

1 Comprehensive resource and integrative model for 2 functional genomics of the adult brain

3 Abstract

4 Understanding how genomic variation influences brain phenotypes remains a key challenge in
5 neuroscience, one where the potential of functional genomic approaches has not yet been fully
6 realized. To this end, the psychENCODE consortium developed a comprehensive, population-
7 level resource that includes thousands of samples processed for healthy controls and
8 neuropsychiatric disorders. Available online, the resource comprises genotyping, RNA-seq,
9 ChIP-seq, and single-cell data, in addition to analytic summaries of quantitative trait loci
10 (>5,000,000 expression QTLs and >5,000 chromatin QTLs), brain-active enhancers,
11 differentially expressed genes and transcripts, and novel non-coding RNAs. Leveraging and
12 comparing this resource with other data, we show that the brain has distinct expression and
13 epigenetic profiles as evident from spectral analysis and more non-coding transcription from
14 most other tissues. Also, using single cell data, we deconvolved the tissue-level gene
15 expression of this resource to find the populations of different cell types corresponding to
16 particular phenotypes. Finally, we developed and built an integrative epigenome- and
17 transcriptome-wide association model (eTWAS) to predict the brain phenotypes using high-
18 dimensional functional genomics data with genotype-phenotype associations in this resource to
19 highlight key brain genes and modules and relate the mechanisms on how variants in these
20 affect gene expression. This model allows us to quantitatively impute missing transcriptional and
21 epigenetic information for samples with genotypes only. This model also shows that the
22 integrated data has significantly improved the prediction accuracy over individual genomic data
23 types and relates these predictions to well characterized functions and pathways in the brain.

24 Introduction

25 Disorders of the brain affect nearly 20% of the world's population (ref). Unlike cardiac disease,
26 where lifestyle and pharmacological modification of environmental risk factors has had a
27 profound effect on disease morbidity and mortality (ref), or cancer, which is now understood to
28 be a disorder of the genomic functions (ref), until recently, little progress has been made in our
29 fundamental understanding of the molecular cause of the brain disorders. This recent progress
30 has come in the form of genetic association signals from large GWAS studies of the psychiatric
31 and neurological disorders and currently hundreds of genomic locations that alter the disease
32 risk are known (ref of review, or list disorders in text below, depending on space).
33 Unfortunately, for most of these locations, we have little to no understanding of which base pairs
34 alterations constitute the functional genomic alteration, which transcripts and networks are
35 altered, and what are the molecular mechanisms that cause those alterations. It is presumed
36 that these changes in transcription modify the proteome, which leads to changes in brain

37 structure and function, and these changes interact with environmental factors to change the
38 probability of developing a brain disorder.

39
40 To this end, a variety of genomic elements have been found by many GWAS studies [refs] to
41 associate with psychiatric behaviors such as ones in mental diseases. [[JK: Add in details of
42 other GWASs we have in the paper, once we know which ones they are.]] For example, the
43 Psychiatric Genomics Consortium (PGC) identified a set of genomic variants including SNPs
44 and CNVs associated with psychiatric disorders; e.g., 108 GWAS loci associated with
45 schizophrenia (SCZ) , which explained ~20% liability across major disorders \cite{23933821}. In
46 addition to genotype, a number of genes have been reported to have specific transcriptional
47 activities in mental diseases; e.g., the specific gene expression in mental diseases \cite{xx}. In
48 another context, recent large consortia such as GTEx, ENCODE and Epigenomics Roadmap
49 have generated large-scale RNA-seq and ChIP-seq data for dozens of brain tissues and cell
50 lines (N=xxx) in order to systematically identify brain specific genes, transcripts and regulatory
51 elements [[JK: Maybe more details here, such as samples size]]. However, these studies were
52 limited to healthy brains, so their data is unable to be used to find genomic elements for mental
53 health. For neuropsychiatric-specific analysis, the CommonMind Consortium and others have
54 generated gene expression and genotyping data for both healthy and schizophrenia samples
55 (N=279 vs. 258), identifying ~693 differentially expressed genes in schizophrenia. However,
56 their results still suggested that thousands of samples would be required to achieve statistical
57 power of 0.8 for detecting differential expression of eQTL-associated genes [refs]. Moreover,
58 recent studies show that specific chromatin activity of the regulatory elements such as
59 enhancers has been found to potentially control gene expression in brain [ref], and that single
60 cell techniques can detect gene expression and epigenetic patterns for neuronal and non-
61 neuronal cell types from brain tissues [ref]. Given the complexity of adult brain, we need a
62 variety of additional samples to gain the statistical power necessary for discovering a complete
63 set of genomic elements for neuropsychiatric disorders and other phenotypes. In addition,
64 individual molecules do not independently affect brain, and instead interact with each other in a
65 network. Thus, effort is needed to model and analyze the molecular interactions that drive the
66 phenotypes of adult brain including neuropsychiatric disorders.

67
68 In fact, understanding the molecular mechanisms on how these genomic elements affect
69 various brain functions and phenotypes is still a key challenge in neuroscience. To address it,
70 the PsychENCODE Consortium integrates a group of projects to produce a public resource of
71 multi-dimensional genomic data from thousands of high quality healthy and diseased human
72 post-mortem brains (PEC ref) (6). Particularly, it has generated and assembled a robust large-
73 scale dataset on the adult human brain to address this challenge, including genotyping, RNA-
74 seq, ChIP-seq and single-cell transcriptomic data on the brain tissue samples of 1931 adult
75 individuals with different phenotypes and these data are housed in a central, publically available
76 depository (xxxx). In addition, for these analyses, we have supplemented the PEC data with the
77 primary data at both tissue and single cell levels from other related genomic resources, such as:
78 ENCODE, CommonMind, GTEx, Epigenomics Roadmap, recent neuronal and non-neuronal
79 single cells [refs], and uniformly processed all the data together and performed integrated
80 analyses with up to X,XXX samples. Using single cell data, we also calculate the fractions of

Deleted: ~2000 (or 1945) adult

Deleted: etc,

Deleted: We have also supplemented the PEC

Deleted: -

Deleted: with the primary data from recent publications (refs), reprocessed and analyzed all the data jointly to find gene expression signatures and

88 neuronal and non-neuronal cell types in normal and disease states, [for individual tissue samples](#).
89 We provide all the PEC data and [integrative](#) analyses in an online resource which contains all
90 possible functional genomic elements for adult brain including the brain-active enhancers,
91 transcripts, expression models, imputed regulatory networks, eQTLs and cQTLs for various
92 phenotypes, and an integrated deep-learning model, Deep Structured Phenotype Networks
93 (DSPN) for predicting and imputing brain phenotypes. We then use this resource to discover the
94 properties of brain gene expression, non-coding transcription and enhancers, and to build this
95 [model](#), to describe how interactions between genomic variants, gene expression, enhancers
96 might work together molecularly to alter disease risk.

Deleted: .

Deleted: mega-

Deleted: DSPN

97 Comprehensive resource for adult brain functional genomics

98 The PsychENCODE consortium has generated and assembled a large-scale dataset of
99 genotypes, RNA-seq, ChIP-seq, ATAC-seq, Hi-C and single-cell transcriptomic data from adult
100 brains of [1931](#) individuals, with and without several mental illnesses (Figure 1, Assay summary
101 in Methods). To harmonize and integrate the datasets across multiple consortia, we processed
102 these datasets using standard bioinformatic pipelines in common use (Methods). For instance,
103 we adopted the ENCODE processing pipelines for the bulk and single cell RNA-seq and ChIP-
104 seq data. Likewise, we used the GTEx eQTL pipeline and associated parameters, to allow
105 comparison to previously published eQTL maps. All these uniformly processed datasets are
106 available in our XXXX resource (URL here). Finally, we also compared the resource data
107 against various phenotypes, and identified the brain specific data (derived data type). For
108 example, this resource includes the regulatory variants such as QTLs, brain active enhancers,
109 differentially expressed genes and transcripts, novel transcribed regions and non-coding RNAs,
110 and putative genome-wide regulatory networks. It is also publicly accessible and available on
111 the PyschENCODE website (<http://adult.psychencode.org/pec/>).

Deleted: 1945

Deleted: xxxx), such as using the interactive web app.

112
113 Overall, this resource is structured in a pyramid shape (Figure 1), with the largest scale and raw
114 data at the bottom level and the lightest and most interpretive data at the top level.

116 Next generation sequencing data for brain functional genomics

117 [At the bottom](#), we have the large scale raw data and the phenotype information for [1931](#)
118 individuals, much of which is private and under controlled access. Based on this, we have then
119 uniformly processed raw datasets from PyschENCODE and other consortia (ENCODE,
120 CommonMind, GTEx, Epigenomics Roadmap, etc), including RNA-seq expression
121 quantifications, ChIP-seq signal track qualifications and peak identifications using ENCODE
122 standard pipelines, and private imputed genotypes. The processed functional genomic data is
123 much easier to interpret but still rather large scale. In details, they include the following major
124 [data](#) types:

Deleted: .

Deleted: 1945

125
126 *Phenotypes* - the PsychENCODE data covers a number of phenotypes on mental health. They
127 are normal control (n=[1445](#)), SCZ (n=[270](#)), BP (n=[160](#)), ASD (n=[65](#)), AFF (n=[8](#)), Male (n=[1244](#)),
128 Female (n=[700](#)), Age (distribution), etc. (Supplement).

Deleted: xxx

Deleted: xxx

Deleted: xxx

Deleted: xxx

Deleted: xxx

Deleted: xxx

143 *Epigenomics* - we used the ENCODE standard ChIP-seq pipeline and uniformly processed the
144 ChIP-seq data of available [tissue](#) samples in PsychENCODE and Roadmap Epigenomics, [and](#)
145 [neu+ and neu- cell samples](#) for the signal track qualifications and peak identifications.

146
147 *Transcriptomics* - we also used the ENCODE standard RNA-seq pipeline to uniformly process
148 the RNA-seq data of available samples from a number of PsychENCODE-related studies,
149 ENCODE and GTEx to quantify the expression levels for the protein coding genes, transcripts,
150 noncoding RNA and novel transcribed regions.

151
152 *Chromatin [interactomics](#)* - we generated and processed the Hi-C data for adult brain [including](#)
153 [three reference brains](#), and identified the xxx regions on which the enhancers and promoters
154 interact. [Using this full Hi-C data for adult brain, we identified xxx Topologically Associating](#)
155 [Domains \(TADs\) of adult brain. These TADs provide the regions at which the enhancers interact](#)
156 [with target gene promoters in adult brain, and](#) enable the systems identification of potential [cis-](#)
157 regulatory enhancers [of the genes](#). [more from HJ&DH]

Deleted: structures

Deleted: These

Deleted: (e.g., TADs)

Deleted: if they interact the target genes' promoters

159 System identification of the specific transcriptomic and epigenomic 160 elements in adult brain

161 Given the large-scale transcriptomic and epigenomic data in resource, we further integrated
162 them and identified the genomic elements that have specific activities in adult brain. We used
163 the uniformly processed data and compared against various phenotypes to have even more
164 interpreted functional elements such as sets of differentially expressed genes characterizing
165 various brain regions and phenotypes, sets of aggregated brain enhancers from merging the
166 K27 peaks on the ENCODE regulatory elements. And then above these individual elements, we
167 even identified more interpreted association relationship data such as the QTLs affecting gene
168 expression and enhancers, and imputed the regulatory networks consisting of QTLs,
169 transcriptional factors (TFs), enhancers and genes. This includes:

170
171 *Brain active enhancers* - we identified the brain enhancers from the uniformly processed ChIP-
172 seq data and related them with the regulatory elements in ENCODE and Epigenomics
173 Roadmap, and summarized a list of PsychENCODE brain enhancers which are activated on
174 major brain regions [such as](#) ~88,800 enhancers in pre-frontal cortex [including xxx ones in adult](#)
175 [brain TADs](#) (Supplement).

Deleted: including

Formatted: Font color: R,G,B (34,34,34)

Deleted:

Deleted: .

... 11

176
177 *Differentially expressed genes, transcripts and brain splicing patterns* - we compared expression
178 changes in uniformly processed RNA-seq data from brain samples across PsychENCODE-
179 related studies, ENCODE, and GTEx, and found xxx expressed genes and ~79,000 transcripts
180 in pre-frontal cortex, ~11k eGenes associated with eQTLs (Methods), and xxx non-coding
181 RNAs and novel transcribed regions. We also derived phenotype-specific genes and transcripts.
182 In addition, we calculated the alternative splicing patterns at the transcript level; i.e., the
183 percentage of the transcript abundance over its gene abundance, and found the brain-specific
184 spliced transcripts. Our resource contains differentially expressed and spliced genes and
185 transcripts across a number of biological variables, including neuropsychiatric disorders and
186 developmental stages.

195
196
197
198
199
200
201
202
203
204
205
206
207
208
209

Gene co-expression modules - Also, the brain specific gene expression is likely driven by a group of genes, rather than individual genes, so we constructed the gene co-expression network using all PsychENCODE and GTEx samples, and clustered it into gene co-expression modules using WGCNA [Methods]. The genes clustered in a same module are highly likely co-regulated by similar mechanisms. Our co-expression analysis indeed found several modules whose eigengenes show very different expression levels between brain and non-brain samples (Figure Sxxx, Supplement), which suggests that there exist brain specific regulatory mechanisms drive these brain co-expression modules.

Moved (insertion) [1]

We should emphasize that our comparative analysis is consistent for finding various brain elements including brain enhancers, genes and transcripts. More specifically, we compared them against a same set of brain and non-brain tissues; e.g., the RNA-seq gene expression data from GTEx and the ChIP-seq binding signal data from Epigenomics Roadmap for brain pre-frontal cortex vs. other non-brain tissues including liver, lung, blood, etc.

Deleted: H3K27AC

210 System identification of the QTLs and gene regulatory networks associated 211 with adult brain transcriptomics and epigenomics

212 To understand how the genotype affects the transcriptomic and epigenetic activities in adult
213 brain, we first used the resource data as above to identify more interpreted association
214 relationship data such as the quantitative trait loci (QTLs) affecting gene expression and
215 chromatin activity. In particular, we merged genotype and gene expression and chromatin data
216 of Brain DFC region from a number of studies relating to PyschENCODE. We calculated the
217 association of imputed SNPs with normalized gene expression and chromatin states (Methods)
218 to find the quantitative trait loci associating with gene expression and epigenomic activities in
219 adult brain, including three major categories: expression QTLs (eQTLs), chromatin QTLs
220 (cQTLs), splicing QTLs (sQTLs) and even cell fractions (fQTLs, more details from the single-cell
221 analysis as below). We used the GTEx standard pipeline for discovering eQTLs to find the
222 associations, which is based on an additive linear model from QTLtools. Given the complex
223 relationships between genotype and phenotype, potentially driven by batch effects and biases
224 (e.g., merging different chromatin datasets), this linear model was also adjusted by covariates
225 like PEER factors of gene expression, genotype PCs and disease diagnosis. Among these
226 SNPs, we identified a great number of the regulatory variants significantly associated with brain
227 transcriptional and epigenomic activity: >1 million expression QTLs (eQTLs) with ~11k
228 eGenes, >5 thousand chromatin QTLs (cQTLs) for histone modification signals, and xxx splicing
229 QTLs for alternative splicing patterns. The distributions of detailed QTL annotations on genomic
230 regions are shown in Figure xxx.

231
232 Given a great number of QTLs we identified, we are further interested to see how they relate to
233 the known variants for brain. In particular, we compared them with existing QTLs databases and
234 subdivided our QTLs into different functional categories, mainly including the disease GWAS
235 SNPs, the SNPs breaking the TF binding sites, etc (Table/Figure xxx). Collectively, these QTLs
236 annotate a larger fraction of GWAS SNPs involving the brain (e.g., 6% in schizophrenia, 10% in

238 bipolar) than previously observed, providing leads on which genes are affected in disease. We
239 also evaluated the overlap of eQTLs with cQTLs and found that XX% of cQTLs are overlapped
240 with eQTLs. The SNPs in cis-eQTL list(Cis-eSNPs) were enriched within XXXX, and depleted
241 XXXXXX (Fig. X). We examined the enrichment of most significant eQTLs per gene in
242 Roadmap Epigenomics Consortium and ENCODE enhancers across XX human tissues and cell
243 lines. Cis-eQTL were enriched for enhancer sequences present in brain tissues and the
244 strongest enrichment is observed in DLPFC enhancers. We also calculate the enrichment of
245 cis-QTLs on GWAS SNPs of brain related disorders (schizophrenia, bipolar disorders and
246 parkinson's disease) and non-brain related disorders (CAD, asthma and type 2 diabetes). Cis-
247 QTLs have more significant enrichment for GWAS SNPs of brain related disorders than the
248 ones of non-brain related disorders. In addition, we link the QTLs that overlap the enhancers
249 and promoters in the resource to reveal the potential regulatory activities. We thus classified the
250 QTLs into subgroups in terms of their gene regulatory characteristics including the regulatory
251 QTLs (rQTLs) that break TF binding sites on promoters and/or enhancers, and the modular
252 QTLs (mQTLs) that highly associate with a set of co-expressed genes. Finally, we found that
253 the eQTLs/eGenes number can be predicted from the sample size using a fitted curve (Figure
254 xxx).

255
256 *Gene regulatory networks* - we also integrated and imputed the regulatory relationships in
257 brain such as the enhancers, transcription factors (TFs), miRNAs and target genes [refs] in this
258 resource (Methods). For example, we found the TF binding motifs using ENCODE data and
259 inferred the TF-target gene relationships if TFs have enriched binding motifs on the target
260 gene's regulatory regions such as promoters and enhancers. We also used Hi-C data to filter
261 the enhancers that are not in the TAD regions for given target genes. In total, we included xxx
262 enhancer-gene, xxx TF-gene, and xxx miRNA-gene regulatory linkages, providing a reference
263 wiring network on gene regulation in brain. It should be noted that activations of these regulatory
264 wires are highly attributed to the genotypes of QTLs, leading to various phenotypes. Thus, using
265 these "wiring" regulatory relationships, we inferred the gene regulatory networks that identify the
266 regulatory relationships on how QTLs, enhancers, and transcription factors relate to target gene
267 expression (Methods). In particular, given a target gene, we found its related regulatory
268 elements from the resource including the eQTLs, the enhancers that control its gene expression
269 [JEME] plus their cQTLs, and predicted the transcription factors (TFs) that have enriched
270 binding sites on these enhancers and its promoter. We then used RNA-seq and ChIP-seq data
271 based on the Elastic Net model with regularization that combines the L1 and L2 penalties of the
272 lasso and ridge regressions to predict the regression coefficients of genotypes of various QTLs,
273 the chromatin stages of enhancers, splicing patterns and TFs gene expression to the target
274 gene expression, and identified the highly predictive relationships (i.e., large coefficients). We
275 repeated this for all genes and found how various subgroups of QTLs affect gene expression;
276 e.g., a significantly number of predictive QTLs break the TFBSs on the enhancers or promoters
277 (xx%, Figure xxx). We thus constructed a gene regulatory networks consisting of the QTLs,
278 enhancers, TFs and target genes with high predictive relationships (Methods), revealing the
279 biological mechanisms on how QTLs regulate the target gene expression in the adult brain.

Deleted: coeff. > xxx,

280

282 In summary, the establishment of this comprehensive resource enables the modeling and
283 analysis for the biological processes in adult brain and helps understand the molecular
284 mechanisms between genotypes and phenotypes. Therefore, we later analyzed and modeled
285 the data from this resource to further reveal the brain specific genomic and transcriptomic
286 activities, and the biological mechanisms explaining how the brain specific elements affect the
287 phenotypes and diseases in the adult brain.

288 Comparative analysis reveals the brain related transcriptomic and 289 epigenomic activity

290 We leveraged this resource to compare the human brain with other tissues. To reveal potential
291 brain specific genomic activities, particularly relating to transcriptomic and epigenomic activities,
292 we performed a consistent spectral analysis and compared the similarities of [RNA-seq](#) gene
293 expression and [ChIP-seq](#) binding signals on enhancers and found that the brain has more
294 distinct expression patterns compared to most other tissues, including a greater amount of non-
295 coding transcription. However, the differences in epigenetics are relatively smaller.

296
297 For gene expression, we compared the adult brain samples from our resource with the other
298 tissue samples from GTEx, using uniformly reprocessed RNA-seq data. [We tested three well
299 established dimensionality reduction methods to identify structures of gene expression. Principal
300 Component Analysis \(PCA\) was able to capture some, but not all structure of human tissues.
301 On the other hand, tSNE is too sensitive to batch effects and exposed structures that have not
302 originated from biological differences. We finally tested Reference Component Analysis \(RCA\),
303 that projects the gene expression into a reference panel of tissues and genes and shows
304 highlights intermediate structures in the data. Using the reference component RCA, we show
305 that the brain samples, though from different studies are clustered together in a major cluster,
306 significantly separated from the other major cluster consisting of non-brain samples from their
307 leading reduced dimension \(left vs. right clusters in Figure xxx\). This suggests that the brain has
308 unique and distinctive gene expression programs, which are involved by the brain elements
309 including brain expressed genes, transcripts and non-coding RNAs in our resource. In addition,
310 the samples of PsychENCODE that include psychiatric disorders have larger variations than
311 other tissue clusters \(Figure xxx\). The cluster radiuses were estimated by fitting the two main
312 principal components into a multivariate normal model and finding a 0.95 confidence interval
313 \(Methods\). This suggests that the psychiatric diseases still have larger variations of gene
314 expression, and different gene regulatory programs from the normal, though even more distant
315 from other organs. \[Thus, we then want to check all unified transcriptional activities on the
316 genome scale in brain including potentially novel transcribed non-annotated regions.
317 Specifically,\]\(#\) to understand where the human brain sits in regards of its the transcription diversity
318 compared to other tissues, we estimated the proportion of genome that is transcriptionally active
319 across hundreds of samples. We first found that transcript diversity is mostly saturated at the
320 scale of hundreds of individuals \(Figure xxx\). The saturation is observed for both the \[annotated\]\(#\)
321 and non-\[annotated\]\(#\) portions of the genome. The human brain does not stand as a highly diverse
322 in protein coding regions. For example, the tissues such as the testis is highly diverse \[Ref\];
323 however, we found that the brain has more transcriptional activity at the non-\[annotated\]\(#\) and](#)

Deleted: H3K27AC

Deleted: It shows

Deleted: the reference brain samples and

Deleted:

Deleted: Additionally,

Deleted: coding

Deleted: coding

Deleted: coding

332 novel transcribed regions than most other tissues (Figure xxx). Which implies that the non-
333 coding transcription is highly likely another factor to make the brain tissues unique.

334
335 As shown above, the brain samples have different chromatin and gene expression activities
336 from other organs, implying that the brain also has specific gene regulatory activities. Therefore,
337 we are further interested to compare the enhancers between brain and other tissues to see any
338 brain epigenomic activities. In particular, we integrated the H3K27Ac ChIP-seq signal data of
339 enhancers in the resource and performed dimensional reduction analysis consistent to the for
340 gene expression RCA to compare the similarities of epigenetic profiles of PsychENCODE
341 samples with Epigenomic Roadmap data. It is also interesting to find dissimilar patterns with the
342 gene expression comparison; e.g., while the brain samples separates from other tissues when
343 using genes expression data, the active enhancers are not able to separate brain from other
344 tissues (Figure xxx). This result suggests that the brain has less specific and distinct epigenomic
345 activities, involving the brain active enhancers from our resource. Thus, there may exist more
346 complex regulatory mechanisms among the brain enhancers with low signal variability than
347 other tissues to drive the brain distinct gene expression. One important mechanism is that the
348 brain active enhancers or gene expression patterns are intermediate phenotypes, potentially
349 driven by particular large set of brain regulatory variants such as our QTLs as previously
350 described.

351
352 Our comparative analysis reveals that the brain is different from other organs in gene
353 expression. Thus, we are then interested to identify the functional genomic elements in brain
354 that give rise to the uniqueness of brain. To systematically find the specific expressed functional
355 elements in brain, we identified the differentially expressed genes for phenotypes such as
356 gender (Methods and Figure XX) for the resource. For example, we identified a group of genes
357 that differentially express across different ages (Figure xxx). In particular, the gene involved in
358 early growth response is down-regulated at elder samples whereas the gene with ceruloplasmin
359 is down-regulated around the middle ages. Finally, we report the DEX genes for all phenotypes
360 in our resource along with their enriched functions and pathways in supplement.

361 Single cell analysis and deconvolution explain gene expression 362 changes across adult phenotypes

363 The brain tissues have been found to comprise a variety of cell types including neuronal and
364 non-neuronal cells such as astrocytes [refs]. One issue with the changes of gene expression in
365 our brain tissue samples is whether the changes are driven by gene expression in a particular
366 cell type or different cell-type populations. To address this tissue, we integrated the single cell
367 gene expression data to discover how the gene expression from various cell types including
368 both neuronal and non-neuronal contribute to the gene expression at the tissue level. In
369 particular, we used the biomarker genes with strong expression signals in single cell to
370 deconvolve the gene expression data of individual tissues over both novel and known cell types
371 to find the cell fractions for individuals, and relate to the individual phenotypes. We found that
372 the gene expression changes across individual tissue samples can be largely explained by the

Deleted: the
Deleted: spectral analysis
Deleted: as above
Deleted: Epigenomics

Deleted: and non-coding RNAs
Deleted: various
Deleted: including mental disease,
Deleted: , regions
Deleted: xxx genes have been found to differentially express between SCZ and normal samples; i.e., SCZ DEX genes, and they are also enriched with the pathways and functions relating to SCZ (Figure Sxxx). Moreover,
Deleted: For example

Moved up [1]: Also, the brain specific gene expression is likely driven by a group of genes, rather than individual genes, so we constructed the gene co-expression network using all PsychENCODE and GTEx samples, and clustered it into gene co-expression modules using WGCNA [Methods]. The genes clustered in a same module are highly likely co-regulated by similar mechanisms. Our co-expression analysis indeed found several modules whose eigengenes show very different expression levels between brain and non-brain samples (Figure Sxxx, Supplement), which suggests that there exist brain specific regulatory mechanisms drive these brain co-expression modules.

401 single cell gene expression, and the changes of single cell fractions are also associated with the
402 individual phenotypes.

403
404 Specifically, we integrated and used the same pipeline to uniformly process the single cell RNA-
405 seq data for the neuronal and non-neuronal cell types from PsychENCODE and recent
406 publications [Lake&quaker]. In total, we included 23 single cell types (Supplement) and found
407 that the same-type cells generally can be clustered together (Figure Sxxx) using our uniformly
408 processed data. We also include these single cell data as well as their cell-type biomarker
409 genes in the resource. Moreover, we found that a group of psychiatric disorder related genes
410 indeed show the expression dynamic changes among cells. For example, the dopamine
411 receptor genes (DRD) that associate with SCZ, are significantly more highly expressed in
412 neuronal cells than others (Figure Sxxx), and their expression levels across cells vary
413 significantly larger than tissue samples, suggesting that the cell fraction changes potentially
414 equalize the tissue expression variability. Therefore, we are further interested to see if the brain
415 gene expression at the tissue level in our resource is contributed by the above cell types and
416 affected by the cell fractions.

417
418 To this end, we decomposed the gene expression data across individuals at the tissue level
419 from our resource using non-negative matrix factorization (NMF, see Methods). Indeed, we
420 found that three groups of top principal components of NMF (NMF-PCs) capturing the most
421 covariance of brain gene expression across individual tissues, highly correlate with the
422 biomarker gene expression signatures of neuronal, non-neuronal and fetal cell types as above,
423 respectively. For example, the NMF-PCs shown in Figure xxx. This suggests that the large
424 portion of tissue's gene expression changes is a linear combination of these cell types' gene
425 expression. Thus, we want to further identify the cell fractions showing how individual single
426 cells contribute the tissue's gene expression, using the deconvolution.

427
428 Therefore, we deconvolved the tissue-level gene expression data of all 1931 individuals' tissue
429 samples using single-cell gene expression data of 450 biomarker genes to find the fraction of
430 different cell types corresponding, and compare cell fractions across different phenotypes
431 (Supplement). The single cells used in deconvolution cover all 16 neuronal types, five non-
432 neuronal types and xxx additional fetal types from PsychENCODE single cell data [ref:
433 brainspan]. It is very interesting that the linear combinations of single cell expression of 23 cell
434 types, where combinational coefficients, can explain >80% of the gene expression variations
435 across 1931 individual tissues (Figure xx). The coefficients of cell types for linear combination
436 are estimated from our deconvolution analysis (Methods in supplement), and proportional to the
437 cell fractions of individuals. In addition, we found that the cell fractions of individuals (i.e.,
438 deconvolution coefficients) vary, and a number of cell population changes highly associate with
439 different phenotypes and disorders (Figure xxx). For example, the fraction(s) of neuronal type(s)
440 (Inhibitory X) is significantly anti-correlated with Age ($r = xxx$), and Inhibitory X cells have
441 functions of XXX involving the differentially expressed genes in Age from our resource (Figure
442 xxx). The excitatory neuronal cell populations (e.g., EX1) increase significantly in ASD samples
443 ($p < xxx$) while the non-neuronal cells decreasing (e.g., oligodendrocytes). Finally, we report the

Deleted: ~3000

Deleted: cells with 8 excitatory

Deleted: 8 inhibitory types [Lake's 2016 paper], and ~400 cells including 5

Deleted: types, astrocytes, endothelial, microglia, oligodendrocytes and Oligodendrocyte progenitor

Deleted: (OPC), and ~800 cells

Deleted: PsychENCODE for potentially additional cell types in embryonic and fetal brain tissues [ref brainspan].

Deleted:). We first compared these single cells based on the (biomarker) gene expression similarity using tSNE,

Deleted:). This suggests that

Deleted: integration has removed the batch effects of single cell

Deleted: from different studies. In particular, xx% PsychENCODE cells have been found to cluster together with known cell types (xx% neuronal, xx% non-neuronal, details in supplement). In addition, xx% PsychENCODE cells form their own clusters, away from known cell types, suggesting that the potential novel cell types found by PsychENCODE for brain tissues.

Deleted: and

Deleted: for those differentially expressed genes at the tissue level from our resource,

Deleted: further checked their expression changes across various single cells, and

Deleted: then

Deleted: (three blocks in Figure xxx). For example, No. 22 and 23 NMF-PCs of the non-neuronal group highly correlate with astrocytes, No. 2 NMF-PC correlate with fetal cells, and No. 1, 5, 10, 24 and 25 NMF-PCs of the neuronal group correlate with excitatory neuronal cell types.

Deleted: 1945

Deleted: 1945

Deleted:).

483 individual cell populations along with significantly associated relationships between particular
484 cell type fractions and phenotypes (Supplement).

485

486 Furthermore, we are interested to see if any genotype is also associated with two single cell
487 features: (1) the cell fractions and (2) the gene expression changes that can't be explained by
488 the cell fractions. In particular, we used our QTL pipeline and identified xxx SNPs whose
489 genotypes are significantly associated with yyy neuronal cell fractions across individuals, (or zzz
490 non-neuronal cell types); i.e., cell fraction QTLs (fQTLs). This suggests that these fQTLs
491 potentially can be used to predict the yyy cell fractions in adult brain. Moreover, we identified
492 xxx SNPs significantly associated with the gene expression changes across individual tissues
493 unexplained by our single cell deconvolution; i.e., Y-WX (Methods). These SNPs are likely
494 causing certain gene expression changes driven by unknown cell types in adult brain.

495 Integrative modeling to explain the molecular mechanisms for 496 genotype-phenotype relationships in adult brain

497 The interaction between genotype and phenotype is a very complex process, involving multiple
498 intermediate stages including gene expression, signaling, modulation and so on. Thus, to
499 understand [this](#) and [merge all these stages in one model](#), we introduce an interpretable deep-
500 learning framework, Deep Structured Phenotype Networks (DSPN), which provides insight into
501 how the brain genomic variants affect gene expression and regulation, and eventually predict
502 phenotypes; i.e., [the DSPN pathways from genotype to phenotype](#) (Figure xxx). This model
503 combines a Deep Boltzmann Machine architecture with conditional and lateral connections
504 derived from the QTLs and regulatory networks estimated in our resource. [On the resource](#)
505 [website, we provide a list of DSPN pathways for each phenotype and disease. We also make](#)
506 [the model downloadable as a set of simplified files summarizing represented genotype-](#)
507 [phenotype pathways. In particular, this model](#) integrates all high dimensional functional data
508 types in this resource including genomics, transcriptomics, epigenetics and regulatomics, and
509 genotype-phenotype relationships, and also allows us to quantitatively impute missing
510 transcriptional and epigenetic information for samples with genotypes only. The model is trained
511 as a deep generative model to represent the conditional distribution of all variables given the
512 genotype. Unlike a feed-forward network architecture, the undirected form of the Boltzmann
513 machine allows information to flow in top-down, bottom-up and lateral directions during
514 inference, so that intermediate and high-level phenotypes may be jointly inferred while
515 respecting their mutual dependencies. This allows us for instance to impute transcriptome and
516 epigenome data when it is missing. [In particular, our inference](#) is performed using a mean-field
517 approximation, and training is performed using a Persistent Markov Chain Monte Carlo
518 algorithm [which is able to ensemble multi-dimensional datasets \(Supplement\)](#).

519

520 As shown in Figure xxx, the DSPN consists of four layers: 1) genotypes such as QTLs; 2)
521 molecules and genomic elements, including genes and enhancers; 3) functional modules and
522 other mid-level phenotypes at a series of intermediate layers; i.e., the hidden nodes of deep
523 learning modeling; 4) high-level phenotypes such as brain traits. In addition, we enforce the
524 DSPN to have sparse connectivity (Supplement). Specifically, we built each layer of our model

Deleted: the entire process of how genotype

Deleted: phenotype relate to each other

Deleted: It

Deleted:

Deleted: Inference

Deleted: (see supplement)

531 as follows. We first used the imputed gene regulatory networks that identify the regulatory
532 connectivities on how QTLs, enhancers, and transcription factors relate to target gene
533 expression (Supplement). We then connected the nodes on the molecular layer of our model to
534 follow the inferred gene regulatory network structures; i.e., embedding the gene regulatory
535 network. In particular, many intermediate-layer modules (i.e., strongly predictive features on
536 Layer 3) that correspond to known gene sets associated with well-characterized pathways and
537 functions in the brain; e.g., the module xxx is connected to genes enriched in the dopaminergic
538 and glutamatergic synapse (GSEA enrichment score > xxx, Figure xx). Also, some modules are
539 used to capture the information on single cell populations; e.g., the module yyy is connecting to
540 Age, and represents the neuronal cell fractions (Figure xxx). Furthermore, we used this model to
541 recapitulate the pathways comprising the cross-layer nodes and predictive edges for particular
542 phenotypes. For example, as highlighted in Figure xxx, the schizophrenia (SCZ) trait is activated
543 by two modules on the layer of hidden nodes corresponding to glutamatergic signaling and
544 excitatory synapse, respectively. The modules are connected by a set of genes including
545 GRIN1, which are regulated by corresponding QTLs (e.g., rs1146020) and enhancers (e.g.,
546 GH09H137166) as shown in the blowup gene regulatory mechanism. In addition, we discovered
547 additional molecular mechanisms for SCZ such as module(s) corresponding to dopamine-
548 related pathways and complement pathways (Figure xxx). These modules are connected to the
549 C4 family genes, regulated by eQTLs and enhancers ($p < 1e-4$).

Deleted: Layer 2

550
551 Moreover, the model also enables practical imputation of a subset of the transcriptome and
552 epigenome, with an accuracy of ~70% (Figure xxx). We use the model to improve prediction of
553 biological variables and psychiatric diseases by the addition of transcriptomic data to genotype,
554 as compared to genotype alone. In particular, we can predict bipolar disease and schizophrenia
555 with much higher accuracy from the transcriptome than from genotype alone; i.e., three times
556 improvements (+18% vs. +6%) from the random prediction 50% for schizophrenia, Figure XXX).
557 The imputed transcriptome also clearly adds predictive value, as we can predict schizophrenia
558 with an accuracy of 61% using our model and an imputed transcriptome compared to 56% with
559 genotype alone. This result demonstrates the usefulness of even a limited amount of functional
560 genomics information for unraveling gene-disease relationships.

Deleted: novel

Deleted: On the resource website, we provide a list of DSPN pathways for each endophenotype and disease. We also make the model available as distributive software and as a set of simplified files summarizing represented genotype-phenotype pathways.

561 Discussion

562 We integrated the genomic, transcriptomic and regulatomic PsychENCODE datasets from
563 ~2000 samples and developed this comprehensive resource consisting of various functional
564 genomic elements for the adult brain. Developing this resource and integrated model to a
565 population-level scale serves as an important step in gaining meaningful biological insights from
566 functional genomics studies in neuroscience. In particular, we compared it with other tissues
567 such as GTEx data and identified the genotypes and QTLs, the specific expressed genes,
568 transcripts and noncoding RNAs, active chromatin regions, the regulatory networks that
569 significantly relate with different brain phenotypes at both cellular and tissue levels. For
570 example, the QTLs allow one to potentially interpret most of the known brain-associated GWAS
571 SNPs in terms of perturbations to specific genes. Thus, the neuroscientist can use this resource
572 as a reference to compare with their data, generate hypotheses and help design experimental

580 validations. In addition, this resource is publicly available online and can be extendable and
581 scalable to integrate additional data types and phenotypes. For example, it can add the
582 individual's fMRI image features measuring functional neuro-connectivity, and use our model to
583 identify the genotypes that associated with image features such as image-QTLs (iQTLs) [xx].
584 Also, our resource can incorporate with the neurodegenerative diseases like Alzheimer or
585 developmental stages.
586
587 Moreover, we built an integrative epigenome- and transcriptome-wide association model
588 (eTWAS), built on the Deep Boltzmann Machine (RBM) and integrates the high dimensional
589 functional genomic and phenotypic data at multiple layers, using the hierarchical structures in
590 deep learning. The model reveals the relationships among various data types from a number of
591 directions for genotype to phenotype. In particular, this model also incorporates the derived data
592 types into its hierarchical structure such as imputed gene regulatory networks and QTLs, and
593 provides the additional statistical powers to better predict the genotype to phenotype. This
594 model allows us to quantitatively impute missing transcriptional and epigenetic information for
595 samples with genotypes only. More importantly, it integrates high-dimensional functional
596 genomics data with genotype-phenotype associations to highlight key brain genes and modules
597 and relate how variants in these regulate gene expression. This integrative model is also
598 available online as a general purpose platform. The users can apply it to impute missing data ,
599 predict the genotype-phenotype relationships, and reveal potentially novel gene regulatory
600 mechanisms and modules for additional phenotypes. Also, the model can be used to make in-
601 silico predictions for the perturbation outcomes. For example, we can identify the module X that
602 have the extremely highest connection weights to Autism, and thus knocking down the genes
603 connecting to the module highly likely will deactivate Autism. Furthermore, while the model does
604 provide better predictive performance, some of these correlations are deliberately set to be
605 interpreted simplifications, such as the known enhancers, or gene regulatory network structure,
606 to make the model more interpretable and easier to use. Thus, another major goal of the model
607 is to provide a compression of larger amount of functional genomic datasets for brain; e.g., XXX
608 KB of model files vs. XXX TB of total resource data, beyond a purely predictive network from
609 genotype to phenotype.
610
611 Though single cell remains challenging to reliably quantify the low-abundant transcripts/genes
612 and interrogate the biological variations using single-cell sequencing technology, it is still
613 worthwhile using the biomarker genes with strong expression signals in single cell to
614 deconvolve the gene expression data of individual tissues over both novel and known cell types
615 to find the cell populations for individuals, and relate to the individual phenotypes. With
616 increasing amount of single cell data in near future, we could deconvolve the resource data at
617 tissue level to find potential new cell types and obtain more complete cell populations. The
618 current single-cell sequencing technology suffers from the low capture efficiency [PMCID:
619 PMC4758375, PMCID: PMC4132710]. Due to this reason, the single-cell sequencing will only
620 measure a small fraction of cellular transcriptome as the final sequencing library only contains a
621 subset of input materials. Furthermore, the limited amount of RNA molecules in single cell
622 makes it even harder to capture the weak signals, which makes the data sensitive to technical
623 noise. Thus, given that the RNA decaying issues in single cell RNA-seq, we could also relate

624 this resource to the in situ transcriptomic data such as optogenetic techniques measuring the
625 spatial gene expression, and find the consistent expressed gene for the brain phenotypes at the
626 tissue level.

Formatted Table

References

1. Fromer M, Roussos P, Sieberts SK, Johnson JS, Kavanagh DH, Perumal TM, Ruderfer DM, Oh EC, Topol A, Shah HR *et al*: **Gene expression elucidates functional impact of polygenic risk for schizophrenia**. *Nat Neurosci* 2016, **19**(11):1442-1453.
2. Consortium GT: **Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans**. *Science* 2015, **348**(6235):648-660.
3. Psych EC, Akbarian S, Liu C, Knowles JA, Vaccarino FM, Farnham PJ, Crawford GE, Jaffe AE, Pinto D, Dracheva S *et al*: **The PsychENCODE project**. *Nat Neurosci* 2015, **18**(12):1707-1712.
4. Neale BM, Sklar P: **Genetic analysis of schizophrenia and bipolar disorder reveals polygenicity but also suggests new directions for molecular interrogation**. *Curr Opin Neurobiol* 2015, **30**:131-138.
5. Schizophrenia Working Group of the Psychiatric Genomics C: **Biological insights from 108 schizophrenia-associated genetic loci**. *Nature* 2014, **511**(7510):421-427.
6. Consortium GT, Laboratory DA, Coordinating Center -Analysis Working G, Statistical Methods groups-Analysis Working G, Enhancing Gg, Fund NIHC, Nih/Nci, Nih/Nhgri, Nih/Nimh, Nih/Nida *et al*: **Genetic effects on gene expression across human tissues**. *Nature* 2017, **550**(7675):204-213.
7. Waszak SM, Delaneau O, Gschwind AR, Kilpinen H, Raghav SK, Witwicki RM, Orioli A, Wiederkehr M, Panousis NI, Yurovsky A *et al*: **Population Variation and Genetic Control of Modular Chromatin Architecture in Humans**. *Cell* 2015, **162**(5):1039-1050.
8. Roshara NR, Horn K, Kirsten H, Ahnert P, Scholz M: **Comparing performance of modern genotype imputation methods in different ethnicities**. *Sci Rep* 2016, **6**:34386.
9. McCarthy S, Das S, Kretschmar W, Delaneau O, Wood AR, Teumer A, Kang HM, Fuchsberger C, Danecek P, Sharp K *et al*: **A reference panel of 64,976 haplotypes for genotype imputation**. *Nat Genet* 2016, **48**(10):1279-1283.
10. Won H, de la Torre-Ubieta L, Stein JL, Parikhshak NN, Huang J, Opland CK, Gandal MJ, Sutton GJ, Hormozdiari F, Lu D *et al*: **Chromosome conformation elucidates regulatory relationships in developing human brain**. *Nature* 2016, **538**(7626):523-527.
11. Geschwind DH, Flint J: **Genetics and genomics of psychiatric disease**. *Science* 2015, **349**(6255):1489-1494.
12. Ongen H, Buil A, Brown AA, Dermizakis ET, Delaneau O: **Fast and efficient QTL mapper for thousands of molecular phenotypes**. *Bioinformatics* 2016, **32**(10):1479-1485.
13. What constitutes the prefrontal cortex? *Science* 2017, DOI: 10.1126/science.aan8868

Supplement

Please edit

<https://docs.google.com/document/d/1vJ1PIW1AVwkMSpR036AOJbrNCnIFrd5RImXSX3rRNTk/edit?usp=sharing>

627

Topologically associating domains – we used a full Hi-C data for adult brain and identified xxx in xxx Topologically Associating Domains (TADs) of adult brain. These TADs provide the regions at which the enhancers interact with target gene promoters in adult brain. [more from HJ&DH]