

Journal of Visualized Experiments

Mapping Alzheimer's Disease Variants to Their Target Genes Using Computational Analysis of Chromatin Configuration --Manuscript Draft--

Article Type:	Invited Methods Article - JoVE Produced Video
Manuscript Number:	JoVE60428R1
Full Title:	Mapping Alzheimer's Disease Variants to Their Target Genes Using Computational Analysis of Chromatin Configuration
Section/Category:	JoVE Genetics
Keywords:	Hi-C, GWAS, non-coding variants, gene mapping, functional genomics, Alzheimer's disease
Corresponding Author:	Hyejung Won University of North Carolina at Chapel Hill Chapel Hill , NC UNITED STATES
Corresponding Author's Institution:	University of North Carolina at Chapel Hill
Corresponding Author E-Mail:	hyejung_won@med.unc.edu
Order of Authors:	Hyejung Won Nana Matoba Doug H Phanstiel Yoseli Quiroga
Additional Information:	
Question	Response
Please indicate whether this article will be Standard Access or Open Access.	Standard Access (US\$2,400)
Please indicate the city, state/province, and country where this article will be filmed . Please do not use abbreviations.	Chapel Hill, NC, USA

TITLE:**Mapping Alzheimer's Disease Variants to Their Target Genes Using Computational Analysis of Chromatin Configuration****AUTHORS AND AFFILIATIONS:**

Nana Matoba^{1,2}, Ivana Y. Quiroga³, Douglas H. Phanstiel^{3,4,*}, Hyejung Won^{1,2,*}

¹Department of Genetics, University of North Carolina, Chapel Hill, NC, USA

²Neuroscience Center, University of North Carolina, Chapel Hill, NC, USA

³Thurston Arthritis Research Center, University of North Carolina, Chapel Hill, NC, USA

⁴Department of Cell Biology and Physiology, University of North Carolina, Chapel Hill, NC, USA

*These authors contributed equally.

Email addresses of co-authors:

Nana Matoba (nana_matoba@med.unc.edu)

Ivana Y. Quiroga (yoseli@live.unc.edu)

Corresponding author:

Hyejung Won (hyejung_won@med.unc.edu)

Douglas H. Phanstiel (douglas_phanstiel@med.unc.edu)

KEYWORDS:

Hi-C, GWAS, non-coding variants, gene mapping, functional genomics, Alzheimer's disease

SUMMARY:

We present a protocol to identify functional implications of non-coding variants identified by genome-wide association studies (GWAS) using three-dimensional chromatin interactions.

ABSTRACT:

Genome-wide association studies (GWAS) have successfully identified hundreds of genomic loci that are associated with human traits and disease. However, because the majority of the genome-wide significant (GWS) loci fall onto the non-coding genome, the functional impact of many remain unknown. Three-dimensional chromatin interactions identified by Hi-C or its derivatives can provide useful tools to annotate these loci by linking non-coding variants to their actionable genes. Here, we outline a protocol to map GWAS non-coding variants to their putative genes using Alzheimer's disease (AD) GWAS and Hi-C datasets from human adult brain tissue. Putative causal single-nucleotide polymorphisms (SNPs) are identified by application of fine-mapping algorithms. SNPs are then mapped to their putative target genes using enhancer-promoter interactions based on Hi-C. The resulting gene set represents AD risk genes, as they are potentially regulated by AD risk variants. To garner further biological insights into molecular mechanisms underlying AD, we characterize AD risk genes using developmental brain expression data and brain single-cell expression profiles. This protocol can be expanded to any GWAS and

Hi-C datasets to identify putative target genes and molecular mechanisms underlying various human traits and diseases.

INTRODUCTION:

Genome-wide association studies (GWAS) have played a pivotal role in revealing the genetic basis of a range of human traits and diseases. This large-scale genotyping has uncovered thousands of genomic variants associated with phenotypes ranging from height to schizophrenia risk. However, despite the enormous success of GWAS in identifying disease and trait associated loci, a mechanistic understanding of how these variants contribute to phenotype has been challenging because most phenotype associated variants reside in the non-coding fraction of the human genome. Since these variants often overlap with predicted regulatory elements, they are likely to alter transcriptional control of a nearby gene. However, non-coding loci can influence transcription of genes at linear distances exceeding one megabase, making the genes affected by each variant difficult to identify. Three-dimensional (3D) chromatin structure plays an important role in mediating connections between distant regulatory loci and gene promoters and can be used to identify genes affected by phenotype associated single-nucleotide polymorphisms (SNPs).

Gene regulation is mediated by a complex process, which involves enhancer activation and chromatin loop formation that physically connect enhancers to gene promoters to which the transcriptional machinery can be directed¹⁻³. Because chromatin loops often span several hundred kilobases (kb), detailed maps of 3D chromatin architecture are required to decipher gene regulatory mechanisms. Multiple chromatin conformation capture technologies have been invented to identify the 3D chromatin architecture⁴. Among these technologies, Hi-C provides the most comprehensive architecture, as it captures genome-wide 3D chromatin interaction profiles. Hi-C datasets have been quickly adapted to interpret non-coding genome-wide significant (GWS) loci⁵⁻¹³, as it can link non-coding variants to their putative target genes based on chromatin interaction profiles.

In this article, we outline a protocol to computationally predict putative target genes of GWAS risk variants using chromatin interaction profiles. We apply this protocol to map AD GWS loci¹⁴ to their target genes using Hi-C datasets in the adult human brain⁹. The resulting AD risk genes are characterized by other functional genomic datasets that include single cell transcriptomic and developmental expression profiles.

PROTOCOL:

1. Workstation setup

1.1. Install R (version 3.5.0) and RStudio Desktop. Open RStudio.

1.2. Install the following libraries in R by typing the following code into the console window in RStudio.

```
if (!"BiocManager" %in% rownames(installed.packages()))
```

```
88 install.packages("BiocManager", repos="https://cran.r-project.org")
89 BiocManager::install("GenomicRanges")
90 BiocManager::install("biomaRt")
91 BiocManager::install("WGCNA")
92 install.packages("reshape")
93 install.packages("ggplot2")
94 install.packages("corrplot")
95 install.packages("gProfileR")
96 install.packages("tidyverse")
97 install.packages("ggpubr")
98
```

99 1.3. Download files.

100
101 NOTE: In this protocol, all files are required to be downloaded to ~/work directory.

102 1.3.1. Download the following files by clicking the links provided in **Table of Materials**.

103 1.3.1.1. Download fine-mapped credible SNPs for AD (Supplementary Table 8 from Jansen et al.¹⁴).

104
105
106
107 NOTE: Before analysis, open sheet eight in 41588_2018_311_MOESM3_ESM.xlsx, remove the
108 first three rows and save the sheet as Supplementary_Table_8_Jansen.txt with tab separated
109 format.

110 1.3.1.2. Download 10 kb resolution Hi-C interaction profiles in the adult brain from psychencode 111 (described as *Promoter-anchored_chromatin_loops.bed* below).

112
113
114 NOTE: This file has the following format: chromosome, TSS_start, TSS_end, Enhancer_start, and
115 Enhancer_end. In case other Hi-C datasets are used, this protocol requires Hi-C datasets
116 processed at high resolution (5–20 kb).

117 1.3.1.3. Download single cell expression datasets from the PsychENCODE (described as 118 *singlecell.rda* below).

119
120
121 NOTE: These are from neurotypical control samples.

122 1.3.1.4. Download developmental expression datasets from the BrainSpan (described as 123 *devExpr.rda* below).

124
125
126 NOTE: 267666527 is a zipped file, so unzip the 267666527 to extract "columns_metadata.csv",
127 "expression_matrix.csv", and "rows_metadata.csv" to generate devExpr.rda (see section 3).

128 1.3.2. Download exonic coordinates (see **Supplementary Files**, described as 129 *Gencode19_exon.bed* and *Gencode19_promoter.bed* below) from Gencode version 19.

NOTE: Promoters are defined as 2 kb upstream of transcription start site (TSS). These files have the following format: chromosome, start, end, and gene.

1.3.3. Download gene annotation file (see **Supplementary Files**, described as *geneAnno.rda* below) from biomaRt.

NOTE: This file can be used to match genes based on Ensembl gene IDs and HUGO Gene Nomenclature Committee (HGNC) symbol.

2. Generation of a GRanges object for credible SNPs

2.1. Set up in R by typing the following code into the console window in RStudio.

```
library(GenomicRanges)
options(stringsAsFactors = F)
setwd("~/work") # This is the path to the working directory.
credSNP = read.delim("Supplementary_Table_8_Jansen.txt", header=T)
credSNP = credSNP[credSNP$Credible.Causal=="Yes", ]
```

2.2. Make a GRanges object by typing the following code into the console window in RStudio.

```
credranges = GRanges(credSNP$Chr, IRanges(credSNP$bp, credSNP$bp), rsid=credSNP$SNP,
P=credSNP$P)
save(credranges, file="AD_credibleSNP.rda")
```

3. Positional mapping

NOTE: For each step, type the corresponding code into the console window in RStudio.

3.1. Set up in R.

```
options(stringsAsFactors=F)
library(GenomicRanges)
load("AD_credibleSNP.rda") # (see 2)
```

3.2. Positional mapping of promoter/exonic SNPs to genes

3.2.1. Load promoter and exonic region and generate a GRange object.

```
exon = read.table("Gencode19_exon.bed")
exonranges = GRanges(exon[,1],IRanges(exon[,2],exon[,3]),gene=exon[,4])
promoter = read.table("Gencode19_promoter.bed")
promoterranges = GRanges(promoter[,1], IRanges(promoter[,2], promoter[,3]),
gene=promoter[,4])
```

```

176
177 3.2.2. Overlap credible SNPs with exonic regions.
178
179 olap = findOverlaps(credranges, exonranges)
180 credexon = credranges[queryHits(olap)]
181 mcols(credexon) = cbind(mcols(credexon), mcols(exonranges[subjectHits(olap)]))
182
183 3.2.3. Overlap credible SNPs with promoter regions.
184
185 olap = findOverlaps(credranges, promoterranges)
186 credpromoter = credranges[queryHits(olap)]
187 mcols(credpromoter) = cbind(mcols(credpromoter), mcols(promoterranges[subjectHits(olap)]))
188
189 3.3. Link SNPs to their putative target genes using chromatin interactions.
190
191 3.3.1. Load Hi-C dataset and generate a GRange object.
192
193 hic = read.table("Promoter-anchored_chromatin_loops.bed ", skip=1)
194 colnames(hic) = c("chr", "TSS_start", "TSS_end", "Enhancer_start", "Enhancer_end")
195 hicranges = GRanges(hic$chr, IRanges(hic$TSS_start, hic$TSS_end),
196 enhancer=hic$Enhancer_start)
197 olap = findOverlaps(hicranges, promoterranges)
198 hicpromoter = hicranges[queryHits(olap)]
199 mcols(hicpromoter) = cbind(mcols(hicpromoter), mcols(promoterranges[subjectHits(olap)]))
200 hicenhancer = GRanges(seqnames(hicpromoter), IRanges(hicpromoter$enhancer,
201 hicpromoter$enhancer+10000), gene=hicpromoter$gene)
202
203 3.3.2. Overlap credible SNPs with Hi-C GRange object.
204
205 olap = findOverlaps(credranges, hicenhancer)
206 credhic = credranges[queryHits(olap)]
207 mcols(credhic) = cbind(mcols(credhic), mcols(hicenhancer[subjectHits(olap)]))
208
209 3.4. Compile AD candidate genes defined by positional mapping and chromatin interaction
210 profiles.
211
212 ### The resulting candidate genes for AD:
213 ADgenes = Reduce(union, list(credhic$gene, credexon$gene, credpromoter$gene))
214 ### to convert Ensembl Gene ID to HGNC symbol
215 load("geneAnno.rda")
216 ADhgnc = geneAnno1[match(ADgenes, geneAnno1$ensembl_gene_id), "hgnc_symbol"]
217 ADhgnc = ADhgnc[ADhgnc!=""]
218 save(ADgenes, ADhgnc, file="ADgenes.rda")
219 write.table(ADhgnc, file="ADgenes.txt", row.names=F, col.names=F, quote=F, sep="\t")

```

220

221 4. Developmental expression trajectories

222

223 NOTE: For each step, type the corresponding code into the console window in RStudio.

224

225 4.1. Set up in R.

226

```
227 library(reshape); library(ggplot2); library(GenomicRanges); library(biomaRt)
```

```
228 library("WGCNA")
```

```
229 options(stringsAsFactors=F)
```

230

231 4.2. Process expression and meta data.

232

```
233 datExpr = read.csv("expression_matrix.csv", head = FALSE)
```

```
234 datExpr = datExpr[,-1]
```

```
235 datMeta = read.csv("columns_metadata.csv")
```

```
236 datProbes = read.csv("rows_metadata.csv")
```

```
237 datExpr = datExpr[datProbes$ensembl_gene_id!="",]
```

```
238 datProbes = datProbes[datProbes$ensembl_gene_id!="",]
```

```
239 datExpr.cr= collapseRows(datExpr, rowGroup = datProbes$ensembl_gene_id, rowID=
240 rownames(datExpr))
```

```
241 datExpr = datExpr.cr$datETcollapsed
```

```
242 genome = data.frame(datExpr.cr$group2row)
```

```
243 rownames(datExpr) = genome$group
```

244

245 4.2.1. Specify developmental stages.

246

```
247 datMeta$Unit = "Postnatal"
```

```
248 idx = grep("pcw", datMeta$age)
```

```
249 datMeta$Unit[idx] = "Prenatal"
```

```
250 idx = grep("yrs", datMeta$age)
```

```
251 datMeta$Unit[idx] = "Postnatal"
```

```
252 datMeta$Unit = factor(datMeta$Unit, levels=c("Prenatal", "Postnatal"))
```

253

254 4.2.2. Select cortical regions.

255

```
256 datMeta$Unit = "Postnatal"
```

```
257 datMeta$Region = "SubCTX"
```

```
258 r = c("A1C", "STC", "ITC", "TCx", "OFC", "DFC", "VFC", "MFC", "M1C", "S1C", "IPC", "M1C-S1C",
259 "PCx", "V1C", "Ocx")
```

```
260 datMeta$Region[datMeta$structure_acronym %in% r] = "CTX"
```

```
261 datExpr = datExpr[,which(datMeta$Region=="CTX")]
```

```
262 datMeta = datMeta[which(datMeta$Region=="CTX"),]
```

```
263 save(datExpr, datMeta, file="devExpr.rda")
```

264

265 4.3. Extract developmental expression profiles of AD risk genes.

266

267 load("ADgenes.rda")

268 exprdat = apply(datExpr[match(ADgenes, rownames(datExpr)),],2,mean,na.rm=T)

269 dat = data.frame(Region=datMeta\$Region, Unit=datMeta\$Unit, Expr=exprdat)

270

271 4.4. Compare prenatal versus postnatal expression levels of AD risk genes.

272

273 pdf(file="developmental_expression.pdf")

274 ggplot(dat,aes(x=Unit, y=Expr, fill=Unit, alpha=Unit)) + ylab("Normalized expression") +

275 geom_boxplot(outlier.size = NA) + ggtitle("Brain Expression") + xlab("") +

276 scale_alpha_manual(values=c(0.2, 1)) + theme_classic() + theme(legend.position="na")

277 dev.off()

278

279 5. Cell-type expression profiles

280

281 NOTE: For each step, type the corresponding code into the console window in RStudio.

282

283 5.1. Set up in R.

284

285 load("ADgenes.rda")

286 load("singlecell.rda")

287 load("geneAnno.rda")

288 targetname = "AD"

289 targetgene = ADhgnc

290 datExpr = scale(cellexp,center=T, scale=F)

291

292 5.2. Extract cellular expression profiles of AD risk genes.

293

294 exprdat = apply(datExpr[match(targetgene, rownames(datExpr)),],2,mean,na.rm=T)

295 dat = data.frame(Group=targetname, cell=names(exprdat), Expr=exprdat)

296

297 dat\$celltype = unlist(lapply(strsplit(dat\$cell, split="["),'["),1))

298 dat = dat[-grep("Ex|In",dat\$celltype),]

299 dat\$celltype = factor(dat\$celltype, levels=c("Neurons","Astrocytes","Microglia","Endothelial",
300 "Oligodendrocytes","OPC","Fetal"))

301

302 pdf(file="singlecell_expression_ADgenes.pdf")

303 ggplot(dat,aes(x=celltype, y=Expr, fill=celltype)) +

304 ylab("Normalized expression") + xlab("") + geom_violin() +

305 theme(axis.text.x=element_text(angle = 90, hjust=1)) + theme(legend.position="none") +

306 ggtitle(paste0("Cellular expression profiles of AD risk genes"))

307 dev.off()

6. Gene annotation enrichment analysis of AD risk genes

6.1. Download and configure HOMER by typing the commands below in terminal.

```
mkdir homer
cd homer
wget http://homer.ucsd.edu/homer/configureHomer.pl
perl ./configureHomer.pl -install
perl ./configureHomer.pl -install human-p
perl ./configureHomer.pl -install human-o
```

6.2. Run HOMER by typing the commands below in terminal.

```
export PATH=$PATH:~/work/homer/bin
findMotifs.pl ~/work/ADgenes.txt human ~/work/
```

6.3. Plot the enriched terms by typing the following code into the console window in RStudio.

```
library(ggpubr)
options(stringsAsFactors=F)
pdf("GO_enrichment.pdf",width=15,height=8)
plot_barplot = function(dbname,name,color){
  input = read.delim(paste0(dbname,".txt"),header=T)
  input = input[,c(-1,-10,-11)]
  input = unique(input)
  input$FDR = p.adjust(exp(input$logP))
  input_sig = input[input$FDR < 0.1,]
  input_sig$FDR = -log10(input_sig$FDR)
  input_sig = input_sig[order(input_sig$FDR),]

  p = ggbarplot(input_sig, x = "Term", y = "FDR", fill = color, color = "white", sort.val = "asc", ylab =
expression(-log[10](italic(FDR))), xlab = paste0(name," Terms"), rotate = TRUE, label =
paste0(input_sig$Target.Genes.in.Term,"/",input_sig$Genes.in.Term), font.label = list(color =
"white", size = 9), lab.vjust = 0.5, lab.hjust = 1)
  p = p+geom_hline(yintercept = -log10(0.05), linetype = 2, color = "lightgray")
  return(p)
}

p1 = plot_barplot("biological_process","GO Biological Process","#00AFBB")
p2 = plot_barplot("kegg","KEGG","#E7B800")
p3 = plot_barplot("reactome","Reactome","#FC4E07")

ggarrange(p1, p2, p3, labels = c("A", "B", "C"), ncol = 2, nrow = 2)
```

dev.off()

REPRESENTATIVE RESULTS:

The process described here was applied to a set of 800 credible SNPs that were defined by the original study¹⁴. Positional mapping revealed that 103 SNPs overlapped with promoters (43 unique genes) and 42 SNPs overlapped with exons (27 unique genes). After positional mapping, 84% (669) SNPs remained unannotated. Using Hi-C datasets in the adult brain, we were able to link an additional 208 SNPs to 64 genes based on physical proximity. In total, we mapped 284 AD credible SNPs to 112 AD risk genes (**Figure 1A**). AD risk genes were associated with amyloid precursor proteins, amyloid-beta formation, and immune response, reflecting the known biology of AD^{15–18} (**Figure 1B–D**). Developmental expression profiles of AD risk genes showed marked postnatal enrichment, indicative of the age-associated elevated risk of AD (**Figure 2A**). Finally, AD risk genes were highly expressed in microglia, primary immune cells in the brain (**Figure 2B**). This is in agreement with the recurrent findings that AD has a strong immune basis and microglia are the central player in AD pathogenesis^{14,19,20}.

FIGURES LEGENDS:

Figure 1: Defining putative target genes of AD GWS loci. (A) Credible SNPs derived from the top 29 AD loci were categorized into promoter SNPs, exonic SNPs, and unannotated non-coding SNPs. Promoter and exonic SNPs were directly assigned to their target genes by positional mapping, while chromatin interaction profiles in the adult brain were additionally used to map SNPs based on physical interactions. (B–D) Enrichment of GO (B), KEGG (C), and Reactome (D) terms in AD risk genes was performed using HOMER as described in protocol section 6. The x axis represents the false discovery rate (FDR) corrected $-\log_{10}$ (P-value). Enriched terms with $FDR < 0.1$ were plotted. Grey vertical lines represent $FDR = 0.05$. APP amyloid precursor protein. Numerator, the number of AD risk genes represented in each term; denominator, the number of genes in each term.

Figure 2: Characterization of AD risk genes. (A) AD risk genes are highly expressed in the postnatal cortex compared to the prenatal cortex. (B) Violin plots depicting distributions of gene expression values (normalized expression) in different cell types from the cortex. These results show that AD risk genes are highly expressed in microglia, consistent with previous studies¹⁴.

DISCUSSION:

Here we describe an analytic framework that can be used to functionally annotate GWS loci based on positional mapping and chromatin interactions. This process involves multiple steps (for more details see this review¹³). First, given that chromatin interaction profiles are highly cell-type specific, Hi-C data obtained from the appropriate cell/tissue types that best capture underlying biology of the disorder needs to be used. Given that AD is a neurodegenerative disorder, we used adult brain Hi-C data⁹ to annotate GWS loci. Second, each GWS locus often has up to hundreds of SNPs that are associated with the trait because of linkage disequilibrium (LD), so it is important to obtain putative causal ('credible') SNPs by computationally predicting the causality through the use of fine-mapping algorithms^{21,22} or experimentally testing regulatory activities using high-throughput approaches such as massively parallel reporter assays (MPRA)²³ or self-transcribing

active regulatory region sequencing (STARR-seq)²⁴. For the work described here, we used credible SNPs reported in Jansen et al.¹⁴. Third, promoter and exonic SNPs are annotated based on positional mapping. We used a simple positional mapping strategy in which SNPs were mapped to the genes when they overlapped with promoters (defined as 2 kb upstream of transcription start site) or exons. However, this approach can be further elaborated by assessing the functional consequences of exonic SNPs, such as whether the SNP induces nonsense mediated decay, missense variation, or nonsense variation. Fourth, chromatin interaction profiles from the appropriate tissue/cell type can be used to assign SNPs to their putative target genes based on physical proximity. We used interaction profiles anchored to promoters, but we can further refine or expand the interaction profiles by taking enhancer activities (guided by histone H3 K27 acetylation or chromatin accessibility) or exonic interactions into account. One important consideration in this process is to use consistent human genome build. For example, if genomic positions of summary statistics are not based on hg19 (i.e., hg18 or hg38), an appropriate version of the reference genome should be obtained or the summary statistics need to be converted to hg19 using liftover²⁵.

We applied this framework to identify putative target genes for AD GWAS, assigning 284 SNPs to 112 AD risk genes. Using developmental expression profiles²⁶ and cell-type specific expression profiles⁹, we then demonstrated that this gene set was consistent with what is known about AD pathology, revealing the cell types (microglia), biological functions (immune response and amyloid beta), and elevated risk upon age.

While we presented a framework that delineates potential target genes of AD and its underlying biology, it is of note that Hi-C based annotation can be expanded to annotate any non-coding variation. As more whole-genome sequencing data becomes available and our understanding about the non-coding rare variation grows, Hi-C will provide a key resource for interpretation of disease-associated genetic variants. A compendium of Hi-C resources obtained from multiple tissue and cell types will be therefore critical to facilitating a wide application of this framework to garner biological insights into various human traits and disease.

ACKNOWLEDGMENTS:

This work was supported by the NIH grant R00MH113823 (to H.W.) and R35GM128645 (to D.H.P.), NARSAD Young Investigator Award (to H.W.), and SPARK grant from the Simons Foundation Autism Research Initiative (SFARI, to N.M. and H.W.).

DISCLOSURES:

The authors have nothing to disclose.

REFERENCES:

1. Dekker, J., Misteli, T. Long-Range Chromatin Interactions. *Cold Spring Harbor Perspectives in Biology*. **7** (10), a019356 (2015).
2. Sanyal, A., Lajoie, B.R., Jain, G., Dekker, J. The long-range interaction landscape of gene promoters. *Nature*. **489** (7414), 109–113 (2012).
3. Plank, J.L., Dean, A. Enhancer function: mechanistic and genome-wide insights come together.

440 *Molecular Cell*. **55** (1), 5–14 (2014).

441 4. Dekker, J., Marti-Renom, M.A., Mirny, L.A. Exploring the three-dimensional organization of

442 genomes: interpreting chromatin interaction data. *Nature Reviews Genetics*. **14** (6), 390–403

443 (2013).

444 5. Martin, P. et al. Capture Hi-C reveals novel candidate genes and complex long-range

445 interactions with related autoimmune risk loci. *Nature Communications*. **6**, 10069 (2015).

446 6. Won, H. et al. Chromosome conformation elucidates regulatory relationships in developing

447 human brain. *Nature*. **538** (7626), 523–527 (2016).

448 7. Jäger, R. et al. Capture Hi-C identifies the chromatin interactome of colorectal cancer risk loci.

449 *Nature Communications*. **6**, 6178 (2015).

450 8. Chen, J.A.A. et al. Joint genome-wide association study of progressive supranuclear palsy

451 identifies novel susceptibility loci and genetic correlation to neurodegenerative diseases.

452 *Molecular Neurodegeneration*. **13** (1), 41–41 (2018).

453 9. Wang, D. et al. Comprehensive functional genomic resource and integrative model for the

454 adult brain. *Science*. **362** (6420) eaat8464 (2018).

455 10. Demontis, D. et al. Discovery of the first genome-wide significant risk loci for attention

456 deficit/hyperactivity disorder. *Nature Genetics*. **51** (1), 63–75 (2019).

457 11. Grove, J. et al. Identification of common genetic risk variants for autism spectrum disorder.

458 *Nature Genetics*. **51** (3), 431–444 (2019).

459 12. Lee, P.H. et al. Genome wide meta-analysis identifies genomic relationships, novel loci, and

460 pleiotropic mechanisms across eight psychiatric disorders. *bioRxiv*. 528117–528117 (2019).

461 13. Mah, W., Won, H. The three-dimensional landscape of the genome in human brain tissue

462 unveils regulatory mechanisms leading to schizophrenia risk. *Schizophrenia Research*. In press

463 (2019).

464 14. Jansen, I.E. et al. Genome-wide meta-analysis identifies new loci and functional pathways

465 influencing Alzheimer’s disease risk. *Nature Genetics*. **51** (3), 404–413 (2019).

466 15. Viola, K.L., Klein, W.L. Amyloid β oligomers in Alzheimer’s disease pathogenesis, treatment,

467 and diagnosis. *Acta Neuropathologica*. **129** (2), 183–206 (2015).

468 16. Mroczko, B., Groblewska, M., Litman-Zawadzka, A., Kornhuber, J., Lewczuk, P. Amyloid β

469 oligomers (A β Os) in Alzheimer’s disease. *Journal of Neural Transmission*. **125** (2), 177–191 (2018).

470 17. Heneka, M.T. et al. Neuroinflammation in Alzheimer’s disease. *Lancet Neurology*. **14** (4), 388–

471 405 (2015).

472 18. Minter, M.R., Taylor, J.M., Crack, P.J. The contribution of neuroinflammation to amyloid

473 toxicity in Alzheimer’s disease. *Journal of Neurochemistry*. **136** (3), 457–474 (2016).

474 19. Hansen, D.V., Hanson, J.E., Sheng, M. Microglia in Alzheimer’s disease. *The Journal of Cell*

475 *Biology*. **217** (2), 459–472 (2018).

476 20. Gjoneska, E. et al. Conserved epigenomic signals in mice and humans reveal immune basis of

477 Alzheimer’s disease. *Nature*. **518** (7539), 365–369 (2015).

478 21. Benner, C. et al. FINEMAP: efficient variable selection using summary data from genome-wide

479 association studies. *Bioinformatics*. **32** (10), 1493–1501 (2016).

480 22. Hormozdiari, F., Kostem, E., Kang, E.Y., Pasaniuc, B., Eskin, E. Identifying causal variants at loci

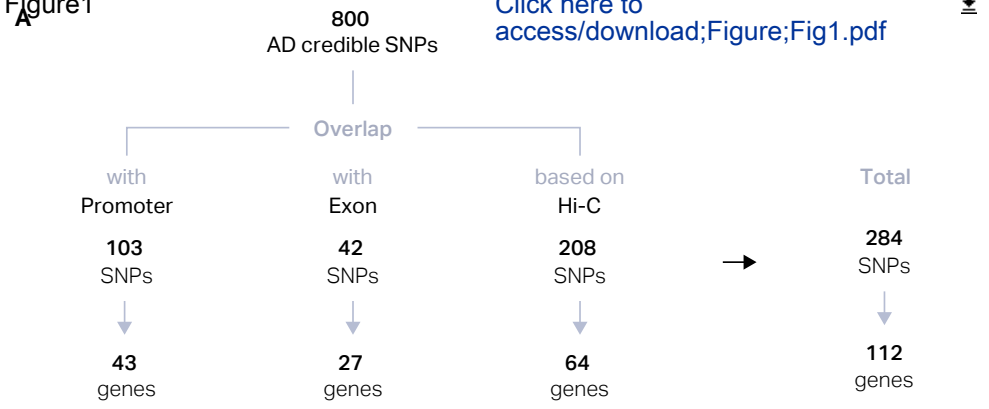
481 with multiple signals of association. *Genetics*. **198** (2), 497–508 (2014).

482 23. Tewhey, R. et al. Direct Identification of Hundreds of Expression-Modulating Variants using a

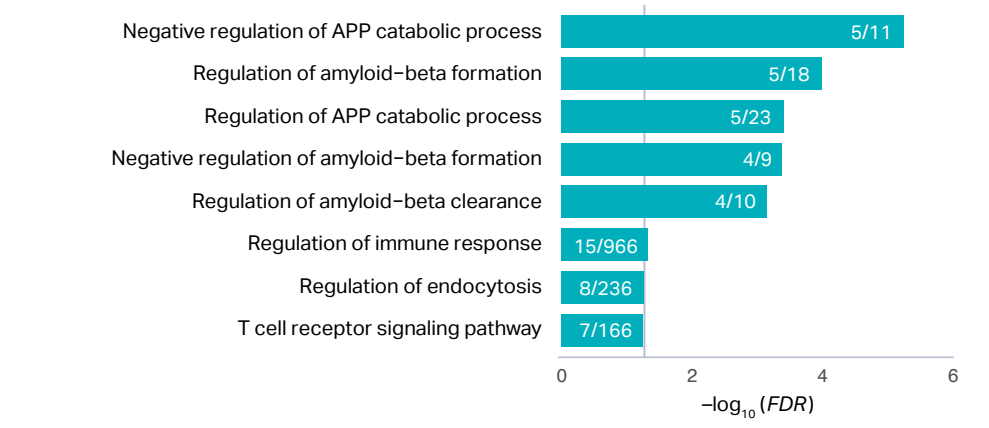
483 Multiplexed Reporter Assay. *Cell*. **165** (6), 1519–1529 (2016).

- 484 24. Arnold, C.D. et al. Genome-wide quantitative enhancer activity maps identified by STARR-seq.
485 *Science*. **339** (6123), 1074–1077 (2013).
- 486 25. Kent, W.J. et al. The human genome browser at UCSC. *Genome Research*. **12** (6), 996–1006
487 (2002).
- 488 26. Kang, H.J. et al. Spatio-temporal transcriptome of the human brain. *Nature*. **478** (7370), 483–
489 489 (2011).
- 490

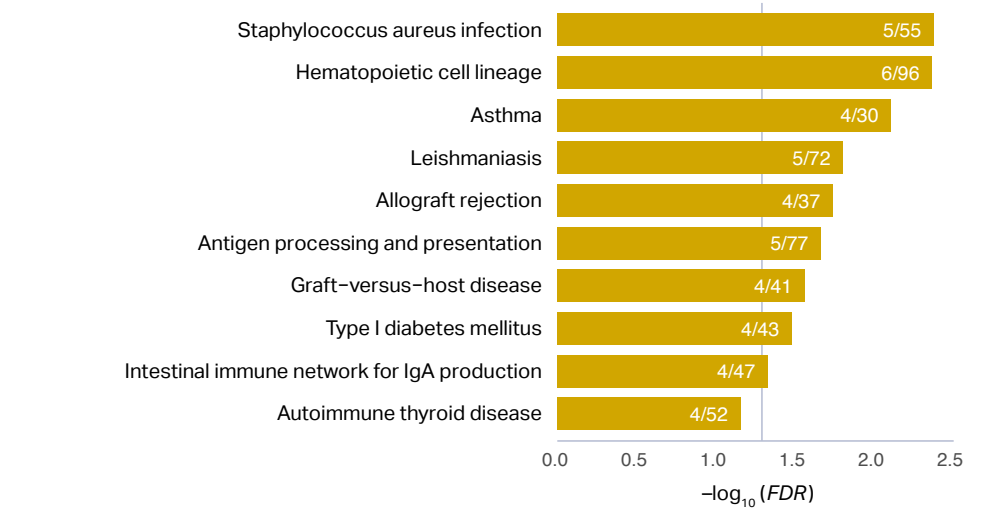
Figure1



B GO Biological Process Terms



C KEGG Terms



D Reactome Terms

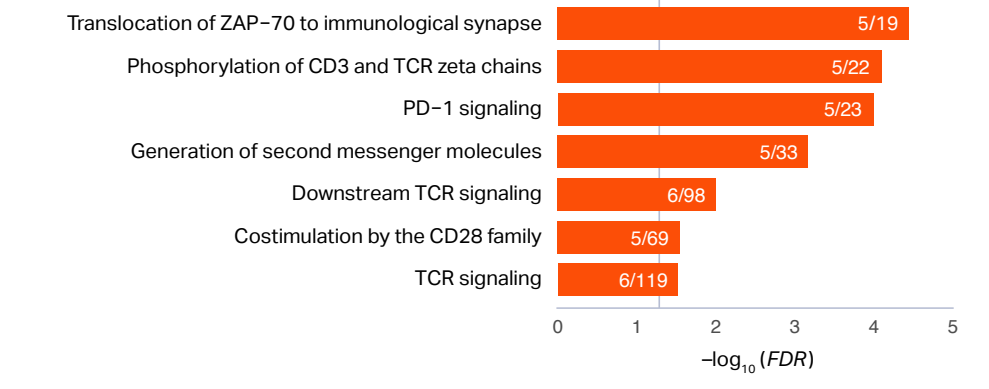
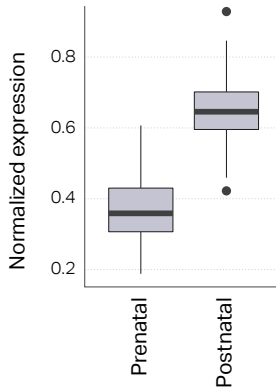


Figure2

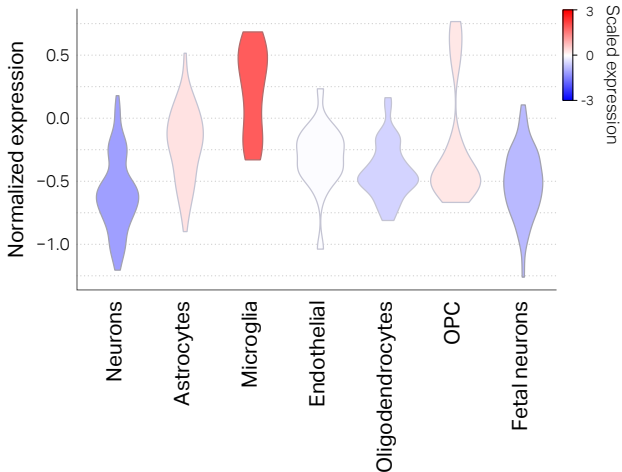
A

Brain Expression



B

Cellular expression profiles of AD risk genes



[Click here to access/download/figure/figure2.pdf](#)



Name of Material/Equipment/Files	Company	Catalog Number
10 kb resolution Hi-C interaction profiles in the adult brain from psychencode		
Developmental expression datasets		
Fine-mapped credible SNPs for AD		
(Supplementary Table 8 from Jansen et al. ¹⁴)		
Single cell expression datasets		
R (version 3.5.0)		
RStudio Desktop		
HOMER		

Comments/Description

http://adult.psychencode.org/Datasets/Integrative/Promoter-anchored_chromatin_loops.bed

http://www.brainspan.org/api/v2/well_known_file_download/267666527

https://static-content.springer.com/esm/art%3A10.1038%2Fs41588-018-0311-9/MediaObjects/41588_2018_311_MOESM3_ESM.xlsx

http://adult.psychencode.org/Datasets/Derived/SC_Decomp/DER-19_Single_cell_markergenes_TPM.xlsx

<https://www.r-project.org/>

<https://www.rstudio.com/products/rstudio/download/>

<http://homer.ucsd.edu/homer/configureHomer.pl>

Editorial comments:

1. Please take this opportunity to thoroughly proofread the manuscript to ensure that there are no spelling or grammar issues. The JoVE editor will not copy-edit your manuscript and any errors in the submitted revision may be present in the published version.

We proofread the manuscript as suggested.

2. Title: Please revise to avoid the use of punctuation (colon, dash, etc.).

We now changed the title into “Mapping Alzheimer’s disease variants to their target genes using computational analysis of chromatin configuration”

3. Please revise the Protocol to contain only action items that direct the reader to do something (e.g., “Do this,” “Ensure that,” etc.). The actions should be described in the imperative tense in complete sentences wherever possible. Avoid usage of phrases such as “could be,” “should be,” and “would be” throughout the Protocol. Any text that cannot be written in the imperative tense may be added as a “NOTE.” Please include all safety procedures and use of hoods, etc. However, notes should be used sparingly and actions should be described in the imperative tense wherever possible. Please move the discussion about the protocol to the Discussion.

We ensured that the protocol does not have any phrases such as “could be”, “should be”, and “would be.”

4. Please add more details to your protocol steps. There should be enough detail in each step to supplement the actions seen in the video so that viewers can easily replicate the protocol. Please ensure you answer the “how” question, i.e., how is the step performed?

This request is similar to request #5. Please see our response to request #5.

5. For actions involving software usage, please provide all specific details (e.g., button clicks, software commands, any user inputs, etc.) needed to execute the actions. Please include a step-wise description of software usage; mention what button is clicked on in the software, or which menu items need to be selected, and provide user input commands, etc.

We made the following changes in response to the request #4 and #5.

1. We directed the readers to install RStudio (Line 93: Install RStudio Desktop: <https://www.rstudio.com/products/rstudio/download/>), which will help them run the code provided.

2. We added a sentence “Type the following code into the console window in RStudio (e.g. Line 96)” in front of each section of code for added clarity.
3. We have revised the Download files section (section **1.3**) to include links to all files and explicit directions for downloading data.

6. Line 283: Figure 2C does not exist. Please revise.

Thanks for pointing this out. We changed it to **Figure 2B**. We also found out that Figures were not labeled properly within the manuscript, which is now corrected.

7. Please remove the embedded figure(s) from the manuscript.

We removed the figures from the manuscript, and submitted separately.

Reviewers' comments:

Reviewer #1:

The authors in this paper developed a computational pipeline for linking GWAS loci to genes using Hi-C data, and applied to Alzheimer's disease for discovering AD risk genes. The pipeline mapped credible SNPs from AD GWAS to various regions including enhancers, gene bodies and promoters, and then linked SNPs to genes if the mapped regions have potential interactions from Hi-C data. Overall, the paper was well organized and provided a complete set of R codes for pipeline implementation. Before recommending for publication, I have the following minor concerns that authors need to address:

We thank the reviewer for his/her positive and constructive comments.

1. Hi-C data description such as protocol, resolution & tissue/cell type is missing. Can the pipeline be scalable to the Hi-C datasets with different resolutions, which authors discussed at the end?

We thank the reviewer for this comment. We already described that Hi-C interaction profiles were generated from the adult brain. Thanks to the comment, we also added Line 125: “NOTE: In case other Hi-C datasets are used, this protocol requires Hi-C datasets processed at high resolution (5-20kb).”

2. Fig. 2A seems mix multiple regions together. Is there any particular developmental expression pattern for particular regions? Authors didn't introduce how to normalize gene expression either.

We used cortical expression data from brainspan. We now updated **Figure 2A legend** to describe the brain region: “AD risk genes are highly expressed in the postnatal cortex compared to the prenatal cortex”. We used the expression data provided by BrainSpan (<https://www.brainspan.org/>) and did not perform any additional normalization or processing. We now updated section **1.3.1** with detailed instructions for how to download the expression data file and process it.

3. Was the single cell data from healthy or AD brain? If healthy, the cell type specific expression might not represent AD cells. Will cell type Hi-C improve gene linking over tissue Hi-C? Also, details on gene expression normalization is missing.

Single cell data was from healthy (neurotypical) brain. We used the expression data provided by the original paper (Wang et al., Science 2018) which has been already normalized. We also updated section **1.3.1** to describe that the single cell data comes from healthy brains: Line 127: “Single cell expression datasets from the PsychENCODE (Described as *singlecell.rda* below). This is from neurotypical control samples.”

Reviewer #2:

Manuscript Summary:

The manuscript describes a method to annotate non-coding variations to the candidate genes for GWAS SNPs. Furthermore, the candidate genes have been subjected to enrichment analysis and cross-cell expression comparisons. Addressing the following comments will improve the quality of the manuscript.

We thank the reviewer for his/her critical and insightful comments.

Major Concerns:

1. The explanation about figure 1B is completely missing in the manuscript. It's important to mention and explain what resource was being used to get the figure and what purpose does it serve as the results of the work done. Additionally, the information about the total number of genes involved in each of the GO terms should also be written. Do all the genes take part in each of these GO terms? I don't think so. Also, what does the $-\log_{10}(\text{FDR})$ mean? This needs explanation as well. Overall, the results have been represented in a very abstract way. They need to be explained in detail.

We are sorry that **Figure 1B** was labelled incorrectly. We corrected the figure captions for Figures 1 and 2: “AD risk genes were associated with amyloid precursor proteins, amyloid-beta formation, and immune response, reflecting the known biology of AD^{15–18} (**Figure 1B**).”

Based on the reviewer's suggestion, we

(1) Updated **Figure 1B-D** with the # of genes in each term and # of genes represented in our list.

(2) Updated **Figure 1B-D legend** as below:

(B-D) Enrichment of GO (**B**), KEGG (**C**), and Reactome (**D**) terms in AD risk genes was performed using HOMER as described in step 5. The x-axis represents the FDR corrected -log₁₀ (P-value). Enriched terms with FDR<0.1 were plotted. Grey vertical lines represent FDR=0.05. APP amyloid precursor protein. Numerator, the number of AD risk genes represented in each term; denominator, the number of genes in each term.

2. This is a follow up to the 1st comment: The authors have done a sort of enrichment analysis depicting top 10 GO terms. It's highly recommended to perform and include results from Gene Set Enrichment Analysis (GSEA) using Broad Institute's Molecular signature database. This will give a better overview of associated pathways from KEGG and Reactome for the genes the authors have shortlisted in their work. In addition to this, it is also recommended to use NeuroMMsig (<https://neurommsig.scai.fraunhofer.de/>) to find out what Alzheimer Disease (AD) mechanisms are represented by the list of genes. NeuroMMsig is a mechanism enrichment analysis platform developed for neurodegenerative diseases, especially for AD and Parkinson disease (PD).

Thanks to this insightful suggestion, we used HOMER to analyze KEGG and Reactome for AD risk genes and reported the result in Figure 1. We think this greatly improved the overall findings of the manuscript, and gave more detailed descriptions about the potential function of AD risk genes.

3. It's difficult to comprehend what do the different shapes of different cellular expression profiles refer to.

Sorry for the confusion. We now updated **Figure 2B legend** as below:

“(B) Violin plots depicting distributions of gene expression values (normalized expression) in different cell types from the cortex. These results show that AD risk genes are highly expressed in microglia, consistent with previous studies¹⁴.”

4. The authors have reported about highly expressed genes in microglia (Figure 2B) as a whole. It is good to make overall comparison of gene expressions across different cell types but when it comes to gene expression analysis, the down-regulated genes are also of interests. Therefore, the author might have to report about individual genes that are either over-expressed or under-expressed to improve the quality of the manuscript.

While it is often of interest to investigate genes with different patterns of change, the intent of this plot is to show that the genes implicated by our analysis are expressed in cell types with a known role in AD. This offers support to the validity of our analyses.

While the reviewer raised an important point, we believe that the analysis of other genes that are differentially expressed across these cell types is outside the scope of this manuscript.

5. The R data files for *devExpr.rda*, *ADgenes.rda* and *singlecell.rda* should be made available.

Because *devExpr.rda* and *singlecell.rda* use datasets that have been made publicly available by the original paper, we do not feel comfortable providing these files. Instead, we updated section 1.3.1 in which we described the download procedure in detail so that readers can easily follow and retrieve the same data. We are providing *ADgenes.rda* with this paper.

Minor Concerns:

Line 60: There are double full stops.

Line 70: First mention of GWS, therefore needs to be written with the full name and the abbreviation.

Line 75: Double "the".

Line 326-328: Needs to be checked again.

We thank the reviewer for catching these minor errors. We now corrected all the points brought up by the reviewer.

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function

Addins

Project: (None)

Console

Jobs x

~/

R version 3.5.0 (2018-04-23) -- "Joy in Playing"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[workspace loaded from ~/.RData]

> |

Type codes here

R Console

You can open a code editor by clicking here

Environment

History

Connections

Import Dataset

List

Global Environment

Environment is empty

Data and Command History

Files

Plots

Packages

Help

Viewer

Zoom

Export

Once you generate plot, the plot will be appeared here

File Edit Code View Plots Session Build Debug Profile Tools Help

Go to file/function

Addins

Project: (None)

Untitled1*

Source on Save

Run

Source

```
9 install.packages("gProfiler")
10 install.packages("tidyverse")
11 install.packages("ggpubr")
12 library(GenomicRanges)
13 options(stringsAsFactors = F)
14 setwd("~/work") # This is the path to the working directory.
15 credSNP = read.delim("Supplementary_Table_8_Jansen.txt", header=T)
16 credSNP = credSNP[credSNP$Credible.Causal=="Yes", ]
17
18 credranges = GRanges(credSNP$Chr, IRanges(credSNP$bp, credSNP$bp), rsid=cr
19 save(credranges, file="AD_credibleSNP.rda")
20
```

Code Editor

Environment History Connections

Import Dataset

List

Global Environment		
datMeta	345 obs. of 10 variables	
datProbes	17085 obs. of 5	
exon	324422 obs. of 4 variables	
exonranges	Large GRanges (324422 ele...	
genome	16767 obs. of 2 variables	
geneAnno1	57317 obs. of 8 variables	
hic	149097 obs. of 5 variables	

Data

Console Jobs

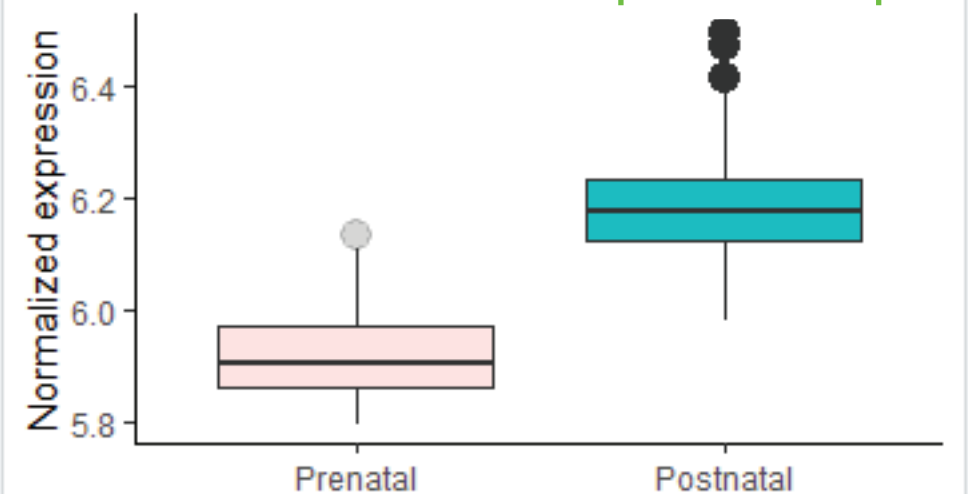
```
~/work/
V4 CTX Prenatal 5.803337
V7 CTX Prenatal 5.810975
V9 CTX Prenatal 5.827192
V11 CTX Prenatal 5.883260
V12 CTX Prenatal 5.980518
V13 CTX Prenatal 5.827781
> ggplot(dat,aes(x=Unit, y=Expr, fill=Unit, alpha=Unit)) + ylab("Normalized expr
ession") + geom_boxplot(outlier.size = NA) + ggtitle("Brain Expression") + xlab
("") + scale_alpha_manual(values=c(0.2, 1)) + theme_classic() + theme(legend.pos
ition="na")
>
>
```

R Console

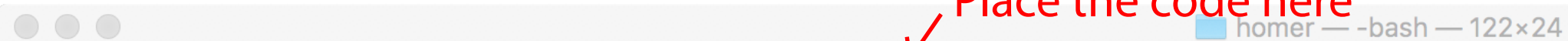
Files Plots Packages Help Viewer

Zoom Export

Brain Expression Graphical output



Type codes here

The terminal window has a light gray title bar with three window control buttons (red, yellow, green) on the left. The title text is "homer — -bash — 122x24".

homer — -bash — 122x24

[phanstiel3@:~/work/homer\$ export PATH=\$PATH:~/work/homer/bin]

phanstiel3@:~/work/homer\$ findMotifs.pl ~/work/ADgenes.txt human ~/work/

Place the code here

A red arrow points from the text "Place the code here" to the terminal prompt "phanstiel3@:~/work/homer\$".

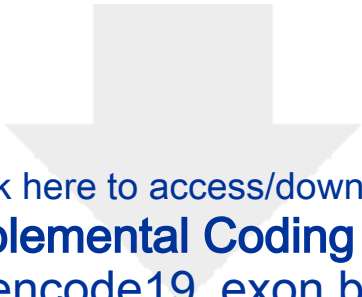


[Click here to access/download](#)

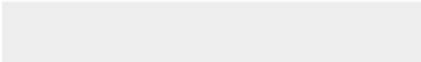

Supplemental Coding Files

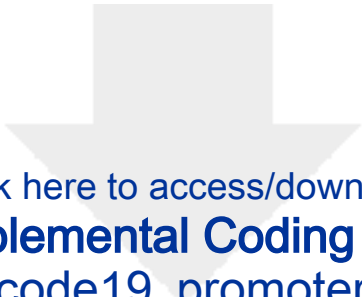
b2c23739-67a1-4695-ac6f-cbd4657c4c5f



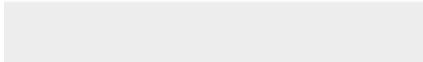



Click here to access/download
Supplemental Coding Files
Gencode19_exon.bed





[Click here to access/download](#)
Supplemental Coding Files
[Gencode19_promoter.bed](#)





Click here to access/download
Supplemental Coding Files
ADgenes.rda