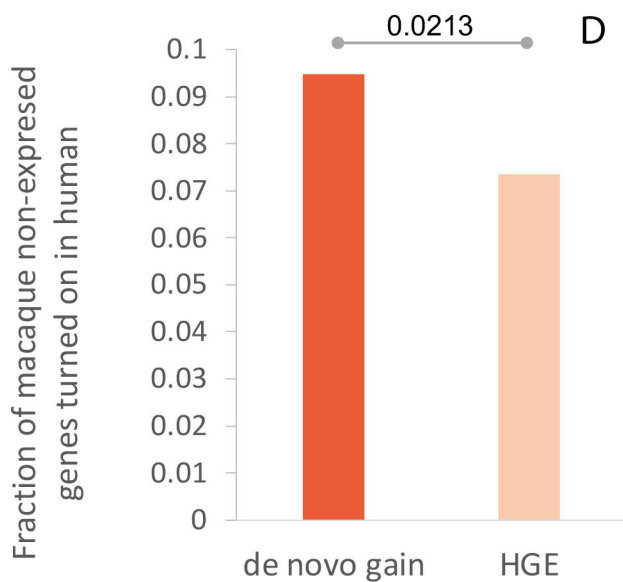
A



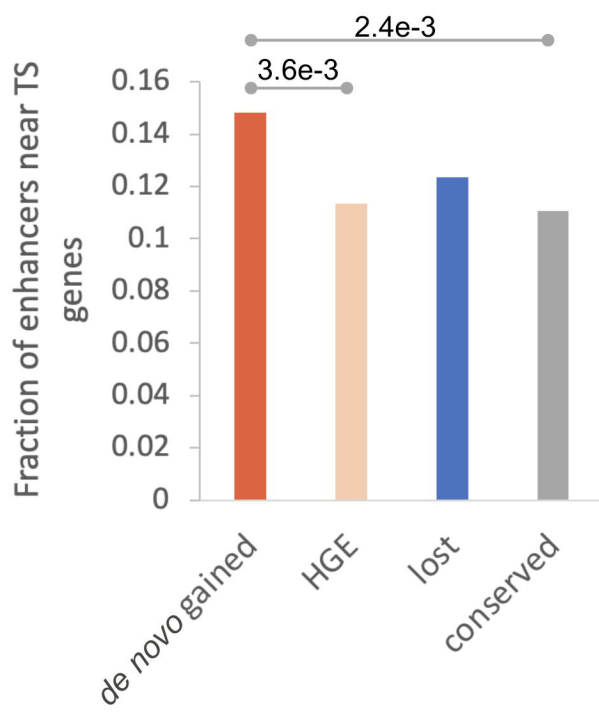
B



C



D



Approximation based