

**Computational methods to dissect the genetic basis  
of human disease**

by

Samuel Sungil Kim

B.E.E., Auburn University (2013)

S.M., Massachusetts Institute of Technology (2019)

Submitted to the Department of Electrical Engineering and Computer  
Science

in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2021

© Massachusetts Institute of Technology 2021. All rights reserved.

Author .....

Department of Electrical Engineering and Computer Science

May 10, 2021

Certified by .....

Alkes L. Price

Professor of Statistical Genetics

Harvard T.H. Chan School of Public Health

Thesis Supervisor

Certified by .....

Manolis Kellis

Professor of Electrical Engineering and Computer Science

Thesis Supervisor

Accepted by .....

Leslie A. Kolodziejcki

Professor of Electrical Engineering and Computer Science

Chair, Department Committee on Graduate Students



# Computational methods to dissect the genetic basis of human disease

by

Samuel Sungil Kim

Submitted to the Department of Electrical Engineering and Computer Science  
on May 10, 2021, in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy

## Abstract

Genome-wide association studies (GWAS) have been successful in identifying disease-associated genetic variants. However, the path from GWAS to biological insight remains challenging, notably in identifying relevant biological pathways, explaining mechanistic links between diseases, and nominating disease-critical tissues and cell types. In this thesis, I introduce computational methods to dissect the genetic basis of human disease by integrating GWAS with functional data. In the first chapter, I integrate the GWAS with biological pathways and gene networks to elucidate biological mechanisms. I identify significantly associated pathways and highlight the importance of accounting for regulatory annotations in pathway enrichment and gene network analyses. In the second chapter, I investigate the shared genetic architecture between Mendelian disease and common disease by developing a machine learning framework to impute and denoise Mendelian disease-derived pathogenicity scores. I assess the informativeness of Mendelian pathogenicity scores for common disease and improve upon existing scores. In the third chapter, I prioritize disease-critical cell types by integrating GWAS with single-cell gene expression and chromatin accessibility profiling of fetal and adult brains. I show that identified disease-cell type associations recapitulates known biology while informing future analyses of disease mechanisms.

Thesis Supervisor: Alkes L. Price  
Title: Professor of Statistical Genetics  
Harvard T.H. Chan School of Public Health

Thesis Supervisor: Manolis Kellis  
Title: Professor of Electrical Engineering and Computer Science

## Acknowledgments

I would like to thank my thesis committee member: Alkes Price, Manolis Kellis, and Soumya Raychaudhuri for scientific feedback, project ideas, and valuable career advice. Alkes, I vividly recall when we first met in your office and talked about the omnigenic model. That moment was the beginning of my research in statistical genetics. Thank you for introducing me to the field, teaching me how to become a better scientist, writer, and thinker, and supporting my journey with great enthusiasm. I would also like to thank my thesis co-advisor Manolis Kellis. I am grateful for your continued support and scientific insights. Moreover, I would like to thank current and past members of Price group: Steven Gazal, Kushal Dey, Farhad Hormozdiari, Bryce van de Geijn, Omer Weissbrod, Karthik Jagadeesh, Yakir Reshef, Hilary Finucane, Po-Ru Loh, Carla Marquez-Luna, Armin Schoech, Luke O'Connor, Xuanyao Liu, Margaux Hujoel, Huwenbo Shi, Wouter Peyrot, Martin Zhang, Katie Siewert, Tiffany Amariuta. I learned greatly from you and truly enjoyed our scientific discussions. I would also like to thank current and past members of Kellis group: Yongjin Park, Carles Boix, Khoi Nguyen, Xinchun Wang, and Abhishek Sarkar. I would also express my appreciation to the NIH for the generous support.

To my family, partner, and friends, without your support, this thesis would not be possible.

# Contents

<b>1 Genes with high network connectivity are enriched for disease heritability</b>	<b>45</b>
1.1 Preface . . . . .	45
1.2 Abstract . . . . .	45
1.3 Introduction . . . . .	46
1.4 Methods . . . . .	48
1.5 Results . . . . .	61
1.6 Discussion . . . . .	77
1.7 Web resources . . . . .	79
<b>2 Improving the informativeness of Mendelian disease-derived pathogenicity scores for common disease</b>	<b>81</b>
2.1 Preface . . . . .	81
2.2 Abstract . . . . .	81
2.3 Introduction . . . . .	82
2.4 Methods . . . . .	83
2.5 Results . . . . .	96
2.6 Discussion . . . . .	117
2.7 Web resources . . . . .	120
<b>3 Leveraging single-cell ATAC-seq to identify disease-critical fetal and adult brain cell types</b>	<b>121</b>
3.1 Preface . . . . .	121

3.2	Abstract . . . . .	121
3.3	Introduction . . . . .	122
3.4	Methods . . . . .	123
3.5	Results . . . . .	127
3.6	Discussion . . . . .	138
3.7	Web resources . . . . .	140
<b>4</b>	<b>Conclusion</b>	<b>142</b>
<b>A</b>	<b>Supplementary tables for Chapter 1</b>	<b>143</b>
<b>B</b>	<b>Supplementary figures for Chapter 1</b>	<b>150</b>
<b>C</b>	<b>Supplementary tables for Chapter 2</b>	<b>154</b>
<b>D</b>	<b>Supplementary figures for Chapter 2</b>	<b>161</b>
<b>E</b>	<b>Supplementary tables for Chapter 3</b>	<b>174</b>
<b>F</b>	<b>Supplementary figures for Chapter 3</b>	<b>179</b>

# List of Figures

- 1-1 **Enriched pathways for three representative traits.** For (A) Crohn’s disease, (B) rheumatoid arthritis and (C) schizophrenia, we report the proportion of heritability explained and proportion of SNPs for each of 18,119 pathways analyzed. Red points indicate significantly enriched pathways ( $FDR < 5\%$ ), and grey points indicate non-significant pathways. Numerical results for all 42 diseases and complex traits are reported in Table A.4. . . . . 62
- 1-2 **Comparison of closeness centrality to other metrics that quantify the biological importance of each gene.** (A) For each of 7 gene sets, we report the excess overlap of genes in each decile bin of closeness centrality for the Greene thyroid network. Error bars represent 95% confidence intervals. (B) For each of 7 gene sets, we report the distribution of closeness centrality for the Greene thyroid network. Colored dots denote genes, white dots denote medians, and grey lines denote box plots. Numbers in parentheses below each gene set denote the number of genes. Results for all four networks and all 18 gene sets analyzed are reported in Figure B-3 and Table A.13 (for panel A) and Figure B-4 and A.14 (for panel B). Lists of genes for each of the 18 gene sets are provided in A.3. . . . . 67

1-3 **Heritability enrichment of network annotations.** We report (A) excess (fold) overlap between network annotations and baseline-LD functional categories; (B) heritability enrichment of network annotations, meta-analyzed across 42 independent traits; and (C)  $\tau^*$  values of network annotations, conditioned on either just the all-genes annotation, or the all-genes annotation and the baseline-LD model, meta-analyzed across 42 independent traits. In (A-C), the percentage under each bar indicates the proportion of SNPs in each annotation (defined for probabilistic annotations as the average value of the annotation), and error bars represent 95% confidence intervals. Numerical results for (A) are reported in Table A.16, and numerical results for (B) and (C) are reported in A.22. The S-LDSC results for the complete set of 168 network-trait pairs are reported in A.22. . . . . 69

1-4 **Heritability enrichment of pathway+network annotations.** We report (A) excess (fold) overlap between pathway+network annotations (averaged across up to 156 pathway-trait pairs); (B) heritability enrichment of pathway+network annotations, meta-analyzed across up to 156 pathway-trait pairs; and (C)  $\tau^*$  values of pathway+network annotations, conditioned on either just the all-genes annotation and the corresponding pathway and network annotations, or the baseline-LD model as well, meta-analyzed across up to 156 pathway-trait pairs. In (A-C), the percentage under each bar indicates the proportion of SNPs in each annotation (defined for probabilistic annotations as the average value of the annotation), and error bars represent 95% confidence intervals. Numerical results for (A) are reported in A.16, and numerical results for (B) and (C) are reported in Table A.34. The S-LDSC results for the complete set of 590 pathway-trait pairs are reported in A.34. . . . . 74



1-5	<b>Genes in enriched pathways have high network connectivity.</b>	
	For each of four networks, we report the sum of edge weights in the network between genes in the pathway, averaged across 141 enriched pathways. For comparison purposes, we report the same quantity averaged across 10,000 null pathways with the same number of genes. Error bars represent 95% confidence intervals (smaller than data points for null pathways). Numerical results, and analogous results for network connectivity between a pathway and interacting genes outside the pathway, are reported in Table A.36. . . . .	76

2-1 **Informativeness for common disease of binary annotations derived from 11 Mendelian disease-derived missense scores and corresponding boosted scores.** We report (A) heritability enrichment of binary annotations derived from published and boosted Mendelian disease-derived missense scores, meta-analyzed across 41 independent traits; (B) conditional  $\tau^*$  values, conditioning on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores). We report results for 10 Mendelian disease-derived missense scores (of 11 analyzed) for which annotations derived from published and/or boosted scores were conditionally significant; the published M-CAP score spanned too few SNPs to be included in the S-LDSC analysis. The percentage under each bar denotes the proportion of SNPs in the annotation; the proportion of top SNPs included in each annotation was optimized to maximize informativeness (largest  $|\tau^*|$  among Bonferroni-significant annotations, or most significant p-value if no annotation was Bonferroni-significant). Error bars denote 95% confidence intervals. In panel (B), \* denotes conditionally significant annotations. Numerical results are reported in Table C.2. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25. . . . . 100

2-2 **Informativeness for common disease of binary annotations derived from 6 genome-wide Mendelian disease-derived scores and corresponding boosted scores.** We report (A) heritability enrichment of binary annotations derived from published and boosted genome-wide Mendelian disease-derived scores, meta-analyzed across 41 independent traits; (B)  $\tau^*$  values, conditioning on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores). We report results for 6 genome-wide Mendelian disease-derived scores (of 6 analyzed) for which annotations derived from published and/or boosted scores were conditionally significant. The percentage under each bar denotes the proportion of SNPs in the annotation; the proportion of top SNPs included in each annotation was optimized to maximize informativeness (largest  $|\tau^*|$  among Bonferroni-significant annotations, or top 5% if no annotation was Bonferroni-significant; top 5% was the average optimized proportion among significant annotations). Error bars denote 95% confidence intervals. In panel (B), \* denotes marginally conditionally significant annotations. Numerical results are reported in Table C.13. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25. . . . . 105

2-3 **Informativeness for common disease of binary annotations derived from 18 additional genome-wide scores + 47 baseline-LD model annotations and corresponding boosted scores.** We report (A) heritability enrichments of binary annotations derived from published and boosted additional genome-wide scores, meta-analyzed across 41 independent traits; (B) conditional  $\tau^*$  values, conditioning on the baseline-LD model and 8 Roadmap annotations (for annotations derived from published scores) or the baseline-LD model, 8 Roadmap annotations, and corresponding published annotations (for annotations derived from boosted scores); (C) heritability enrichments of binary annotations derived from published and boosted baseline-LD model annotations; and (D) conditional  $\tau^*$  values of binary annotations derived from published and boosted baseline-LD model annotations. In (A) and (B), we report results for the 10 most informative additional genome-wide scores (of 18 analyzed). In (C) and (D), we report results for the 10 most informative baseline-LD model annotations (of 47 analyzed). The percentage under each bar denotes the proportion of SNPs in the annotation; the proportion of top SNPs included in each annotation was optimized to maximize informativeness (largest  $|\tau^*|$  among Bonferroni-significant annotations, or top 5% if no annotation was Bonferroni-significant; top 5% was the average optimized proportion among significant annotations). Error bars denote 95% confidence intervals. In panels (B) and (D), \* denotes conditionally significant annotations. Numerical results are reported in Table C.16. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25. . . . . 109

2-4 **Informativeness for common disease of 11 jointly significant binary annotations from combined joint model.** We report (A) heritability enrichment of 11 jointly significant binary annotations, meta-analyzed across 41 independent traits; (B) joint  $\tau^*$  values, conditioned on the baseline-LD model, 8 Roadmap annotations, and each other. We report results for the 11 jointly conditionally informative annotations in the combined joint model (S-LDSC  $\tau^*$   $p < 0.0001$  and  $|\tau^*| \geq 0.25$ ). The percentage under each bar denotes the proportion of SNPs in the annotation. Error bars denote 95% confidence intervals. Numerical results are reported in Table C.21. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25. . . . . 114

2-5 **Evaluation of improvement in heritability model fit.** We report (A) average  $\Delta \log l_{SS}$  (an approximate model likelihood metric<sup>1</sup>) across 30 UKBB traits; (B)  $\Delta \log l_{SS}$  of the baseline-LD and baseline-LD+marginal models for each trait.  $\Delta \log l_{SS}$  is computed as  $\log l_{SS}$  of a given model -  $\log l_{SS}$  of a model with no functional annotations (baseline-LD-nofunct model: MAF/LD annotations only). In panel (A),  $k$  denotes the number of new annotations beyond the baseline-LD model. Numerical results are reported in Table C.22. . . . . 115

3-1 **Disease enrichments of cell-type annotations derived from fetal brain.** We report (A)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 14) fetal brain scATAC-seq cell type annotations; (B)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 34) fetal brain scRNA-seq cell type annotations; and (C) comparison of results for 13 cell types included in both fetal brain scATAC-seq and scRNA-seq data. In A-B, only statistically significant results (FDR > 5%) are colored ( $-\log_{10}(\text{p-value}) \geq 1.67$  for scATAC-seq,  $\geq 2.70$  for scRNA-seq). In A-B, cell types appearing in both datasets are denoted in red font. Numerical results for all diseases/traits and cell types are reported in Table E.4, Table E.6, and Table E.7. ADHD: attention deficit hyperactivity disorder, SCZ: schizophrenia, MDD: major depressive disorder, BMI: body mass index. 130

3-2 **Disease enrichments of cell-type annotations derived from adult brain.** We report (A)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 18) adult brain scATAC-seq cell type annotations; (B)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 17) adult brain scRNA-seq cell type annotations; (C) comparison of results for 8 cell types included in both adult brain scATAC-seq and scRNA-seq data; and (D) comparison of results for 3 cell types included in both fetal and adult brain scATAC-seq data. In A-B, only statistically significant results (FDR > 5%) are colored ( $-\log_{10}(\text{p-value}) \geq 1.79$  for scATAC-seq,  $\geq 2.04$  for scRNA-seq). In A-B, cell types appearing in both datasets are denoted in red font. Numerical results for all diseases/traits and cell types are reported in Table E.14, Table E.15, Table E.16, and Table E.17. ADHD: attention deficit hyperactivity disorder, SCZ: schizophrenia, MDD: major depressive disorder, BMI: body mass index. . . . . 135

**B-1 Summary of pathways analyzed.** We analyzed 18,119 pathways, each consisting of at least 10 and at most 500 coding genes. Number in the boxplot represents the median number of genes for each database. Number of pathways in each database is shown on the top of the figure. BS: NCBI BioSystems. PC: PathwayCommons. MGI: Mouse Genome Informatics. GO: Genome Ontology. InWeb: InWeb protein-protein interactions. See Table A.1 for a description of 18,119 pathways analyzed. 150

**B-2 Closeness centrality is independent from the gene size and the proportion of exon.** For each of four network annotations, we computed a Pearson correlation between probabilistic annotation values and (1) gene size and (2) proportion of exon. We calculated the proportion of exon as the size of coding regions (bp) that lie inside the gene divided by the size of the gene, defined as (transcription stop - transcription start). . . . . 151

**B-3 Excess overlap of 18 genes in each decile bin of closeness centrality.** For each of 18 gene sets, we report the excess overlap of genes in each decile bin of closeness centrality for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. Error bars represent 95% confidence intervals. There is no olfactory receptor gene in Saha network. . . . . 152

**B-4 Distribution of closeness centrality in 18 gene sets.** For each of 18 gene sets, we show the distribution of closeness for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. Box plot is shown as a grey line inside the violin plot where white dot represents median. Each colored dot represents a gene. Outliers are not displayed for Saha and InWeb (see Web resources for closeness scores) . . . . . 153

**D-1 Overview of AnnotBoost framework.** We describe the AnnotBoost model training. AnnotBoost requires only one input, a pathogenicity score to boost, and generates a genome-wide (probabilistic) boosted pathogenicity score. From the input pathogenicity score (e.g. CADD as shown here), we built a classification model, each for even and odd chromosome SNPs using 75 baseline-LD annotations<sup>2,3</sup> as features. We assessed informativeness of annotations derived from published scores (input) and boosted scores (output) for common disease using S-LDSC<sup>4</sup>.161

**D-2 Feature importance of boosted Mendelian disease-derived missense pathogenicity scores.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD drives the prediction of 11 boosted missense scores by AnnotBoost. We report the signed impact of top 20 features for each of 11 predictive models: (A) PolyPhen-2, (B) PolyPhen-2-HVAR, (C) MetaLR, (D) MetaSVM, (E) PROVEAN, (F) SIFT 4G, (G) REVEL, (H) M-CAP, (I) Primate-AI, (J) MPC, and (K) MVP. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources). . . . . 162

**D-3 Excess overlap between gene scores derived from input pathogenicity scores and 165 reference gene sets of biological importance.** We report the excess overlap of genes linked to published and boosted scores in existing gene sets of biological importance (summarized in Table C.10): (A) PolyPhen-2<sup>6,7</sup> gene quintiles from published and boosted scores, (B) CADD<sup>8,9</sup> gene quintiles from published and boosted scores, and (C) CCR<sup>10</sup> gene quintiles from published and boosted scores. Error bars represent 95% confidence intervals. Numeric results for excess overlap and correlaton among gene scores are shown in Table C.11. Numeric results for odds ratios and p-values from Fisher’s exact test between published gene quintiles and boosted gene quintiles are reported in Table C.12. . . . . 163



**D-4 Feature importance of boosted genome-wide Mendelian disease-derived pathogenicity scores.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD drives the prediction of 11 boosted missense scores by AnnotBoost. We report the signed impact of top 20 features for each of 6 genome-wide Mendelian disease-derived pathogenicity scores: (A) CADD, (B) Eigen, (C) Eigen-PC, (D) ReMM, (E) NCBoost, and (F) ncER. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources). . . . . 164

**D-5 Feature importance of boosted scores derived from 18 additional genome-wide scores and 47 baseline-LD model annotations.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD drives the prediction of 18 boosted additional scores by AnnotBoost. We report the signed impact of top 20 features for each of 18 additional scores: (A) CDTS, (B) CCR, (C-I) DeepSEA-CTCF, -DNase, -H3K27ac, -H3K4me1, -H3K4me2, -H3K4me3, -H3K9ac, (J-K) DIS-DNA, -RNA, (L) pLI, (M) LIMBR, (N-Q) Gene network connectivity-Saha, Greene, InWeb, Sonawane, (R) EDS. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see SHAP results of 47 boosted baseline-LD scores online; see Web resources). . . . . 165

**D-6 Classification of fine-mapped disease SNPs using aggregated scores.** We report the true positive rate, false positive rate, precision, and recall along with the classification accuracy (AUPRCs and AUROCs) of four aggregated scores on classifying 5 different independent SNP sets: (A, F) 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, (B, G) 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, (C, H) 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup>, (D, I) 1,379 fine-mapped SNPs without functional data for 47 traits from UK Biobank<sup>13</sup>, and (E, J) 14,807 GWAS significant SNPs<sup>14,15</sup>, from 10 LD-, MAF-, and genomic element-matched control SNPs. We report the average AUPRCs and AUROCs of even/odd-chromosome classifiers. Differences for AUROCs and AUPRCs attained between (1) baseline-LD and baseline-LD+joint model, (2) baseline-LD and baseline-LD+marginal model, and (3) baseline-LD+joint and baseline-LD+marginal were largely significant ( $p\text{-val} < 0.008$ ). Numerical results, including results using the most matched control SNPs (instead of 10), are reported in Table C.23. 166

**D-7 Feature importance of baseline-LD+joint model in predicting fine-mapped or GWAS significant SNPs.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD+joint drive the prediction of fine-mapped or GWAS significant SNPs from 10 matched control SNPs for each positive SNP. We report the signed impact of top 20 feature for each of 4 fine-mapped SNPs and GWAS significant SNPs: (A) Farh et al., (B) Huang et al., (C) Weissbrod et al., (D) Weissbrod et al. (fine-mapped without functional data), (E) GWAS significant SNPs. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources). . . . . 167

**D-8 Feature importance of baseline-LD+marginal model in predicting fine-mapped or GWAS significant SNPs.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD+marginal drive the prediction of fine-mapped or GWAS significant SNPs from 10 matched control SNPs for each positive SNP. We report the signed impact of top 20 feature for each of 4 fine-mapped SNPs and GWAS significant SNPs: (A) Farh et al., (B) Huang et al., (C) Weissbrod et al., (D) Weissbrod et al. (fine-mapped without functional data), (E) GWAS significant SNPs. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources). . . . . 168

**D-9 Classification of fine-mapped disease SNPs in single-score analysis using published and boosted scores.** We report the classification accuracy (AUROC and AUPRC) of each of the 82 boosted scores compared to the corresponding published score. We report AUROC (resp. AUPRC) on (A, F) 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, (B, G) 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, (C, H) 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup>, (D, I) 1,379 fine-mapped SNPs without functional data for 47 traits from UK Biobank<sup>13</sup>, and (E, J) 14,807 GWAS significant SNPs<sup>14,15</sup> from 10 LD-, MAF-, and genomic element-matched control SNPs. Numerical results, including results using the most matched control SNPs (instead of 10), are reported in Table C.25 . . . . . 169

**D-10 Informativeness of the baseline-LD model before and after adding 11 jointly significant binary annotations** We report meta-analyzed  $\tau^*$  of the baseline-LD model annotations, across 41 independent traits, from two different S-LDSC analyses: (1) the baseline-LD model + 8 Roadmap annotations and (2) the baseline-LD model + 8 Roadmap annotations + 11 jointly significant annotations. Error bars represent 95% confidence intervals. Numerical results are reported in Table C.27. . . . . 170

**D-11 Informativeness for common disease of binary annotations derived from boosted CDTS scores with or without MAF features.** We applied AnnotBoost to CDTS annotation using all baseline-Ld features and all features excluding MAF bins. Then, we applied S-LDSC, conditioning on published binary CDTS annotations (five thresholds from 90th percentile to 99.9th percentile) and baseline-LD model annotations; and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . Error bars represent 95% confidence intervals. . . . . 171

**D-12 Informativeness for common disease of binary annotations derived from boosted CDTS scores using imbalanced data.** We applied AnnotBoost to CDTS annotation of varying training data. Then, we applied S-LDSC, conditioning on published binary CDTS annotations (five thresholds from 90th percentile to 99.9th percentile) and baseline-LD model annotations; and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . Error bars represent 95% confidence intervals. . . . . 172

D-13	<b>Informativeness for common disease of binary and probabilistic annotations derived from published CDTS scores.</b>	We constructed binary and probabilistic annotations for published CDTS <sup>16</sup> . We applied S-LDSC, conditional on baseline-LD model annotations and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and $\tau^*$ . To construct probabilistic annotations of varying proportion of SNPs, we performed the following transformation to upweight the upper percentile and downweight the lower percentile SNPs: $\frac{e^{\alpha \cdot annot}}{e^\alpha}$ with $\alpha$ varied from 3 to 2000. Error bars represent 95% confidence intervals. . . . .	173
F-1	<b>Overview of methods and analyses.</b>	We describe the overview of methods building cell-type annotations from single-cell sequencing datasets and evaluating disease informativeness applying S-LDSC across GWAS summary statistics. ABC+Roadmap S2G refers to the brain-specific SNPs-to-Genes linking strategy using enhancer-gene links <sup>17-20</sup> . We separately analyzed fetal and adult brain data. . . . .	179
F-2	<b>Comparison between fetal brain scRNA-seq and adult brain scRNA-seq cell-type annotations.</b>	We report $-\log_{10}(\tau^*$ p-values) of fetal brain scRNA-seq and adult brain scRNA-seq annotations for 6 matched cell types (astrocytes, endothelial cells, microglia, oligodendrocytes, excitatory neurons, inhibitory neurons), conditioning on the baseline model, union of open chromatin regions, and each other. Numeric results are found in Table E.18. Correlation among cell-type annotations is found in Table E.8. . . . .	180

# List of Tables

- 1.1 **Top enriched pathway for each trait.** We report the top enriched pathway (most significant  $\tau^*$ ) for each of 34 traits with at least one significantly enriched pathway. The first 23 traits (above the line) are UK Biobank traits. Enrichment of the "melanin biosynthetic process" pathway for skin color is consistent with previous studies<sup>21,22</sup>, and enrichment of the "absent corpus callosum" pathway for years of education was reported in a previous genetic study<sup>23</sup>. The complete set of 156 significant pathway-trait pairs is reported in Table A.5. . . 63
- 2.1 **11 Mendelian disease-derived missense and 6 genome-wide Mendelian disease-derived pathogenicity scores.** For each of 17 Mendelian disease-derived pathogenicity scores analyzed, we provide a description and report the coverage (% of SNPs scored) and corresponding reference. The first 11 annotations are scores for missense variants, and the last 6 annotations are genome-wide scores. Annotations are ordered first by type and then by the year of publication. . . 99

2.2 **Summary of informativeness for common disease of annotations derived from 82 published scores and corresponding boosted scores** For each category of scores, we report the number of scores; the number of marginally conditionally informative annotations (S-LDSC  $\tau^*$   $p < 0.0001$ , conditional on the baseline-LD model) (baseline-LD+marginal model); and the number of jointly conditionally informative annotations in a combined joint model (S-LDSC  $\tau^*$   $p < 0.0001$  and  $|\tau^*| \geq 0.25$ , conditional on the baseline-LD model and each other) (baseline-LD+joint model). \*Based on 9/11 published Mendelian missense scores analyzed, as binarized annotations were too small to analyze for the remaining 2 published Mendelian missense scores. \*\*Based on 16/18 published additional scores analyzed, as binarized annotations were too small to analyzed for the remaining 2 published additional scores. . . . . 101

3.1 **Summary of findings.** For each of 4 single-cell chromatin and gene expression data sets analyzed, we report the number of brain cell types analyzed, the total number of disease-cell type pairs analyzed (based on 28 diseases/traits), the number of significant disease-cell type pairs (FDR  $< 5\%$  for positive  $\tau^*$ ), and the number of diseases/traits with a significant disease-cell type pair. . . . . 129

3.2 **Notable disease-cell type associations.** We report the disease/trait, cell type, data source, standardized effect size ( $\tau^*$ ), p-value for positive  $\tau^*$ , and FDR q-value for selected results. Full results are reported in Table E.4, Table E.6, Table E.14, Table E.15. A description of diseases/traits analyzed is provided in Table E.1. MDD: major depressive disorder, SCZ: schizophrenia, BMI: body mass index, ADHD: attention deficit hyperactivity disorder. . . . . 131

A.1	<b>List of 18,119 pathways analyzed.</b> For each pathway, we report a pathway ID, pathway description, database, Entrez IDs for genes, and the number of protein-coding genes. . . . .	143
A.2	<b>List of 42 independent diseases and traits analyzed.</b> For each trait, we report a trait identifier, trait description, reference, sample size, and heritability z-score. We selected these 42 traits based on a heritability z-score > 6 (see Background). We further indicated brain or blood-related traits. . . . .	143
A.3	<b>Lists of genes in 18 gene sets compared with closeness centrality.</b> We compiled lists of genes for 18 metrics (gene sets) that we compared with closeness centrality. We report Entrez IDs for genes in the 18 gene sets. . . . .	143
A.4	<b>S-LDSC results of all pathway-trait pairs</b> We applied S-LDSC to 760,869 pathway-trait pairs, conditioning on all-genes annotation and the baseline-LD model. For each pathway-trait pair, we report a proportion of SNPs, enrichment, and $\tau$ . . . . .	144
A.5	<b>S-LDSC results for 156 significantly enriched pathway-trait pairs.</b> For each significantly enriched pathway-trait pair, we report a proportion of SNPs, enrichment, and a $\tau^*$ . The 8 significant pathway-trait pairs were reported in previous genetic studies: "pathways in cancer" for height <sup>24</sup> ; "neuropeptide hormone activity" for BMI <sup>25</sup> ; "immune response" for both Crohn's disease and ulcerative colitis <sup>26</sup> ; "T-cell receptor," "abnormal T-cell physiology," and "cytokine-mediated signaling pathway" for rheumatoid arthritis <sup>27</sup> ; "absent corpus callosum" for years of education <sup>23</sup> . . . . .	144
A.6	<b>Average gene size of annotations.</b> For all-genes annotation, all pathways, pathway, network, and pathway+network, and Quack annotations, we report an average size of genes (and its standard deviation) and an average number of genes. . . . .	144



A.7	<b>Heritability enrichment of enriched pathway-trait pairs.</b> We meta-analyzed (A) 156 enriched pathway-trait pairs; (B) 13 enriched pathway-trait pairs for ExAC, Cassa, and Samocha gene sets; (C) 169 enriched pathway-trait pairs (a and b combined). In each case, we report meta-analyzed enrichments and $\tau^*$ . . . . .	144
A.8	<b>S-LDSC results of 156 enriched pathway-trait pairs excluding genes implicated by GWAS.</b> We removed genes implicated by previous GWAS studies (see 3.1 Methods; average of 5% of genes ( 2 genes) removed) and applied S-LDSC conditional on all-genes and baseline-LD model annotations. For each pathway-trait pair, we report a proportion of SNPs, enrichment, $\tau^*$ , and the number of genes in a pathway excluding GWAS significant genes. . . . .	144
A.9	<b>S-LDSC results of 195 pathway-trait pairs from previous pathway enrichment studies.</b> We applied S-LDSC to 95 pathway-trait pairs from five previous genetic studies <sup>28-32</sup> and 100 from a recent study <sup>33</sup> (A) conditioning on the baseline-LD model and all-genes annotation and (B) conditioning on all-genes annotation only. In each pathway-trait pairs for each case, we report a proportion of SNPs, enrichment, and $\tau$ . We assessed the statistical significance based on global FDR < 5% across 18,119 pathways tested ( $\tau^* < 0.000989$ ). . .	144
A.10	<b>S-LDSC results of 13 enriched pathway-trait pairs for ExAC, Cassa, Samocha gene sets.</b> For each of 13 enriched pathway-trait pairs, we report a proportion of SNPs, enrichment, and $\tau$ . . . . .	145

A.11	<b>Correlation of network annotations with baseline-LD model annotations.</b> We report the Pearson correlation between baseline-LD model annotations and (A) Saha network annotations of different centralities, (B) Saha network annotations of different tissues, (C) Greene network annotations of different centralities, (D) Greene network annotations of different tissue, (E) InWeb network annotations of different centralities, (F) Sonawane network annotations of different centralities, and (G) Sonawane network annotations of different tissue. . . . .	145
A.12	<b>Summary of gene networks analyzed.</b> We report the number of genes, the number of edges, and the distribution of edge weights for each of four networks (InWeb, Saha, Sonawane, Greene). . . . .	145
A.13	<b>Excess overlap of 18 genes in each decile bin of closeness centrality.</b> For each of 18 gene sets, we report the excess overlap (and standard error) of genes in each decile bin of closeness centrality for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. . . . .	145
A.14	<b>Per-gene closeness centrality scores and gene membership in 18 gene sets</b> We report the closeness centrality for all protein-coding genes that exist in each of (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. We indicate gene membership in each of 18 gene sets, marking ‘1’ if in the corresponding gene set. . . . .	145
A.15	<b>Correlation between closeness centrality and 18 gene sets.</b> We report Pearson correlations between closeness centrality and 18 gene sets analyzed for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. . . . .	145
A.16	<b>Excess fold overlap of network and pathway+network annotations with baseline-LD model annotations.</b> We report the excess fold overlap between baseline-LD model annotations and network, pathway+network, Quack, and all-genes annotations. . . . .	145

A.17	<b>Correlation of network and pathway+network annotations with baseline-LD model annotations.</b> We report the Pearson correlation between baseline-LD model annotations and network, pathway+network, Quack, and all-genes annotations. . . . .	146
A.18	<b>Excess fold overlap / correlation among functional annotations from the baseline-LD model.</b> We report (A) excess fold overlap and (B) correlation among functional annotations from the baseline-LD model. . . . .	146
A.19	<b>TFs enriched in high closeness centrality genes.</b> For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, for high closeness centrality genes (top decile), we report significantly enriched TFs (Benjamini-Hochberg adjusted p-value < 0.05). The description of TFs is provided in the ENCODE Chip-Seq Significance Tool (see Web resources). . . .	146
A.20	<b>Correlation of deciles of closeness centrality with baseline-LD model annotations.</b> For four networks, we constructed binarized network annotations based on deciles of closeness centrality. We report the Pearson correlation between these annotations and baseline-LD model annotations. . . . .	146
A.21	<b>Enriched GO terms of high closeness centrality genes.</b> For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we report significantly enriched GO terms in the following GO categories: biological process (BP), cellular component (CC), and molecular function (MP) (Benjamini-Hochberg adjusted p-value < 0.05). . . . .	146
A.22	<b>Heritability enrichment of network annotations.</b> For each of 4 network annotations, we report meta-analyzed enrichments and $\tau^*$ across 42 independent traits. We highlight the network attaining highest enrichment for each trait. For the three tissue-specific networks (Saha, Greene, Sonawane), we also report meta-analyzed enrichments and $\tau^*$ of network annotations constructed using the tissue that maximized the excess overlap with the High pLI (ExAC) gene set. . . . .	146

**A.23 Heritability enrichment of network annotations conditioning on one annotation from the baseline-LD at a time.** For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we meta-analyzed network annotations across 42 independent traits, conditioning on one annotation from the baseline-LD model at a time. We report meta-analyzed enrichments and  $\tau^*$  across 42 independent traits. We highlighted annotations that significantly reduced  $\tau^*$  (using Bonferroni-corrected p-val). . . . . 146

**A.24 Heritability enrichment of deciles of closeness centrality.** For four networks, we constructed binarized network annotations based on deciles of closeness centrality. We applied S-LDSC and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . For the three tissue-specific networks (Saha, Greene, Sonawane), we also report meta-analyzed enrichments and  $\tau^*$  of binarized network annotations constructed using the tissue that maximized the excess overlap with the High pLI (ExAC) gene set. . . . . 147

**A.25 Heritability enrichment of network annotations from network perturbation analysis.** We randomly removed 10% to 90% of edges from the original networks and computed closeness centrality on networks with edges removed; we performed five separate perturbation analyses for each value of the proportion of edges removed. We applied S-LDSC and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . . . . . 147

**A.26 Heritability enrichment of DSD-network annotations.** We applied the diffusion state distance (DSD) algorithm<sup>34</sup> to transform gene networks' edge weights with a random walk ( $k = 5$ ). Then, we constructed network annotations by re-computing closeness on DSD-transformed networks and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . . . . . 147

A.27	<b>Heritability enrichment of networks annotations from consensus networks.</b> We constructed consensus networks and made (A) probabilistic annotations based on closeness centrality and (B) binary annotations based on deciles of closeness centrality. We applied S-LDSC and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and $\tau^*$ . . . . .	147
A.28	<b>Correlation between closeness and gene expression.</b> For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we report the correlation between closeness centrality and gene expression across 53 GTEx tissues. . . . .	147
A.29	<b>Heritability enrichment of Saha TSN and TWN gene sets.</b> We constructed gene sets based on membership of genes in Saha tissue-specific networks (TSN) and transcriptome-wide networks (TWN). We applied S-LDSC to 36 TSN and 16 TWN gene sets across 42 independent traits. We report meta-analyzed enrichments and $\tau^*$ . . . . .	147
A.30	<b>Excess overlap of top deciles of closeness centrality of tissue-specific networks with High pLI (ExAC) genes.</b> For each tissue-specific network (Saha, Greene, Sonawane), for each tissue, we report the excess overlap between High pLI (ExAC) genes and the top decile of closeness centrality. . . . .	148
A.31	<b>Heritability enrichment of network annotations using relevant tissues for brain-related and blood-related traits.</b> For (A) 8 brain-related traits and (B) 10 blood-related traits, we report meta-analyzed enrichments and $\tau^*$ . . . . .	148
A.32	<b>Heritability enrichment of network annotations using 6 other network centrality metrics.</b> For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we constructed network annotations based on 6 other network centrality metrics and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and $\tau^*$ . . .	148

<p><b>A.33 Heritability enrichment of network annotations with different window sizes.</b> Instead of 100kb windows around genes, we added (A) 10kb or (B) 1Mb windows around genes when constructing network annotations. We report meta-analyzed enrichments and <math>\tau^*</math> across 42 independent traits. . . . .</p>	148
<p><b>A.34 Heritability enrichment of pathway+network annotations.</b> We meta-analyzed 156 pathway-trait pairs (122 for Saha, which has less pairs as all genes in some pathways do not exist in the Saha network). For each 590 pathway-trait pair, we report a proportion of SNPs, enrichment, and <math>\tau</math>. We also report meta-analyzed enrichments and <math>\tau^*</math> across 156 pathway-trait pairs (122 for Saha). We highlighted the network attaining highest enrichment for each pathway-trait pairs. . . .</p>	148
<p><b>A.35 Heritability enrichment of pathway+network annotations conditioning on one annotation from the baseline-LD at a time.</b> For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we constructed an average annotation across 156 pathway+network annotations and meta-analyzed averaged pathway+network annotations across 42 independent traits, conditioning on one annotation from the baseline-LD model at a time. We report meta-analyzed enrichments and <math>\tau^*</math> across 42 independent traits. . . . .</p>	148
<p><b>A.36 Network connectivity of enriched pathways and null pathways.</b> Using (A) sum of edge weights or (B) number of edges as network connectivity metrics, we report the number of interacting genes and network connectivity between genes in a pathway and neighboring genes outside the pathway. We constructed null pathways in two ways: (1) each gene sampled from a randomly chosen pathway or (2) each gene randomly sampled from all protein-coding genes. . . . .</p>	149

**A.37 Heritability enrichment of Quack annotations.** We applied the Quack random-forest classifier algorithm<sup>35</sup>. We used 18,119 pathways as a training data and applied Quack to four gene networks. We used the output of Quack to construct Quack pathway+network annotations (with 100kb window), applied S-LDSC, and meta-analyzed across 156 pathway-trait pairs (122 for Saha). We report meta-analyzed enrichments and  $\tau^*$ . . . . . 149

**A.38 Heritability enrichment of 53 pathway+network annotation using pathways excluding genes implicated by GWAS.** From 53 significant pathway-trait pairs after excluding GWAS significant genes, we constructed pathway+network annotations and meta-analyzed results across 53 pathway-trait pairs (40 for Saha). We report meta-analyzed enrichments and  $\tau^*$ . . . . . 149

**C.1 List of 41 independent diseases and complex traits analyzed.** Analogous to a previous study<sup>36</sup>, we considered 89 GWAS summary association statistics, including 34 traits from publicly available sources and 55 traits from the UK Biobank (up to N = 459K); summary association statistics were computed using BOLT-LMM v2.3<sup>37,38</sup>. We obtained 41 independent traits (average N = 320K) with genetic correlation less than 0.9 (computed using cross-trait LDSC<sup>39</sup>). For 6 traits, we analyzed two different sources (both publicly available and UK Biobank), resulting in total 47 summary statistics analyzed. For each trait, we report a trait identifier, trait description, reference, sample size, and heritability z-score. . . . . 154

C.2	<b>Informativeness for common disease of binary annotations derived from 11 Mendelian disease-derived missense scores and corresponding boosted scores.</b> We applied S-LDSC, conditional on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and $\tau^*$ .	154
C.3	<b>Correlations between published scores and genome-wide boosted scores</b> We report Pearson correlation ( $r$ ) between published scores and genome-wide boosted scores for (A) Mendelian disease-derived missense scores, (B) genome-wide Mendelian disease-derived scores, (C) other scores, and (D) all scores. We regard originally unscored SNPs as having a score of zero.	154
C.4	<b>Predictive accuracy of AnnotBoost in out-of-sample predictions of input published scores.</b> We evaluated the AnnotBoost model classification performance based on AUROC and AUPRC on a 20% held-out testing set (consisted of odd (resp. even) chromosome SNPs) through the 5-fold cross-validation. For each boosted score, we reported the average AUROCs of odd-chr and even-chr model.	155
C.5	<b>Expected heritability enrichments of binary annotations derived from boosted Mendelian disease-derived missense scores with significantly negative <math>\tau^*</math>.</b> We report the heritability enrichment that is expected based on the boosted annotation's overlap with the baseline-LD model and corresponding published annotations, by assuming that the $\tau$ of the annotation is zero. In general, expected enrichments were significantly less than observed enrichments.	155



C.6	<b>Informativeness for common disease of binary annotations derived from boosted Mendelian disease-derived missense scores, restricted to non-coding regions.</b> We restricted boosted Mendelian missense annotations (whose top X% variants are already optimized) to non-coding SNPs and applied S-LDSC (conditioning on the baseline-LD model and corresponding published annotations). We report meta-analyzed enrichments, $\tau$ , and $\tau^*$ across 41 independent traits; for meta-analysis, we used the same weights from the meta-analysis of the unrestricted boosted annotations. . . . .	155
C.7	<b>Trait-specific informativeness for common disease of binary annotations.</b> Instead of meta-analyzing S-LDSC results across traits, we report the trait-specific results for 5,311 annotation-trait pairs (113 annotations x 47 traits; 113 annotations whose proportions of top SNPs were optimized). We indicated 152 annotation-trait pairs that were significant at $FDR < 5\%$ on $\tau^*$ p-value. . . . .	155
C.8	<b>Heterogeneity of heritability enrichment and <math>\tau^*</math> across traits.</b> We report the heterogeneity p-value of heritability enrichment and $\tau^*$ across 41 traits, for each of annotations analyzed. . . . .	155
C.9	<b>Gene scores derived from published and boosted pathogenicity scores.</b> We report genes scores (maximum pathogenicity score of linked SNPs) for genes linked to each of 203 scores analyzed (consisting of 35 published, 82 boosted, 86 baseline-LD annotations). . . . .	155
C.10	<b>165 reference gene sets of biological importance.</b> We report (A) gene set's short description, number of genes, and reference; (B) ENSG ID for each of 165 gene sets (limited to protein-coding genes). . . . .	156

**C.11 Excess overlap of 165 reference gene sets in each quintile bin of gene scores derived from input pathogenicity scores.** We report (A) the excess fold overlap of genes in each quintile bin of both published and boosted gene scores; (B) the standard error of the excess overlap; (C) the Pearson correlation ( $r$ ) among the gene scores; and (D) the Spearman correlation ( $\rho$ ) among the gene scores. . . . . 156

**C.12 Odds ratios for fold overlaps by gene quintiles from boosted scores and corresponding gene quintiles from published scores in 165 reference gene sets.** We computed the significance of the association between gene score quintiles from each of 82 published-boosted score pairs using Fisher’s exact test. We report the odds ratios and p-value for (A) fifth quintile (top 20%) gene sets, (B) fourth quintile gene sets, (C) third quintile gene sets, (D) second quintile gene sets, and (E) first quintile (bottom 20%) gene sets. . . . . 156

**C.13 Informativeness for common disease of binary annotations derived from 6 genome-wide Mendelian disease-derived scores and corresponding boosted scores.** We applied S-LDSC, conditional on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . . . . . 156

**C.14 Informativeness for common disease of binary annotations derived from boosted genome-wide Mendelian pathogenicity scores, restricted to previously unscored variants.** We restricted boosted Mendelian annotations (whose coverage is  $< 100\%$ ; see Table 2.1) to previously unscored variants and applied S-LDSC (conditioning on the baseline-LD model and corresponding published annotations). We report meta-analyzed enrichments,  $\tau$ , and  $\tau^*$  across 41 independent traits; for meta-analysis, we used the same weights from the meta-analysis of the unrestricted boosted annotations. . . . . 156

**C.15 Summary of 18 additional genome-wide scores and 47 baseline-LD model annotations.** We provide a description of 18 additional genome-wide scores and 47 baseline-LD model annotations. These include 2 constraint-related annotations, 9 predicted regulatory annotations (7 epigenetic states and 2 predicted transcriptional/post-transcriptional regulatory effects), 7 gene-based annotations, and 47 annotations from the baseline-LD model. We report the coverage, out of 9,997,231 common SNPs we considered; we note that the baseline-LD model annotations consider genome-wide SNPs and do not have any missingness (see Methods). For binary annotations, the proportion of SNPs with value 1 is considered as coverage. Annotations are first ordered by the type and then by the year of publication. base (all SNPs), 10 MAF bins, and 6 LD-related annotations constitute baseline-LD-nofunct model (model with no functional annotations). . . 157

**C.16 Informativeness for common disease of binary annotations derived from 18 additional genome-wide scores + 47 baseline-LD model annotations and corresponding boosted scores.** We applied S-LDSC, conditional on the baseline-LD model and 8 Roadmap annotations (for annotations derived from published scores) or the baseline-LD model, 8 Roadmap annotations, and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . . . . . 157

**C.17 Informativeness for common disease of binary annotations derived from boosted versions of additional genome-wide scores and baseline-LD model annotations, restricted to previously unscored variants.** We restricted boosted additional annotations (whose coverage is  $< 100\%$ ; see Table C.15) to previously unscored variants and applied S-LDSC (conditioning on the baseline-LD model and corresponding published annotations). We report meta-analyzed enrichments,  $\tau$ , and  $\tau^*$  across 41 independent traits; for meta-analysis, we used the same weights from the meta-analysis of the unrestricted boosted annotations. . . . . 157

**C.18 Informativeness for common disease of binary annotations derived from 18 additional genome-wide scores + 47 baseline-LD model annotations and corresponding boosted scores without conditioning on 8 Roadmap annotations.** We applied S-LDSC, conditional on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . . . . . 158

**C.19 Informativeness for common disease of marginally conditionally significant annotations derived from Mendelian disease-derived missense scores and genome-wide Mendelian disease-derived scores and corresponding boosted scores, when additionally conditioning on 8 Roadmap annotations.** We applied S-LDSC to marginally significant Mendelian missense and genome-wide Mendelian annotations, conditioning on the baseline-LD model, 8 Roadmap annotations, and published annotations (for boosted annotations), and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . . . . . 158

**C.20 Genome-wide correlations among functional annotations.** We report Pearson correlation ( $r$ ) (A) among 11 jointly significant annotations, (B) between 11 jointly significant annotations and functional annotations from the baseline-LD model, and (C) among all 198 annotations (85 baseline-LD annotations (excluding ‘base’ annotation; all SNPs annotated with 1), 31 annotations derived from published scores, and 82 boosted annotations). . . . . 158

**C.21 Informativeness for common disease of 11 jointly significant binary annotations in a combined joint model.** We applied S-LDSC, conditioned on the baseline-LD model, 8 Roadmap annotations, and each other, and meta-analyzed results across 41 independent traits. We report proportion of SNPs, meta-analyzed enrichments, and  $\tau^*$ . . . 158

**C.22 Evaluation of heritability model fit using  $logl_{SS}$ .** We report  $logl_{SS}$ ,  $\Delta logl_{SS}$ , and AIC for GCTA (number of annotations = 1), baseline-LD-nofunct (17), baseline-LD (86), baseline-LD+joint (97), baseline-LD+marginal $\tau^* \geq 0.5$  (105), and baseline-LD+marginal model (150), after applying SumHer<sup>1</sup> across 30 UKBB summary statistics.  $\Delta logl_{SS}$  is computed as  $logl_{SS}$  of a given model -  $logl_{SS}$  of baseline-LD-nofunct model. AIC (to account for increase in model complexity) is computed as  $2 \times$  number of annotations -  $2logl_{SS}$ , as described in ref. speed2020. We also report the p-value on the differences in  $\Delta logl_{SS}$ . . 159

**C.23 Classification of fine-mapped disease SNPs: multi-score analysis.** We report AUROCs and AUPRCs of different combined scores in classifying 5 different independent SNP sets: 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup> (stringently defined by causal posterior probability  $\geq 0.95$ ; both functionally-informed and non functionally-informed), and 14,807 GWAS significant SNPs<sup>14,15</sup> (from 10 LD-, MAF-, and genomic element-matched control SNPs as well as the most matched low frequency and common SNPs in the reference panel<sup>40</sup>). . . . . 159

**C.24 Summary of fine-mapped SNPs analyzed and results of single-score classification analysis.** We report the reference SNP ID (RSID) of four fine-mapped SNP sets and GWAS significant SNPs and their corresponding control SNPs that are found in the reference panel (MAF  $\geq 0.5\%$ ). We also report their short description and the number of unique SNPs of 10 matched control SNPs and most matched control SNPs. We also report the summary of single-score classification analysis, including the number of scores, where boosted scores perform better than the corresponding published scores, and the average AUROC and AUPRC improvement. Numeric results for the single-score analysis can be found in Table C.25. . . . . 159

**C.25 Classification of fine-mapped disease SNPs: single-score analysis.** We report AUROCs and AUPRCs of published and boosted scores in classifying 5 different independent SNP sets: 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup> (stringently defined by causal posterior probability  $\geq 0.95$ ; both functionally-informed and non functionally-informed), and 14,807 GWAS significant SNPs<sup>14,15</sup> ((A) from 10 LD-, MAF-, and genomic element-matched control SNPs as well as (B) the most matched low frequency and common SNPs in the reference panel<sup>40</sup>). We report the S-LDSC results (including standardized enrichments) of derived annotations based on top prioritized SNPs side by side. The summary of single-score analysis can be found in Table C.24. . . . . 160

<b>C.26 Correlation between S-LDSC metrics and AUROCs from single-score analysis of classifying fine-mapped SNPs.</b>	
For (A) 75 published annotations (excluding small annotations with proportion of SNPs < 0.02%) and (B) 82 boosted annotations, we computed Pearson correlation ( $r$ ) between different S-LDSC metrics (standardized enrichment, $\tau^*$ ) and AUROCs on classifying fine-mapped SNPs. We note AUROCs are computed on all SNPs, and S-LDSC results correspond to results on top prioritized SNPs. We further note it is 75 published annotations, not 82, because 7 scores were excluded in S-LDSC analysis due to small annotation size for 4 scores and not comparable enrichments for 3 continuous-valued baseline-LD annotations. . . . .	160
<b>C.27 Informativeness of the baseline-LD model before and after adding 11 jointly significant binary annotations</b>	
We applied S-LDSC to (A) the baseline-LD model + 8 Roadmap annotations and (B) the baseline-LD model + 8 Roadmap annotations + 11 jointly significant annotations and meta-analyzed results across 41 independent traits. We report proportion of SNPs, meta-analyzed enrichments, and $\tau^*$ . . . . .	160
<b>E.1 List of 28 brain-related traits and 6 non-brain-related traits analyzed.</b>	
For each trait, we report a trait identifier, trait description, reference, sample size, and heritability z-score. . . . .	174
<b>E.2 Description of the baseline model.</b>	
For each of 53 functional annotations in the baseline model, we provide the description, proportion of SNPs annotated, and reference. . . . .	174
<b>E.3 Tissues and cell types analyzed.</b>	
We summarize 83 cell types (for 4 single-cell datasets) and 26 tissues (for bulk chromatin data) analyzed.	174



E.4	<b>Disease enrichments of cell-type annotations derived from fetal brain scATAC-seq data.</b> For each of 392 trait-cell type pairs (28 brain traits * 14 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	174
E.5	<b>Matched cell types across datasets.</b> We report the matched cell types, defined as either exactly or closely matching cell types appearing in (A) fetal brain (scATAC-seq and scRNA-seq), (B) adult brain (scATAC-seq and scRNA-seq), (C) scATAC-seq (fetal and adult brain), and (D) scRNA-seq (fetal and adult brain). . . . .	174
E.6	<b>Disease enrichments of cell-type annotations derived from fetal brain scRNA-seq data.</b> For each of 952 trait-cell type pairs (28 brain traits * 34 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	175
E.7	<b>Comparison between fetal brain scATAC-seq and fetal brain scRNA-seq cell-type annotations.</b> For 13 cell types appearing in both fetal brain scATAC-seq and scRNA-seq datasets, we report correlation between two corresponding annotations and $\tau^*$ (conditioning on the baseline model, union of open chromatin regions, brain-specific enhancer-gene links of all genes, and to each other). . . . .	175
E.8	<b>Correlation among cell-type annotations.</b> We report the genome-wide Pearson correlation ( $r$ ) across cell-type annotations from 4 single-cell datasets. . . . .	175
E.9	<b>Disease enrichments of cell-type annotations derived from fetal brain scATAC-seq data in non-brain traits.</b> For 84 trait-cell type pairs (6 non-brain traits * 14 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	175
E.10	<b>Disease enrichments of cell-type annotations derived from fetal brain scRNA-seq data in non-brain traits.</b> For 204 trait-cell type pairs (6 non-brain traits * 34 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	175

E.11	<b>GREAT gene set enrichment of cell-type annotations derived from fetal brain datasets.</b> For each cell-type annotation (A) from scATAC-seq and (B) from scRNA-seq, we report the top enriched gene ontology (GO) terms, fold enrichment, binomial test p-value, and FDR-corrected q-value. . . . .	175
E.12	<b>Disease enrichments of cell-type annotations derived from fetal brain scRNA-seq data (using window-based approach).</b> We annotated SNPs +/- 100kb around genes, instead of brain-specific enhancer-gene links, and assessed the heritability enrichments. For each of 952 trait-cell type pairs (28 brain traits * 34 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . .	175
E.13	<b>Disease-critical tissues using bulk chromatin data of fetal epigenomes.</b> For each of 980 disease-tissue-chromatin mark triplets (28 brain traits * 5 fetal tissues * 7 chromatin marks), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . We considered 833 high-quality epigenomesboix2021 and selected 25 fetal biosamples related to brain tissues. We considered tier 1 assays (DNase-seq, H3K4me1, H3K4me3, H3K27ac, H3K36me3, H3K9me3 and H3K27me3), as most of observed data comes from tier 1 assays. Consistent with previous study <sup>2,4,17</sup> , we built tissue annotations by annotating chromatin peaks using the MACS2 <sup>41</sup> . . . . .	176
E.14	<b>Disease enrichments of cell-type annotations derived from adult brain scATAC-seq data.</b> For each of 504 trait-cell type pairs (28 brain traits * 18 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	176
E.15	<b>Disease enrichments of cell-type annotations derived from adult brain scRNA-seq data.</b> For each of 476 trait-cell type pairs (28 brain traits * 17 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	176

E.16	<b>Comparison between adult brain scATAC-seq and adult brain scRNA-seq annotations.</b> For 6 cell types appearing in both adult brain scATAC-seq and scRNA-seq datasets, we report correlation between two corresponding annotations and $\tau^*$ (conditioning on the baseline model, union of open chromatin regions, brain-specific enhancer-gene links of all genes, and to each other). . . . .	176
E.17	<b>Comparison between fetal brain scATAC-seq and adult brain scATAC-seq annotations.</b> For 3 cell types appearing in both fetal brain scATAC-seq and adult brain scATAC-seq datasets, we report correlation between two corresponding annotations and $\tau^*$ (conditioning on the baseline model, union of open chromatin regions, brain-specific enhancer-gene links of all genes, and to each other). . . . .	176
E.18	<b>Comparison between fetal brain scRNA-seq and adult brain scRNA-seq annotations.</b> For 6 cell types appearing in both fetal brain scRNA-seq and adult brain scRNA-seq datasets, we report correlation between two corresponding annotations and $\tau^*$ (conditioning on the baseline model, brain-specific enhancer-gene links of all genes, and to each other). . . . .	176
E.19	<b>Disease enrichments of cell-type annotations derived from adult brain scATAC-seq data in non-brain traits.</b> For 108 trait-cell type pairs (6 non-brain traits * 18 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	177
E.20	<b>Disease enrichments of cell-type annotations derived from adult brain scATAC-seq data in non-brain traits.</b> For 102 trait-cell type pairs (6 non-brain traits * 17 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	177

E.21	<b>Disease enrichments of cell-type annotations derived from scATAC-seq data, conditioning on only baseline model.</b> We assessed heritability enrichments of (A) adult scATAC-seq and (B) fetal scATAC-seq cell-type annotations, conditioning on only baseline model (instead of conditioning on the baseline model and the union of chromatin marks across cell types). We report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . . .	177
E.22	<b>GREAT gene set enrichment of cell-type annotations derived from adult brain datasets.</b> For each cell-type annotation (A) from scATAC-seq and (B) from scRNA-seq, we report the top enriched gene ontology (GO) terms, fold enrichment, binomial test p-value, and FDR-corrected q-value. . . . .	177
E.23	<b>Disease enrichments of cell-type annotations derived from adult brain scRNA-seq data (using window-based approach).</b> We annotated SNPs +/- 100kb around genes, instead of brain-specific enhancer-gene links, and assessed the heritability enrichments. For each of 476 trait-cell type pairs (28 brain traits * 17 cell types), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . . . .	177
E.24	<b>Disease-critical tissues using bulk chromatin data of adult epigenomes.</b> For each of 4,116 disease-tissue-chromatin mark triplets (28 brain traits * 21 adult tissues * 7 chromatin marks), we report the proportion of SNPs, heritability enrichments, and $\tau^*$ . We considered 32 adult biosamples related to brain tissues and tier 1 assays (DNase-seq, H3K4me1, H3K4me3, H3K27ac, H3K36me3, H3K9me3 and H3K27me3). Donor information can be found in ref. boix2021. . .	177
E.25	<b>Genetic correlation among 28 brain-related traits.</b> We report the genetic correlation as estimated by S-LDSC across 28 GWAS summary association statistics. . . . .	178

# Chapter 1

## Genes with high network connectivity are enriched for disease heritability

### 1.1 Preface

This chapter is previously published in *American Journal of Human Genetics* as "Genes with high network connectivity are enriched for disease heritability" (2019; see ref.<sup>42</sup>) with coauthors Chengzhen Dai, Farhad Hormozdiari, Bryce van de Geijn, Steven Gazal, Yongjin Park, Luke O'Connor, Tiffany Amariuta, Po-Ru Loh, Hilary Finucane, Soumya Raychaudhuri, and Alkes L Price. I am also grateful to Manolis Kellis, Yakir Reshef, Evan Boyle, April Kim, Taibo Li, Marieke Kuijjer, and Xinchun Wang for helpful discussions.

### 1.2 Abstract

Recent studies have highlighted the role of gene networks in disease biology. To formally assess this, we constructed a broad set of pathway, network, and pathway+network annotations and applied stratified LD score regression to 42 diseases and complex traits (average  $N=323K$ ) to identify enriched annotations. First, we analyzed 18,119 biological pathways. We identified 156 pathway-trait pairs whose disease enrichment was statistically significant ( $FDR < 5\%$ ) after conditioning on all

genes and 75 known functional annotations (from the baseline-LD model), a stringent step that greatly reduced the number of pathways detected; most significant pathway-trait pairs were previously unreported. Next, for each of four published gene networks, we constructed probabilistic annotations based on network connectivity. For each gene network, the network connectivity annotation was strongly significantly enriched. Surprisingly, the enrichments were fully explained by excess overlap between network annotations and regulatory annotations from the baseline-LD model, validating the informativeness of the baseline-LD model and emphasizing the importance of accounting for regulatory annotations in gene network analyses. Finally, for each of the 156 enriched pathway-trait pairs, for each of the four gene networks, we constructed pathway+network annotations by annotating genes with high network connectivity to the input pathway. For each gene network, these pathway+network annotations were strongly significantly enriched for the corresponding traits. Once again, the enrichments were largely explained by the baseline-LD model. In conclusion, gene network connectivity is highly informative for disease architectures, but the information in gene networks may be subsumed by regulatory annotations, emphasizing the importance of accounting for known annotations.

### 1.3 Introduction

Human diseases and complex traits are heritable and highly polygenic, potentially involving a large number of disease-associated genes connected by dense cellular networks<sup>28,43–46</sup>. Recent work has employed several approaches to infer gene interaction networks, including protein-protein interaction networks<sup>47–49</sup>, tissue-specific co-expression networks<sup>50,51</sup>, and tissue-specific regulatory networks<sup>52–54</sup>. An appealing extension of traditional genome-wide association studies (GWAS) is to identify genes and gene pathways associated with disease by leveraging gene networks and network connectivity between disease-associated genes<sup>50,52,55–63</sup>. However, despite considerable progress on inferring gene interaction networks and applying newly developed methods for network connectivity-informed GWAS to identify specific genes and gene

pathways associated to disease, an overall assessment and interpretation of the contribution of gene networks to the genetic architecture of disease has remained elusive. In particular, the extent to which this contribution can be explained by disease enrichments of known functional annotations<sup>2,3,17,64–66</sup> is unknown.

Here, we sought to answer three questions. First, what is the contribution of disease-associated gene pathways<sup>28–33,55,67–70</sup> to disease heritability, irrespective of network connectivity. Second, what is the contribution of genes with high network connectivity in known gene networks<sup>49–51,53</sup> to disease heritability. Third, what is the contribution of genes with high network connectivity to disease-associated gene pathways to disease heritability.

To answer these questions, we constructed a broad set of pathway, network, and pathway+network annotations. The pathway annotations were constructed from known gene pathways by including 100kb windows around each gene; the network annotations were constructed by quantifying network connectivity using closeness centrality, a measure of how close a gene is to other genes in the network<sup>35,71</sup>; and the pathway+network annotations were constructed by annotating genes with high network connectivity to the input pathway, again quantified using closeness centrality. We applied stratified LD score regression<sup>2,3</sup> to quantify the contribution of the pathway, network, and pathway+network annotations to disease heritability. We conditioned our analyses on all genes and on the baseline-LD model, which includes a broad set of coding, conserved, regulatory and LD-related annotations<sup>3</sup>. In each case, we compared results before and after conditioning on the baseline-LD model, to assess the extent to which the disease enrichments that we identified could be explained by known functional annotations.

## 1.4 Methods

### Genomic annotations and the baseline-LD model

We define a genomic annotation as an assignment of a numeric value to each SNP in a predefined reference panel (e.g. 1000 Genomes<sup>40</sup>; see Web Resources). Continuous-valued annotations can have any real value. Probabilistic annotations can have any real value between 0 and 1. Binary annotations can have value 0 or 1 only. A binary annotation can be viewed as a subset of SNPs (the set of SNPs with annotation value 1). Annotations that correspond to known or predicted function are referred to as functional annotations.

The baseline-LD model<sup>3</sup> (v1.1) contains 75 functional annotations (see Web Resources). These annotations include binary coding, conserved, and regulatory annotations (e.g. promoter, enhancer, histone marks, transcription factor (TF) binding sites), and continuous-valued linkage disequilibrium (LD)-related annotations.

### Excess overlap between binary and/or probabilistic annotations

To study which annotation(s) from the baseline-LD model captures information of focal annotation, excess (fold) overlap of SNPs in annotations can be computed for any pair of binary and/or probabilistic annotations. We define excess (fold) overlap between a pair of annotations as the fraction of overlap between the annotations divided by the amount of overlap expected by chance, computed as the following:

$$\text{excess overlap}(\text{annot1}, \text{annot2}) = M * \frac{\sum_{j=1}^M \text{annot1}_j * \text{annot2}_j}{\sum_{j=1}^M \text{annot1}_j * \sum_{j=1}^M \text{annot2}_j} \quad (1.1)$$

where M is the total number of SNPs (5,961,159). When there is excess overlap, the excess fold overlap is  $> 1$ ; when there is depletion, the excess fold overlap is  $< 1$ . We did not compute excess overlap for continuous-valued annotations that are not probabilistic (e.g. those that can have negative values). A separate definition of excess overlap between gene sets is provided below.



## Effect size ( $\tau^*$ ) and enrichment metrics estimated by S-LDSC

We used stratified LD score regression (S-LDSC<sup>2,3</sup>) to estimate the enrichment and the standardized effect size ( $\tau^*$ ) of an annotation. Let  $a_{cj}$  represent the annotation value of the SNP  $j$  for the annotation  $c$ . S-LDSC assumes the variance of per-normalized genotype effect sizes is a linear additive contribution to the annotation  $c$ :

$$\text{Var}(\beta_j) = \sum_c a_{cj} \tau_c \quad (1.2)$$

where  $\tau_c$  is the per-SNP contribution of the annotation  $c$ . S-LDSC estimates  $\tau_c$  using the following equation:

$$\text{E}[\chi_j^2] = N \sum_c \ell(j, c) \tau_c + 1 \quad (1.3)$$

where  $N$  is the sample size of the GWAS and  $\ell(j, c)$  is the LD score of the SNP  $j$  to the annotation  $c$ . The LD score is computed as follow:  $\ell(j, c) = \sum_k a_{ck} r_{jk}^2$  where  $r_{jk}$  is the correlation between the SNPs  $j$  and  $k$ . The standardized effect size ( $\tau^*$ ), the proportionate change in per-SNP heritability associated with a one standard deviation increase in the value of the annotation (conditional on all the other annotations in the model), is defined as the following by previous study<sup>3</sup>:

$$\tau_{c^*} = \frac{\tau_c \text{sd}(C)}{h_g^2 / M} \quad (1.4)$$

where  $\text{sd}(C)$  is the standard deviation of the annotation  $c$ ,  $h_g^2$  is the estimated SNP-heritability, and  $M$  is the number of variants used to compute  $h_g^2$  (in our experiment,  $M$  is equal to 5,961,159). Unlike enrichment,  $\tau^*$  quantifies effects that are unique to the focal annotation.

Enrichment of the annotation is the fraction of heritability explained by SNPs in the annotation divided by the proportion of SNPs in the annotation. The definition

of enrichment could be extended to continuous annotations<sup>3,36</sup> as the following:

$$\text{Enrichment} = \frac{\%h_g^2(C)}{\%\text{SNP}(C)} = \frac{\frac{h_g^2(C)}{h_g^2}}{\frac{\sum_j a_{jc}}{M}} \quad (1.5)$$

where  $h_g^2(C)$  is the heritability captured by the  $c$ th annotation. When the annotation is enriched for trait heritability, the enrichment is  $> 1$ ; the overlap is greater than one would expect given the trait heritability and the size of the annotation.

The significance for enrichment for each annotation is computed using the block jackknife as mentioned in previous studies<sup>2,4,36</sup>. The significance for the effect size for each annotation, as mentioned in previous studies<sup>3,36</sup>, is computed as  $(\frac{\tau^*}{\text{se}(\tau^*)} \sim N(0, 1))$ , assuming that  $\frac{\tau^*}{\text{se}(\tau^*)}$  follows a normal distribution with zero mean and unit variance.

In all our analyses, we used the European samples in 1000G<sup>40</sup> (see Web Resources) as reference SNPs. Regression SNPs were obtained from HapMap 3<sup>72</sup> (see Web Resources). SNPs with marginal association statistics  $> 80$  and SNPs in the major histocompatibility complex (MHC) region were excluded. Unless stated otherwise, we included the baseline-LD model<sup>3</sup> in all primary analyses using S-LDSC, both to minimize the risk of bias in enrichment estimates due to model mis-specification<sup>2,3</sup> and to estimate effect sizes ( $\tau^*$ ) conditional on known functional annotations.

## Pathway annotations

We define a pathway as a set of genes (gene set). We considered 18,119 pathways from five sources: 2,118 biological pathways from the BioSystem (BS) database<sup>73</sup> (which includes pathways from BioCyc<sup>74</sup>, Kyoto Encyclopedia of Genes and Genomes (KEGG)<sup>75</sup>, Pathway Interaction Database (PID)<sup>76</sup>, REACTOME<sup>77</sup>, WikiPathways<sup>78</sup>), 1,927 biological pathways from the Pathway Commons (PC) database<sup>79</sup> (which includes pathways from HumanCyc<sup>80</sup>, Integrating Network Objects with Hierarchies (INOH)<sup>81</sup>, KEGG<sup>75</sup>, PANTHER<sup>82</sup>, PID<sup>76</sup>, REACTOME<sup>77</sup>, SMPDB<sup>83</sup>, NetPath<sup>84</sup>), 7,209 protein-protein interaction gene sets from the InWeb database<sup>49</sup>, 3,903 mouse

phenotype gene sets from the Mouse Genome Informatics (MGI; sets of genes whose orthologs are associated to mouse phenotypes) database<sup>85</sup>(i.e. sets of genes whose orthologs are associated to mouse phenotypes), and 2,961 gene ontology gene sets from the Genome Ontology (GO) database<sup>86</sup>. This set of pathways, which contain at least 10 genes and at most 500 genes (consistent with previous studies<sup>33,69</sup>, significantly overlap with mSigDB<sup>87</sup> and pathways analyzed in a previous study<sup>69</sup>. The complete list of pathways is provided in Table S1, and a histogram of the number of genes in pathways from each of the five sources is provided in Figure S1.

For each of 18,119 pathways, we constructed a binary pathway annotation by annotating a value of 1 for variants around the protein-coding genes in a given pathway ( $+/-$  100kb as in previous work<sup>4,33</sup>) and 0 for all other variants. To evaluate the contribution of each pathway to disease/trait heritability, we applied S-LDSC<sup>2,3</sup> to 18,119 pathway annotations across 42 independent diseases and complex traits (average  $N = 323K$ ; including 30 UK Biobank traits; see Table S2). We conditioned on the 75 functional annotations from the baseline-LD model<sup>3</sup> (which includes a broad set of coding, conserved, regulatory and LD-related annotations), as well as an “all-genes” annotation representing the set of all 19,031 protein-coding genes ( $\pm 100kb$ ). We removed 129 pathway-trait pairs whose annotated SNPs are less than 0.02% of the reference genome (European samples from the 1000 Genomes Project<sup>40</sup>; see Web Resources) as S-LDSC is not well-equipped for annotations that span very small proportion of the genome.

For each of 760,869 pathway-trait pairs (roughly 18,119 pathways x 42 traits), we assessed the statistical significance of the pathway annotations based on global FDR  $< 5\%$  on the pathway annotation’s standardized effect size ( $\tau^*$ ) p-value, defined as the proportionate change in per-SNP heritability associated to a one standard deviation increase in the value of the annotation, conditioned on other annotations included in the model. We note that controlling FDR for each trait and for all traits did not make a major difference in the number of identified enriched pathway-trait pairs. We further note that our choice of two-tailed test on the significance is conservative, partially attributing to the reduced number of significantly associated pathway-trait

pairs. Among significantly enriched pathway-trait pairs, we calculated a pairwise correlation for every pair of annotations and retained the more significant pathway for correlated pathways with  $r \geq 0.5$  (as in a previous study<sup>35</sup>). If correlated pathways were enriched for different traits, we retained both pathways.

For each of 156 significantly enriched pathway-trait pairs, we also constructed pathway annotations excluding genes implicated by GWAS. First, we downloaded all GWAS associations from the GWAS Catalog (see Web Resources); we restricted to significant associations (p-value  $\leq 5e-8$ ). For each of the 34 traits, we defined genes implicated by GWAS by including genes mapped to the lead SNP; if the SNP was intergenic, we included the nearest upstream and downstream genes (We note that the nearest gene might not be the correct target gene.<sup>88</sup>) Then, for each of the 156 enriched pathway-trait pairs, we removed trait-specific GWAS significant genes (5% of the genes, on average across 156 pathway-trait pairs) and rebuilt our pathway annotations with a 100kb window. For pathways significant for multiple traits, we built unique pathway annotations excluding genes implicated by GWAS for the corresponding traits.

## Gene networks and data processing

A gene network is defined by an edge weight (which we normalized to lie between 0 and 1) for each pair of genes, representing their connectivity in the network. We considered four gene networks: two co-expression networks (Saha<sup>51</sup> and Greene<sup>50</sup>), one protein-protein interaction network (InWeb<sup>49</sup>), and one regulatory network (Sonawane<sup>53</sup>). We note that the Sonawane regulatory network contains nonzero edge weights only for pairs of genes in which at least one of the genes is a known TF.

We processed gene networks into a uniform format. Each gene networks is represented by an  $N \times 3$  matrix, where  $N$  is the number of edges in the network. Three columns represent gene 1, gene 2, and the edge weight between gene 1 and gene 2. We used Entrez IDs as gene identifiers. We used the ensembl biomart tool (GRCh37 assembly; see Web Resources) for gene identifier conversion. We note that the Saha and Sonawane networks can have negative edge weights. These were transformed to

positive edge weights (see below) to avoid calculating connectivity metrics on negative edge weights. In all analyses, we considered protein-coding genes only.

*Saha network:* we downloaded Saha transcriptome-wide networks (TWN) for each of 16 tissues (see Web Resources). (We did not consider the Saha tissue-specific networks (TSN) in our primary analyses, as they contained very few edges and small numbers of genes for computing connectivity metrics.) We collapsed isoforms by converting Ensembl transcript to Entrez gene ID; no conflicting edges (multiple non-zero edges between a pair of mapped genes) were found. We treated all edges equally regardless edge types (TE-TE (total expression), TE-IR (isoform ratio), or IR-IR). We transformed negative edge weights to positive edge weights by taking the absolute value.

*Greene network:* we downloaded “top edges” gene networks for each of 144 tissues (see Web Resources). These networks contain edges with evidence supporting a tissue-specific functional interactions. We used Greene networks as downloaded without any further modification.

*InWeb network:* we downloaded the InWeb protein-protein interaction network (v.20160912; see Web Resources). We re-format the network to follow the consistent format as other gene networks with three columns (Entrez ID 1, Entrez ID 2, edge weight) and used the confidence score between a pair of genes as the edge weight.

*Sonawane network:* we downloaded Sonawane networks for each of 38 tissues (see Web Resources). We transformed edge weights using the formula  $\ln(e^{\text{edge weight}} + 1)$  to avoid negative edge weights as stated in Sonawane et al.<sup>53</sup>, downweighting large negative edge weights. We intersected ‘net’ and ‘netTS’ to obtain edge weights of tissue-specific edges and concatenated with a set of edges to construct three-column matrix (TF, gene, weight). Then, we converted genes to Entrez gene ID.

## Network annotations

For each network, we constructed seven different probabilistic annotations based on the following network connectivities (centralities, i.e. how connected the gene is to other genes in the network), which we computed using Graph-tool (see Web Re-

sources):

Let  $G$  be a weighted, undirected graph with a set of vertices ( $V$ ) and a set of edges ( $E$ ).

*Closeness centrality ( $v$ ):* how close gene  $v$  is to all other genes in the network. It is defined as  $1/\sum_{v,v \neq w} d_{vw}$  where  $d_{vw}$  is the weighted distance from  $v$  to  $w$ . If there is no path between two genes, the distance of zero is used.

*Degree ( $v$ ):* the number of vertices connected to  $v$ ; that is the number of neighboring genes for each gene in the graph.

*Maximum edge weight ( $v$ ):* the maximum of the weights of all edges connected to  $v$ .

*Sum edge weight ( $v$ ):* the sum of the weights of all edges connected to  $v$ .

*Betweenness centrality ( $v$ ):* The number of shortest paths that pass through gene  $v$ . It is defined as  $\sum \sigma_{uw}(v)/\sigma_{uw}$  where  $\sigma_{uw}(v)$  is the number of shortest paths from node  $u$  to  $w$  that pass through  $v$ , and  $\sigma_{uw}$  is the total number of shortest paths from  $u$  to  $w$ .

*Eigenvector centrality ( $v$ ):* Intuitively, the eigencentality of  $v$  is proportional to the sum of the centralities of its neighbors<sup>89</sup>. It is defined as the solution of  $Ax = \lambda x$ , where  $x$  is the eigenvector of the weighted adjacency matrix  $A$  with the largest eigenvalue  $\lambda$ .

*Pagerank ( $v$ ):* Pagerank is similar to eigenvector centrality, except it contains a damping factor, the probability that the person who randomly visit genes will continue, under the assumption that more important genes are more likely to receive more visits. It is defined as  $(1 - d)/N + d \sum_{u \in N(v)} \text{Pagerank}(u)/d^+(u)$  where  $d$  is a damping factor,  $N(v)$  are the in-neighbors of  $v$ ,  $d^+(u)$  is the out-degree of  $u$ . Because  $G$  is an undirected graph, pagerank treats it as a directed graph by making edges bidirectional. After computing network connectivities for all genes that exist in a network, we linearly transformed scores to lie between 0 to 1 and annotated variants around genes (+/- 100kb). When a variant is spanned by multiple genes with the 100kb window, we assigned the maximum connectivity score. For each of 168 network-trait pairs (4 networks x 42 traits), we applied S-LDSC and assessed the statistical significance of the network annotation's  $\tau^*$  conditioned on other annotations. We also

computed the enrichment, which naturally extends to probabilistic annotations<sup>36</sup>. As a secondary analysis, we constructed network annotations with different window sizes (+/- 10kb and 1Mb).

## Comparison of closeness centrality to 18 gene sets

To compare closeness centrality to other metrics that quantify the biological importance of each gene, we considered 18 gene sets that reflect a broad range of gene essentiality metrics<sup>90</sup> (see Web Resources). The 18 gene sets are provided in Table S3 and briefly described below; the number of genes corresponds to protein-coding genes with an Entrez ID:

*All genes*: 19,031 genes with protein product according to HGNC<sup>91</sup> (HUGO Gene Nomenclature Committee) that have an Entrez ID.

*MGI essential genes*<sup>92-94</sup>: 2,371 genes for which homozygous knockout in mice results in pre-, peri- or post-natal lethality.

*Autosomal dominant genes*<sup>95,96</sup>: 698 genes among OMIM disease genes that are deemed to follow autosomal dominant inheritance.

*Haploinsufficient genes*<sup>97</sup>: 174 genes of severe, moderate, and mild haploinsufficiency, where having only a single functioning copy of a gene is not enough for normal function. *High pLI genes*<sup>97</sup>: 3,104 loss-of-function (LoF) genes with pLI > 0.9, i.e. strongly depleted for protein-truncating variants.

*High  $s_{het}$  genes*<sup>98</sup>: 2,853 constrained genes with  $s_{het} > 0.1$ , i.e. strong selection against protein-truncating variants.

*High Phi genes*<sup>99</sup>: 588 LoF-constrained genes with probability of haploinsufficiency (Phi) > 0.95.

*High missense Z genes*<sup>100</sup>: 1,440 constrained genes strongly depleted for missense mutations, with  $\text{exp\_syn} \geq 5$ ,  $\text{syn\_z\_sign} < 3.09$ , and  $\text{mis\_z\_sign} > 3.09$ , as retrieved in ref.<sup>97</sup>.

*ClinVar genes*<sup>101</sup>: 5,428 genes with a pathogenic or likely pathogenic variant with no conflict among studies.

*OMIM disease genes*<sup>102</sup>: 2,266 genes deposited in the Online Mendelian Inheritance

in Man (OMIM), as retrieved in ref<sup>103</sup>

*GWAS nearest genes*<sup>14</sup>: 6,271 genes nearest to peak GWAS significant loci (p-value  $\leq 5e-8$ ) in the GWAS Catalog.

*Transcription factors*<sup>104</sup>: 1,610 human transcription factors.

*DrugBank genes*<sup>105</sup>: 373 genes whose protein products are human targets of FDA-approved drugs with known mechanisms of action.

*High EDS genes*<sup>106</sup>: 2,664 genes among top 3,000 genes highly scored in enhancer domain score (EDS).

*Olfactory receptors*<sup>107</sup>: 369 olfactory receptor genes.

*eQTL-deficient genes*<sup>108</sup>: 604 genes with no significant variant-gene association in all 48 tissues in GTEx v.7 single-tissue cis-eQTL data.

*Genes with more independent SNPs*: 1,884 genes, defined as the top 10% of genes ranked based on the number of independent ( $r^2 < 0.1$ ) SNPs near the gene ( $\pm 100$ kb) relative to the length of the gene. We created PLINK files from 1000 Genomes<sup>40</sup> European Phase 3 reference genome individuals (see Web Resources) and SNPs with minor allele frequency (MAF)  $\geq 5\%$ . We filtered SNPs by applying LD-pruning to retain SNPs with  $r^2 < 0.1$  using PLINK (see Web Resources) with the window size of 50kp and the step size of 10kp. We excluded SNPs in the major histocompatibility complex (MHC) region, in our other analyses. We computed the number of pruned independent SNPs near each protein-coding gene ( $\pm 100$ kb) and divided by the length of the gene. We obtained highly correlated gene sets with different  $r^2$  thresholds ( $< 0.3$  and  $0.5$ ) and using a 10kb window.

*Genes with more SNPs*: 1,884 genes, defined as the top 10% of genes ranked based on the total number of SNPs near the gene ( $\pm 100$ kb) relative to the length of the gene. We obtained highly correlated gene sets using a 10kb window.

We extend the definition of excess overlap of annotations provided earlier defined earlier to a definition of excess overlap between gene sets. Let “gene set 1” denote one of the 18 gene sets defined above and “gene set 2” denote a given decile bin of



closeness centrality for a given network. We define:

$$\text{excess overlap}(\text{gene set 1, gene set2}) = P_d/P_{tot} \quad (1.6)$$

where  $P_d = \frac{|\text{gene set 1} \cap \text{gene set 2}|}{|\text{gene set 2}|}$  and  $P_{tot} = \frac{|\text{gene set 1} \cap \text{genes in network}|}{|\text{genes in network}|}$ . The standard error for the excess overlap is similarly scaled:

$$SE = \sqrt{\frac{P_d(1 - P_d)}{|\text{gene set 2}|}}/P_{tot} \quad (1.7)$$

When there is excess overlap, the excess fold overlap is  $> 1$ ; when there is depletion, the excess fold overlap is  $< 1$ . More generally, excess overlap can be computed for any pair of gene sets.

## **Assessing whether genes with high closeness centrality are heavily regulating other genes or are heavily regulated by other genes**

To assess whether genes with high closeness centrality are heavily regulating other genes or are heavily regulated by other genes, we performed three analyses. First, we assessed the excess overlap of 1,610 known human TFs<sup>104</sup> in the top decile of closeness centrality for each network.

Second, we used the ENCODE ChIP-Seq Significance Tool<sup>109</sup> (see Web Resources) to assess the excess overlap of TF binding sites in the promoters of high closeness centrality genes (top 10% of genes for each gene network analyzed). We considered 220 TFs from the union of all 91 available ENCODE<sup>66</sup> cell lines. We used all protein-coding genes with an Entrez ID as background regions and defined promoters based on  $\pm 1$  kb from the TSS. We report TFs that have significant excess overlap with the promoters of high closeness centrality genes (Benjamini-Hochberg adjusted p-value  $< 0.05$ ) and calculated the significance of the excess overlap using the hypergeometric test. In addition, for each network, we computed excess overlap of the ENCODE TF

binding sites annotations<sup>66</sup> (from the baseline-LD model) in 10 deciles of closeness centrality (+/- 100kb).

Third, we used DAVID<sup>110</sup> (see Web Resources) to assess GO enrichments of high closeness centrality genes (top decile) for each of four gene networks. We used all protein-coding genes with an Entrez ID as a background set. We considered three GO categories: biological process, cellular component, and molecular function. We reported significant GO terms (Benjamini-Hochberg adjusted p-value < 0.05).

## Assessing the impact of noise in gene networks

We performed three analyses to assess the impact of noise in gene networks on our results. First, we performed a network perturbation analysis by randomly removing a subset of edges (from 10% to 90% with 10% increment) using an established protocol<sup>111</sup>. We performed five separate perturbation analyses for each value of the proportion of edges removed. For each networks with edges removed, we computed network connectivity metrics and applied S-LDSC to estimate disease heritability enrichment and  $\tau^*$ .

Second, we applied the diffusion state distance algorithm (DSD<sup>34</sup>; see Web Resources) to de-noise networks. (We note that the DSD algorithm has been shown to be effective in de-noising networks; see ref.<sup>112</sup> for performance in application to networks from Greene et al. networks<sup>50</sup>.) We ran DSD using default parameters. Because smaller values in the DSD output correspond to stronger interactions, we used the inverse of DSD output as new edge weights. We computed network connectivity metrics on the transformed gene networks and applied S-LDSC to estimate disease heritability enrichment and  $\tau^*$ .

Third, we constructed two consensus networks by intersecting edges in either (i) the Greene and InWeb networks or (ii) the Greene, InWeb, and Sonawane networks. We note that the intersection of all four networks did not contain any edges, as the Saha network is very sparse. We computed consensus networks both including and excluding the Sonawane network, which contains nonzero edge weights only for pairs

of genes in which at least one of the genes is a known TF. For tissue-specific networks, we used the same tissue as in our primary analysis (Greene thyroid and Sonawane testis). For edges weights, we used the mean edge weight. As each consensus network contained disjoint connected components, we computed closeness centrality using the connected component with the largest number of edges. We applied S-LDSC to estimate disease heritability enrichment and  $\tau^*$ .

## Correlation between closeness and gene expression

To assess the correlation between closeness centrality and gene expression, we used gene expression (TPM) across 53 tissues from GTEx RNA sequencing data<sup>108</sup> (see Web Resources), restricted to protein-coding genes that are present in a given network. We computed the Pearson correlation between a vectors of gene expression and closeness centrality, for each of 53 GTEx tissues.

## Pathway + Network annotations

In our pathway+network analyses, we considered 156 significant pathway-trait pairs (from our pathway analyses) and four gene networks for a total of 590 (pathway-trait, network) pairs (122 pathway-trait pairs for Saha + 156 pathway-trait pairs x 3 other networks). We constructed pathway+network annotations, specific to an input pathway and a gene network. For each of 141 enriched pathways, we first constructed an adjacency matrix by mapping genes in a pathway to a gene network and identifying the set of neighboring genes outside the pathway, using Graph-tool (see Web Resources). In the adjacency matrix representing a pathway-specific subnetwork, a vector represents an Entrez ID of a gene in a pathway, an Entrez ID of a neighboring gene, and the interaction score (e.g. posterior probability of these two genes). We computed closeness centrality, using the inverse of the interaction score as the cost. Using the set of genes with the linearly transformed closeness scores that lie from 0 to 1, we annotated variants around genes (+/- 100kb). When a variant is spanned by multiple genes within the 100kb window, we assigned the maximum closeness score.

For each of 590 (pathway-trait, network) pairs, we applied S-LDSC and assessed the statistical significance of the pathway+network annotation’s  $\tau^*$  conditioned on other annotations.

## **Network connectivity of a pathway**

Given a pathway and a network, we quantified how tightly connected the pathway is in the network, compared to a null pathway. Let P be a set of genes in the pathway that exist in the network. Let Q be a set of neighboring genes; that is, genes in P are connected with genes in Q with at least one edge, in the network. We calculated three metrics: (1) size of Q (number of neighboring genes), (2) sum of edges among genes in P, and (3) sum of edges between genes in P and genes in Q. We did not consider genes in the pathway that are not coding or do not appear in the network.

For each of 141 enriched pathways, we constructed a corresponding null pathway as follows: (1) randomly choose a pathway from the full set of 18,119 pathways (as shown in Table S1), (2) randomly choose a gene from the sampled pathway, and (3) repeat (1) and (2) N times, where N is the number of genes in the pathway. For each of four gene networks analyzed, we repeated this procedure 10,000 times and reported three connectivity metrics of enriched and null pathways (for null, the mean of 10,000 permutations is reported).

## **Set of 42 independent traits**

Analogous to a previous study<sup>36</sup>, we considered 89 GWAS summary association statistics, including 34 traits from publicly available sources and 55 traits from the UK Biobank (up to  $N = 459K$ ); summary association statistics were computed using BOLT-LMM v2.3<sup>37,38</sup>. Among 47 summary statistics with z-scores of total SNP heritability of at least 6 (computed using S-LDSC with the baseline-LD model), we further removed 5 summary statistics that have genetic correlation of at least 0.9 (computed using cross-trait LDSC<sup>39</sup>). Whenever applicable, meta-analysis across 42 independent traits (Table S2), whose GWAS summary statistics are publicly available

(see Web Resources), was performed using a random-effect meta-analysis using the R package *rmeta*.

## Set of 10 blood-related and 8 brain-related traits

We analyzed ten independent blood- and autoimmune-related traits: Crohn’s disease<sup>26</sup>, rheumatoid arthritis<sup>27</sup>, and ulcerative colitis<sup>26</sup> from publicly available datasets and eczema, autoimmune diseases. eosinophil count, platelet count, red blood cell count, white blood cell count, and red blood cell distribution width from the UK Biobank.

We analyzed eight independent brain-related traits: autism spectrum<sup>113</sup>, depressive symptoms<sup>23</sup>, and schizophrenia<sup>114</sup> from publicly available datasets, and age at menarche, body mass index, neuroticism, smoking status, and years of education from the UK Biobank. We selected these traits from 42 independent traits we analyzed, as inferred from heritability enrichment of tissue-specific gene expression and chromatin annotations<sup>2,4</sup> and eQTL annotations<sup>36</sup>. We additionally considered autism spectrum based on previous studies<sup>115,116</sup> that the brains of subjects with autism have altered expression.

## 1.5 Results

### Enrichment of disease heritability in pathway annotations

We sought to identify pathways that are enriched for disease heritability. We applied S-LDSC to 760,869 pathway-trait pairs, spanning 18,119 pathways from five sources (Table A.1 and Figure B-1) and 42 independent diseases and complex traits, including 30 UK Biobank traits (average  $N = 323K$ ; Table A.2). We identified 156 pathway-trait pairs that were significantly enriched after conditioning on the baseline-LD model and the all-genes annotation (FDR<5% for positive  $\tau^*$ ). The 156 pathway-trait pairs spanned 141 pathways and 34 traits, implying that most pathway enrichments are trait-specific. Complete results for three representative traits<sup>28</sup> (Crohn’s disease (IBD

[MIM: 266600]), rheumatoid arthritis (RA [MIM: 180300]), and schizophrenia (SCZ [MIM: 181500])) are reported in Figure 2-1, and complete results for all traits are reported in Table A.4. The top pathway (most significant  $\tau^*$ ) for each of the 34 traits is reported in Table 2.1, and the complete set of 156 significant pathway-trait pairs is reported in Table A.5. Genes in the 141 enriched pathways had a larger gene size (92kb on average) compared to all protein-coding genes (58kb) and genes in all pathways (76kb) (see Table A.6). We meta-analyzed the 156 pathway-trait pairs using random-effect meta-analysis (analogous to previous work<sup>2,3,36</sup>). Both the enrichment (4.13, s.e. = 0.12;  $P = 4.74e-158$ ) and  $\tau^*$  (0.15, s.e. = 0.0061;  $P=3.41e-131$ ) were large and highly statistically significant (Table A.7); we caution that these p-values are slightly inflated because we meta-analyzed across significant pathway-trait pairs only. When we repeated our analysis excluding from each associated pathway all genes harboring genome-wide significant associations for the corresponding trait (average of 2 genes removed per pathway-trait pair; see Material and Methods), only 53 pathway-trait pairs remained significant (Table A.8). For each of 156 pathway-trait pairs, genes after excluding GWAS significant genes are provided in Table A.5.

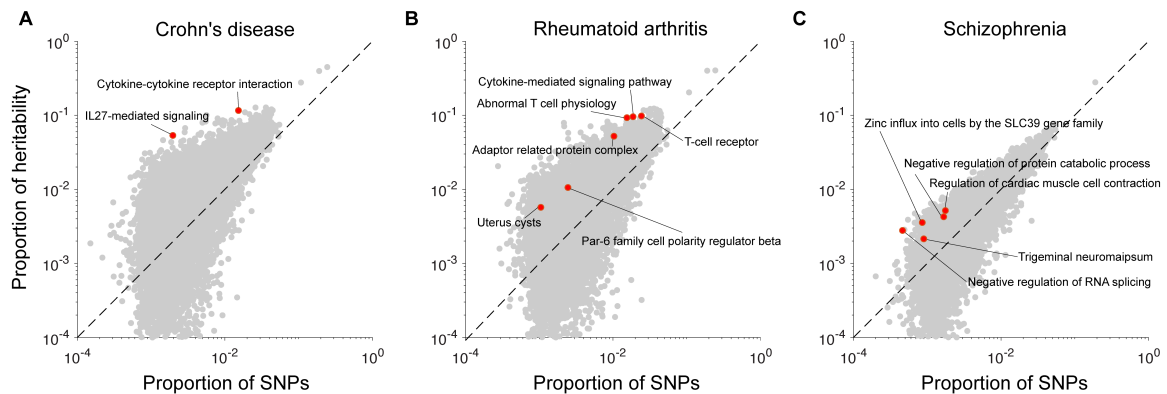


Figure 1-1: **Enriched pathways for three representative traits.** For (A) Crohn’s disease, (B) rheumatoid arthritis and (C) schizophrenia, we report the proportion of heritability explained and proportion of SNPs for each of 18,119 pathways analyzed. Red points indicate significantly enriched pathways ( $FDR < 5\%$ ), and grey points indicate non-significant pathways. Numerical results for all 42 diseases and complex traits are reported in Table A.4.

Our results include 8 pathway-trait pairs reported in previous genetic studies (see Table 2.1 and Table A.5). These include “immune response” for both Crohn’s dis-

Trait	Top enriched pathway	Database	# Genes	Enr. (s.e.)	$\tau^*$ (s.e.)
Age at menarche	SIX homeobox 6	InWeb	11	3.06 (0.45)	0.06 (0.02)
Auto immune traits	Mismatch repair directed by MSH2:MSH6	BS	14	8.53 (2.05)	0.28 (0.07)
BMI	Positive regulation of synapse maturation	GO	11	3.34 (0.25)	0.07 (0.01)
Dermatologic diseases	Aldo-keto reductase family 1 member E2	InWeb	12	7.51 (1.80)	0.19 (0.05)
Eczema	Th17 cell differentiation	BS	102	9.01 (1.56)	0.64 (0.15)
Eosinophil count	Jak-STAT signaling pathway	BS	152	7.67 (1.11)	0.46 (0.11)
Forced vital capacity	Absent acrosome	MGI	11	2.21 (0.31)	0.04 (0.01)
Heel T Score	Skeletal system development	GO	142	4.43 (0.60)	0.31 (0.07)
Height	Abnormal trabecular bone morphology	MGI	172	2.93 (0.39)	0.20 (0.05)
Hypothyroidism	Cytokine-mediated signaling pathway	GO	235	4.88 (0.78)	0.42 (0.11)
Morning person	N-acetylglucosamine metabolic process	GO	13	2.07 (0.26)	0.03 (0.01)
Neuroticism	HSF1 activation	BS	11	2.54 (0.27)	0.04 (0.01)
Platelet count	Platelet activation	GO	208	4.32 (0.62)	0.30 (0.09)
Red blood cell count	Exogenous drug catabolic process	GO	10	3.62 (0.63)	0.08 (0.02)
Red blood cell distribution width	Decreased erythrocyte cell number	MGI	206	4.61 (0.74)	0.39 (0.11)
Respiratory/ear-nose-throat diseases	Abnormal T cell activation	MGI	108	4.28 (0.75)	0.26 (0.07)
Skin color	Melanin biosynthetic process	GO	11	192.14 (55.37)	5.90 (1.73)
Smoking status	Fibromodulin	InWeb	17	2.53 (0.48)	0.06 (0.02)
Systolic blood pressure	cGMP-PKG signaling pathway	BS	164	3.41 (0.49)	0.25 (0.07)
Type 2 diabetes	Glucuronidation	BS	30	2.74 (0.32)	0.06 (0.01)
Waist-hip ratio	Negative regulation of transcription	GO	496	2.60 (0.22)	0.19 (0.05)
White blood cell count	ERK cascade	PC	15	4.65 (0.69)	0.11 (0.03)
Years of education	Absent corpus callosum	MGI	45	2.50 (0.33)	0.08 (0.02)
Age first birth	Receptor signaling protein tyrosine kinase	GO	10	3.62 (0.58)	0.10 (0.02)
Anorexia	Activation of the AP-1 family of TF	BS	10	5.64 (1.04)	0.15 (0.03)
Autism spectrum	GABA-A receptor activity	GO	15	4.01 (0.73)	0.12 (0.02)
Coronary artery disease	Sodium-independent organic anion transport	GO	12	7.05 (1.49)	0.22 (0.04)
Crohn's disease	Cytokine-cytokine receptor interaction	BS	266	7.59 (1.50)	0.65 (0.18)
Depressive symptoms	Glycoprotein metabolic process	GO	12	5.22 (1.10)	0.15 (0.04)
LDL	Increased erythroblast number	MGI	21	5.63 (1.15)	0.18 (0.05)
Number children even born	Lipoxygenase pathway	GO	12	9.81 (1.59)	0.27 (0.04)
Rheumatoid arthritis	Par-6 family cell polarity regulator beta	InWeb	16	4.31 (1.01)	0.17 (0.04)
Schizophrenia	Trigeminal neuroma	MGI	10	2.39 (0.31)	0.03 (0.01)
Ulcerative colitis	FGFR2b ligand binding and activation	BS	10	8.00 (1.97)	0.18 (0.05)

Table 1.1: **Top enriched pathway for each trait.** We report the top enriched pathway (most significant  $\tau^*$ ) for each of 34 traits with at least one significantly enriched pathway. The first 23 traits (above the line) are UK Biobank traits. Enrichment of the "melanin biosynthetic process" pathway for skin color is consistent with previous studies<sup>21,22</sup>, and enrichment of the "absent corpus callosum" pathway for years of education was reported in a previous genetic study<sup>23</sup>. The complete set of 156 significant pathway-trait pairs is reported in Table A.5.

ease and ulcerative colitis<sup>26</sup>; “T-cell receptor”, “cytokine-mediated signaling pathway”, and “abnormal T-cell physiology” for rheumatoid arthritis<sup>27</sup>; and “Absent corpus callosum” for years of education<sup>23</sup>. In addition, “melanin biosynthetic process” was overwhelmingly enriched for skin color (Table 2.1), consistent with the fact that genetic variants in constituent genes are strongly associated with skin pigmentation (e.g. *MC1R*) and other pigmentation traits (e.g. *TYR* [MIM: 606933], *OCA2* [MIM: 611409], *SLC45A2* [MIM: 606202])<sup>21,22</sup>.

Surprisingly, most pathway-trait pairs reported in recent studies<sup>28–33</sup> using genome-wide polygenic methods<sup>2,33,68</sup> were not significant in our analysis. Specifically, we considered 95 significant pathway-trait pairs for the 6 traits (schizophrenia, Crohn’s disease, rheumatoid arthritis, neuroticism, intelligence, depressive symptoms), analyzed in five previous studies<sup>28–32</sup>, restricting to at most the top 20 pathways per trait per study. We assessed the significance of these 95 pathway-trait pairs in our analysis, based on global FDR < 5% across 18,119 pathways tested ( $\tau^* < 0.000989$ ). Only 15/95 pathway-trait pairs were significant in our primary analysis, after conditioning on the baseline-LD model and all-genes annotation (Table A.9A). However, 67/95 were significant when we repeated the S-LDSC analysis conditioning on just the all-genes annotation and not the baseline-LD model (Table A.9B). We obtained similar results for a pathway-trait pairs reported in a very recent study<sup>33</sup> (see Web Resources; Table A.9A-B). Enriched pathways that were fully explained by the baseline-LD model could potentially be due to factors that do not play a direct role in trait biology<sup>117</sup>, although we caution that our analyses do not resolve which factors are causal (see Discussion).

Our results also highlight pathway-trait pairs that have not previously been reported but are consistent with known biology. These include “GABA-A receptor activity” with autism (Table 2.1), consistent with the finding that the brains of subjects with autism have altered expression of GABA receptors<sup>115,116</sup>; “Oncostatin M (OSM [MIM: 165095])” for ulcerative colitis (Table A.5), consistent with the finding that inflamed intestinal tissues from patients with inflammatory bowel diseases contained higher expression of OSM, and OSM-deficient mice displayed significantly attenuated



colitis<sup>118</sup>; and "zinc influx into cells by the SLC39 gene family" with schizophrenia (Table A.5), consistent with the finding that *SLC39A12A* expression in dorsolateral prefrontal cortex is associated with schizophrenia<sup>119,120</sup>.

We also analyzed three additional gene sets (distinct from the 18,119 pathways) reflecting genes under strong selection: ExAC<sup>97</sup> (high pLI genes; genes strongly depleted for protein-truncating variants), Cassa<sup>98</sup> (high  $s_{\text{het}}$  genes; genes with strong selection against protein-truncating variants), and Samocha<sup>100</sup> (high missense Z scores genes; genes strongly depleted for missense mutations), which were previously shown to be enriched for heritability in a meta-analysis across traits<sup>36</sup>. We identified 13 significantly enriched gene set-trait pairs (7 for ExAC, 2 for Cassa, and 4 for Samocha) spanning 9 traits, after conditioning on the baseline-LD model and the all-genes annotation (FDR < 5% for positive  $\tau^*$ ; Table A.10). In a meta-analysis of these 13 gene set-trait pairs, both the enrichment (1.57, s.e. = 0.053; P = 7.52e-33) and  $\tau^*$  (0.13, s.e. = 0.0090; P=3.30e-48) were highly statistically significant (Table A.7).

## Enrichment of disease heritability in network annotations

We sought to assess the hypothesis that genes with high network connectivity are enriched for disease heritability. We constructed probabilistic annotations based on closeness centrality for each of the Saha, Greene, InWeb, and Sonawane networks (see Material and Methods). For three networks that include tissue-specific networks, we selected the Saha-skin(sun exposed lower leg), Greene-thyroid and Sonawane-testis networks for our primary analyses, as these tissue-specific networks maximized the correlation of the resulting network annotations with H3K27ac (Table A.11); we also considered other criteria for selecting tissue-specific networks (see below). We determined that closeness centrality was independent of gene size ( $r = -0.015$  to  $0.019$ ) and exon proportion ( $r = -0.18$  to  $0.008$ ; Figure B-2). We note that different tissue-specific networks from the same source were only weakly correlated ( $r = 0.027$  to  $0.076$  for Saha,  $0.17$ - $0.21$  for Greene,  $0.062$ - $0.24$  for Sonawane; see Table A.11). The number of genes, number of edges, and distribution of edge weights for each network are reported in Table A.12. The Greene network is very dense (25,825 genes with

mean degree 6,484), the Saha network is very sparse (2,381 genes with mean degree 7.9), and the InWeb and Sonawane networks have intermediate density.

To compare closeness centrality to other metrics that quantify the biological importance of each gene, we computed the excess overlap between genes in each decile bin of closeness centrality (for each network) and 18 other gene sets (Table A.3; see Material and Methods). We first describe results for the Greene network. We determined that high closeness centrality genes (top decile of closeness centrality) had high excess overlap with gene sets defined by constraint metrics (e.g. high pLI (ExAC) genes<sup>97</sup>, high  $s_{\text{het}}$  (Cassa) genes<sup>98</sup>) and essentiality metrics (e.g. MGI essential genes<sup>92-94</sup>) as compared to other decile bins (Figure 2-2A and Table A.13), consistent with larger values of closeness centrality in these gene sets (Figure 2-2B and Table A.14;  $r = 0.22$  to  $0.27$  between closeness centrality and gene set, Table A.15). On the other hand, high closeness centrality genes were highly depleted for olfactory receptor genes and eQTL-deficient genes (Figure 2-2A and Table A.13), consistent with smaller values of closeness centrality in these gene sets (Figure 2-2B and Table A.14;  $r = -0.19$  to  $-0.21$ , Table A.15). High closeness centrality genes were also depleted for genes with more independent SNPs and more total SNPs relative to the length of the gene (which are less functionally important as they are less likely to be under negative selection<sup>121</sup>) (Figure S3 and Table A.13), consistent with smaller values of closeness centrality in these gene sets (Figure S4 and Table A.14;  $r = -0.208$  to  $-0.211$ , Table A.15); however, this was not the case for the Saha, InWeb and Sonawane networks (Figure B-3, Figure B-4, Table A.13, Table A.14;  $r = -0.047$  to  $0.037$ , Table A.15). For other gene sets, results for the InWeb network, but not Saha or Sonawane networks, were generally similar to those for the Greene network (Figure B-3, Figure B-4, Table A.13, Table A.14, Table A.15).

We computed the excess overlap between each of the four network annotations (probabilistic annotations based on closeness centrality) and representative annotations from the baseline-LD model and determined that each of the network annotations had substantial excess overlap with regulatory annotations (e.g. 1.23-1.36 with H3K27ac, 1.35-1.66 with H3K9ac; Figure 2-3A and Table A.16), significantly stronger

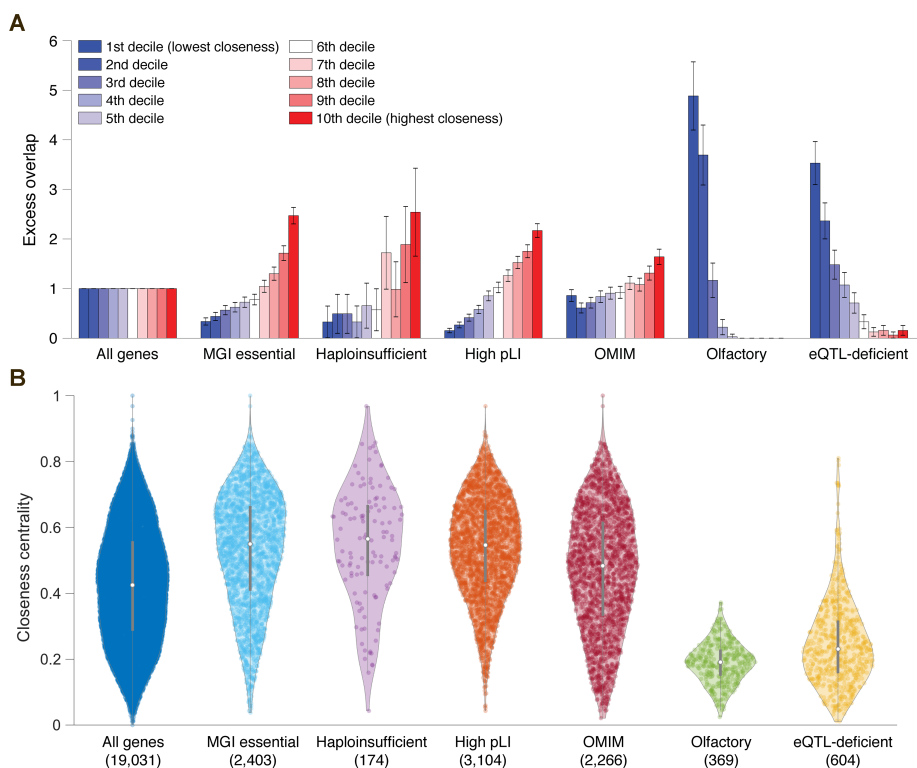


Figure 1-2: **Comparison of closeness centrality to other metrics that quantify the biological importance of each gene.** (A) For each of 7 gene sets, we report the excess overlap of genes in each decile bin of closeness centrality for the Greene thyroid network. Error bars represent 95% confidence intervals. (B) For each of 7 gene sets, we report the distribution of closeness centrality for the Greene thyroid network. Colored dots denote genes, white dots denote medians, and grey lines denote box plots. Numbers in parentheses below each gene set denote the number of genes. Results for all four networks and all 18 gene sets analyzed are reported in Figure B-3 and Table A.13 (for panel A) and Figure B-4 and A.14 (for panel B). Lists of genes for each of the 18 gene sets are provided in A.3.

than the excess overlap between the all-genes annotation and regulatory annotations (e.g. 1.11 with H3K27ac, 1.14 with H3K9ac; Figure 2-3A); this implies that, for each network, genes with high connectivity to other genes in the network are enriched for the presence of nearby regulatory marks, perhaps because they are regulated by many other genes (see below). Correlations between these annotations produced similar conclusions (Table A.17; larger correlation for H3K27ac than H3K9ac), although we consider excess overlap to be a more robust metric because the size of the network annotations varies from 1%-34% of SNPs. We also report excess overlap and correlations between baseline-LD model annotations in Table A.18.

To investigate whether genes with high closeness centrality are heavily regulated by other genes, or are heavily regulating other genes, we performed three analyses; we did not reach a consistent conclusion. First, we assessed the excess overlap of 1,610 known human TFs<sup>104</sup> in the top decile of closeness centrality for each network (see Material and Methods). We determined that TFs (which regulate other genes) were depleted in high closeness centrality genes for the Greene and Sonawane networks, but enriched in high closeness centrality genes for the Saha and InWeb networks (Figure B-3). Thus, this analysis did not reach a consistent conclusion. Second, we assessed the excess overlap of TF binding sites in the promoters ( $\pm$  1kb from TSS) of high closeness centrality genes, by applying the ENCODE ChIP-seq significance tool<sup>109</sup> to ENCODE ChIP-seq data<sup>66</sup> spanning 220 TFs and 91 cell lines (see Material and Methods). We determined that, for each gene network, binding sites for the majority of TFs (132 to 206 out of 220 TFs) had significant excess overlap (FDR < 0.05) in the promoters of high closeness centrality genes (Table A.19); we also observed significant excess overlap throughout high closeness centrality genes ( $\pm$  100kb) (Table A.20). Thus, this analysis supports the hypothesis that genes with high closeness centrality are heavily regulated by other genes. Third, we assessed the GO enrichment of high closeness centrality genes using DAVID<sup>110</sup> (see Material and Methods). We used all protein-coding genes as a background set, and evaluated gene sets corresponding to three GO categories: biological process, cellular component, and molecular function. We determined that high closeness centrality genes from all four networks were often

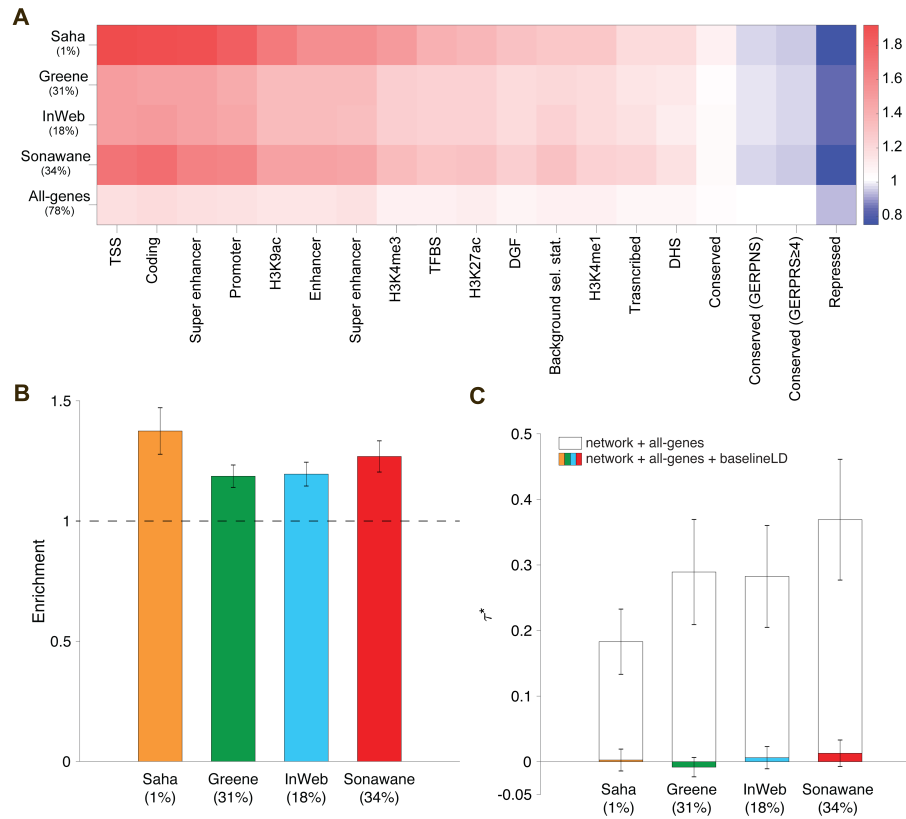


Figure 1-3: **Heritability enrichment of network annotations.** We report (A) excess (fold) overlap between network annotations and baseline-LD functional categories; (B) heritability enrichment of network annotations, meta-analyzed across 42 independent traits; and (C)  $\tau^*$  values of network annotations, conditioned on either just the all-genes annotation, or the all-genes annotation and the baseline-LD model, meta-analyzed across 42 independent traits. In (A-C), the percentage under each bar indicates the proportion of SNPs in each annotation (defined for probabilistic annotations as the average value of the annotation), and error bars represent 95% confidence intervals. Numerical results for (A) are reported in Table A.16, and numerical results for (B) and (C) are reported in A.22. The S-LDSC results for the complete set of 168 network-trait pairs are reported in A.22.

significantly enriched ( $FDR < 0.05$ ) in protein/TF/DNA/regulatory region binding (e.g. “transcription regulatory region DNA binding (GO:0044212)”, “transcription factor binding (GO:0008134)”) and transcriptional regulation (e.g. “positive regulation of transcription (GO:0045944)”) gene sets (Table A.21). Thus, this analysis supports the hypothesis that genes with high closeness centrality are heavily regulating other genes. Overall, our analyses did not reach a consistent conclusion on the regulatory role of high closeness centrality genes.

For each of the four network annotations, we applied S-LDSC to the 42 independent diseases and complex traits, conditioning on the baseline-LD model and the all-genes annotation, and meta-analyzed the results across traits using random-effects meta-analysis. We identified strongly significant enrichments for each network annotation: 1.19 (s.e. = 0.024;  $P = 2.5e-30$ ) to 1.37 (s.e. = 0.049;  $P = 1.6e-17$ ) (Figure 2-3B and Table A.22). However, estimates of  $\tau^*$ , quantifying effects unique to the network annotations, were not significant ( $P = 0.21$  to  $0.77$ ) (Figure 2-3C and Table A.22). This implies that the enrichment signal in the network annotations (Figure 2-3B) is entirely explained by the excess overlap between the network and baseline-LD model annotations (Figure 2-3A); accordingly, when we repeated the S-LDSC analysis conditioning only on the all-genes annotation and not on the baseline-LD model,  $\tau^*$  estimates were large and highly significant (Figure 2-3C). We repeated the S-LDSC analysis conditional on one annotation from the baseline-LD model at a time and confirmed that regulatory annotations (primarily histone marks and transcription factor binding sites) reduced estimates of  $\tau^*$  by 74-100% (Table A.23). We note that other annotations did not have this effect (e.g. 16-20% for conservation annotations). We also partitioned all genes in a given network into 10 deciles of closeness centrality and applied S-LDSC to estimate the heritability enrichments for each decile. For the Greene, InWeb and Sonewane networks, we determined that the top decile of closeness centrality had the highest enrichment (1.68 (s.e. 0.062) - 1.76 (s.e. 0.070)), and the bottom decile of closeness centrality had the lowest enrichment (0.92 (s.e. 0.031) - 1.24 (0.038); see Table A.24). For the Saha network, the top decile had the highest enrichment (1.98 (s.e. 0.099)), but the bottom decile had the fourth lowest

enrichment (1.464 (s.e. 0.071)). We also observed that the top decile had significantly larger correlations with regulatory annotations than the bottom decile (Table A.20). On the one hand, these findings represent a negative result for efforts to improve upon the baseline-LD model. On the other hand, these findings provide a strong validation of the baseline-LD model, in that the information about diseases/traits in annotations from other sources that broadly reflect the action of gene regulation are fully captured by the baseline-LD model.

We performed three analyses to assess the impact of noise in gene networks on our results. First, we performed a network perturbation analysis by repeating the S-LDSC analysis on networks with 10-90% of edges randomly removed, following an established protocol<sup>111</sup> (see Material and Methods). We obtained similar results, including significant enrichments and non-significant  $\tau^*$  (Table A.25). Second, we de-noised each network by applying the diffusion state distance (DSD) algorithm<sup>34</sup>, computed closeness centrality on the transformed networks (see Material and Methods), and repeated the S-LDSC analysis. We determined that the DSD-transformed network annotations were 1.9-25.4% more enriched, but  $\tau^*$  remained non-significant (Table A.26). Third, we constructed two “consensus” networks by intersecting the edges in either (i) the Greene and InWeb networks or (ii) the Greene, InWeb, and Sonawane networks (the intersection of all four networks did not contain any edges, as the Saha network is very sparse), and repeated the S-LDSC analysis. We obtained similar results, including significant enrichments and non-significant  $\tau^*$  (Table A.27). These analyses support the robustness of our results.

We performed several secondary analyses. First, we examined the correlation between closeness centrality (for each network) and gene expression, using GTEx RNA-seq data<sup>108</sup> (see Material and Methods). We determined that closeness centrality in the Saha and Sonawane networks was independent of gene expression ( $r = -0.045$  to  $0.030$ ) for each of 53 GTEx tissues), whereas closeness centrality in the Greene and InWeb networks was moderately correlated with gene expression, particularly in brain tissues (e.g.  $r = 0.17$  to  $0.27$  for cerebellum) (Table A.28). Second, we performed heritability analyses on gene sets defined by membership in gene networks

(analogous to pathway annotations). We focused our analyses on the Saha networks, which are much sparser than the other networks. We defined gene sets using the 16 Saha transcriptome-wide networks (TWN) used in our primary analyses, as well as the 36 Saha tissue-specific networks (TSN), and applied S-LDSC to these gene sets. All 16 gene sets derived from TWNs and 4 of 36 gene sets derived from TSNs were significantly enriched for disease/trait heritability, but none had a significant  $\tau^*$  conditional on the baseline-LD model (Table A.29). Third, we repeated the S-LDSC analysis for the 3 tissue-specific networks (Saha, Greene, Sonawane) using the tissue that maximized the excess overlap of the High pLI (ExAC) gene set with the top decile of closeness centrality of the tissue-specific network (Table A.30; see Material and Methods). We obtained similar results (as compared to our primary analysis in which tissue-specific networks were selected so as to maximize correlation with H3K27ac), including significant enrichments for all three tissue-specific networks but non-significant  $\tau^*$  conditional on the baseline-LD model (Table A.22 and Table A.24). Fourth, for 10 blood-related traits and 8 brain-related traits, we repeated the S-LDSC analysis for the 3 tissue-specific networks using the most biologically relevant tissue (as inferred from heritability enrichment of tissue-specific specifically expressed gene and chromatin annotations<sup>4</sup>; see Table A.2 for list of traits and tissues). For the 10 blood-related traits, the resulting tissue-specific networks (blood) produced slightly larger enrichments (significant enrichment for all 3 networks; significant difference for Sonawane only), but  $\tau^*$  conditional on the baseline-LD model remained non-significant (Table A.31; consistent with Table A.22). For the 8 brain-related traits, both the resulting tissue-specific networks (brain) and the default tissue-specific networks produced non-significant enrichments (Table A.31; consistent with Table A.22). Finally, we also obtained similar results when we repeated the S-LDSC analysis using six connectivity metrics other than closeness centrality (see Material and Methods; Table A.32), using window sizes other than 100kb (10kb and 1Mb; Table A.33).



## Enrichment of disease heritability in integrated pathway+network annotations

We sought to assess the hypothesis that genes with high network connectivity to genes in enriched pathways are also enriched for disease heritability; these enriched pathways serve as relevant, disease-specific signals that network analysis often lacks. For each of four networks and 141 enriched pathways, we constructed probabilistic pathway+network annotations based on closeness centrality within the subnetwork consisting of input genes and their one-degree neighbors. As most of the pathway enrichments that we identified are disease-specific, the pathway+network annotations are expected to harbor disease-specific signals. As before, for three networks that include tissue-specific networks, we selected the Saha-skin(sun exposed lower leg), Greene-thyroid and Sonawane-testis networks for our primary analyses.

For each network, we computed the excess overlap between the 141 pathway+network annotations and representative annotations from the baseline-LD model and averaged results across the 141 pathways. We observed higher excess overlap with regulatory annotations (e.g. 1.27-1.35 with H3K27ac, 1.42-1.70 with H3K9ac; Figure 1-4A and Table A.16) than the analogous excess overlaps for network annotations (Figure 2-3A); this implies that, for each network, genes with high network connectivity to genes in input pathways are enriched for the presence of nearby regulatory marks, perhaps because they are regulated by many other genes. Correlations between network, pathway+network and baseline-LD model annotations are reported in Table A.17. We determined that pathway+network annotations were often correlated with network annotations, particularly for the corresponding network (from  $r = 0.12$  for Saha to  $r = 0.93$  for Greene). We further determined that pathway+network annotations constructed using the Greene, InWeb and Sonawane networks were moderately correlated ( $r = 0.27$  to  $0.42$ ), whereas pathway+network annotations constructed using the Saha network were more distinct ( $r = 0.05$ - $0.06$ ), primarily due to its small annotation size (Table A.17).

For each of 590 (pathway-trait,network) pairs (122 pathway-trait pairs for Saha

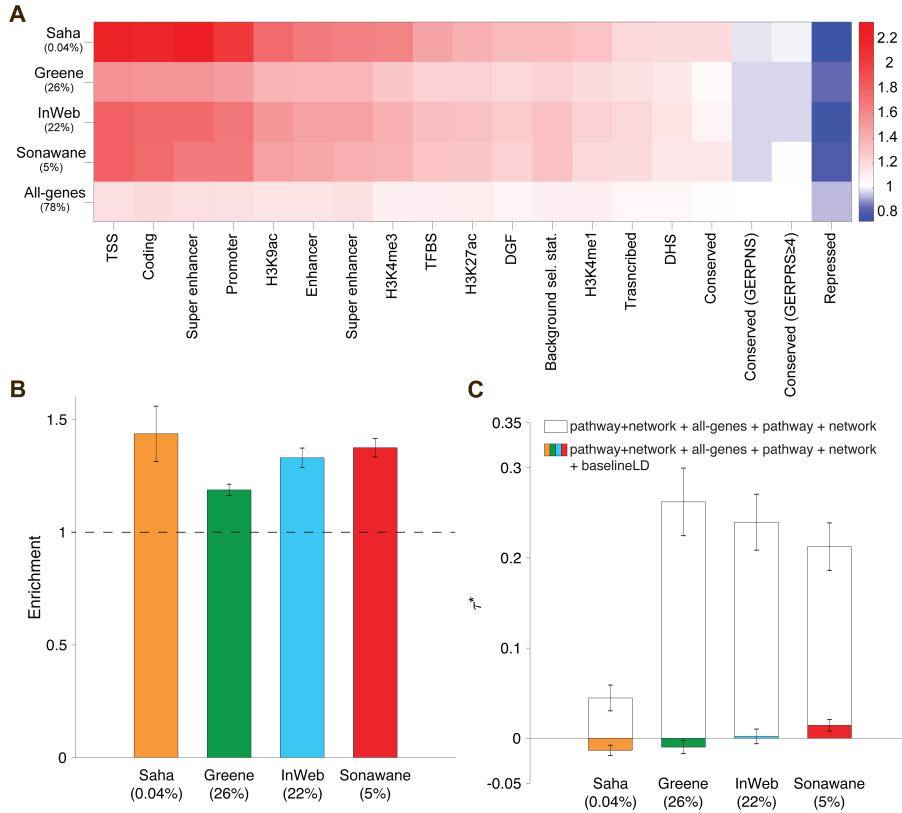


Figure 1-4: **Heritability enrichment of pathway+network annotations.** We report (A) excess (fold) overlap between pathway+network annotations (averaged across up to 156 pathway-trait pairs); (B) heritability enrichment of pathway+network annotations, meta-analyzed across up to 156 pathway-trait pairs; and (C)  $\tau^*$  values of pathway+network annotations, conditioned on either just the all-genes annotation and the corresponding pathway and network annotations, or the baseline-LD model as well, meta-analyzed across up to 156 pathway-trait pairs. In (A-C), the percentage under each bar indicates the proportion of SNPs in each annotation (defined for probabilistic annotations as the average value of the annotation), and error bars represent 95% confidence intervals. Numerical results for (A) are reported in A.16, and numerical results for (B) and (C) are reported in Table A.34. The S-LDSC results for the complete set of 590 pathway-trait pairs are reported in A.34.

+ 156 pathway-trait pairs x 3 other networks; see Material and Methods), we applied S-LDSC to the resulting pathway+network annotation and the corresponding trait, conditioning on the baseline-LD model, the all-genes annotation, and the corresponding pathway and network annotations, and meta-analyzed the results for each network using random-effects meta-analysis. We identified strongly significant enrichments for all of our pathway+network annotations: 1.19 (s.e. = 0.01; P = 1.5e-49) to 1.44 (s.e. = 0.06; P = 3.8e-12) (Figure 1-4B and Table A.34). On average, the pathway+network annotations most enriched for trait heritability are those derived from the Saha network (Figure 1-4B). However, estimates of  $\tau^*$ , quantifying effects unique to the network annotations, were not significant or only weakly significant (P = 0.62 to 4.5e-6) (Figure 1-4C and Table A.34). Once again, this implies that the enrichment signal in the pathway+network annotations is entirely explained by the excess overlap between the pathway+network and baseline-LD model annotations; accordingly, when we repeated the S-LDSC analysis conditioning only on the all-genes annotation and the corresponding pathway and network annotations, and not on the baseline-LD model,  $\tau^*$  estimates were large and highly significant, except for the Saha network (Figure 1-4C). We repeated the S-LDSC analysis conditional on one annotation from the baseline-LD model at a time, and determined that inclusion of regulatory annotations (primarily histone marks and transcription factor binding sites) reduced estimates of  $\tau^*$  by 17-89% (Table A.35).

We assessed whether genes in enriched pathways have higher network connectivity than other genes. For each of the four gene networks, for each of the 141 enriched pathways, we assessed both the network connectivity within the pathway and the network connectivity between the pathway and interacting genes (one-degree neighbors) outside the pathway, as compared to 10,000 null pathways with the same number of genes, each randomly sampled from a randomly chosen pathway from the full set of 18,119 pathways (see Material and Methods). We assessed network connectivity using the sum of edge weights between genes. For each network, we averaged results across pathways. We determined that genes in enriched pathways have higher network connectivity within the pathway (1.43x-7.60x more edges; Figure 1-5 and Table

A.36), but do not necessarily have higher network connectivity with interacting genes outside the pathway (0.69x-1.56x; Table A.36); we note that there is no significant difference in the number of interacting genes between the 141 enriched pathways and the 10,000 null pathways (Table A.36).

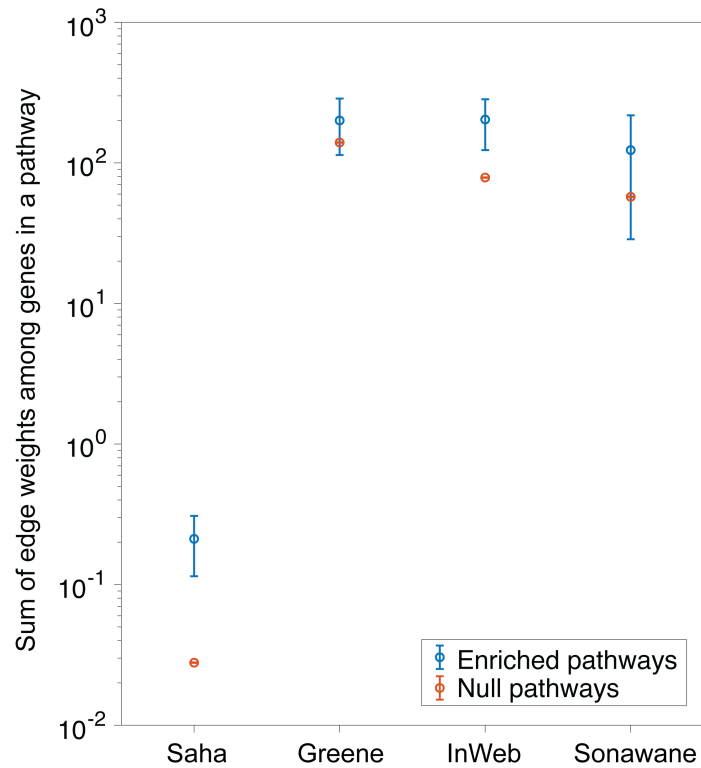


Figure 1-5: **Genes in enriched pathways have high network connectivity.** For each of four networks, we report the sum of edge weights in the network between genes in the pathway, averaged across 141 enriched pathways. For comparison purposes, we report the same quantity averaged across 10,000 null pathways with the same number of genes. Error bars represent 95% confidence intervals (smaller than data points for null pathways). Numerical results, and analogous results for network connectivity between a pathway and interacting genes outside the pathway, are reported in Table A.36.

We repeated the S-LDSC analysis (Figure 1-4B-C) using new pathway+network annotations constructed using a random-forest classifier, Quack<sup>35</sup>, that identifies new candidate genes that have similar topological patterns and network centrality metrics as genes in the input pathway. Because genes in enriched input pathways have high network connectivity (Figure 1-5), this is closely related to our primary strategy of defining pathway+network annotations based on genes with high network connec-

tivity to genes in enriched input pathways. Indeed, Quack annotations were highly correlated with our main pathway+network annotations ( $r = 0.50-0.67$ ; Table A.17), and produced S-LDSC results similar to our main analysis (Figure 1-4B-C), including significant enrichments for all four networks but non-significant  $\tau^*$  conditional on the baseline-LD model (Table A.37).

## 1.6 Discussion

We analyzed 42 diseases and complex traits (average  $N=323K$ ) to show that genes with high network connectivity are enriched for disease heritability, but that it is critical for gene network and pathway analyses to account for known functional annotations, such as those from our baseline-LD model<sup>3</sup>. First, in analyses of pathway annotations, we identified 156 pathway-trait pairs with significant heritability enrichment after conditioning on the baseline-LD model, a stringent step that caused a majority of pathway-trait pairs reported in recent studies to become non-significant in our analyses. Second, we determined that network annotations based on closeness centrality, a measure of network connectivity, are strongly enriched for disease heritability, but that these enrichments were fully explained by annotations from the baseline-LD model. Third, for each of the 156 significant pathway-trait pairs, we determined that pathway+network annotations constructed from genes with network connectivity to the input pathway were strongly enriched for the corresponding traits, but that once again these enrichments were largely explained by annotations from the baseline-LD model.

Our findings have important ramifications for studies connecting gene networks and pathways to disease<sup>28-33,50,52,55-58,60-63,67-70</sup>. Specifically, it is important to account for known functional annotations when seeking to elucidate biological mechanisms. For some methods, such as S-LDSC<sup>2,28,30</sup>, it is straightforward to incorporate known functional annotations such as those from the baseline-LD model<sup>3</sup>, and we emphasize the importance of doing so. For other methods, it is of high interest to investigate how functional annotations could be incorporated. More generally, it is

of broad interest to re-assess previously reported results while accounting for known functional annotations; for example, this could be achieved by running S-LDSC both with and without incorporating functional annotations from the baseline-LD model.

We note several limitations of our work. First, S-LDSC is not well-suited to analysis of annotations spanning a very small proportion of the genome<sup>2</sup> and does not model sparsity in trait effect sizes, potentially explaining why we did not identify enriched pathways for 8 traits that are less polygenic<sup>122</sup> (e.g. age at menopause, balding, hair color, sunburn). Nonetheless, our main results attained high statistical significance. Second, we did not explicitly compare S-LDSC to other methods. However, previous work suggests that S-LDSC compares favorably to other gene set enrichment methods, both in simulations and in analyses of real traits<sup>4</sup>. Third, interpretation of pathway-trait enrichments is complicated by the possibility that enrichment signals may be driven by a small number of highly significant genes<sup>33</sup>. However, we verified that repeating our main pathway+network analyses using the remaining significant pathway-trait pairs (after excluding genes implicated by GWAS) produced similar conclusions (Table A.38). Fourth, gene networks may include false-positive interactions, even after correcting for technical confounding<sup>51,71,123</sup>. However, our network perturbation analysis, DSD-transformed network analysis, and consensus network analysis all support the robustness of our results. Fifth, inferences about components of heritability can potentially be biased by failure to account for LD-dependent architectures<sup>3,124-126</sup>. All of our main analyses used the baseline-LD model, which includes 6 LD-related annotations<sup>3</sup>. The baseline-LD model is supported by formal model comparisons using likelihood and polygenic prediction methods, as well as analyses using a combined model incorporating alternative approaches<sup>127</sup>; however, there can be no guarantee that the baseline-LD model perfectly captures LD-dependent architectures. Sixth, although we showed that many pathway-trait pairs reported in recent studies were fully explained by the baseline-LD model and thus could potentially be due to factors that do not play a direct role in trait biology<sup>117</sup>, our analyses do not resolve which factors are causal. Nonetheless, because it is plausible that the regulatory annotations of the baseline-LD model may be the causal factors, ac-

counting for known functional annotations is an appropriate conservative measure to avoid incorrect conclusions. Finally, although we identified many significantly enriched pathways conditional on the baseline-LD model, our results for network and network+pathway annotations represent a negative result for efforts to improve upon the baseline-LD model<sup>3</sup>, further emphasizing the importance of accounting for known functional annotations in network and pathway analyses.

## 1.7 Web resources

Pathway, network, and pathway+network annotations: [https://alkesgroup.broadinstitute.org/LDSCORE/Kim\\_pathwaynetwork/](https://alkesgroup.broadinstitute.org/LDSCORE/Kim_pathwaynetwork/)

Source codes, gene scores, 18 gene sets: <https://github.com/samskim/networkconnectivity>

S-LDSC software: <https://github.com/bulik/ldsc>

baseline-LD annotations: <https://alkesgroup.broadinstitute.org/LDSCORE/>

Saha transcriptome-wide networks: [https://storage.googleapis.com/gtex\\_analysis\\_v6p/coexpression\\_networks/coexpression\\_networks\\_v6p.zip](https://storage.googleapis.com/gtex_analysis_v6p/coexpression_networks/coexpression_networks_v6p.zip)

Greene tissue-specific co-expression networks: <http://hb.flatironinstitute.org/>

InWeb protein-protein interaction network (inBio\_Map\_core\_2016\_09\_12.zip): <https://www.intomics.com/inbio/map>

Sonawane gene regulatory networks: <https://zenodo.org/record/838734>

DAVID (Database for Annotation, Visualization and Integrated Discovery, v.6.8): <https://david.ncifcrf.gov>

DSD (Diffusion State Distance) algorithm: <http://dsd.cs.tufts.edu/capdsd/>

Ensembl biomart: <https://www.ensembl.org/biomart>

ENCODE ChIP-Seq Significance Tool: <http://encodeqt.simple-encode.org/>

HapMap: <ftp://ftp.ncbi.nlm.nih.gov/hapmap/>

Graph-tool: <https://graph-tool.skewed.de>

Subset of 18 gene sets: [https://github.com/macarthur-lab/gene\\_lists](https://github.com/macarthur-lab/gene_lists)

Zhu & Stephens pathway studies: <https://zenodo.org/record/838734#.W89JDx>

NKiAw

GTEEx (Release v7): <http://www.gtexportal.org>

GWAS Catalog (Release v1.0): <https://www.ebi.ac.uk/gwas>

1000 Genomes Project Phase 3 data: <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502>

PLINK software: <https://www.cog-genomics.org/plink2>

BOLT-LMM software: <https://data.broadinstitute.org/alkesgroup/BOLT-LMM>

BOLT-LMM summary statistics for UK Biobank traits: <https://data.broadinstitute.org/alkesgroup/UKBB>

UK Biobank: <http://www.ukbiobank.ac.uk/>

UK Biobank Genotyping and QC Documentation: [http://www.ukbiobank.ac.uk/wp-content/uploads/2014/04/UKBiobank\\_genotyping\\_QC\\_documentation-web.pdf](http://www.ukbiobank.ac.uk/wp-content/uploads/2014/04/UKBiobank_genotyping_QC_documentation-web.pdf)

Online Mendelian Inheritance in Man: <http://www.omim.org>



## Chapter 2

# Improving the informativeness of Mendelian disease-derived pathogenicity scores for common disease

### 2.1 Preface

This chapter is previously published in *Nature Communications* as "Improving the informativeness of Mendelian disease-derived pathogenicity scores for common disease" (2020; see ref.<sup>128</sup>) with coauthors Kushal K. Dey, Omer Weissbrod, Carla Márquez-Luna, Steven Gazal, and Alkes L. Price. I am also grateful to Bryce van de Geijn, Farhad Hormozdiari, Armin Schoech, and Huwenbo Shi for helpful discussion.

### 2.2 Abstract

Despite considerable progress on pathogenicity scores prioritizing variants for Mendelian disease, little is known about the utility of these scores for common disease. Here, we assess the informativeness of Mendelian disease-derived pathogenicity scores for common disease and improve upon existing scores. We first apply stratified link-

age disequilibrium (LD) score regression to evaluate published pathogenicity scores across 41 common diseases and complex traits (average  $N=320K$ ). Several of the resulting annotations are informative for common disease, even after conditioning on a broad set of functional annotations. We then improve upon published pathogenicity scores by developing AnnotBoost, a machine learning framework to impute and denoise pathogenicity scores using a broad set of functional annotations. AnnotBoost substantially increases the informativeness for common disease of both previously uninformative and previously informative pathogenicity scores, implying that Mendelian and common disease variants share similar properties. The boosted scores also produce improvements in heritability model fit and in classifying disease-associated, fine-mapped SNPs. Our boosted scores may improve fine-mapping and candidate gene discovery for common disease.

## 2.3 Introduction

Despite considerable progress on pathogenicity scores prioritizing both coding and non-coding variants for Mendelian disease<sup>6,8,129–136</sup> (reviewed in ref.<sup>137</sup>), little is known about the utility of these pathogenicity scores for common disease. The shared genetic architecture between Mendelian disease and common disease has been implicated in studies reporting the impact of genes underlying monogenic forms of common diseases on the corresponding common diseases<sup>138</sup>, significant comorbidities among Mendelian and complex diseases<sup>139</sup>, and gene-level overlap between Mendelian diseases and cardiovascular diseases<sup>140–142</sup>, neurodevelopmental traits<sup>143,144</sup>, and other complex traits<sup>145</sup>. However, variant-level assessment of shared genetic architecture using Mendelian disease-derived pathogenicity scores has not been explored. Thus, our current understanding of the genetic relationship between Mendelian disease and common disease remains limited.

Here, we assess the informativeness of Mendelian disease-derived pathogenicity scores for common disease and improve upon existing scores. We focus our attention on polygenic common and low-frequency variant architectures, which explain the bulk

of common disease heritability<sup>33,122,146–148</sup>. We assess the informativeness of annotations defined by top variants from published Mendelian disease-derived pathogenicity scores by applying stratified linkage disequilibrium (LD) score regression<sup>2</sup> (S-LDSC) with the baseline-LD model<sup>3,127</sup> to 41 independent common diseases and complex traits (average  $N = 320\text{K}$ ). We assess informativeness conditional on the baseline-LD model, which includes a broad set of coding, conserved, regulatory and LD-related annotations. Then, we improve upon the published pathogenicity scores by developing AnnotBoost, a gradient boosting-based machine learning framework to impute and denoise pathogenicity scores using functional annotations from the baseline-LD model. We assess the informativeness of annotations defined by top variants from the boosted scores by applying S-LDSC to quantify conditional informativeness after considering annotations from the baseline-LD model as well as annotations derived from the corresponding published scores. We also assess the informativeness of the published and boosted pathogenicity scores in producing improvements in heritability model fit and in predicting disease-associated, fine-mapped SNPs.

We find that several annotations derived from published pathogenicity scores are informative for common disease, even after conditioning on annotations from the baseline-LD model. Furthermore, AnnotBoost substantially increases the informativeness for common disease of both previously uninformative and previously informative pathogenicity scores, implying that Mendelian and common disease variants share similar properties. We conclude that our boosted scores have high potential to improve fine-mapping and candidate gene discovery for common disease.

## 2.4 Methods

### Genomic annotations and the baseline-LD model

We define a genomic annotation as an assignment of a numeric value to each SNP above a specified minor allele frequency (e.g.  $\text{MAF} \geq 0.5\%$ ) in a predefined reference panel (e.g. 1000 Genomes<sup>40</sup>). Continuous-valued annotations can have any real value.

Probabilistic annotations can have any real value between 0 and 1. A binary annotation can be viewed as a subset of SNPs (the set of SNPs with annotation value 1); we note all annotations analyzed in this work are binary annotations. Annotations that correspond to known or predicted function are referred to as functional annotations.

The baseline-LD model<sup>3</sup> (v2.1) contains 86 functional annotations (see Web resources). We use these annotations as features of AnnotBoost (see below). These annotations include genomic elements (e.g. coding, enhancer, promoter), conservation (e.g. GERP, PhastCon), regulatory elements (e.g. histone marks, DNaseI-hypersensitive sites (DHS), transcription factor (TF) binding sites), and linkage disequilibrium (LD)-related annotations (e.g. predicted allele age, recombination rate, SNPs with low levels of LD).

## Enrichment and $\tau^*$ metrics

We used stratified LD score regression (S-LDSC<sup>2,3</sup>) to assess the contribution of an annotation to disease heritability by estimating the enrichment and the standardized effect size ( $\tau^*$ ) of an annotation.

Let  $a_{cj}$  represent the (binary or probabilistic) annotation value of the SNP  $j$  for the annotation  $c$ . S-LDSC assumes the variance of per normalized genotype effect sizes is a linear additive contribution to the annotation  $c$ :

$$\text{Var}(\beta_j) = \sum_c a_{cj} \tau_c \quad (2.1)$$

where  $\tau_c$  is the per-SNP contribution of the annotation  $c$ . S-LDSC estimates  $\tau_c$  using the following equation:

$$\text{E}[\chi_j^2] = N \sum_c \ell(j, c) \tau_c + 1 \quad (2.2)$$

where  $N$  is the sample size of the GWAS and  $\ell(j, c)$  is the LD score of the SNP  $j$  to the annotation  $c$ . The LD score is computed as follow  $\ell(j, c) = \sum_k a_{ck} r_{jk}^2$  where  $r_{jk}$  is the correlation between the SNPs  $j$  and  $k$ .

We used two metrics to assess the informativeness of an annotation. First, the standardized effect size ( $\tau^*$ ), the proportionate change in per-SNP heritability associated with a one standard deviation increase in the value of the annotation (conditional on all the other annotations in the model), is defined as follows:

$$\tau_{c^*} = \frac{\tau_c sd(C)}{h_g^2/M} \quad (2.3)$$

where  $sd(C)$  is the standard deviation of the annotation  $c$ ,  $h_g^2$  is the estimated SNP-heritability, and  $M$  is the number of variants used to compute  $h_g^2$  (in our experiment,  $M$  is equal to 5,961,159, the number of common SNPs in the reference panel). The significance for the effect size for each annotation, as mentioned in previous studies<sup>3,36,42</sup>, is computed as  $(\frac{\tau^*}{se(\tau^*)} \sim N(0,1))$ , assuming that  $\frac{\tau^*}{se(\tau^*)}$  follows a normal distribution with zero mean and unit variance.

Second, enrichment of the binary and probabilistic annotation is the fraction of heritability explained by SNPs in the annotation divided by the proportion of SNPs in the annotation, as shown below:

$$\text{Enrichment} = \frac{\%h_g^2(C)}{\%\text{SNP}(C)} = \frac{\frac{h_g^2(C)}{h_g^2}}{\frac{\sum_j a_{jc}}{M}} \quad (2.4)$$

where  $h_g^2(C)$  is the heritability captured by the  $c$ th annotation. When the annotation is enriched for trait heritability, the enrichment is  $> 1$ ; the overlap is greater than one would expect given the trait heritability and the size of the annotation. The significance for enrichment is computed using the block jackknife as mentioned in previous studies<sup>2,4,36,42</sup>). The key difference between enrichment and  $\tau^*$  is that  $\tau^*$  quantifies effects that are unique to the focal annotation after conditioning on all the other annotations in the model, while enrichment quantifies effects that are unique and/or non-unique to the focal annotation.

In all our analyses, we used the European samples in 1000G<sup>40</sup> (see Web resources) as reference SNPs. Regression SNPs were obtained from HapMap 3<sup>72</sup> (see Web resources). SNPs with marginal association statistics  $> 80$  and SNPs in the major

histocompatibility complex (MHC) region were excluded. Unless stated otherwise, we included the baseline-LD model<sup>3</sup> in all primary analyses using S-LDSC, both to minimize the risk of bias in enrichment estimates due to model mis-specification<sup>2,3</sup> and to estimate effect sizes ( $\tau^*$ ) conditional on known functional annotations.

## Published Mendelian disease-derived pathogenicity scores

We considered a total 35 published scores: 11 Mendelian disease-derived missense pathogenicity scores, 6 genome-wide Mendelian disease-derived pathogenicity scores, and 18 additional scores (see Table 2.1 and Table C.15). Here, we provide a short description for Mendelian missense and genome-wide Mendelian disease-derived pathogenicity scores. Details for 18 additional scores and the baseline-LD annotations are provided in Table C.15. Our curated pathogenicity scores are available online (see Web resources).

For all scores, we constructed annotations using GRCh37 (hg19) assembly limited to all 9,997,231 low-frequency and common SNPs (with  $\text{MAF} \geq 0.5\%$ ) found in 1000 Genomes<sup>40</sup> European Phase 3 reference genome individuals (see Web resources). Mendelian missense scores were readily available from dbNSFP database<sup>149,150</sup> using a rankscore (a converted score based on the rank among scored SNPs); genome-wide Mendelian disease-derived scores were individually downloaded and used with no modification to original scores (see Web resources). For each pathogenicity score, we constructed a binary annotation based on optimized threshold (See below). Short descriptions for each pathogenicity score (excluding 18 additional scores and the baseline-LD annotations; provided in Table C.15) are provided below:

Mendelian disease-derived missense pathogenicity scores:

*PolyPhen-2*<sup>6,7</sup> (*HDIV and HVAR*): Higher scores indicate higher probability of the missense mutation being damaging on the protein function and structure. The default predictor is based on a naive Bayes classifier using HumDiv (*HDIV*), and the other is trained using HumVar (*HVAR*), using 8 sequence-based and 3 structure-based features.

*MetaLR/MetaSVM*<sup>151</sup>: An ensemble prediction score based on logistic regression (LR) or support vector machine (SVM) to classify pathogenic mutations from background SNPs in whole exome sequencing, combining 9 prediction scores and one additional feature (maximum minor allele frequency).

*PROVEAN*<sup>152,153</sup>: An alignment-based score to predict the damaging single amino acid substitutions.

*SIFT 4G*<sup>131</sup>: Predicted deleterious effects of an amino acid substitution to protein function based on sequence homology and physical properties of amino acids.

*REVEL*<sup>132</sup>: An ensemble prediction score based on a random forest classifier trained on 6,182 missense disease mutations from HGMD<sup>154</sup>, using 18 pathogenicity scores as features.

*M-CAP*<sup>133</sup>: An ensemble prediction score based on a gradient boosting classifier trained on pathogenic variants from HGMD<sup>154</sup> and benign variants from ExAC data set<sup>97</sup>, using 9 existing pathogenicity scores, 7 base-pair, amino acid, genomic region, and gene-based features, and 4 features from multiple sequence alignments across 99 species.

*PrimateAI*<sup>134</sup>: A deep-learning-based score trained on the amino acid sequence flanking the variant of interest and the orthologous sequence alignments in other species and eliminating common missense variants identified in 6 non-human primate species.

*MPC*<sup>155</sup> (*missense badness, PolyPhen-2, and constraint*): Logistic regression-based score to identify regions within genes that are depleted for missense variants in ExAC data<sup>97</sup> and incorporating variant-level metrics to predict the impact of missense variants. Higher MPC score indicates increased deleteriousness of amino acid substitutions once occurred in missense-constrained regions.

*MVP*<sup>156</sup>: A deep-learning-based score trained on 32,074 pathogenic variants from ClinVar<sup>157</sup>, HGMD<sup>154</sup>, and UniProt<sup>158</sup>, using 38 local context, constraint, conservation, protein structure, gene-based, and existing pathogenicity scores as features.

Genome-wide Mendelian disease-derived pathogenicity scores:

*CADD*<sup>8,9</sup>: An ensemble prediction score based on a support vector machine classifier

trained to differentiate 14.7 million high-frequency human-derived alleles from 14.7 million simulated variants, using 63 conservation, regulatory, protein-level, and existing pathogenicity scores as features. We used PHRED-scaled CADD score for all possible SNVs of GRCh37.

*Eigen/Eigen-PC*<sup>129</sup>: Unsupervised machine learning score based on 29 functional annotations and leveraging blockwise conditional independence between annotations to differentiate functional vs. non-functional variants. Eigen-PC uses the lead eigenvector of the annotation covariance matrix to weight the annotations. For both Eigen and Eigen-PC, we used PHRED-scaled scores and combined coding and non-coding regions to make it as a single genome-wide score. Higher score indicates more important (predicted) functional roles.

*ReMM*<sup>130</sup> (*regulatory Mendelian mutation*): An ensemble prediction score based on a random forest classifier to distinguish 406 hand-curated Mendelian mutations from neutral variants using conservation scores and functional annotations. Higher ReMM score indicate greater potential to cause a Mendelian disease if mutated.

*NCBoost*<sup>136</sup>: An ensemble prediction score based on a gradient boosting classifier trained on 283 pathogenic non-coding SNPs associated with Mendelian disease genes and 2830 common SNPs, using 53 conservation, natural selection, gene-based, sequence context, and epigenetic features.

*ncER*<sup>135</sup> (*non-coding essential regulation*): An ensemble prediction score based on a gradient boosting classifier trained on 782 non-coding pathogenic variants from ClinVar<sup>157</sup> and HGMD<sup>154</sup>, using 38 gene essentiality, 3D chromatin structure, regulatory, and existing pathogenicity scores as features.

## AnnotBoost framework

AnnotBoost is based on gradient boosting, a machine learning method for classification; the AnnotBoost model is trained using the XGBoost gradient boosting software<sup>159</sup> (see Web resources). AnnotBoost requires only one input, a pathogenicity score to boost, and generates a genome-wide (probabilistic) pathogenicity score (as



described in Supplementary Figure D-1). During the training, AnnotBoost uses decision trees, where each node in a tree splits SNPs into two classes (pathogenic and benign) using 75 coding, conserved, regulatory, and LD-related features from the baseline-LD model<sup>3</sup> (excluding 10 MAF bins features; we obtained similar results with or without MAF bins features; see Supplementary Figure D-11). We note that the baseline-LD annotations considered all low-frequency and common SNPs thus do not have unscored SNPs. The method generates training data from the input pathogenicity scores without using external variant data; top 10% SNPs from the input pathogenicity score are labeled as a positive training set, and bottom 40% SNPs are labeled as a control training set; we obtained similar results with other training data ratios (see Supplementary Figure D-12). The prediction is based on  $T$  additive estimators (we use  $T = 200$  to  $300$ ; see below), minimizing the following loss objective function  $L^t$  at the  $t$ -th iteration:

$$L^t = \sum_{i=1}^n l(y_i, \hat{y}_i^{t-1} + f_t(x_i)) + \gamma(f_t) \quad (2.5)$$

where  $l$  is a differentiable convex loss function (which measures the difference between the prediction ( $\hat{y}_i$ ) and the target  $y_i$  at the  $i$ -th instance),  $f_t$  is an independent tree structure, and last term  $\gamma(f_t)$  penalizes the complexity of the model, helping to avoid over-fitting. The prediction ( $\hat{y}_i$ ) is made by  $\sum_{t=1}^T f_t(x_i)$  by ensembling outputs of multiple weak-learner trees. Odd (resp. even) chromosome SNPs are used for training to score even (resp. odd) chromosome SNPs. The output of the classifier is the probability of being similar to the positive training SNPs and dissimilar to the control training SNPs.

We used the following model parameters: the number of estimators (200, 250, 300), depth of the tree (25, 30, 35), learning rate (0.05), gamma (minimum loss reduction required before additional partitioning on a leaf node; 10), minimum child weight (6, 8, 10), and subsample (0.6, 0.8, 1); we optimized parameters with hyperparameters tuning (a randomized search) with five-fold cross-validation. Two important parameters to avoid over-fitting are gamma and learning rate; we chose these values consistent

with previous studies<sup>135,136</sup>. The model with the highest AUROCs on the held-out data was selected and used to make a prediction.

To identify which feature(s) drives the prediction output with less bias, AnnotBoost uses Shapley Addictive Explanation (SHAP<sup>5</sup>), a widely used tool to interpret complex non-linear models, instead of built-in feature importance tool because of SHAP’s property of satisfying symmetry, dummy player, and additivity axioms. SHAP uses the training matrix (features x SNP labels) and the trained model to generate a signed impact of each baseline-LD features on the AnnotBoost prediction.

To evaluate the performance of classifiers, we plotted receiver operating characteristic (ROC) and precision-recall (PR) curves. As we train AnnotBoost by splitting SNPs into odd and even chromosomes, we report the average out-of-sample area under the curve (AUC) of the odd and even chromosomes classifier. We used the threshold of 0.5 to define a class; that is, class 1 includes SNPs with the output probability  $> 0.5$ . We caution that high classification accuracy does not necessarily translate into conditional informativeness for common disease<sup>160</sup>.

## **Constructing binary annotations using top variants from published and boosted scores**

For published Mendelian disease-derived missense pathogenicity scores, we considered five different thresholds to construct binary annotations: top 50%, 40%, 30%, 20% or 10% of scored variants. For published scores that produce Bonferroni-significant binary annotations, we report results for the binary annotation with largest  $|\tau^*|$  among those that are Bonferroni-significant. For published scores that do not produce Bonferroni-significant binary annotations, we report results for the threshold with most significant  $\tau^*$  (even though not Bonferroni-significant).

For all other published pathogenicity scores, we considered the top 10%, 5%, 1%, 0.5% or 0.1% of scored variants to construct binary annotations; we used more inclusive thresholds for published Mendelian disease-derived missense pathogenicity scores due to the small proportion of variants scored ( $\sim 0.3\%$ ; see Table 2.1). For

published scores that produce Bonferroni-significant binary annotations, we report results for the binary annotation with largest  $|\tau^*|$  among those that are Bonferroni-significant. For published scores that do not produce Bonferroni-significant binary annotations, we report results for the top 5% of variants (the average optimized proportion among Bonferroni-significant binary annotations); we made this choice because (in contrast to published Mendelian missense scores) for many other published scores the most significant  $\tau^*$  was not even weakly significant.

For boosted pathogenicity scores, we considered the top 10%, 5%, 1%, 0.5% or 0.1% of scored variants, as well as variants with boosted scores  $\geq 0.5$ ; we note that top 10% of SNPs does not necessarily translate to 10% of SNPs, as some SNPs share the same score, and some genomic regions (e.g. MHC) are excluded when running S-LDSC (see below). For boosted scores that produce Bonferroni-significant binary annotations, we report results for the binary annotation with largest  $|\tau^*|$  among those that are Bonferroni-significant. For boosted scores that do not produce Bonferroni-significant binary annotations, we report results for the top 5% of variants.

In all analyses, we excluded binary annotations with proportion of SNPs  $< 0.02\%$  (the same threshold used in ref.<sup>42</sup>), because S-LDSC does not perform well for small annotations<sup>2</sup>. We analyzed 155 annotations derived from published scores (31 published scores (Table 2.2), 5 thresholds for top x% of variants,  $31 \times 5 = 155$ ), such that 500 hypotheses is a conservative correction in the analysis of published scores. We also analyzed 492 annotations derived from boosted scores (82 underlying published scores including 47 baseline-LD model annotations (Table 2.2), 6 thresholds for top x% of variants,  $82 \times 6 = 492$ ), such that 500 hypotheses is a roughly appropriate correction in the analysis of boosted scores. For simplicity, we corrected for  $\max(155, 492) \approx 500$  hypotheses throughout. We note that, in the meta-analysis  $\tau^*$  p-values, a global FDR  $< 5\%$  corresponds to  $P < 0.0305$ ; thus, our choice of  $P < 0.05/500 = 0.0001$  is conservative.

In all primary analyses, we analyzed only binary annotations. However, we verified in a secondary analysis of the CDTS score<sup>16</sup> that probabilistic annotations produced results similar to binary annotations (see Supplementary Figure D-13).

## Heterogeneity of enrichment and $\tau^*$ across traits

For a given annotation, we assessed the heterogeneity of enrichment and  $\tau^*$  (across 41 independent traits) by estimating the standard deviation of the true parameter value across traits, as analogous to ref.<sup>148</sup>. We calculated the cross-trait  $\tau^*$  as the inverse variance weighted mean across the traits. Then, we compared  $\sum_{i=1}^n (\hat{\tau}_i - \hat{\tau}_{cross-trait})^2 / (std.error_i^2)$  to a  $\chi_n^2$  null statistic, where  $n = 47$  (41 independent traits; 47 summary statistics; see Supplementary Table C.1). We repeated analysis for heritability enrichment by using enrichments and standard errors of enrichment estimates from S-LDSC.

## Overlap between gene score quintiles informed by input pathogenicity scores and 165 reference gene sets

For a given pathogenicity score, we scored genes based on the maximum pathogenicity score of linked SNPs, where SNPs were linked to a unique nearest gene using ANNOVAR<sup>161</sup>: 9,997,231 SNP-gene links, decreasing to 5,059,740 S2G links after restricting to 18,117 genes with a protein product (according to HGNC<sup>91</sup>) that have an Ensembl gene identifier (ENSG ID). Gene scores are reported in Table C.9. We constructed quintiles of genes scores and assessed gene-level excess fold overlap with 165 reference gene sets of biological importance (see below; summarized in Table C.10). We note that this analysis used continuous-valued pathogenicity scores, instead of binary annotations.

The 165 reference gene sets (Table C.10) reflected a broad range of gene essentiality<sup>90</sup> metrics, as outlined in ref.<sup>42</sup>. They included known phenotype-specific Mendelian disease genes<sup>145</sup>, constrained genes<sup>97-100</sup>, essential genes<sup>92,93,162</sup>, specifically expressed genes across GTEx tissues<sup>4</sup>, dosage outlier genes across GTEx tissues<sup>163</sup>, genes with a ClinVar pathogenic or likely pathogenic variants<sup>157</sup>, genes in the Online Mendelian Inheritance in Man (OMIM<sup>102</sup>), high network connectivity genes in different gene networks<sup>42,164</sup>, genes with more independent SNPs<sup>42</sup>, known drug targets<sup>165</sup>, human targets of FDA-approved drugs<sup>105</sup>, eQTL-deficient genes<sup>108,166</sup>, and

housekeeping genes<sup>167</sup>; a subset of these gene sets were previously analyzed in ref.<sup>42</sup>.

As defined in our previous study<sup>42</sup>, the excess fold overlap of gene set 1 and gene set 2 is defined as follows:

$$\text{excess overlap}(\text{gene set 1, gene set 2}) = P_d/P_{tot} \quad (2.6)$$

where  $P_d = \frac{|\text{gene set 1} \cap \text{gene set 2}|}{|\text{gene set 2}|}$  and  $P_{tot} = \frac{|\text{gene set 1} \cap \text{all protein-coding genes}|}{|\text{all protein-coding genes}|}$ . The standard error for the excess overlap is similarly scaled:

$$SE = \sqrt{\frac{P_d(1 - P_d)}{|\text{gene set 2}|}}/P_{tot} \quad (2.7)$$

When there is excess overlap, the excess fold overlap is  $> 1$ ; when there is depletion, the excess fold overlap is  $< 1$ . We assessed the odds ratio and significance in the difference between the excess overlap between the boosted gene quintile and the published gene quintile by the Fisher’s exact test (two-sided).

## Evaluating heritability model fit using $\log l_{SS}$

Given a heritability model (e.g. the baseline-LD model, combined joint model, or combined marginal model), we define the  $\Delta \log l_{SS}$  of that heritability model as the  $\log l_{SS}$  of that heritability model minus the  $\log l_{SS}$  of a model with no functional annotations (baseline-LD-nofunct; 17 LD and MAF annotations from the baseline-LD model<sup>3</sup>), where  $\log l_{SS}^1$  is an approximate likelihood metric that has been shown to be consistent with the exact likelihood from restricted maximum likelihood (REML; see Web resources). We compute p-values for  $\Delta \log l_{SS}$  using the asymptotic distribution of the Likelihood Ratio Test (LRT) statistic:  $-2 \log l_{SS}$  follows a  $\chi^2$  distribution with degrees of freedom equal to the number of annotations in the focal model, so that  $-2 \Delta \log l_{SS}$  follows a  $\chi^2$  distribution with degrees of freedom equal to the difference in number of annotations between the focal model and the baseline-LD-nofunct model. We used UK10K as the LD reference panel and analyzed 4,631,901 HRC (haplotype reference panel<sup>168</sup>) well-imputed SNPs with  $MAF \geq 0.01$  and  $INFO \geq 0.99$  in the

reference panel; We removed SNPs in the MHC region, SNPs explaining  $> 1\%$  of phenotypic variance and SNPs in LD with these SNPs.

We computed  $\Delta \log l_{SS}$  for 4 heritability models:

- baseline-LD: annotations from the baseline-LD model<sup>2,3</sup> (86 annotations)
- baseline-LD+joint: baseline-LD model + 11 jointly significant annotations (3 published, 8 boosted; 97 annotations)
- baseline-LD+marginal-stringent: baseline-LD model + 19 marginally significant annotations with conditional  $\tau^* > 0.5$  (105 annotations)
- baseline-LD+marginal: baseline-LD model + 64 marginally significant annotations (11 published, 53 boosted; 150 annotations)

## Classification of fine-mapped disease SNPs

We assessed the classification accuracy of fine-mapped disease SNPs. Here, we consider only low-frequency and common SNPs ( $MAF \geq 0.5\%$ ) and report the total number of unique SNPs (regardless MAF). we assessed the accuracy of classifying five different SNP sets (summarized in Table C.24): (1) 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup> (of 7,747 total SNPs; 95% credible sets), (2) 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup> (of 4,311 total SNPs; 95% credible sets), (3) 1,851 SNPs (of 2,225 SNPs, spanning 3,025 SNP-trait pairs; stringently defined by causal posterior probability  $\geq 0.95$ ) functionally-informed fine-mapped for 49 UK Biobank traits from Weissbrod et al.<sup>13</sup>, (4) 1,379 (of 1,853 total SNPs with causal posterior probability  $\geq 0.95$ ) non-functionally-informed fine-mapped SNPs for 49 UK Biobank traits from Weissbrod et al.<sup>13</sup>, and (5) 14,807 SNPs from the NHGRI GWAS catalog<sup>14,15</sup> (2019-07-12 version; p-value  $< 5e-8$ ; we note only about 5% of GWAS SNPs are expected to be causal<sup>11</sup>).

For each of these five SNP sets, we matched 10 control SNPs for each positive fine-mapped SNP by matching LD, MAF, and genomic element, as in previous studies<sup>16,135,136</sup>; we note that these studies emphasized the need for matching the relative genomic region distribution in performance evaluation. MAF was based on the same reference panel (European samples from 1000 Genomes Phase 3<sup>40</sup>), and LD was es-

timated by applying S-LDSC on the all SNPs annotation ('base'). To identify the genomic element of each SNPs and the nearest gene, we annotated these five sets of SNPs using ANNOVAR<sup>161</sup> using the gene-based annotation. For assigning the genomic element to each SNP, we used the default ANNOVAR prioritization rule for gene-based annotations: exonic = splicing (defined by 10bp of a splicing junction) > ncRNA > UTR5 = UTR3 > intronic > upstream = downstream > intergenic. When SNP (in the intergenic or intronic region) is associated with overlapping genes, the nearest protein-coding gene (based on the distance to the TSS or TSE) is retained. To obtain 10 control SNPs for each positive fine-mapped SNP, we first searched the control SNPs within the same genomic element and the same chromosome of that positive SNP; then kept the 10 control SNPs with the most similar LD and MAF (based on the average of rank(LD difference from the positive SNP) and rank(MAF difference from the positive SNP)). In secondary analyses, we instead retained a unique most closely matched control SNP.

Given positive and control SNP sets, we computed the AUPRCs (and AUROCs) by an individual score (each of 82 published and 82 boosted scores). We refer this as a single-score analysis. We used AUPRC as a primary metric, as AUPRC is more robust for imbalanced data<sup>169</sup>. We assessed the significance of the difference between two AUPRCs using 1,000 samples bootstrapped standard errors then performed 2-sample t-test; variance of AUPRCs (and AUROCs) from 1,000 samples was sufficiently small. We note this single-score analysis measures an improvement between two scores, where one score is derived from the other (e.g. our boosted score from published score). Also, we performed a multi-score analysis. For each heritability model, we aggregated scores by training a gradient boosting model (features: aggregated scores, positive label: each of five sets of SNPs, control label: LD-, MAF-, and genomic-element-matched control sets of SNPs); we used odd (resp. even) chromosomes as training data to make predictions for even (resp. odd) chromosomes. We used the same training parameters as AnnotBoost (carefully selected to avoid over-fitting, consistent with the previous study<sup>135,136</sup>) with hyperparameters tuned using a randomized search method with five-fold cross-validation. We report the average AUPRC and AUROC of odd and

even chromosome classifiers. We also computed  $\Delta$ AUPRC as AUPRC of a given model minus AUPRC of baseline-LD-nofunct model. We note that no disease data (five sets of SNPs used as labels) was re-used in these analyses, as AnnotBoost uses only the input pathogenicity scores to generate positive and negative sets of training data. We assessed the significance of the difference as described above. To identify which feature(s) drives the prediction output, we applied SHAP<sup>5</sup> to generate a signed impact of each baseline-LD, published, and boosted score features on classifying fine-mapped disease SNPs.

## 2.5 Results

### Overview of methods

We define a binary annotation as an assignment of a binary value to each of low-frequency ( $0.5\% \leq \text{MAF} < 5\%$ ) and common ( $\text{MAF} \geq 5\%$ ) SNP in a 1000 Genomes Project European reference panel<sup>40</sup>, as in our previous work<sup>2,127</sup>. We define a pathogenicity score as an assignment of a numeric value quantifying predicted pathogenicity, deleteriousness, and/or protein function to some or all of these SNPs; we refer to these score as Mendelian disease-derived pathogenicity scores, as these scores have predominantly been developed and assessed in the context of Mendelian disease (e.g. using pathogenic variants from ClinVar<sup>157</sup> and HGMD<sup>154</sup>). We analyze 11 Mendelian disease-derived missense scores, 6 genome-wide Mendelian disease-derived scores, and 18 additional scores. Our primary focus is on binary annotations defined either using top variants from published (missense or genome-wide) Mendelian disease-derived pathogenicity scores, or using top variants from boosted scores that we constructed from those pathogenicity scores using AnnotBoost, a gradient boosting-based framework that we developed to impute and denoise pathogenicity scores using 75 coding, conserved, regulatory and LD-related annotations from the baseline-LD model<sup>3,127</sup> (Figure D-1; see Methods). AnnotBoost uses decision trees to distinguish pathogenic variants (defined using the input pathogenicity score) from benign variants; the An-



notBoost model is trained using the XGBoost gradient boosting software<sup>159</sup> (see Web resources). AnnotBoost uses odd (resp. even) chromosomes as training data to make predictions for even (resp. odd) chromosomes; the output of AnnotBoost is the predicted probability of being pathogenic. We note that Mendelian disease-derived pathogenicity scores may score a subset of SNPs, but every baseline-LD model annotation scores all SNPs. Further details are provided in the Methods section; we have publicly released open-source software implementing AnnotBoost, as well as all pathogenicity scores and binary annotations analyzed in this work (see Web resources).

We assessed the informativeness of the resulting binary annotations for common disease heritability by applying S-LDSC<sup>2</sup> to 41 independent common diseases and complex traits<sup>36</sup> (average  $N = 320\text{K}$ ; Table C.1; see Web resources), conditioned on coding, conserved, regulatory and LD-related annotations from the baseline-LD model<sup>3,127</sup> and meta-analyzing results across traits. We assessed informativeness for common disease using standardized effect size ( $\tau^*$ ), defined as the proportionate change in per-SNP heritability associated to a one standard deviation increase in the value of the annotation, conditional on other annotations<sup>3</sup> (see Methods). We also computed the heritability enrichment, defined as the proportion of heritability divided by the proportion of SNPs. Unlike enrichment,  $\tau^*$  quantifies effects that are unique to the focal annotation; annotations with significantly positive or negative  $\tau^*$  are informative after considering all other annotations in the model, whereas annotations with  $\tau^* = 0$  contain no unique information, even if they are enriched for heritability (see Methods). While S-LDSC models linear combinations of functional annotations, AnnotBoost constructs (linear and) non-linear combinations of baseline-LD model annotations to provide unique information.

## Informativeness of Mendelian disease-derived missense scores for common disease

We assessed the informativeness for common disease of binary annotations derived from 11 Mendelian pathogenicity scores for missense variants<sup>6,7,131–134,151,152,155,156</sup> (see Table 2.1). These scores reflect the predicted impact of missense mutations on Mendelian disease; we note that our analyses of common disease are focused on common and low-frequency variants, but these scores were primarily trained using very rare variants from ClinVar<sup>157</sup> and Human Gene Mutation Database (HGMD)<sup>154</sup>. For each of the 11 missense scores, we constructed binary annotations based on top missense variants using 5 different thresholds (from top 50% to top 10% of missense variants) and applied S-LDSC<sup>2,3</sup> to 41 independent common diseases and complex traits (Table C.1), conditioning on coding, conserved, regulatory and LD-related annotations from the baseline-LD model<sup>3,127</sup> and meta-analyzing results across traits; proportions of top SNPs were optimized to maximize informativeness (see Methods). We incorporated the 5 different thresholds into the number of hypotheses tested when assessing statistical significance (Bonferroni  $P < 0.05/500 = 0.0001$ , based on a total of  $\approx 500$  hypotheses tested in this study; see Methods). We identified (Bonferroni-significant) conditionally informative binary annotations derived from 2 published missense scores: the top 30% of SNPs from MPC<sup>155</sup> (enrichment = 27x (s.e. 2.5),  $\tau^* = 0.60$  (s.e. 0.07)) and the top 50% of SNPs from PrimateAI<sup>134</sup> (enrichment = 17x (s.e. 2.0),  $\tau^* = 0.42$  (s.e. 0.09)) (Figure 2-1, Table 2.2 and Table C.2). The MPC (Missense badness, PolyPhen-2, and Constraint) score<sup>155</sup> is computed by identifying regions within genes that are depleted for missense variants in ExAC data<sup>97</sup> and incorporating variant-level metrics to predict the impact of missense variants; the PrimateAI score<sup>134</sup> is computed by eliminating common missense variants identified in other primate species (which are presumed to be benign in humans), incorporating a deep learning model trained on the amino acid sequence flanking the variant of interest and the orthologous sequence alignments in other species. The remaining published Mendelian disease-derived missense scores all had derived binary annotations

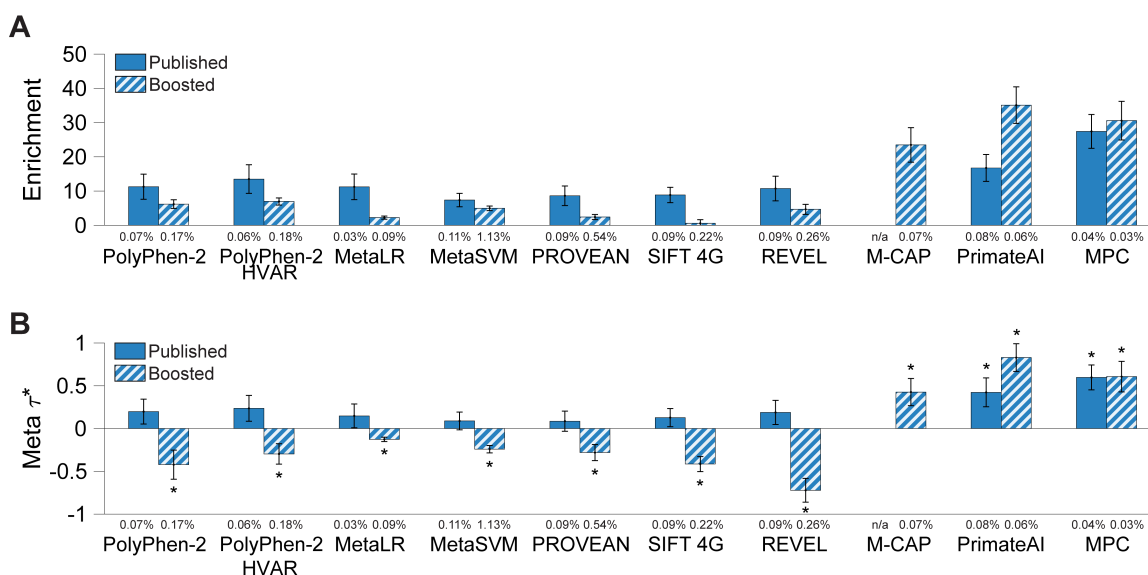
that were significantly enriched for disease heritability (after Bonferroni correction) but not conditionally informative (except for the published M-CAP<sup>133</sup> score, which spanned too few SNPs to be included in the S-LDSC analysis).

Score	Description	Coverage (% SNPs scored)	Ref.
PolyPhen-2	Impact of missense variants using protein sequence and structure using HumDiv	0.28%	ref. <sup>6,7</sup>
PolyPhen-2-HVAR	Impact of missense variants using protein sequence and structure using HumVar	0.28%	ref. <sup>6,7</sup>
MetaLR	Deleterious missense mutations using ensemble scoring (logistic regression)	0.32%	ref. <sup>151</sup>
MetaSVM	Deleterious missense mutations using ensemble scoring (support vector machine)	0.32%	ref. <sup>151</sup>
PROVEAN	Impact of an amino acid change on protein function	0.31%	ref. <sup>152,153</sup>
SIFT 4G	Impact of an amino acid change on protein function	0.31%	ref. <sup>131</sup>
REVEL	Pathogenic missense variants using ensemble scoring	0.32%	ref. <sup>132</sup>
M-CAP	Pathogenic rare missense variants	0.03%	ref. <sup>133</sup>
PrimateAI	Impact of missense variants using deep neural networks	0.26%	ref. <sup>134</sup>
MPC	Regional missense constraint	0.10%	ref. <sup>155</sup>
MVP	Impact of missense variants using deep neural networks	0.29%	ref. <sup>156</sup>
CADD	Predicted deleterious variants using ensemble scoring	100%	ref. <sup>8,9</sup>
Eigen	Putatively causal variants using unsupervised learning	83.79%	ref. <sup>129</sup>
Eigen-PC	Putatively causal variants using unsupervised learning using the lead eigenvector	83.79%	ref. <sup>129</sup>
ReMM	Pathogenic regulatory variants using ensemble scoring	100%	ref. <sup>130</sup>
NCBoost	Pathogenic non-coding variants using ensemble scoring	28.55%	ref. <sup>136</sup>
ncER	Essential regulatory variants using ensemble scoring	61.94%	ref. <sup>135</sup>

Table 2.1: **11 Mendelian disease-derived missense and 6 genome-wide Mendelian disease-derived pathogenicity scores.** For each of 17 Mendelian disease-derived pathogenicity scores analyzed, we provide a description and report the coverage (% of SNPs scored) and corresponding reference. The first 11 annotations are scores for missense variants, and the last 6 annotations are genome-wide scores. Annotations are ordered first by type and then by the year of publication.

We constructed *boosted* scores from the 11 Mendelian disease-derived missense scores using AnnotBoost, a gradient boosting-based machine learning framework that we developed to impute and denoise pathogenicity scores using functional annotations from the baseline-LD model<sup>3</sup> (see Methods). We note that AnnotBoost scores genome-wide (missense and non-missense) variants, implying low genome-wide correlations between input Mendelian disease-derived missense scores and corresponding genome-wide boosted scores (0.02-0.24; Table C.3A). AnnotBoost attained high predictive accuracy in out-of-sample predictions of input missense scores (AUROC= 0.76-0.94, AUPRC= 0.43-0.82; Table C.4), although we caution that high predictive

accuracy does not necessarily translate into conditional informativeness for common disease<sup>160</sup>. We further note that out-of-sample AUROCs closely tracked the genome-wide correlations between input Mendelian disease-derived missense scores and corresponding genome-wide boosted scores ( $r = 0.65$ ), implying that accurately predicting the input pathogenicity scores results in more correlated boosted scores.



**Figure 2-1: Informativeness for common disease of binary annotations derived from 11 Mendelian disease-derived missense scores and corresponding boosted scores.** We report (A) heritability enrichment of binary annotations derived from published and boosted Mendelian disease-derived missense scores, meta-analyzed across 41 independent traits; (B) conditional  $\tau^*$  values, conditioning on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores). We report results for 10 Mendelian disease-derived missense scores (of 11 analyzed) for which annotations derived from published and/or boosted scores were conditionally significant; the published M-CAP score spanned too few SNPs to be included in the S-LDSC analysis. The percentage under each bar denotes the proportion of SNPs in the annotation; the proportion of top SNPs included in each annotation was optimized to maximize informativeness (largest  $|\tau^*|$  among Bonferroni-significant annotations, or most significant p-value if no annotation was Bonferroni-significant). Error bars denote 95% confidence intervals. In panel (B), \* denotes conditionally significant annotations. Numerical results are reported in Table C.2. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25.

Score	# scores	# marginally significant annotations		# significant annotations in a combined joint model	
		published	boosted	published	boosted
Mendelian missense	11	2*	10	1*	2
Genome-wide Mendelian	6	3	6	2	3
Additional scores	18	6**	13	0**	0
Baseline-LD model annotations	47	n/a	24	n/a	3

**Table 2.2: Summary of informativeness for common disease of annotations derived from 82 published scores and corresponding boosted scores** For each category of scores, we report the number of scores; the number of marginally conditionally informative annotations (S-LDSC  $\tau^* p < 0.0001$ , conditional on the baseline-LD model) (baseline-LD+marginal model); and the number of jointly conditionally informative annotations in a combined joint model (S-LDSC  $\tau^* p < 0.0001$  and  $|\tau^*| \geq 0.25$ , conditional on the baseline-LD model and each other) (baseline-LD+joint model). \*Based on 9/11 published Mendelian missense scores analyzed, as binarized annotations were too small to analyze for the remaining 2 published Mendelian missense scores. \*\*Based on 16/18 published additional scores analyzed, as binarized annotations were too small to analyzed for the remaining 2 published additional scores.

For each missense pathogenicity score, after running AnnotBoost, we constructed binary annotations based on top genome-wide variants, using 6 different thresholds (ranging from top 10% to top 0.1% of genome-wide variants, as well as variants with boosted scores  $\geq 0.5$ ; see Methods). We assessed the informativeness for common disease of binary annotations derived from each of the 11 boosted scores using S-LDSC, conditioning on annotations from the baseline-LD model and 5 binary annotations derived from the corresponding published Mendelian disease-derived missense score (using all 5 thresholds) (baseline-LD+5). We identified conditionally informative binary annotations derived from boosted versions of 10 Mendelian disease-derived missense scores, including 8 previously uninformative scores and the 2 previously informative scores (Figure 2-1, Table 2.2 and Table C.2). Letting  $\uparrow$  denote boosted scores, examples include the top 0.1% of SNPs from M-CAP $\uparrow$ <sup>133</sup>, a previously uninformative score (enrichment = 23x (s.e. 2.6),  $\tau^* = 0.43$  (s.e. 0.08); the published M-CAP pathogenicity score spanned too few SNPs to be included in the S-LDSC analysis of Figure 2-1) and the top 0.1% of SNPs from PrimateAI $\uparrow$ <sup>134</sup>, a previously

informative score (enrichment = 35x (s.e. 2.7),  $\tau^* = 0.83$  (s.e. 0.08)). The M-CAP (Mendelian Clinically Applicable Pathogenicity) score<sup>133</sup> is computed by training a gradient boosting tree classifier to distinguish pathogenic variants from HGMD<sup>154</sup> vs. benign variants from ExAC<sup>97</sup> using 9 pathogenicity likelihood scores as features (including PolyPhen-2<sup>6</sup>, MetaLR<sup>151</sup>, CADD<sup>8</sup>; see Table 2.1); the PrimateAI score is described above. Interestingly, binary annotations derived from 7 boosted scores had significantly negative  $\tau^*$  ( $-0.72$  (s.e. 0.07) to  $-0.13$  (s.e. 0.01)). All but one of these binary annotations were significantly enriched for disease heritability, but less enriched than expected based on annotations from the baseline-LD+5 model (Table C.5; see ref.<sup>170</sup> and Methods, resulting in significantly negative  $\tau^*$ . These annotations are thus conditionally informative for disease heritability (analogous to transposable element annotations in ref.<sup>170</sup> that were significantly depleted, but less depleted than expected and thus conditionally informative); as noted above, annotations with significantly positive or negative  $\tau^*$  are conditionally informative. The boosted version of the remaining Mendelian disease-derived missense score (MVP $\uparrow$ ; not included in Figure 2-1) had a derived binary annotation that was significantly enriched for disease heritability (after Bonferroni correction) but not conditionally informative (Table C.2).

We performed 5 secondary analyses. First, we restricted the 10 significant binary annotations derived from our boosted Mendelian disease-derived missense scores to non-coding regions, which were previously unscored by the Mendelian disease-derived missense scores, and assessed the informativeness of the resulting non-coding binary annotations using S-LDSC. We determined that the non-coding annotations retained the bulk of the overall signals (85%-110% of absolute  $\tau^*$ ; Table C.6), implying that AnnotBoost leverages information about pathogenic missense variants to usefully impute scores for non-missense variants. Second, we investigated which features of the baseline-LD model contributed the most to the informativeness of the boosted annotations by applying Shapley Additive Explanation (SHAP)<sup>5</sup>, a widely used tool for interpreting machine-learning models. We determined that conservation-related features drove the predictions of the boosted annotations, particularly (binary and

continuous) GERP scores<sup>171</sup> (Figure D-2). Third, we examined the trait-specific S-LDSC results, instead of meta-analyzing results across 41 traits. We identified 2 (resp. 32) annotation-trait pairs with significant  $\tau^*$  values for annotations derived from published (resp. boosted) scores (Table C.7); significance was assessed using  $FDR < 5\%$ , as Bonferroni correction would be overly conservative in this case. These annotation-trait pairs spanned 1 (resp. 10) different annotations, all of which were also conditionally significant in the meta-analysis across traits (Figure 2-1B). Fourth, we assessed the heterogeneity of heritability enrichment and conditional informativeness ( $\tau^*$ ) across the 41 traits (as in ref.<sup>148</sup>; see Methods). We determined that 10/20 annotations tested had significant heterogeneity in enrichment and 14/20 annotations tested had significant heterogeneity in  $\tau^*$ , implying substantial heterogeneity across traits (Table C.8). Fifth, we investigated the overlap between genes linked to each of our 12 conditionally significant annotations and 165 gene sets of biological importance, including high-pLI genes<sup>97</sup> and known Mendelian genes<sup>145</sup> (see Methods; Table C.9 and Table C.10). We linked SNPs to the nearest gene, scored genes based on the maximum pathogenicity score of linked SNPs, and assessed overlap between top gene score quintiles and each of 165 reference gene sets. We consistently observed excess overlap for genes for which homozygous knockout in mice results in lethality<sup>92,93</sup> and high-pLI genes<sup>97</sup>, and depleted overlap for olfactory receptor genes<sup>107</sup> (Table C.11A-B). In addition, top gene score quintiles derived from our boosted scores often had significantly more overlap (based on Fisher’s exact test) than corresponding top gene quintiles derived from published scores (Table C.12; results for PolyPhen-2<sup>6,7</sup> shown in Figure D-3); gene scores derived from published and boosted pathogenicity scores were moderately correlated (Table C.11C-D). This implies that our new annotations can help identify biologically important genes.

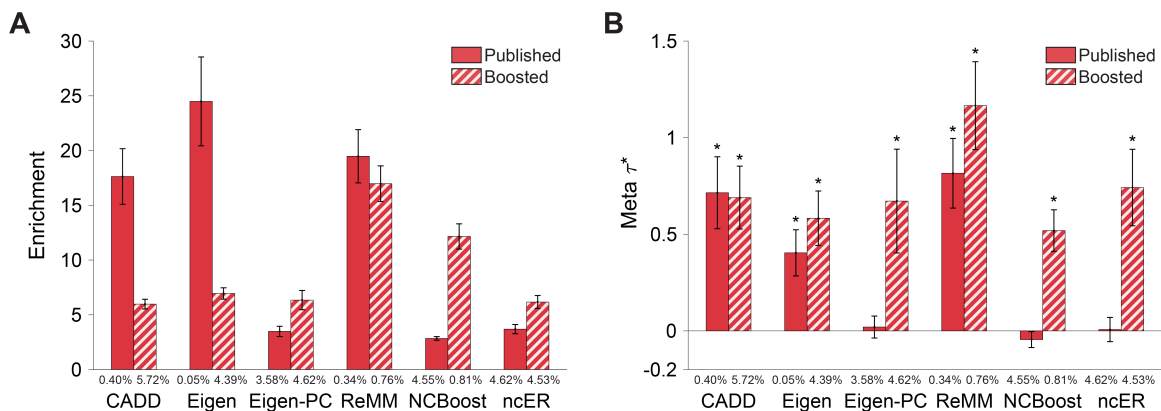
We conclude that 2 Mendelian disease-derived missense annotations and 10 boosted annotations are conditionally informative for common disease, relative to baseline-LD model annotations.

## Informativeness of genome-wide Mendelian disease-derived pathogenicity scores for common disease

We assessed the informativeness for common disease of binary annotations derived from 6 genome-wide Mendelian disease-derived pathogenicity scores<sup>8,129,130,135,136</sup> (see Table 2.1). These scores reflect the predicted impact of (coding and) non-coding variants on Mendelian disease; as above, these scores were primarily trained using very rare variants from ClinVar<sup>157</sup> and HGMD<sup>154</sup>. For each of the 6 genome-wide scores, we constructed binary annotations based on top genome-wide variants using 5 different thresholds (from top 0.1% to top 10% of genome-wide variants) and applied S-LDSC to the 41 traits, conditioning on the baseline-LD model<sup>3</sup> and meta-analyzing results across traits; proportions of top SNPs were optimized to maximize informativeness (see Methods). We identified (Bonferroni-significant) conditionally informative binary annotations derived from 3 genome-wide scores: the top 0.5% of SNPs from ReMM<sup>130</sup> (enrichment = 19x (s.e. 1.2),  $\tau^* = 0.82$  (s.e. 0.09)), the top 0.5% of SNPs from CADD<sup>8,9</sup> (enrichment = 18x (s.e. 1.3),  $\tau^* = 0.71$  (s.e. 0.10)), and the top 0.1% of SNPs from Eigen<sup>129</sup> (enrichment = 24x (s.e. 2.1),  $\tau^* = 0.40$  (s.e. 0.06)) (Figure 2-2, Table 2.2 and Table C.13). The CADD (Combined Annotation Dependent Depletion) score<sup>8,9</sup> is computed by training a support vector machine to distinguish deleterious vs. neutral variants using functional annotations as features; the Eigen score<sup>129</sup> is computed from 29 input functional annotations by using an unsupervised machine learning method (leveraging blockwise conditional independence between annotations) to differentiate functional vs. non-functional variants; the ReMM (Regulatory Mendelian Mutation) score<sup>130</sup> is computed by training a random forest classifier to distinguish 406 hand-curated Mendelian mutations from neutral variants using conservation scores and functional annotations as features. The remaining 3 genome-wide scores all had derived binary annotations that were significantly enriched for disease heritability (after Bonferroni correction) but not conditionally informative (Table C.13).

We applied AnnotBoost to the 6 genome-wide Mendelian disease-derived scores.





**Figure 2-2: Informativeness for common disease of binary annotations derived from 6 genome-wide Mendelian disease-derived scores and corresponding boosted scores.** We report (A) heritability enrichment of binary annotations derived from published and boosted genome-wide Mendelian disease-derived scores, meta-analyzed across 41 independent traits; (B)  $\tau^*$  values, conditioning on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores). We report results for 6 genome-wide Mendelian disease-derived scores (of 6 analyzed) for which annotations derived from published and/or boosted scores were conditionally significant. The percentage under each bar denotes the proportion of SNPs in the annotation; the proportion of top SNPs included in each annotation was optimized to maximize informativeness (largest  $|\tau^*|$  among Bonferroni-significant annotations, or top 5% if no annotation was Bonferroni-significant; top 5% was the average optimized proportion among significant annotations). Error bars denote 95% confidence intervals. In panel (B), \* denotes marginally conditionally significant annotations. Numerical results are reported in Table C.13. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25.

We observed moderate correlations between input genome-wide Mendelian disease-derived scores and corresponding boosted scores ( $r = 0.35-0.66$ ; Table C.3B). Annot-Boost again attained high predictive accuracy in out-of-sample predictions of input genome-wide scores (AUROC = 0.83-1.00, AUPRC = 0.70-1.00; Table C.4); however, out-of-sample AUROCs did not closely track the correlations between input genome-wide scores and corresponding boosted scores ( $r = 0.05$ ).

We again constructed binary annotations based on top genome-wide variants, using 6 different thresholds (ranging from top 0.1% to top 10% of genome-wide variants, as well as variants with boosted scores  $\geq 0.5$ ; see Methods). We assessed the informativeness for common disease of binary annotations derived from each of the 6 boosted scores using S-LDSC, conditioning on annotations from the baseline-LD model and 5 binary annotations derived from the corresponding published genome-wide Mendelian disease-derived score (using all 5 thresholds). We identified conditionally informative binary annotations derived from boosted versions of all 6 genome-wide Mendelian disease-derived scores, including the 3 previously uninformative scores and the 3 previously informative scores (Figure 2-2, Table 2.2 and Table C.2). Examples include the top 5% of SNPs from ncER $\uparrow$ <sup>135</sup> (enrichment = 6.2x (s.e. 0.30),  $\tau^* = 0.74$  (s.e. 0.10)) and the top 0.5% of SNPs from boosted Eigen-PC $\uparrow$ <sup>129</sup> (enrichment = 16x (s.e. 1.1),  $\tau^* = 0.62$  (s.e. 0.12)), both of which were previously uninformative scores, and the top 1% of SNPs from ReMM $\uparrow$ <sup>130</sup> (enrichment = 17x (s.e. 0.8),  $\tau^* = 1.17$  (s.e. 0.12)), a previously informative score. The ncER (non-coding Essential Regulation) score<sup>135</sup> is computed by training a gradient boosting tree classifier to distinguish non-coding pathogenic variants from ClinVar<sup>157</sup> and HGMD<sup>154</sup> vs. benign variants using 38 functional and structural features; the Eigen-PC score<sup>129</sup> (related to the Eigen score) is computed from 29 input functional annotations by using the lead eigenvector of the annotation covariance matrix to weight the annotations; the ReMM score is described above.

We performed 5 secondary analyses. First, for the 4 genome-wide Mendelian disease-derived scores with <100% of SNPs scored (Table 2.1), we restricted the binary annotations derived from our boosted genome-wide Mendelian disease-derived

scores to previously unscored variants and assessed the informativeness of the resulting binary annotations using S-LDSC. We determined that these annotations retained only a minority of the overall signals (17%-54% of absolute  $\tau^*$ ; Table C.14), implying that AnnotBoost usefully denoises previously scored variants. Second, we again investigated which features of the baseline-LD model contributed the most to the informativeness of the boosted annotations by applying SHAP<sup>5</sup>. We determined that both conservation-related features (e.g. GERP scores) and LD-related features (e.g. LLD-AFR; level of LD in Africans) drove the predictions of the boosted annotations (Figure D-4). Third, we examined the trait-specific S-LDSC results, instead of meta-analyzing results across 41 traits. We identified 11 (resp. 20) annotation-trait pairs with significant  $\tau^*$  values for annotations derived from published (resp. boosted) scores (FDR<5%; Table C.7). These annotation-trait pairs spanned 2 (resp. 4) different annotations, all of which were also conditionally significant in the meta-analysis across traits (Figure 2-2B). Fourth, we assessed the heterogeneity of heritability enrichment and  $\tau^*$  across the 41 traits (see Methods). We determined that 10/12 annotations tested had significant heterogeneity in enrichment and 9/12 annotations tested had significant heterogeneity in  $\tau^*$ , implying substantial heterogeneity across traits (Table C.8). Fifth, we investigated the overlap between genes linked to each of our 9 conditionally significant annotations and 165 gene sets of biological importance (see Methods; Table C.9). As above, we consistently observed excess overlap for genes for which homozygous knockout in mice results in lethality<sup>92,93</sup> and high-pLI genes<sup>97</sup>, and depleted overlap for olfactory receptor genes<sup>107</sup> (Table C.11), implying that our new annotations can help identify biologically important genes. For example, the top (bottom) quintile of genes linked to the boosted CADD annotation had 1.7x (0.34x) excess overlap with high-pLI genes<sup>97</sup> vs. excess overlap of 0.8x (0.5x) for the top (bottom) quintile of genes linked to the published CADD<sup>8,9</sup> annotation (OR = 2.4,  $P < 6\text{e-}49$ ) (Figure D-3, Table C.12).

We conclude that 3 genome-wide Mendelian disease-derived annotations and 6 boosted annotations are conditionally informative for common disease, relative to baseline-LD model annotations.

## Informativeness of additional genome-wide scores for common disease

For completeness, we assessed the informativeness for common disease of 18 additional genome-wide scores not related to Mendelian disease, including 2 constraint-related scores<sup>10,16</sup>, 9 scores based on deep learning predictions of epigenetic marks<sup>172-174</sup>, and 7 gene-based scores<sup>42,97,106,175</sup> (see Table C.15). For each of the 18 additional scores, we constructed binary annotations based on top variants using 5 different thresholds and applied S-LDSC to the 41 traits, conditioning on the baseline-LD model<sup>3</sup> and meta-analyzing results across traits; in this analysis, we also conditioned on 8 Roadmap annotations<sup>17</sup> (4 annotations based on the union across cell types and 4 annotations based on the average across cell types, as in ref.<sup>160</sup>), as many of the additional scores pertain to regulatory elements, making this an appropriate conservative step.

We identified (Bonferroni-significant) conditionally binary annotations derived from 6 informative scores, including the top 1% of SNPs from CDTS<sup>16</sup> (enrichment = 9.3x (s.e. 0.75),  $\tau^* = 0.35$  (s.e. 0.06)) and the top 5% of SNPs from DeepSEA-H3K4me3<sup>172,173</sup> (enrichment = 3.9x (s.e. 0.23),  $\tau^* = 0.21$  (s.e. 0.04)) (Figure 2-3, Table 2.2 and Table C.16). CDTS (Context-Dependent Tolerance Score)<sup>16</sup> is a constraint score based on observed vs. expected variation in whole-genome sequence data; DeepSEA-H3K4me3 scores<sup>172,173</sup> are computed by training a deep learning model to predict chromatin marks using proximal reference genome sequence as features and aggregated across different cell types<sup>160</sup> (The DeepSEA annotations in Figure 3 were more significant than those analyzed in ref.<sup>160</sup>, because we optimized binary annotations based on top variants; however, no DeepSEA annotations were included in our combined joint model (see below)). 9 of the remaining 10 scores (excluding two that were not analyzed due to small annotation size) had derived binary annotations that were significantly enriched for disease heritability (after Bonferroni correction) but not conditionally informative (Table C.16).

We applied AnnotBoost to the 18 additional scores, and to the 47 main anno-

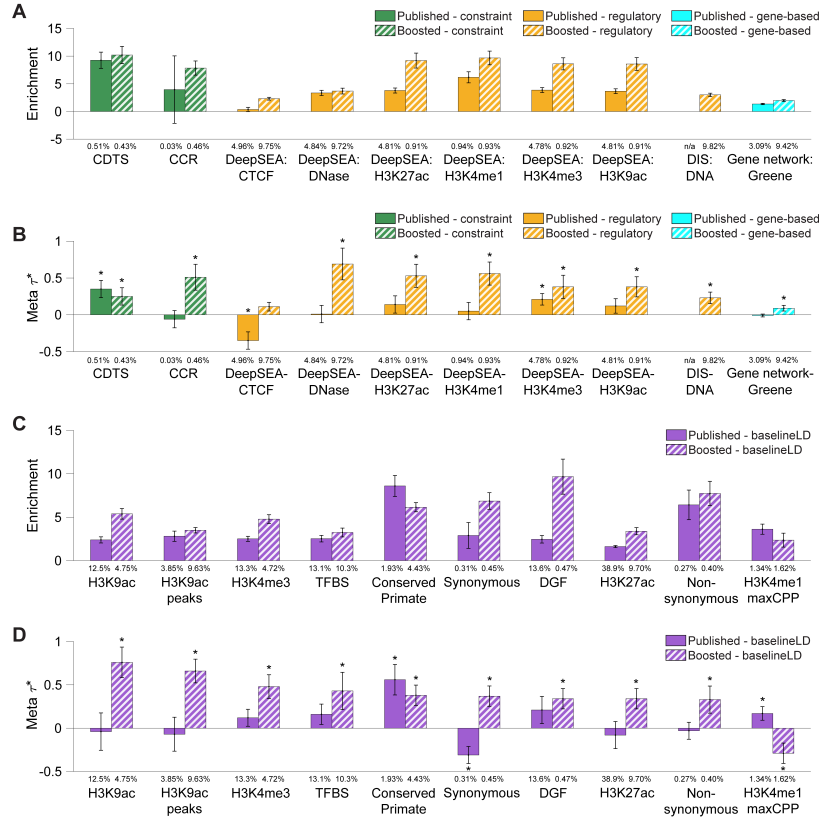


Figure 2-3: **Informativeness for common disease of binary annotations derived from 18 additional genome-wide scores + 47 baseline-LD model annotations and corresponding boosted scores.** We report (A) heritability enrichments of binary annotations derived from published and boosted additional genome-wide scores, meta-analyzed across 41 independent traits; (B) conditional  $\tau^*$  values, conditioning on the baseline-LD model and 8 Roadmap annotations (for annotations derived from published scores) or the baseline-LD model, 8 Roadmap annotations, and corresponding published annotations (for annotations derived from boosted scores); (C) heritability enrichments of binary annotations derived from published and boosted baseline-LD model annotations; and (D) conditional  $\tau^*$  values of binary annotations derived from published and boosted baseline-LD model annotations. In (A) and (B), we report results for the 10 most informative additional genome-wide scores (of 18 analyzed). In (C) and (D), we report results for the 10 most informative baseline-LD model annotations (of 47 analyzed). The percentage under each bar denotes the proportion of SNPs in the annotation; the proportion of top SNPs included in each annotation was optimized to maximize informativeness (largest  $|\tau^*|$  among Bonferroni-significant annotations, or top 5% if no annotation was Bonferroni-significant; top 5% was the average optimized proportion among significant annotations). Error bars denote 95% confidence intervals. In panels (B) and (D), \* denotes conditionally significant annotations. Numerical results are reported in Table C.16. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25.

tations of the baseline-LD model (Table C.15). Correlations between input scores and corresponding boosted scores varied widely ( $r = 0.005-0.93$ ; Table C.3C). AnnotBoost again attained high predictive accuracy in out-of-sample predictions of the input scores (AUROC = 0.55-1.00, AUPRC = 0.23-0.98; Table C.4); out-of-sample AUROCs closely tracked the correlations between input scores and corresponding boosted scores ( $r = 0.65$ ).

We again constructed binary annotations based on top genome-wide variants, using 6 different thresholds (ranging from top 0.1% to top 10% of genome-wide variants, as well as variants with boosted scores  $\geq 0.5$ ; see Methods). We assessed the informativeness for common disease of binary annotations derived from each of the 65 boosted scores using S-LDSC, conditioning on annotations from the baseline-LD model, the 8 Roadmap annotations, and (for the first 18 additional scores only) 5 binary annotations derived from the corresponding input scores (using all 5 thresholds). We identified conditionally informative binary annotations derived from boosted versions of 13/18 additional scores (including 11 previously uninformative scores and 2 previously informative scores) and 24/47 baseline-LD model annotations (Figure 2-3, Table 2.2 and Table C.16). Examples include the top 10% of SNPs from DeepSEA-DNase $\uparrow^{172,173}$  (enrichment = 3.7x (s.e. 0.27),  $\tau^* = 0.69$  (s.e. 0.11)), a previously uninformative score, the top 1% of SNPs from CCR $\uparrow^{10}$  (enrichment = 7.9x (s.e. 0.65),  $\tau^* = 0.51$  (s.e. 0.09)), a previously uninformative score, and the top 5% of SNPs from H3K9ac $\uparrow^{64}$  (enrichment = 5.4x (s.e. 0.31),  $\tau^* = 0.76$  (s.e. 0.09)), a baseline-LD model annotation. The CCR (Constrained Coding Regions) score<sup>10</sup> is a constraint score based on observed vs. expected variation in whole-exome sequence data; DeepSEA scores are described above. We note that the 18 additional scores included 7 gene-based scores, which did not perform well; 3 published gene-based scores and 4 boosted gene-based scores yielded conditionally significant binary annotations, but their  $\tau^*$  were small ( $-0.02$  to  $0.09$ ). Boosted versions of 3 of the remaining 5 additional scores and 20 of the remaining 23 baseline-LD model annotations had derived binary annotations that were significantly enriched for disease heritability (after Bonferroni correction) but not conditionally informative (Table C.16).

We performed 6 secondary analyses. First, for the 31 additional boosted scores which were conditionally significant and for which the underlying published scores had  $< 100\%$  of SNPs scored, we restricted the boosted scores to previously unscored variants and assessed the informativeness of the resulting binary annotations using S-LDSC. We determined that these annotations retained over half of the overall signals (average of 55% of absolute  $\tau^*$ ; Table C.17), implying that AnnotBoost both imputes and denoises existing scores. Second, we again investigated which features of the baseline-LD model contributed the most to the informativeness of the boosted annotations by applying SHAP<sup>5</sup>. We determined that a broad set features contributed to the predictions, including conservation-related features and LD-related features (as above), but also including regulatory features (e.g. H3K4me1, DGF, H3K9ac for boosted DeepSEA $\uparrow$ ) (Figure D-5). Third, we examined the trait-specific S-LDSC results, instead of meta-analyzing results across 41 traits. We identified 9 (resp. 78) annotation-trait pairs with significant  $tau^*$  values for annotations derived from published (resp. boosted) scores (FDR $<5\%$ ; Table C.7). These annotation-trait pairs spanned 5 (resp. 32) different annotations, most of which (3/5 published, 23/32 boosted) were also conditionally significant in the meta-analysis across traits (Figure 2-3B). Fourth, we assessed the heterogeneity of heritability enrichment and  $\tau^*$  across the 41 traits (see Methods). We determined that 64/81 annotations tested had significant heterogeneity in enrichment and 40/81 annotations tested had significant heterogeneity in  $\tau^*$ , implying substantial heterogeneity across traits (Table C.8). Fifth, we repeated the analyses of Figure 2-3 without including the 8 Roadmap annotations. We determined that the number of significant binary annotations increased (Table C.18), confirming the importance of conditioning on the 8 Roadmap annotations as an appropriate conservative step<sup>160</sup>. We further verified that including the 8 Roadmap annotations did not impact results from previous sections (Table C.19). Sixth, we investigated the overlap between genes linked to each of our 43 conditionally significant annotations and 165 gene sets of biological importance (see Methods; Table C.9). As above, we consistently observed excess overlap for genes for which homozygous knockout in mice results in lethality<sup>92,93</sup> and high-pLI genes<sup>97</sup>,

and depleted overlap for olfactory receptor genes<sup>107</sup> (Figure D-3, Table C.11, Table C.12), implying that our new annotations can help identify biologically important genes.

We conclude that 6 additional genome-wide annotations and 37 boosted annotations are conditionally informative for common disease, relative to baseline-LD model annotations.

## Constructing and evaluating combined heritability models

We constructed and evaluated two heritability models incorporating our new functional annotations: (i) a combined marginal model incorporating all binary annotations that were conditionally significant (conditional on the baseline-LD model), and (ii) a combined joint model incorporating only a subset of binary annotations that were jointly and conditionally significant (conditional on each other and the baseline-LD model). We constructed the combined marginal model by merging the baseline-LD model with the 64 conditionally significant annotations (derived from 11 published scores and 53 boosted scores; Table 2.2) (baseline-LD+marginal). We constructed the combined joint model by performing forward stepwise elimination to iteratively remove annotations that had conditionally non-significant  $\tau^*$  values after Bonferroni correction ( $P \geq 0.05/500 = 0.0001$ ) or  $\tau^* < 0.25$  (see Methods) (baseline-LD+joint). We have prioritized this type of combined joint model in previous work<sup>3,36,42,160,170</sup>, hypothesizing that there would be little value in retaining a large number of annotations containing redundant information; we note the substantial correlations between annotations in this study (Table C.20).

The combined joint model included 11 binary annotations derived from 3 published scores and 8 boosted scores (Figure 2-4, Table 2.2 and Table C.21). These 11 annotations are each substantially conditionally informative for common disease and include 5 boosted annotations with  $\tau^* > 0.5$  (e.g. boosted ReMM:  $\tau^* = 1.33$  (s.e. 0.12)); annotations with  $\tau^* > 0.5$  are unusual, and considered to be very important<sup>170</sup>. We note that the top 0.5% of SNPs from REVEL<sup>↑</sup><sup>132</sup> had significantly negative  $\tau^*$  ( $-0.95$  (s.e. 0.08)), as the annotation was significantly enriched for disease



heritability but less enriched than expected based on annotations from the combined joint model; as noted above, annotations with significantly positive or negative  $\tau^*$  are conditionally informative.

We performed two analyses to evaluate the combined joint model and the combined marginal model, compared to the baseline-LD model; these analyses evaluated aggregate heritability models, rather than individual functional annotations. First, we computed the average  $\log l_{SS}^1$  (an approximate likelihood metric) of each model, relative to a model with no functional annotations, across 30 common diseases and complex traits from the UK Biobank<sup>176</sup> (subset of 41 traits; Table C.1) ( $\Delta \log l_{SS}$ ; see Methods). Results are reported in Figure 2-5A and Table C.22. The combined joint model attained a +7.4% larger  $\Delta \log l_{SS}$  than the baseline-LD model ( $P < 6e-38$ ); the combined marginal model attained a +23.9% larger  $\Delta \log l_{SS}$  than the baseline-LD model ( $P < 2e-99$ ), including a significantly larger improvement for 30/30 traits analyzed (Figure 2-5B and Table C.22). The improvements were only slightly smaller when using Akaike Information Criterion (AIC) to account for increases in model complexity<sup>1</sup> (+7.0% and +20.3%; Table C.22).

Second, we assessed each model’s accuracy of classifying three different sets of fine-mapped SNPs (from 10 LD-, MAF-, and genomic-element-matched control SNPs in the reference panel<sup>40</sup>): 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, and 1,851 fine-mapped SNPs for 49 UK Biobank traits from Weissbrod et al.<sup>13</sup>. We note that with the exception of Weissbrod et al. fine-mapped SNPs (stringently defined by causal posterior probability  $\geq 0.95$ ; FDR  $< 0.05$ ), 95% credible fine-mapped SNPs likely include a large fraction of non-causal variants. We computed the AUPRC attained by the combined joint model and the combined marginal model, relative to a model with no functional annotations ( $\Delta \text{AUPRC}$ ), aggregated by training a gradient boosting model (multi-score analysis); we used odd (respectively even) chromosomes as training data to make predictions for even (respectively odd) chromosomes (see Methods). We note that this gradient boosting model uses disease data (fine-mapped SNPs), whereas AnnotBoost does not use disease data to construct

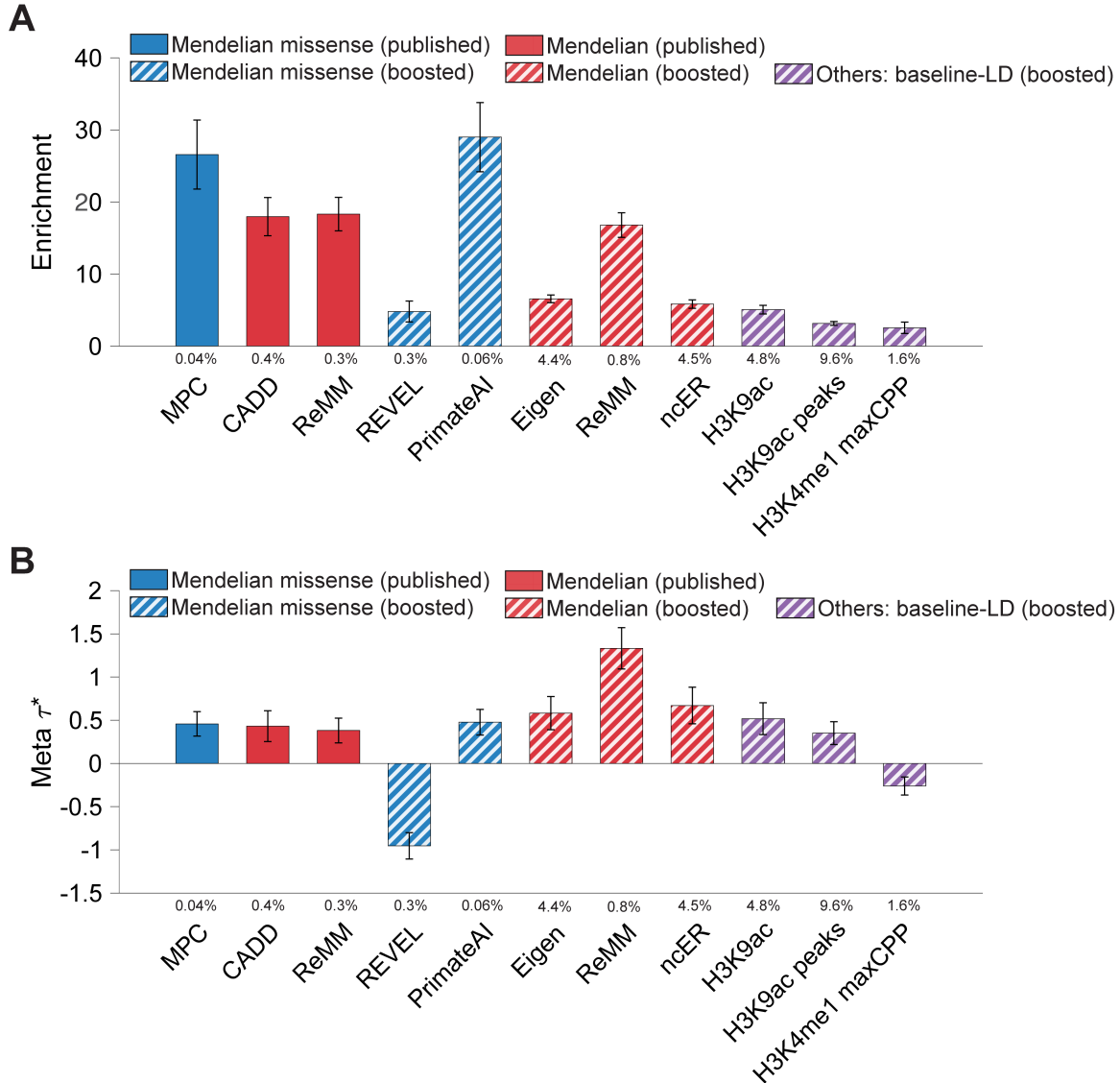


Figure 2-4: **Informativeness for common disease of 11 jointly significant binary annotations from combined joint model.** We report (A) heritability enrichment of 11 jointly significant binary annotations, meta-analyzed across 41 independent traits; (B) joint  $\tau^*$  values, conditioned on the baseline-LD model, 8 Roadmap annotations, and each other. We report results for the 11 jointly conditionally informative annotations in the combined joint model (S-LDSC  $\tau^* p < 0.0001$  and  $|\tau^*| \geq 0.25$ ). The percentage under each bar denotes the proportion of SNPs in the annotation. Error bars denote 95% confidence intervals. Numerical results are reported in Table C.21. Results for standardized enrichment, defined as enrichment times the standard deviation of annotation value (to adjust for annotation size), are reported in Table C.25.

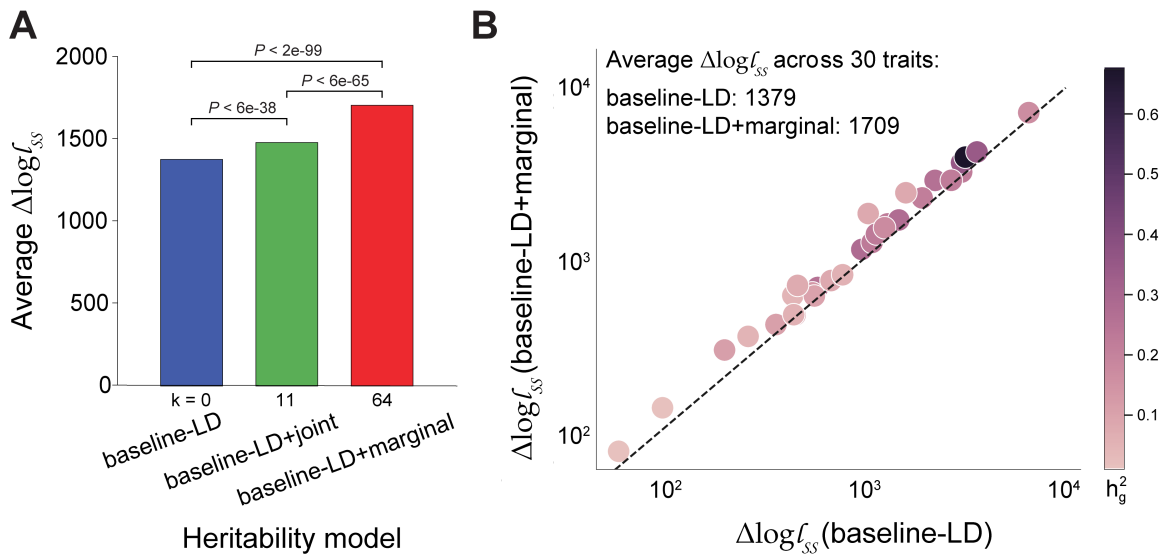


Figure 2-5: **Evaluation of improvement in heritability model fit.** We report (A) average  $\Delta \log \ell_{SS}$  (an approximate model likelihood metric<sup>1</sup>) across 30 UKBB traits; (B)  $\Delta \log \ell_{SS}$  of the baseline-LD and baseline-LD+marginal models for each trait.  $\Delta \log \ell_{SS}$  is computed as  $\log \ell_{SS}$  of a given model -  $\log \ell_{SS}$  of a model with no functional annotations (baseline-LD-nofunct model: MAF/LD annotations only). In panel (A),  $k$  denotes the number of new annotations beyond the baseline-LD model. Numerical results are reported in Table C.22.

boosted pathogenicity scores; specifically, our boosted scores do not use fine-mapped SNPs. The combined joint model attained a +2.5% to 6.9% larger  $\Delta$ AUPRC than the baseline-LD model (each  $P < 3e-28$ ); the combined marginal model attained a +4.9% to 21.3% larger  $\Delta$ AUPRC than the baseline-LD model (each  $P < 7e-100$ ); we obtained similar results using AUROC (Figure D-6, Table C.23). This improvement likely comes from non-linear interactions involving the boosted annotations, published annotations, and the baseline-LD model.

We performed 8 secondary analyses. First, we repeated  $\log l_{SS}$  analysis on the model with 19 new annotations with conditional  $\tau^* > 0.5$ ; we determined this model attained a +10.6% larger  $\Delta \log l_{SS}$  and +9.7% larger AIC than the baseline-LD model ( $P < 2e-50$ ) (Table C.22). Second, we applied SHAP<sup>5</sup> to investigate which features contributed the most to classification of fine-mapped SNPs; we determined that boosted scores often drove the predictions, validating the potential utility of boosted scores in functionally-informed fine-mapping (e.g. H3K9ac $\uparrow$ , CADD) (Figure D-7, Figure D-8). Third, we repeated the classification of fine-mapped SNPs using a single LD-, MAF-, and genomic-element-matched control variant (instead of 10 control variants) for each fine-mapped SNP, and obtained similar results (Table C.23). Fourth, we repeated the classification of fine-mapped disease SNPs analysis of Weissbrod et al. fine-mapped SNPs using 1,379 SNPs that were fine-mapped without using functional information<sup>13</sup> (to ensure that results were not circular), and obtained similar results (Figure D-6, Table C.23). Fifth, we computed the AUPRCs for classifying fine-mapped SNPs individually attained by each of 82 published and 82 boosted scores (single-score analysis), comparing results for boosted scores vs. the corresponding published scores. The boosted scores significantly outperformed the corresponding published scores in each case (66/82 to 80/82 scores; Figure D-9, Table C.24, and Table C.25). We also found that AUPRC and AUROC results for published and boosted scores were moderately correlated with S-LDSC results (up to  $r = 0.67$ ) for binary annotations derived from these scores, validating the S-LDSC results (Table C.26). Sixth, we repeated the single-score and multi-score analysis using 14,807 NHGRI GWAS SNPs<sup>14,15</sup>; we obtained similar results (Figure D-9, Fig-

ure D-6, Table C.24, Table C.23). Seventh, we computed genome-wide correlations between all annotations analyzed including baseline-LD model annotations (Table C.20). Several of the jointly significant annotations were strongly correlated (up to 0.73) with conservation-related annotations from the baseline-LD model, particularly binary GERP scores, consistent with our SHAP results (Figure D-2, Figure D-4 and Figure D-5). Eighth, we compared the informativeness of the baseline-LD model and the combined joint model. We identified the addition of 11 jointly significant annotations greatly reduced the informativeness of several existing baseline-LD annotations, including conservation-related annotations (e.g. conserved primate, binary GERP scores) and other annotations (e.g. coding, CpG content; see Figure D-10 and Table C.27), recapitulating the informativeness of 11 jointly significant annotations.

We conclude that the combined joint model and the combined marginal model both significantly outperformed the baseline-LD model, validating the informativeness of our new annotations for common disease. The improvement was much larger for the combined marginal model; this finding was surprising, in light of our previous work advocating for conservatively restricting to jointly significant annotations when expanding heritability models<sup>3,36,42,160,170</sup>. However, we caution that due to the much larger number of new annotations in the combined marginal model, the combined joint model may still be preferred in some settings.

## 2.6 Discussion

We analyzed the informativeness of a broad set of Mendelian disease-derived pathogenicity scores across 41 independent common diseases and complex traits to show that several annotations derived from published Mendelian disease-derived scores were conditionally informative for common disease after conditioning on the baseline-LD model. We further developed AnnotBoost, a gradient boosting-based machine learning framework to impute and denoise existing pathogenicity scores. We determined that annotations derived from boosted pathogenicity scores were even more informative for common disease, resulting in 64 marginally significant annotations and

11 jointly significant annotations and implying that Mendelian disease variants and common disease variants share similar properties. These variant-level results are substantially different from previous studies of gene-level overlap between Mendelian diseases and complex traits<sup>138–145</sup>. Notably, our new annotations produced significant improvements in heritability model fit and in classifying disease-associated, fine-mapped SNPs. We also detected significant excess overlap between genes linked to our new annotations and biologically important gene sets.

We note three key differences between AnnotBoost and previous approaches that utilized gradient boosting to identify pathogenic missense<sup>133</sup> and non-coding variants<sup>135,136</sup>. First, AnnotBoost uses a pathogenicity score as the only input and does not use disease data (e.g. ClinVar<sup>157</sup> or HGMD<sup>154</sup>). Second, AnnotBoost produces genome-wide scores, even when some SNPs are unscored by the input pathogenicity score. Third, AnnotBoost leverages 75 diverse features from the baseline-LD model<sup>3,127</sup>, significantly more than previous approaches<sup>133,135,136</sup>. Indeed, we determined that AnnotBoost produces strong signals even when conditioned on those approaches.

Our findings have several ramifications for improving our understanding of common disease. First, elucidating specific mechanistic links between Mendelian disease and common disease may yield important biological insights. Second, it is of interest to assess the informativeness for common disease of Mendelian disease pathogenicity scores that may be developed in the future, particularly after imputing and denoising these scores using AnnotBoost; this would further elucidate the shared properties between Mendelian disease variants and common disease variants. Third, annotations derived from published and boosted Mendelian pathogenicity scores can be used to improve functionally informed fine-mapping<sup>13,177–180</sup>, motivating their inclusion in future large-scale fine-mapping studies. (On the other hand, we anticipate that our new annotations will be less useful for improving functionally informed polygenic risk prediction<sup>181,182</sup> and association mapping<sup>183</sup>, because there is pervasive LD between SNPs in an annotation and SNPs outside of an annotation, such that these annotations do not distinguish which LD blocks contain causal signal.) Fourth, the larger

improvement for our combined marginal model versus our combined joint model (Figure 2-5A) advocates for a more inclusive approach to expanding heritability models, as compared to our previous work advocating for conservatively restricting to jointly significant annotations<sup>3,36,42,160,170</sup>. However, the combined marginal model suffers a cost of reduced interpretability (it contains a much larger number of new annotations, and it is unclear which of these annotations are providing the improvement), thus the combined joint model may still be preferred in some settings. Fifth, gene scores derived from published and boosted Mendelian pathogenicity scores can be used to help identify biologically important genes; we constructed gene scores by linking SNPs to their nearest gene, but better strategies for linking regulatory variants to genes<sup>18,184,185</sup> could potentially improve upon our results.

We note several limitations of our work. First, we focused our analyses on common disease (which are driven by common and low-frequency variants) and did not analyze Mendelian diseases (which are driven by very rare variants); the application of AnnotBoost to impute and denoise very rare pathogenic variants for Mendelian disease is a direction for future work. Second, we primarily report results that are meta-analyzed across 41 traits (analogous to previous studies<sup>2,3,36,42,160,170</sup>), but results and their interpretation may vary substantially across traits. Nonetheless, our combined marginal model produced a significant improvement in heritability model fit for 30/30 UK Biobank traits analyzed (Figure 2-5B). Third, S-LDSC is not well-suited to analysis of annotations spanning a very small proportion of the genome, preventing the analysis of a subset of published pathogenicity scores; nonetheless, our main results attained high statistical significance. Fourth, we restricted all of our analyses to European populations, which have the largest available GWAS sample size. However, we expect our results to be generalizable to other populations, as functional enrichments have been shown to be highly consistent across ancestries<sup>178,186,187</sup>; we note that assessing functional enrichments in admixed populations<sup>188</sup> would require the application of an unpublished extension of S-LDSC<sup>189</sup>. Fifth, the gene-based SNP scores that we analyzed did not perform well, perhaps because they were defined using 100kb windows, a crude strategy employed in previous work<sup>4,36,42</sup>; better

strategies for linking regulatory variants to genes<sup>18,184,185</sup> (as shown in above gene scores) could potentially improve upon those results. Despite these limitations, the imputed and denoised pathogenicity scores produced by our AnnotBoost framework have high potential to improve gene discovery and fine-mapping for common disease.

## 2.7 Web resources

All published and boosted pathogenicity scores and binary annotations and SHAP results are available at [https://alkesgroup.broadinstitute.org/LDSCORE/Kim\\_annotboost/](https://alkesgroup.broadinstitute.org/LDSCORE/Kim_annotboost/). GWAS summary statistics are available at [https://alkesgroup.broadinstitute.org/LDSCORE/independent\\_sumstats/](https://alkesgroup.broadinstitute.org/LDSCORE/independent_sumstats/). baseline-LD annotations (v.2.1) are available at <https://alkesgroup.broadinstitute.org/LDSCORE/>. The 1000 Genomes Project Phase 3 data are available at <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502>. 165 reference gene sets are available at <https://github.com/samskim/networkconnectivity>.

AnnotBoost source code is provided here: <https://github.com/samskim/annotboost/>. This chapter primarily uses the S-LDSC software (<https://github.com/bulik/ldsc>). SumHer software for computing  $\log l_{SS}$  is available at <http://dougsped.com/sumher/>.



# Chapter 3

## Leveraging single-cell ATAC-seq to identify disease-critical fetal and adult brain cell types

### 3.1 Preface

This chapter is yet to be published but available as a preprint in *bioRxiv* as "Leveraging single-cell ATAC-seq to identify disease-critical fetal and adult brain cell types" (2021; see ref.<sup>190</sup>) with coauthors Karthik Jagadeesh, Kushal K. Dey, Amber Z. Shen, Soumya Raychaudhuri, Manolis Kellis, and Alkes L. Price. I am also grateful to Tiffany Amariuta, Katie Siewert, and Martin Zhang for helpful discussion.

### 3.2 Abstract

Prioritizing disease-critical cell types by integrating genome-wide association studies (GWAS) with functional data is a fundamental goal. Single-cell chromatin accessibility (scATAC-seq) and gene expression (scRNA-seq) have characterized cell types at high resolution, and early work on integrating GWAS with scRNA-seq has shown promise, but work on integrating GWAS with scATAC-seq has been limited. Here, we identify disease-critical fetal and adult brain cell types by integrating GWAS summary

statistics from 28 brain-related diseases and traits (average  $N=298K$ ) with 3.2 million scATAC-seq and scRNA-seq profiles from 83 cell types. We identified disease-critical fetal (resp. adult) brain cell types for 22 (resp. 23) of 28 traits using scATAC-seq data, and for 8 (resp. 17) of 28 traits using scRNA-seq data. Notable findings using scATAC-seq data included highly significant enrichments of fetal photoreceptor cells for major depressive disorder, fetal ganglion cells for BMI, fetal astrocytes for ADHD, and adult VGLUT2 excitatory neurons for schizophrenia. Our findings improve our understanding of brain-related diseases and traits, and inform future analyses of other diseases/traits.

### 3.3 Introduction

Genome-wide association studies (GWAS) have been successful in identifying disease-associated loci, occasionally producing valuable functional insights<sup>43,191</sup>. Identifying disease-critical cell types is a fundamental goal for understanding disease mechanisms, designing functional follow-ups, and developing disease therapeutics<sup>192</sup>. Several studies have identified disease-critical tissues and cell types using bulk chromatin<sup>4,17,64,65,114,193</sup> and/or gene expression data<sup>2,4,194,195</sup>.

With the emergence of single-cell profiling of diverse tissues and cell types<sup>196–200</sup>, several studies have identified disease-critical cell types at high resolution using and single-cell chromatin accessibility (scATAC-seq)<sup>199–202</sup> and single-cell gene expression (scRNA-seq)<sup>20,194,203</sup>. However, compared to scRNA-seq data, scATAC-seq data has been less well-studied for identifying disease-critical cell types. In addition, while it is widely known that biological processes in the human brain vary with developmental stage<sup>204–208</sup>, the impact on disease risk of cell types in different developmental stages of the brain has not been widely explored.

Here, we sought to infer disease-critical cell types by analyzing scATAC-seq and scRNA-seq data derived from single-cell profiling of over 3 million cells from fetal and adult human brains. We analyzed 83 brain cell types from 4 single-cell datasets<sup>197–200</sup> across 28 brain-related diseases and complex traits (average  $N = 298K$ ). We deter-

mined that both scATAC-seq and scRNA-seq data were highly informative for identifying disease-critical cell types; surprisingly, scATAC-seq data was with somewhat more informative in the data that we analyzed.

## 3.4 Methods

### Set of 28 distinct brain-related traits analyzed

We considered 146 GWAS summary association statistics, including 83 traits from the UK Biobank and 63 traits from publicly available sources. These traits had z-scores of total SNP heritability of at least 6 (computed using S-LDSC with the baseline-LD (v.2.2) model). We selected 31 brain-related traits based on previous studies<sup>4,20,42,200,203</sup>. We removed 3 summary statistics that had a genetic correlation of at least 0.9 with at least one trait among these 31 traits, retaining a final set of 28 distinct brain-related traits (where 7 of whom are from UK Biobank). The genetic correlation among these 28 brain traits are shown in Table E.25.

As a control set of traits, we additionally analyzed 6 distinct non-brain-related traits: coronary artery disease, bone mineral density, rheumatoid arthritis, type 2 diabetes, sunburn occasion, and breast cancer. These 6 traits had a similar sample size and heritability z-scores as above 28 brain-related traits.

### Genomic annotations and the baseline model

We define a binary genomic annotation as a subset of SNPs in a predefined reference panel. We restrict our analysis to SNPs with a minor allele frequency (MAF)  $\geq 0.5\%$  in 1000 Genomes<sup>40</sup> (see Web resources).

The baseline model<sup>3</sup> (v.1.2; see Table E.2) contains 53 functional annotations (see Web resources). These annotations include genomic elements (e.g. coding, enhancer, UTR), regulatory elements (e.g. histone marks), and evolutionary constraint. We included the baseline model, consistent with ref. finucane2018, vandegeijn2020, when assessing the heritability enrichment of the cell-type annotations.

## Single-cell ATAC-seq data analyzed

We considered single-cell ATAC-seq data for fetal brains from Domcke et al.<sup>199</sup> (donor  $N = 26$ ) and adult brains (isocortex, striatum, hippocampus, and substantia nigra) of cognitively healthy individuals from Corces et al.<sup>200</sup> ( $N = 10$ ). We used the cell-type-specific chromatin accessible peaks without modifications (see Web resources). In short, these peaks refer to MACS2<sup>41</sup> peak regions, excluding the ENCODE blacklist regions. For the Domcke et al. data, authors called peaks on each tissue sample and then generated a masterlist of all peaks across all tissue samples. Thus, we used the provided cell-type specificity scores (a higher score indicates a peak that is more specific to a given cell type).

## Single-cell RNA-seq data analyzed

We considered single-cell RNA-seq data for fetal brains from Cao et al.<sup>198</sup> ( $N = 28$ ) and single-cell RNA-seq data for non-fetal brains (prefrontal cortex and anterior cingulate cortex) from Velmeshev et al.<sup>197</sup> ( $N = 31$ ). For Cao et al. data, we processed data from three brain-related organs: cerebellum, cerebrum, and eye. We followed the processing protocol as outlined in ref. jagadeesh2021. For each data set, we downloaded metadata for each cell including the total number of reads and sample ID. We then transformed each expression matrix to  $\log_2(\text{TP10K}+1)$  unit. We performed a dimensional reduction using a principal component analysis with the top 2,000 highly variable genes, batch correction using Harmony<sup>209</sup>, and applied the Leiden graph clustering method<sup>210</sup>.

To obtain differentially expressed (DE) genes for each cell type, we applied a non-parametric Wilcoxon rank-sum test between gene expression from focal cell type vs. gene expression from all the other cell types. We transformed the per-gene DE p-value to lie between 0 and 1, by applying min-max normalization on  $-2\log(\text{p-value})$ , indicating a relative importance of each gene in each cellular process.

To build SNP annotations from gene program, we employed an enhancer-gene linking strategy from the union of the Roadmap<sup>17</sup> and Activity-By-Contact (ABC<sup>18,19</sup>)

strategies. Because we are focusing on the brain-related traits, we used the brain-specific enhancer-gene links. Using these brain enhancer-gene links, we have built a probabilistic SNP annotation representing the relative weight in each cellular processes.

## Enrichment and $\tau^*$ metrics

We used stratified LD score regression (S-LDSC<sup>2,3</sup>) to assess the contribution of an annotation to disease and complex trait heritability.

Let  $a_{cj}$  represent the (binary or probabilistic) annotation value of the SNP  $j$  for the annotation  $c$ . S-LDSC assumes the variance of per normalized genotype effect sizes is a linear additive contribution to the annotation  $c$ :

$$\text{Var}(\beta_j) = \sum_c a_{cj} \tau_c \quad (3.1)$$

where  $\tau_c$  is the per-SNP contribution of the annotation  $c$ . S-LDSC estimates  $\tau_c$  using the following equation:

$$\text{E}[\chi_j^2] = N \sum_c \ell(j, c) \tau_c + 1 \quad (3.2)$$

where  $N$  is the sample size of the GWAS and  $\ell(j, c)$  is the LD score of the SNP  $j$  to the annotation  $c$ . The LD score is computed as follow  $\ell(j, c) = \sum_k a_{ck} r_{jk}^2$  where  $r_{jk}$  is the correlation between the SNPs  $j$  and  $k$ .

We used two metrics to assess the informativeness of an annotation. First, the standardized effect size ( $\tau^*$ ), the proportionate change in per-SNP heritability associated with a one standard deviation increase in the value of the annotation (conditional on all the other annotations in the model), is defined as follows:

$$\tau_{c^*} = \frac{\tau_c \text{sd}(C)}{h_g^2/M} \quad (3.3)$$

where  $\text{sd}(C)$  is the standard deviation of the annotation  $c$ ,  $h_g^2$  is the estimated SNP-

heritability, and  $M$  is the number of variants used to compute  $h_g^2$  (in our experiment,  $M$  is equal to 5,961,159, the number of common SNPs in the reference panel). The significance for the effect size for each annotation, as mentioned in previous studies<sup>3,36,42</sup>, is computed as  $(\frac{\tau^*}{\text{se}(\tau^*)} \sim N(0,1))$ , assuming that  $\frac{\tau^*}{\text{se}(\tau^*)}$  follows a normal distribution with zero mean and unit variance.

Second, enrichment of the binary and probabilistic annotation is the fraction of heritability explained by SNPs in the annotation divided by the proportion of SNPs in the annotation, as shown below:

$$\text{Enrichment} = \frac{\%h_g^2(C)}{\%\text{SNP}(C)} = \frac{\frac{h_g^2(C)}{h_g^2}}{\frac{\sum_j a_{jc}}{M}} \quad (3.4)$$

where  $h_g^2(C)$  is the heritability captured by the  $c$ th annotation. When the annotation is enriched for trait heritability, the enrichment is  $> 1$ ; the overlap is greater than one would expect given the trait heritability and the size of the annotation. The significance for enrichment is computed using the block jackknife as mentioned in previous studies<sup>2,4,36,42</sup>). The key difference between enrichment and  $\tau^*$  is that  $\tau^*$  quantifies effects that are unique to the focal annotation after conditioning on all the other annotations in the model, while enrichment quantifies effects that are unique and/or non-unique to the focal annotation.

We used European samples in 1000G<sup>40</sup> as reference SNPs and HapMap 3<sup>72</sup> SNPs as regression SNPs (see Web resources). We excluded SNPs with marginal association statistics  $> 80$  and SNPs in the major histocompatibility complex region. In all our analyses, we used the p-value of  $\tau^*$  as our primary metric to estimate the effect sizes conditional on known annotations (by including the baseline model as recommended previously<sup>4,211</sup>). We excluded trait-annotation pairs with negative  $\tau^*$ , consistent with previous studies<sup>42,199,212</sup>. We controlled the per-dataset false discovery rate (FDR) across 28 brain traits.

## Gene set enrichment analysis using GREAT<sup>213</sup>

We performed gene set enrichments on each cell-type annotations for the gene ontology (GO) biological process, cellular component, and molecular function. We used GREAT (v.4.0.4) with its default setting, where each gene is assigned a regulatory domain (for proximal: 5kb upstream, 1kb downstream of the TSS; for distal: up to 1Mb). Because annotations from the scRNA-seq were probabilistic, we limited to regions with gene membership probability  $\geq 0.8$  for gene set enrichment analysis. We used all regions for the scATAC-seq annotations as an input. We defined significant results as those with the FDR-corrected binomial test p-value  $< 0.05$ .

## 3.5 Results

### Overview of methods

We define a *cell-type annotation* as an assignment of a binary or probabilistic value between 0 and 1 to each SNP in the 1000 Genomes European reference panel<sup>40</sup>, representing the estimated contribution of that SNP to gene regulation in a particular cell type. Here, we constructed cell-type annotations for 4 datasets: (1) fetal brain scATAC-seq<sup>199</sup> (the number of cell types ( $C$ ) = 14), (2) fetal brain scRNA-seq data<sup>198</sup> ( $C$  = 34), (3) adult brain scATAC-seq<sup>200</sup> ( $C$  = 18), and (4) scRNA-seq data<sup>197</sup> ( $C$  = 17) (see Web resources).

For scATAC-seq cell-type annotations, we used the chromatin accessible peaks (MACS2<sup>41</sup> peak regions) provided by ref. domcke2020, corces2020. These peaks correspond to accessible regions for transcription factor binding, indicative of active gene regulation. For scRNA-seq cell-type annotations, we annotated SNPs linked to specifically expressed genes in a given cell type<sup>4</sup> using brain-specific enhancer-gene links<sup>17-20</sup>.

We assessed the heritability enrichments of the resulting cell-type annotations by applying S-LDSC<sup>2</sup> across 28 distinct brain-related diseases and traits (pairwise genetic correlation  $< 0.9$ ; average  $N=298K$ ; Table E.1) to identify significant disease-cell

type associations (Figure F-1, Table E.2). For each disease-cell type pair, we estimated the heritability enrichment<sup>2</sup> (the proportion of heritability explained divided by the proportion of annotated SNPs) and standardized effect size<sup>3</sup> ( $\tau^*$ , defined as the proportionate change in per-SNP heritability associated to a one standard deviation increase in the value of the annotation, conditional on other annotation). We assessed the statistical significance of disease-cell type associations based on per-dataset FDR  $< 5\%$  (for each of 5 datasets, aggregating diseases and cell types) based on p-values for positive  $\tau^*$ , as  $\tau^*$  quantifies effects that are unique to the cell-type annotation. We conditioned the analyses on a broad set of coding, conserved, and regulatory annotations from the baseline model<sup>2</sup>. For scATAC-seq annotations, we additionally conditioned on the union of open chromatin regions across all cell types (consistent with ref. ziffra2020, freimer2021), a conservative step to ensure cell-type specificity (see Discussion). For scRNA-seq annotations, we additionally conditioned on the union of brain-specific enhancer-gene links across all genes analyzed (consistent with ref. jagadeesh2021). We did not condition on the LD-related annotations included in the baseline-LD model of ref.<sup>3,127</sup>, as these annotations reflect the action of negative selection, which may obscure cell-type-specific signals<sup>211</sup>. Further details are provided in the Methods section. We have publicly released cell-type annotations analyzed in this study and source codes for all primary analyses (see Web resources).

## Identifying disease-critical cell types using fetal brain data

We sought to identify disease-critical cell types using fetal brain data, across 28 distinct brain-related diseases and traits. We analyzed 14 fetal brain cell types from scATAC-seq data<sup>199</sup> (donor size = 26; fetal age of 72-129 days) and 34 fetal brain cell types from scRNA-seq data<sup>198</sup> (donor size = 28; fetal age of 89-125 days) (Table E.3; see Methods).

We first analyzed fetal brain scATAC-seq data spanning 14 cell types<sup>199</sup>. We identified 152 significant disease-cell type pairs (FDR  $< 5\%$  for positive  $\tau^*$  conditional on other annotations; Table 3.1, Table 3.2, Figure 3-1A, Table E.4). Consistent with previous genetic studies<sup>4,200,203</sup>, we identified strong enrichments of excitatory (i.e. gluta-

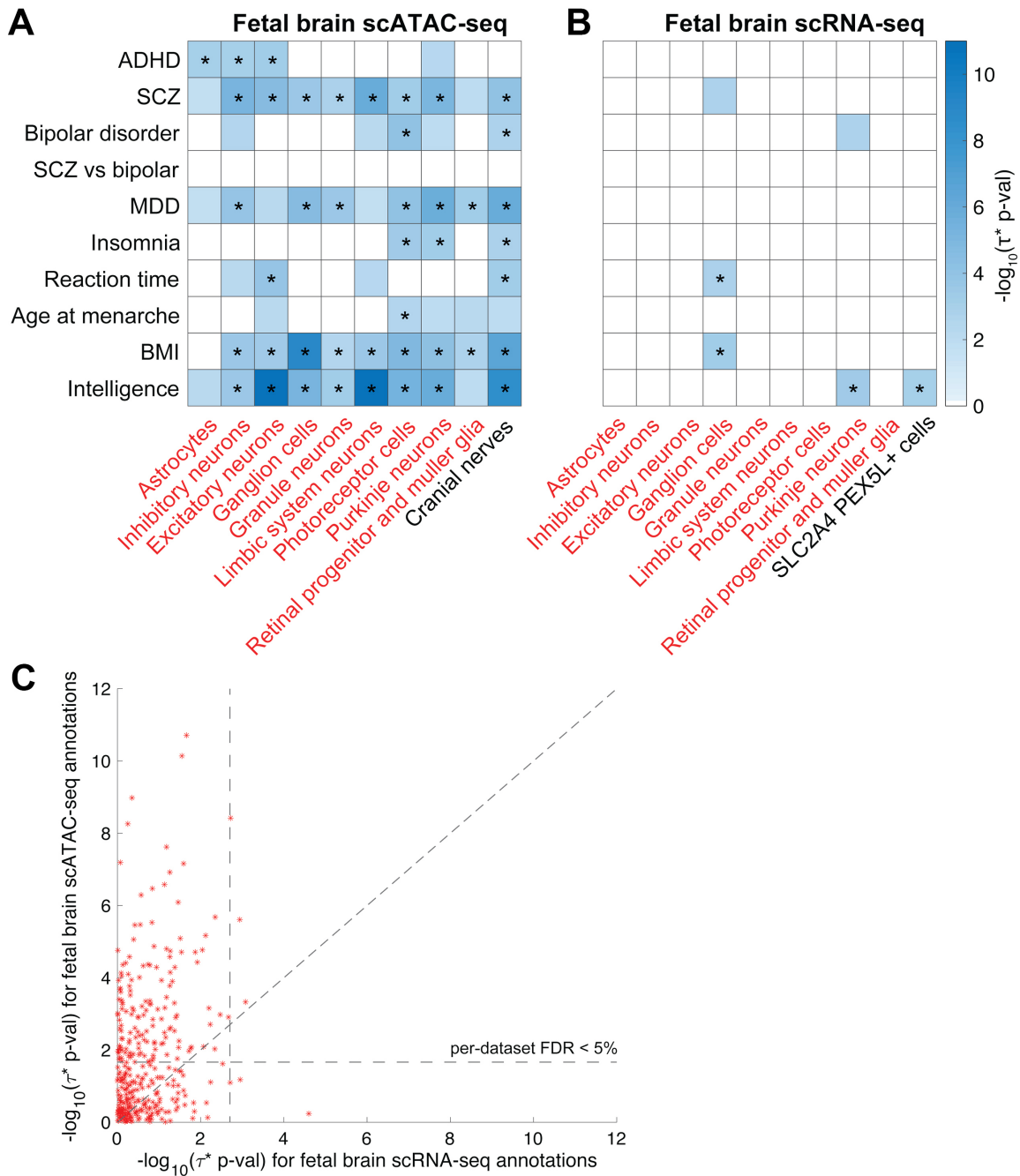


matergic) neurons in psychiatric and neurological disorders, including schizophrenia (SCZ), major depressive disorder (MDD), and attention deficit hyperactivity disorder (ADHD) (Figure 3-1A). In particular, the role of glutamatergic neurons in MDD is well-supported, as evident from decreased glutamatergic neurometabolite levels in subjects with depression<sup>214</sup>.

	Fetal scATAC	Fetal scRNA	Adult scATAC	Adult scRNA
Brain cell types	14	34	18	17
Total disease-cell type pairs	392	952	504	476
Significant disease-cell-type pairs	152	9	168	64
Significant diseases (out of 28)	22	8	23	17
Data source	ref. <sup>199</sup>	ref. <sup>198</sup>	ref. <sup>200</sup>	ref. <sup>197</sup>

Table 3.1: **Summary of findings.** For each of 4 single-cell chromatin and gene expression data sets analyzed, we report the number of brain cell types analyzed, the total number of disease-cell type pairs analyzed (based on 28 diseases/traits), the number of significant disease-cell type pairs (FDR < 5% for positive  $\tau^*$ ), and the number of diseases/traits with a significant disease-cell type pair.

Our results also highlight several disease-cell type associations that have not (to our knowledge) previously been reported in analyses of genetic data (Table 3.2). First, photoreceptor cells were enriched in insomnia. Photoreceptor cells, present in the retina, convert light into signals to the brain, and thus play an essential role in circadian rhythms<sup>215</sup>, explaining their potential role in insomnia. Second, photoreceptor cells were also enriched in MDD, a genetically uncorrelated trait ( $r = -0.01$ ), and neuroticism ( $r = 0.02$ ). Recent studies support the relationship between the degeneration of photoreceptors and anxiety and depression<sup>216</sup>, suggesting that mental illness is not only a consequence but also a cause of visual impairment. Third, inhibitory (GABAergic) neurons were enriched in SCZ, consistent with GABA dysfunction in the cortex of schizophrenics<sup>217</sup>. Fourth, ganglion cells were enriched in BMI. Ganglion cells are the projection neurons of the retina, relaying information from bipolar and amacrine cells to the brain. Patients with morbid obesity display significant differences in retinal ganglion cells, retinal nerve fiber layer thickness, and choroidal thickness<sup>218</sup>. Fifth, purkinje neurons were enriched in insomnia (and two ge-



**Figure 3-1: Disease enrichments of cell-type annotations derived from fetal brain.** We report (A)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 14) fetal brain scATAC-seq cell type annotations; (B)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 34) fetal brain scRNA-seq cell type annotations; and (C) comparison of results for 13 cell types included in both fetal brain scATAC-seq and scRNA-seq data. In A-B, only statistically significant results (FDR > 5%) are colored ( $-\log_{10}(\text{p-value}) \geq 1.67$  for scATAC-seq,  $\geq 2.70$  for scRNA-seq). In A-B, cell types appearing in both datasets are denoted in red font. Numerical results for all diseases/traits and cell types are reported in Table E.4, Table E.6, and Table E.7. ADHD: attention deficit hyperactivity disorder, SCZ: schizophrenia, MDD: major depressive disorder, BMI: body mass index.

netically distinct traits, sleep duration ( $r = -0.03$ ) and chronotype ( $r = -0.03$ ). While purkinje neurons play a major role in controlling motor movement, they also regulate the rhythmicity of neurons, consistent with a role in impacting sleep<sup>219</sup>. Sixth, astrocytes were enriched in ADHD. Astrocytes perform various functions including synaptic support, control of blood flow, and axon guidance<sup>220</sup>. In particular, Nagai et al.<sup>221</sup> highlighted the role of the astrocyte Gi-coupled GABA<sub>B</sub> pathway activation resulting in ADHD-like behaviors in mice.

Disease/trait	Cell type	Data source	$\tau^*$ (SE)	p-value( $\tau^*$ )	q-value
Insomnia <sup>222</sup>	Photoreceptor cells	Fetal brain scATAC	0.81 (0.23)	$4.58 \times 10^{-4}$	$2.02 \times 10^{-3}$
MDD <sup>38</sup>	Photoreceptor cells	Fetal brain scATAC	0.67 (0.17)	$8.45 \times 10^{-5}$	$5.26 \times 10^{-4}$
SCZ <sup>30</sup>	Inhibitory neurons	Fetal brain scATAC	0.98 (0.22)	$6.14 \times 10^{-6}$	$7.08 \times 10^{-5}$
BMI <sup>38</sup>	Ganglion cells	Fetal brain scATAC	0.55 (0.09)	$8.72 \times 10^{-10}$	$6.84 \times 10^{-8}$
Insomnia <sup>223</sup>	Purkinje neurons	Fetal brain scATAC	0.73 (0.21)	$6.01 \times 10^{-4}$	$2.48 \times 10^{-3}$
ADHD <sup>224</sup>	Astrocytes	Fetal brain scATAC	1.05 (0.32)	$9.68 \times 10^{-4}$	$3.72 \times 10^{-3}$
Reaction time <sup>225</sup>	Ganglion cells	Fetal brain scRNA	0.45 (0.14)	$1.26 \times 10^{-3}$	$3.93 \times 10^{-2}$
MDD <sup>226</sup>	BDNF excitatory neurons	Adult brain scATAC	1.31 (0.20)	$1.14 \times 10^{-10}$	$4.10 \times 10^{-9}$
Bipolar disorder <sup>227</sup>	Parvalbumin interneurons	Adult brain scATAC	1.23 (0.28)	$7.35 \times 10^{-6}$	$8.23 \times 10^{-5}$
SCZ <sup>30</sup>	VGLUT2 excitatory neurons	Adult brain scATAC	1.31 (0.24)	$4.14 \times 10^{-8}$	$7.45 \times 10^{-7}$
Intelligence <sup>31</sup>	Corticofugal projection neurons	Adult brain scRNA	0.76 (0.14)	$1.17 \times 10^{-7}$	$1.39 \times 10^{-5}$

**Table 3.2: Notable disease-cell type associations.** We report the disease/trait, cell type, data source, standardized effect size ( $\tau^*$ ), p-value for positive  $\tau^*$ , and FDR q-value for selected results. Full results are reported in Table E.4, Table E.6, Table E.14, Table E.15. A description of diseases/traits analyzed is provided in Table E.1. MDD: major depressive disorder, SCZ: schizophrenia, BMI: body mass index, ADHD: attention deficit hyperactivity disorder.

We next analyzed fetal brain scRNA-seq data spanning 34 cell types<sup>198</sup> (of which 13 were also included in fetal brain scATAC-seq data; Table E.5). We identified 9 significant disease-cell type pairs (FDR < 5% for positive  $\tau^*$  conditional on other annotations; Table 3.1, Table 3.2, Figure 3-1B, Table E.6). When restricting to the 7 significant disease-cell type pairs corresponding to the 13 cell types included in both scATAC-seq and scRNA-seq data, 6 of 7 were also significant in analyses of scATAC-seq data. (Conversely, when restricting to the 134 significant disease-cell type pairs in scATAC-seq data corresponding to cell types included in scRNA-seq data, 7 of 134 were also significant in scRNA-seq data.) In particular, the enrichment

of retinal ganglion cells in reaction time ( $p = 1.26 \times 10^{-3}$  in scRNA-seq data, FDR  $q = 0.039$ ) was non-significant in scATAC-seq data ( $p = 0.028$ , FDR  $q = 0.060$ ). The enrichment of retinal ganglion cells in reaction time has not (to our knowledge) previously been reported in analyses of genetic data. Previous genetic analyses have focused on enrichments of cerebellum and brain cortex in reaction time<sup>225</sup>, but the involvement of retinal ganglion cells in receiving visual information and propagating it to the rest of the brain is consistent with a role in visual reaction time<sup>228</sup>.

We compared the results for 13 fetal brain cell types included in both the scATAC-seq and scRNA-seq datasets (Table E.5 and Table E.7). While scATAC-seq and scRNA-seq cell-type annotations for matched cell types were approximately uncorrelated to each other ( $r = 0.01 - 0.06$ ; Table E.8), the corresponding  $-\log_{10}(\text{p-values})$  for positive  $\tau^*$  were moderately correlated ( $r = 0.24$ ; Figure 3-1C), confirming the shared biological information. We observed more significant p-values for scATAC-seq than for scRNA-seq in these data sets (see Discussion).

We performed 4 secondary analyses. First, we analyzed enrichments of both scATAC-seq and scRNA-seq brain cell types in 6 control, non-brain-related diseases and complex traits. As expected, we did not find significant enrichments of brain cell types in 6 control traits (Table E.9 and Table E.10). Second, we performed gene set enrichment analysis using GREAT<sup>213</sup> for both scATAC-seq and scRNA-seq cell-type annotations. As expected, we identified significant enrichments in relevant gene sets (e.g. “photoreceptor cell differentiation” for photoreceptor cells from scATAC-seq; “negative regulation of cell projection organization” for ganglion cells from scRNA-seq); Table E.11). Third, for the fetal scRNA-seq data<sup>198</sup>, we constructed annotations based on a +/- 100kb window-based strategy (previously used in ref.<sup>4</sup>) instead of brain-specific enhancer-gene links<sup>17-19</sup> (used in ref.<sup>20</sup>). We identified 22 significant disease-cell type pairs (Table E.12), vs. only 9 using brain-specific enhancer-gene links (although we observed a much stronger opposite trend in adult scRNA-seq data; see below). Fourth, we analyzed bulk chromatin data (7 chromatin marks) spanning 5 fetal brain tissues<sup>193</sup> (age 52-142 days). We identified 541 significant disease-tissue-chromatin mark triplets spanning 26 of 28 brain-related traits (Table E.13). These

results are included for completeness, but cannot achieve the same cell-type specificity as analyses of single-cell data.

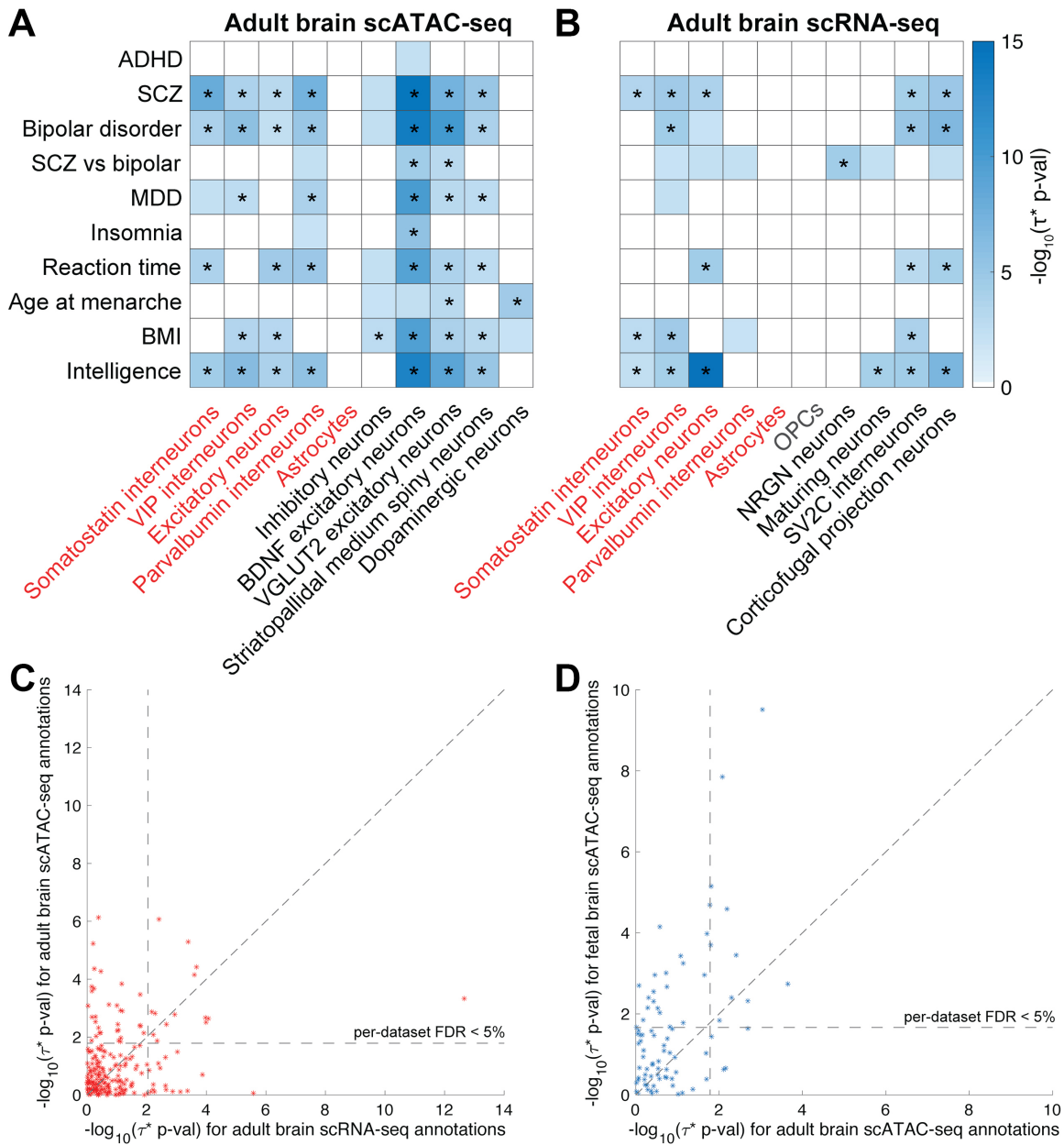
## Identifying disease-critical cell types using adult brain data

We sought to identify disease-critical cell types using adult brain data, across 28 distinct brain-related diseases and traits. Analysis of brains with varying developmental stages might elucidate biological mechanisms, as brains undergo changes in cell type composition and gene expression<sup>208,229</sup>. We analyzed 18 adult brain cell types from scATAC-seq data<sup>200</sup> (donor size = 10; age of 38-95 years) and 17 adult brain cell types from scRNA-seq data<sup>197</sup> (donor size = 31; age of 4-22 years) (Table E.3; see Methods). For brevity, we use the term “adult” to refer to child and adult donors who have surpassed the fetal development stage.

We first analyzed adult brain scATAC-seq data spanning 18 cell types<sup>200</sup>. We identified 168 significant disease-cell type pairs (FDR < 5% for positive  $\tau^*$  conditional on other annotations; Table 3.1, Table 3.2, Figure 3-2A, Table E.14). Consistent with previous genetic studies<sup>4,200,202,212</sup>, we identified strong enrichments of excitatory neurons in bipolar disorder and SCZ (genetic correlation = 0.70) (Figure 3-2A). Although an analysis of mouse scATAC-seq identified a significant enrichment of excitatory neurons in SCZ cases vs. bipolar cases<sup>202</sup>, we did not replicate this finding ( $p = 0.66$  for positive  $\tau^*$ ).

Our results also highlight disease-cell type associations that have not (to our knowledge) previously been reported in analyses of genetic data (Table 3.2). First, brain-derived neurotrophic factor (BDNF) excitatory neurons were highly enriched in MDD (and several other diseases/traits, including bipolar disorder and SCZ). BDNF is involved in supporting survival of existing neurons and differentiating new neurons, and decreased BDNF levels have been observed in untreated MDD<sup>230</sup>, bipolar<sup>231</sup> and SCZ cases<sup>232</sup>. Previous studies identified an enrichment of excitatory neurons in MDD<sup>212</sup>. Second, parvalbumin interneurons were enriched in bipolar disorder (and SCZ). Decreased expression and diminished function of parvalbumin interneurons in regulating balance of excitation and inhibition have been observed in bipolar disorder

and SCZ cases<sup>233,234</sup>. Third, vesicular glutamate transporter (VLUGT2) excitatory neurons were enriched in SCZ (as well as bipolar disorder and intelligence). VLUGT2 knock-out mice display glutamatergic deficiency, diminished maturation of pyramidal neuronal architecture, and impaired spatial learning and memory<sup>235</sup>, supporting a role in SCZ and intelligence.



**Figure 3-2: Disease enrichments of cell-type annotations derived from adult brain.** We report (A)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 18) adult brain scATAC-seq cell type annotations; (B)  $-\log_{10}$  p-values for positive  $\tau^*$  for a subset of 10 (of 28) diseases/traits and 10 (of 17) adult brain scRNA-seq cell type annotations; (C) comparison of results for 8 cell types included in both adult brain scATAC-seq and scRNA-seq data; and (D) comparison of results for 3 cell types included in both fetal and adult brain scATAC-seq data. In A-B, only statistically significant results ( $\text{FDR} > 5\%$ ) are colored ( $-\log_{10}(\text{p-value}) \geq 1.79$  for scATAC-seq,  $\geq 2.04$  for scRNA-seq). In A-B, cell types appearing in both datasets are denoted in red font. Numerical results for all diseases/traits and cell types are reported in Table E.14, Table E.15, Table E.16, and Table E.17. ADHD: attention deficit hyperactivity disorder, SCZ: schizophrenia, MDD: major depressive disorder, BMI: body mass index.

We next analyzed adult brain scRNA-seq data spanning 17 cell types<sup>197</sup> (of which 8 were also included in the fetal brain scATAC-seq data). We identified 64 significant disease-cell type pairs (FDR < 5% for positive  $\tau^*$  conditional on other annotations; Table 3.1, Table 3.2, Figure 3-2B, Table E.15). When restricting to the 33 significant disease-cell type pairs corresponding to 8 cell types included in both scATAC-seq and scRNA-seq data, 20 of 33 were also significant in analyses of scATAC-seq data. (Conversely, when restricting to the 48 significant disease-cell type pairs in scATAC-seq data corresponding to cell types included in scRNA-seq data, 21 of 48 were also significant in scRNA-seq data.)

The most significant enrichment was observed for excitatory neurons in intelligence, consistent with previous genetic studies<sup>203</sup>. We also identified an enrichment of corticofugal projection neurons (CPN) in intelligence, which has not (to our knowledge) previously been reported in analyses of genetic data. CPN connect neocortex and the subcortical regions and transmits axons from the cortex. Imbalance in neuronal activity, particularly regarding excitability of CPNs, has been hypothesized to lead to deficits in learning and memory<sup>236,237</sup>. Recently, Runge et al.<sup>238</sup> reported that NEUROD2 knockout mice display synaptic and physiological defects in CPN along with autism-like behavior abnormalities (where NEUROD2 is a transcription factor involved in early neuronal differentiation). CPN have previously been reported to be enriched in autism spectrum disorder (ASD) genes<sup>239</sup>, we did not detect a significant ASD enrichment for CPN ( $p = 0.056$ ) or any other cell type (see Discussion).

We compared the results for 9 adult brain cell types included in both the scATAC-seq and scRNA-seq datasets (Table E.5 and Table E.16). While scATAC-seq and scRNA-seq cell-type annotations for matched cell types were weakly correlated to each other ( $r = 0.01 - 0.09$ ; Table E.8), the corresponding  $-\log_{10}$ (p-values) for positive  $\tau^*$  were moderately correlated ( $r = 0.25$ ; Figure 3-2C), confirming the shared biological information. We observed more significant p-values for scATAC-seq than for scRNA-seq in these data sets, analogous our analyses of fetal brain data (see Discussion).

We compared the results for 3 cell types (astrocytes, inhibitory neurons, excitatory neurons) included in both fetal brain and adult brain scATAC-seq data sets (Figure



3-2D and Table E.17). While fetal brain and adult brain cell-type annotations for matched cell types were weakly correlated to each other ( $r = 0.00 - 0.01$ ), the corresponding  $-\log_{10}(\text{p-values})$  for positive  $\tau^*$  attained a moderately high correlation ( $r = 0.52$ ), higher than the analogous correlations for scATAC-seq vs. scRNA-seq results ( $r = 0.24$  for fetal brain,  $r = 0.25$  for adult brain; see above). Interestingly, the enrichment in ADHD (as well as neuroticism and worry) for fetal brain astrocytes (see above) was not observed for adult brain astrocytes ( $p = 0.52$  for positive  $\tau^*$ ,  $p = 0.0065$  for difference in  $\tau^*$  of adult brain astrocytes vs. fetal brain astrocytes). While astrocytes participate in defense against stress, energy storage, and tissue repair, they also mediate synaptic pruning (elimination of synaptosomes) during development<sup>240</sup>. Indeed, astrocytes in more mature stages of brain development were found to be less efficient at removing synaptosomes compared to younger, fetal astrocytes<sup>241</sup> (in both *in vitro* in pluripotent stem cells and *in vivo* mice), supporting a fetal brain-specific role of astrocytes in brain-related diseases and traits.

We observed no correlation between fetal brain and adult brain  $-\log_{10}(\text{p-values})$  for positive  $\tau^*$  in analyses of scRNA-seq data ( $r = 0.044$ ; Figure F-2 and Table E.18), possibly due to the lower power of these analyses (particularly for fetal brain scRNA-seq) in the data sets that we analyzed (see Discussion).

We performed 5 secondary analyses. First, we analyzed enrichments of both scATAC-seq and scRNA-seq brain cell types in 6 control (non-brain-related) diseases and complex traits. As expected, we did not find significant enrichments of brain cell types in 6 control traits (Table E.19 and Table E.20). Second, we repeated our disease heritability enrichment analyses of scATAC-seq annotations while conditioning only on the baseline model (and not the union of open chromatin regions across all cell types). We identified 246 significant disease-cell type pairs, as compared to 168 significant disease-cell type pairs in our primary analysis (Table E.21). This underscores the importance of conditioning on the union of open chromatin regions across all cell types, a conservative step to ensure cell-type specificity. (When we repeated the analysis on fetal brain scATAC-seq, we found similar results before or after additionally condition on the union of open chromatin regions.) Third, we performed

gene set enrichment analysis using GREAT<sup>213</sup> for both scATAC-seq and scRNA-seq cell-type annotations from adult brain. As expected, we identified significant enrichments in relevant gene sets (Table E.22). Fourth, for the adult scRNA-seq data<sup>197</sup>, we constructed annotations based on a +/- 100kb window-based strategy (previously used in ref.<sup>4</sup>) instead of brain-specific enhancer-gene links<sup>17-19</sup> (used in ref.<sup>20</sup>). We identified 28 significant trait-cell type pairs (Table E.23), vs. 64 using brain-specific enhancer-gene links. Fifth, we analyzed bulk chromatin data (7 chromatin marks) spanning 21 adult brain tissues<sup>193</sup> (age 27-85 years). We identified 1,710 significant disease-tissue-chromatin mark triplets spanning 26 of 28 brain-related diseases and traits (Table E.24). Once again, these results are included for completeness, but cannot achieve the same cell-type specificity as analyses of single-cell data .

### 3.6 Discussion

We identified a rich set of disease-critical fetal and adult brain cell types by integrating GWAS summary association statistics from 28 brain-related diseases and traits with scATAC-seq and scRNA-seq data from 83 fetal and adult brain cell types<sup>197-200</sup>. We confirmed many previously reported disease-cell type associations, but also identified disease-cell type associations supported by known biology that were not previously reported in analyses of genetic data. We determined that cell-type annotations derived from scATAC-seq were particularly powerful in the data that we analyzed. We also determined that the disease-cell type associations that we identified can be either shared or specific across fetal vs. adult brain developmental stages.

We note 4 key distinctions between our work and previous studies identifying disease-critical tissues and cell types<sup>4,17,20,64,65,114,194,195,201-203</sup>. First, we explicitly compared results from scATAC-seq vs. scRNA-seq data in matched cell types. Although applications of single-cell data to identify disease-critical cell types have largely prioritized analyses of scRNA-seq data<sup>192</sup>, we determined that cell-type annotations derived from scATAC-seq were even more powerful in the data that we analyzed. This findings may be specific to the data that we analyzed, and should not preclude

further prioritization of scRNA-seq data, but does motivate further prioritization of scATAC-seq data. Second, we explicitly compared results for fetal and adult brain in matched cell types. We determined that disease-critical cell types are specific to fetal vs. adult brain developmental stages in some instances, such as the enrichment of fetal astrocytes in ADHD. Third, we rigorously conditioned on a broad set of other functional annotations, a conservative step to ensure cell-type specificity that was included in only a subset of previous studies<sup>212,242</sup>. Fourth, in analyses of scRNA-seq data, we constructed annotations using brain-specific enhancer-gene links<sup>17-19</sup> (used in ref.<sup>20</sup>), an emerging approach that is more powerful than conventional window-based strategies for linking SNPs to genes.

Our findings have implications for improving our understanding of how cell-type specificity impacts disease risk. Better understanding disease-critical cell types is crucial to characterizing disease mechanisms underlying cell type specificity and developing new therapeutics<sup>192</sup>. To this end, the disease-cell type associations that we identified can help guide functional follow-up experiments (e.g. Perturb-seq<sup>243</sup>, saturation mutagenesis<sup>244</sup>, and CRISPR-Cas9 cytosine base editor screen<sup>245</sup>) to study cellular mechanisms of specific loci or genes underlying disease. In addition, our results highlight the benefits of analyzing data from different sequencing platforms and different developmental stages to identify disease-critical cell types. This motivates the prioritization of technologies that simultaneously profile ATAC and RNA expression such as SHARE-seq<sup>246</sup> and continued efforts to profile the developing human brain<sup>212</sup>.

We note several limitations of our work. First, although annotations derived from scATAC-seq generally outperformed annotations derived from scRNA-seq in the data that we analyzed, we caution that we are unable to draw any universal conclusions about which technology is most useful, as our findings may be impacted by the particularities of the data sets. However, we note that for both fetal and adult brain, the scRNA-seq data that we analyzed had larger numbers of donors and nuclei sequenced vs. the scATAC-seq data. Second, our resolution in identifying disease-critical cell types is fundamentally limited by the resolution of the single-

cell data that we analyzed; in particular, rare but biologically important cell types may be poorly represented in these data sets. Third, despite our rigorous efforts to condition on a broad set of functional annotations, we are unable to conclude that the disease-critical cell types that we identify are biologically causal; it may often be the case that they “tag” a biologically causal cell type that is not included in the data that we analyzed. This motivates further research on fine-mapping causal cell types (analogous to research on fine-mapping disease variants<sup>247</sup> and disease genes<sup>248</sup>). Fourth, we failed to identify any significant cell types for 4 disease/traits (autism, anorexia, ischemic stroke, and Alzheimer’s disease), possibly due to limited GWAS power and/or disease heterogeneity. Despite these limitations, the disease-cell type associations that we identified have high potential to improve our understanding of the biological mechanisms of complex disease.

### 3.7 Web resources

Generated cell-type annotations: [https://alkesgroup.broadinstitute.org/LDSC/ORE/Kim\\_ATAC/](https://alkesgroup.broadinstitute.org/LDSC/ORE/Kim_ATAC/)

S-LDSC software: <https://github.com/bulik/ldsc>

GWAS summary statistics: [https://alkesgroup.broadinstitute.org/sumstats\\_formatted/](https://alkesgroup.broadinstitute.org/sumstats_formatted/)

Generated cell-type annotations: [https://alkesgroup.broadinstitute.org/LDSC/ORE/Kim\\_ATAC/](https://alkesgroup.broadinstitute.org/LDSC/ORE/Kim_ATAC/)

Domcke et al. fetal scATAC-seq data: <https://atlas.brotmanbaty.org/bbi/human-chromatin-during-development/>

Cao et al. fetal scRNA-seq data: <https://atlas.brotmanbaty.org/bbi/human-gene-expression-during-development/>

Corces et al. scATAC-seq data: <http://epigenomegateway.wustl.edu/legacy/?genome=hg38&session=drS3o1n4kJ>

Velmeshev et al. scRNA-seq data: <https://autism.cells.ucsc.edu/>

baseline (v.1.2) annotations: <https://data.broadinstitute.org/alkesgroup/LDS>

CORE/

1000 Genomes Project Phase 3 data: <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502>

GREAT (Genomic Regions Enrichment of Annotations Tool): <http://great.stanford.edu/>

# Chapter 4

## Conclusion

In part I of this thesis, I assess the contribution of disease-associated gene pathways and gene sets to disease heritability. I highlight the importance of the conditional pathway enrichment analysis to elucidate biological mechanisms. In part II of the thesis, I introduce a framework, AnnotBoost, to improve the Mendelian disease-derived pathogenicity scores for common disease. Building upon the previous findings on the shared genetic architecture between Mendelian disease and common disease, AnnotBoost generates significant improvements in heritability model fit, suggesting improved power in future large-scale fine-mapping studies. In part III the thesis, I conduct a study to identify disease-critical cell types by leveraging single-cell RNA-seq and ATAC-seq data. I demonstrate that single-cell-derived cell type annotations are powerful to address the biological question of what are the critical cellular processes.

In this thesis, I develop several computational methods and frameworks to interpret GWAS findings to reveal biological insights, a major challenge of past decade. Taken together, this thesis advances our understanding of the genetic basis of common disease and complex trait by dissecting heritability of genome by functional categories, improving existing functional annotations, and identifying disease-critical cell types.

# Appendix A

## Supplementary tables for Chapter 1

See Kim et al.<sup>42</sup> for all supplementary tables. Titles and captions are provided below.

Table A.1: **List of 18,119 pathways analyzed.** For each pathway, we report a pathway ID, pathway description, database, Entrez IDs for genes, and the number of protein-coding genes.

Table A.2: **List of 42 independent diseases and traits analyzed.** For each trait, we report a trait identifier, trait description, reference, sample size, and heritability z-score. We selected these 42 traits based on a heritability z-score  $> 6$  (see Background). We further indicated brain or blood-related traits.

Table A.3: **Lists of genes in 18 gene sets compared with closeness centrality.** We compiled lists of genes for 18 metrics (gene sets) that we compared with closeness centrality. We report Entrez IDs for genes in the 18 gene sets.

Table A.4: **S-LDSC results of all pathway-trait pairs** We applied S-LDSC to 760,869 pathway-trait pairs, conditioning on all-genes annotation and the baseline-LD model. For each pathway-trait pair, we report a proportion of SNPs, enrichment, and  $\tau$ .

Table A.5: **S-LDSC results for 156 significantly enriched pathway-trait pairs.** For each significantly enriched pathway-trait pair, we report a proportion of SNPs, enrichment, and a  $\tau^*$ . The 8 significant pathway-trait pairs were reported in previous genetic studies: "pathways in cancer" for height<sup>24</sup>; "neuropeptide hormone activity" for BMI<sup>25</sup>; "immune response" for both Crohn's disease and ulcerative colitis<sup>26</sup>; "T-cell receptor," "abnormal T-cell physiology," and "cytokine-mediated signaling pathway" for rheumatoid arthritis<sup>27</sup>; "absent corpus callosum" for years of education<sup>23</sup>.

Table A.6: **Average gene size of annotations.** For all-genes annotation, all pathways, pathway, network, and pathway+network, and Quack annotations, we report an average size of genes (and its standard deviation) and an average number of genes.

Table A.7: **Heritability enrichment of enriched pathway-trait pairs.** We meta-analyzed (A) 156 enriched pathway-trait pairs; (B) 13 enriched pathway-trait pairs for ExAC, Cassa, and Samocha gene sets; (C) 169 enriched pathway-trait pairs (a and b combined). In each case, we report meta-analyzed enrichments and  $\tau^*$ .

Table A.8: **S-LDSC results of 156 enriched pathway-trait pairs excluding genes implicated by GWAS.** We removed genes implicated by previous GWAS studies (see 3.1 Methods; average of 5% of genes ( 2 genes) removed) and applied S-LDSC conditional on all-genes and baseline-LD model annotations. For each pathway-trait pair, we report a proportion of SNPs, enrichment,  $\tau^*$ , and the number of genes in a pathway excluding GWAS significant genes.

Table A.9: **S-LDSC results of 195 pathway-trait pairs from previous pathway enrichment studies.** We applied S-LDSC to 95 pathway-trait pairs from five previous genetic studies<sup>28-32</sup> and 100 from a recent study<sup>33</sup> (A) conditioning on the baseline-LD model and all-genes annotation and (B) conditioning on all-genes annotation only. In each pathway-trait pairs for each case, we report a proportion of SNPs, enrichment, and  $\tau$ . We assessed the statistical significance based on global FDR < 5% across 18,119 pathways tested ( $\tau^* < 0.000989$ ).



Table A.10: **S-LDSC results of 13 enriched pathway-trait pairs for ExAC, Cassa, Samocha gene sets.** For each of 13 enriched pathway-trait pairs, we report a proportion of SNPs, enrichment, and  $\tau$ .

Table A.11: **Correlation of network annotations with baseline-LD model annotations.** We report the Pearson correlation between baseline-LD model annotations and (A) Saha network annotations of different centralities, (B) Saha network annotations of different tissues, (C) Greene network annotations of different centralities, (D) Greene network annotations of different tissue, (E) InWeb network annotations of different centralities, (F) Sonawane network annotations of different centralities, and (G) Sonawane network annotations of different tissue.

Table A.12: **Summary of gene networks analyzed.** We report the number of genes, the number of edges, and the distribution of edge weights for each of four networks (InWeb, Saha, Sonawane, Greene).

Table A.13: **Excess overlap of 18 genes in each decile bin of closeness centrality.** For each of 18 gene sets, we report the excess overlap (and standard error) of genes in each decile bin of closeness centrality for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network.

Table A.14: **Per-gene closeness centrality scores and gene membership in 18 gene sets** We report the closeness centrality for all protein-coding genes that exist in each of (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. We indicate gene membership in each of 18 gene sets, marking '1' if in the corresponding gene set.

Table A.15: **Correlation between closeness centrality and 18 gene sets.** We report Pearson correlations between closeness centrality and 18 gene sets analyzed for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network.

Table A.16: **Excess fold overlap of network and pathway+network annotations with baseline-LD model annotations.** We report the excess fold overlap between baseline-LD model annotations and network, pathway+network, Quack, and all-genes annotations.

Table A.17: **Correlation of network and pathway+network annotations with baseline-LD model annotations.** We report the Pearson correlation between baseline-LD model annotations and network, pathway+network, Quack, and all-genes annotations.

Table A.18: **Excess fold overlap / correlation among functional annotations from the baseline-LD model.** We report (A) excess fold overlap and (B) correlation among functional annotations from the baseline-LD model.

Table A.19: **TFs enriched in high closeness centrality genes.** For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, for high closeness centrality genes (top decile), we report significantly enriched TFs (Benjamini-Hochberg adjusted p-value  $< 0.05$ ). The description of TFs is provided in the ENCODE Chip-Seq Significance Tool (see Web resources).

Table A.20: **Correlation of deciles of closeness centrality with baseline-LD model annotations.** For four networks, we constructed binarized network annotations based on deciles of closeness centrality. We report the Pearson correlation between these annotations and baseline-LD model annotations.

Table A.21: **Enriched GO terms of high closeness centrality genes.** For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we report significantly enriched GO terms in the following GO categories: biological process (BP), cellular component (CC), and molecular function (MP) (Benjamini-Hochberg adjusted p-value  $< 0.05$ ).

Table A.22: **Heritability enrichment of network annotations.** For each of 4 network annotations, we report meta-analyzed enrichments and  $\tau^*$  across 42 independent traits. We highlight the network attaining highest enrichment for each trait. For the three tissue-specific networks (Saha, Greene, Sonawane), we also report meta-analyzed enrichments and  $\tau^*$  of network annotations constructed using the tissue that maximized the excess overlap with the High pLI (ExAC) gene set.

Table A.23: **Heritability enrichment of network annotations conditioning on one annotation from the baseline-LD at a time.** For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we meta-analyzed network annotations across 42 independent traits, conditioning on one annotation from the baseline-LD model at a time. We report meta-analyzed enrichments and  $\tau^*$  across 42 independent traits. We highlighted annotations that significantly reduced  $\tau^*$  (using Bonferroni-corrected p-val).

Table A.24: **Heritability enrichment of deciles of closeness centrality.** For four networks, we constructed binarized network annotations based on deciles of closeness centrality. We applied S-LDSC and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . For the three tissue-specific networks (Saha, Greene, Sonawane), we also report meta-analyzed enrichments and  $\tau^*$  of binarized network annotations constructed using the tissue that maximized the excess overlap with the High pLI (ExAC) gene set.

Table A.25: **Heritability enrichment of network annotations from network perturbation analysis.** We randomly removed 10% to 90% of edges from the original networks and computed closeness centrality on networks with edges removed; we performed five separate perturbation analyses for each value of the proportion of edges removed. We applied S-LDSC and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table A.26: **Heritability enrichment of DSD-network annotations.** We applied the diffusion state distance (DSD) algorithm<sup>34</sup> to transform gene networks' edge weights with a random walk ( $k = 5$ ). Then, we constructed network annotations by re-computing closeness on DSD-transformed networks and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table A.27: **Heritability enrichment of networks annotations from consensus networks.** We constructed consensus networks and made (A) probabilistic annotations based on closeness centrality and (B) binary annotations based on deciles of closeness centrality. We applied S-LDSC and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table A.28: **Correlation between closeness and gene expression.** For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we report the correlation between closeness centrality and gene expression across 53 GTEx tissues.

Table A.29: **Heritability enrichment of Saha TSN and TWN gene sets.** We constructed gene sets based on membership of genes in Saha tissue-specific networks (TSN) and transcriptome-wide networks (TWN). We applied S-LDSC to 36 TSN and 16 TWN gene sets across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table A.30: **Excess overlap of top deciles of closeness centrality of tissue-specific networks with High pLI (ExAC) genes.** For each tissue-specific network (Saha, Greene, Sonawane), for each tissue, we report the excess overlap between High pLI (ExAC) genes and the top decile of closeness centrality.

Table A.31: **Heritability enrichment of network annotations using relevant tissues for brain-related and blood-related traits.** For (A) 8 brain-related traits and (B) 10 blood-related traits, we report meta-analyzed enrichments and  $\tau^*$ .

Table A.32: **Heritability enrichment of network annotations using 6 other network centrality metrics.** For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we constructed network annotations based on 6 other network centrality metrics and meta-analyzed results across 42 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table A.33: **Heritability enrichment of network annotations with different window sizes.** Instead of 100kb windows around genes, we added (A) 10kb or (B) 1Mb windows around genes when constructing network annotations. We report meta-analyzed enrichments and  $\tau^*$  across 42 independent traits.

Table A.34: **Heritability enrichment of pathway+network annotations.** We meta-analyzed 156 pathway-trait pairs (122 for Saha, which has less pairs as all genes in some pathways do not exist in the Saha network). For each 590 pathway-trait pair, we report a proportion of SNPs, enrichment, and  $\tau$ . We also report meta-analyzed enrichments and  $\tau^*$  across 156 pathway-trait pairs (122 for Saha). We highlighted the network attaining highest enrichment for each pathway-trait pairs.

Table A.35: **Heritability enrichment of pathway+network annotations conditioning on one annotation from the baseline-LD at a time.** For (A) Saha, (B) Greene, (C) InWeb, (D) Sonawane networks, we constructed an average annotation across 156 pathway+network annotations and meta-analyzed averaged pathway+network annotations across 42 independent traits, conditioning on one annotation from the baseline-LD model at a time. We report meta-analyzed enrichments and  $\tau^*$  across 42 independent traits.

Table A.36: **Network connectivity of enriched pathways and null pathways.** Using (A) sum of edge weights or (B) number of edges as network connectivity metrics, we report the number of interacting genes and network connectivity between genes in a pathway and neighboring genes outside the pathway. We constructed null pathways in two ways: (1) each gene sampled from a randomly chosen pathway or (2) each gene randomly sampled from all protein-coding genes.

Table A.37: **Heritability enrichment of Quack annotations.** We applied the Quack random-forest classifier algorithm<sup>35</sup>. We used 18,119 pathways as a training data and applied Quack to four gene networks. We used the output of Quack to construct Quack pathway+network annotations (with 100kb window), applied S-LDSC, and meta-analyzed across 156 pathway-trait pairs (122 for Saha). We report meta-analyzed enrichments and  $\tau^*$ .

Table A.38: **Heritability enrichment of 53 pathway+network annotation using pathways excluding genes implicated by GWAS.** From 53 significant pathway-trait pairs after excluding GWAS significant genes, we constructed pathway+network annotations and meta-analyzed results across 53 pathway-trait pairs (40 for Saha). We report meta-analyzed enrichments and  $\tau^*$ .

# Appendix B

## Supplementary figures for Chapter 1

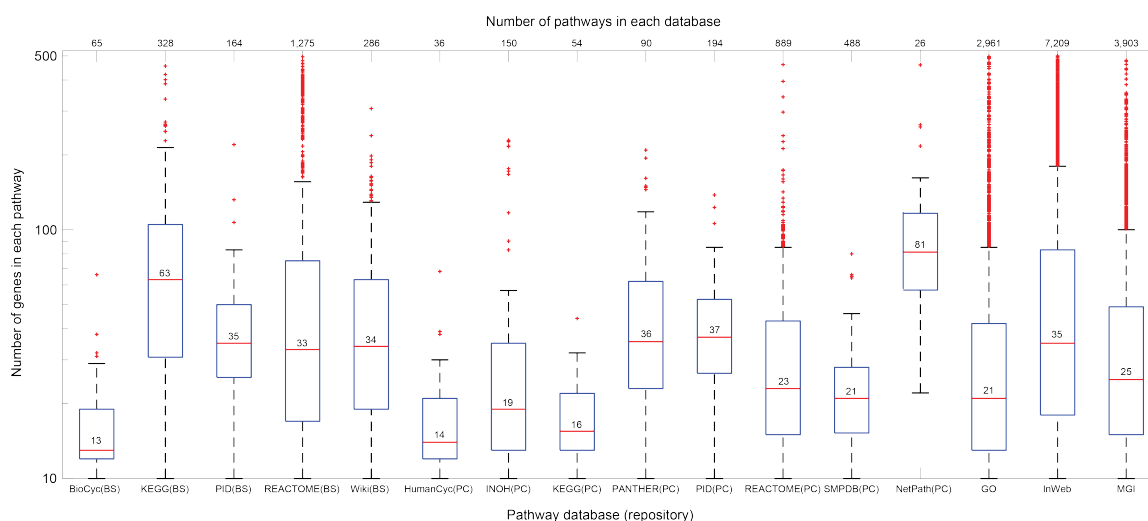


Figure B-1: **Summary of pathways analyzed.** We analyzed 18,119 pathways, each consisting of at least 10 and at most 500 coding genes. Number in the box-plot represents the median number of genes for each database. Number of pathways in each database is shown on the top of the figure. BS: NCBI BioSystems. PC: PathwayCommons. MGI: Mouse Genome Informatics. GO: Genome Ontology. In-Web: InWeb protein-protein interactions. See Table A.1 for a description of 18,119 pathways analyzed.

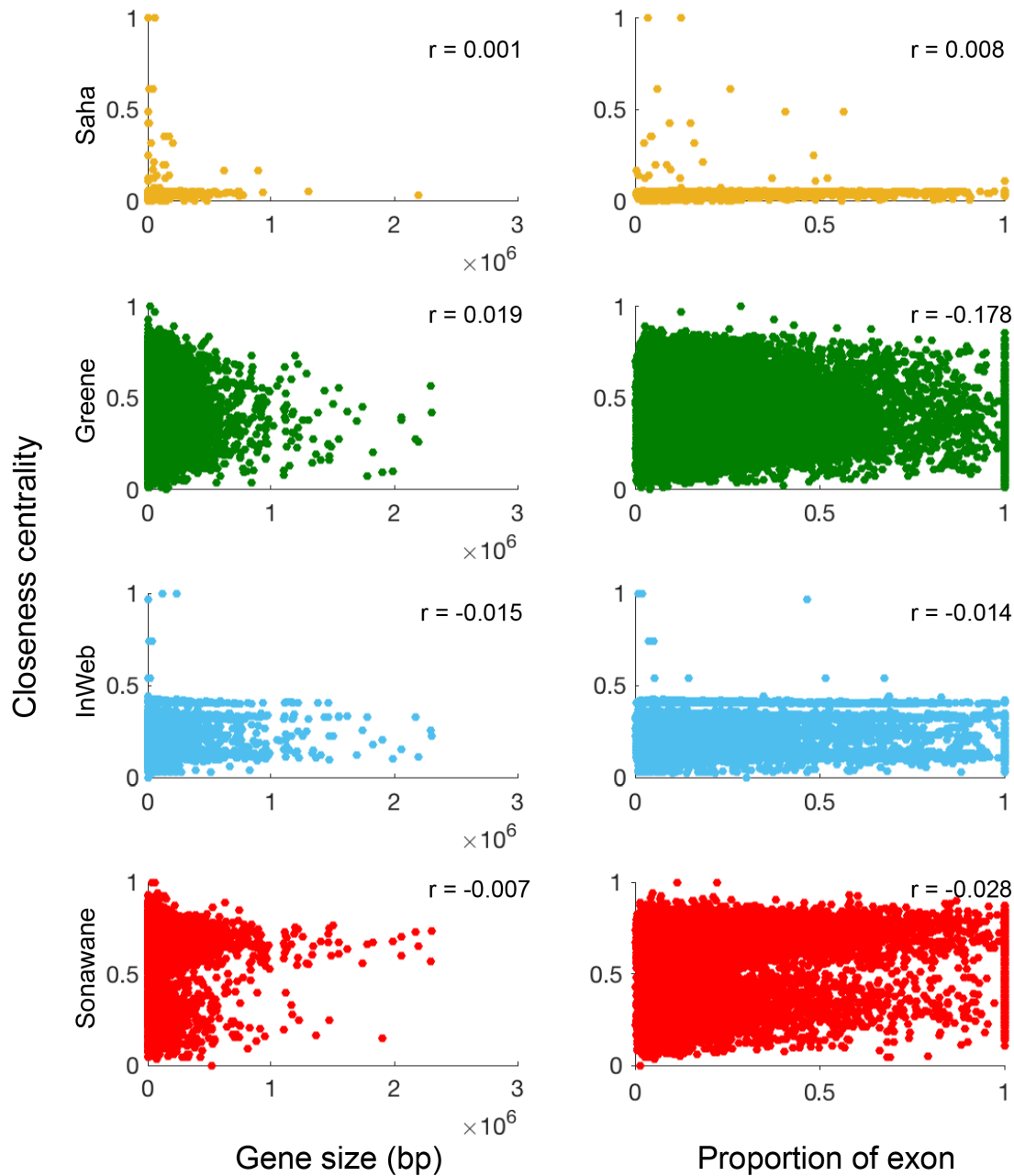


Figure B-2: **Closeness centrality is independent from the gene size and the proportion of exon.** For each of four network annotations, we computed a Pearson correlation between probabilistic annotation values and (1) gene size and (2) proportion of exon. We calculated the proportion of exon as the size of coding regions (bp) that lie inside the gene divided by the size of the gene, defined as (transcription stop - transcription start).

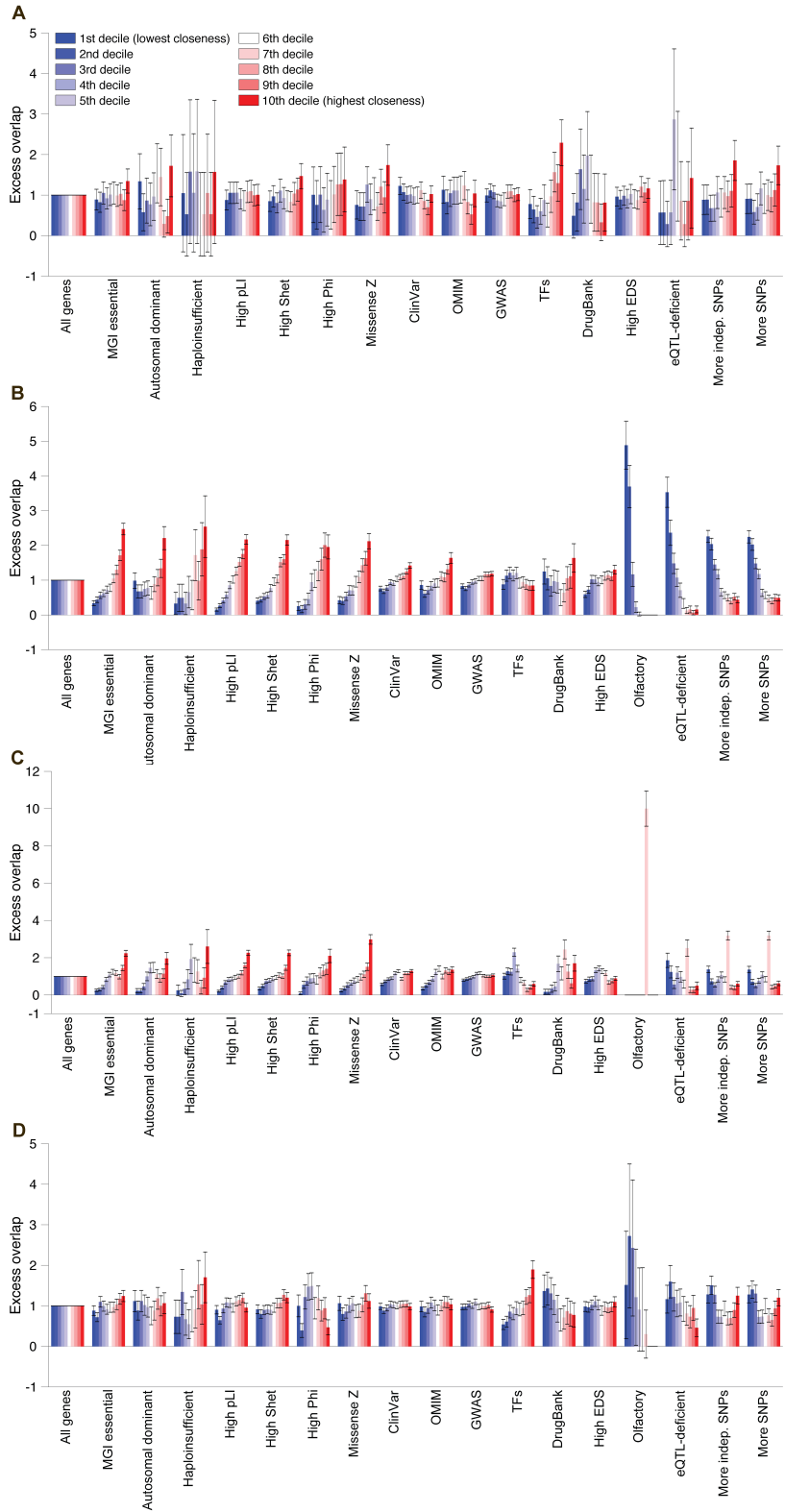


Figure B-3: **Excess overlap of 18 genes in each decile bin of closeness centrality.** For each of 18 gene sets, we report the excess overlap of genes in each decile bin of closeness centrality for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. Error bars represent 95% confidence intervals. There is no olfactory receptor gene in Saha network.



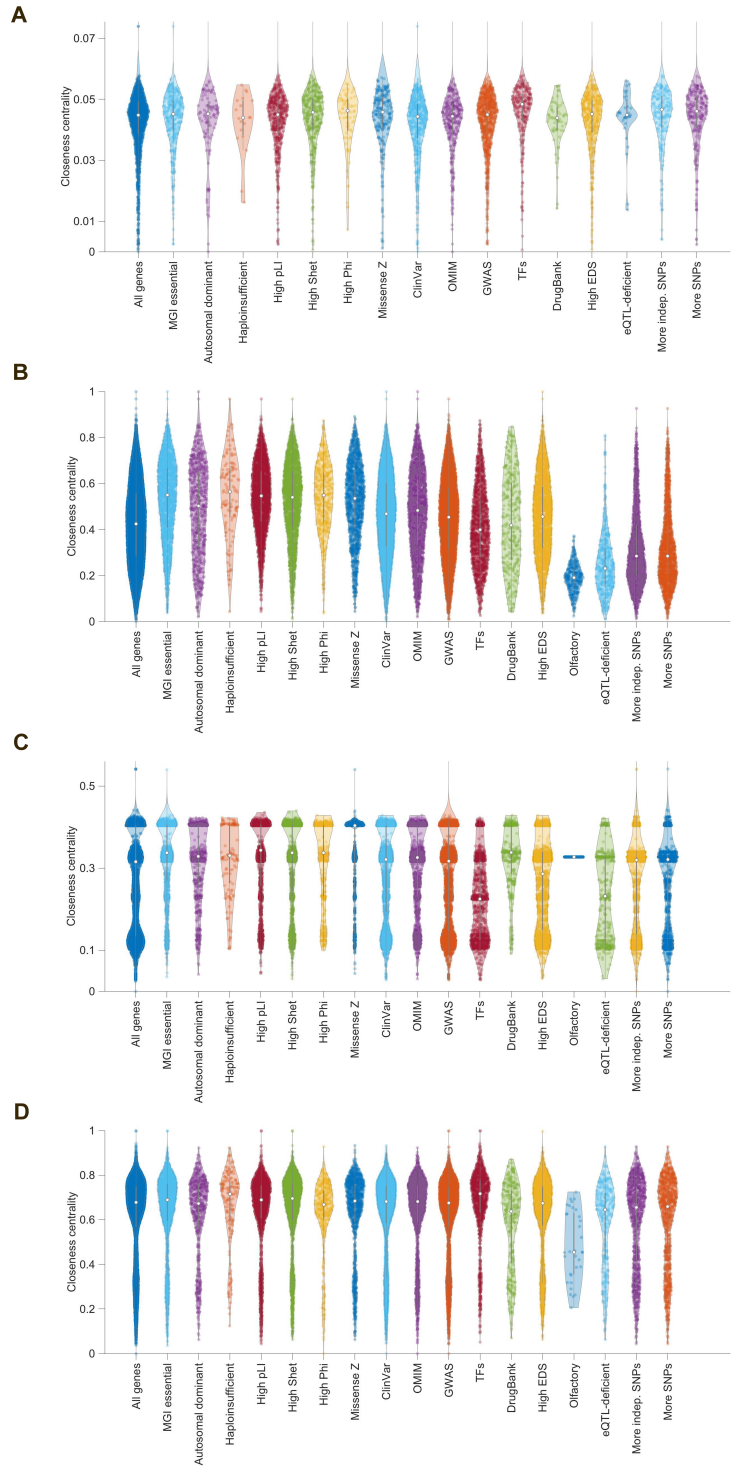


Figure B-4: **Distribution of closeness centrality in 18 gene sets.** For each of 18 gene sets, we show the distribution of closeness for (A) Saha skin network, (B) Greene thyroid network, (C) InWeb network, and (D) Sonawane testis network. Box plot is shown as a grey line inside the violin plot where white dot represents median. Each colored dot represents a gene. Outliers are not displayed for Saha and InWeb (see Web resources for closeness scores)

# Appendix C

## Supplementary tables for Chapter 2

See Kim et al.<sup>128</sup> for all supplementary tables. Titles and captions are provided below.

**Table C.1: List of 41 independent diseases and complex traits analyzed.** Analogous to a previous study<sup>36</sup>, we considered 89 GWAS summary association statistics, including 34 traits from publicly available sources and 55 traits from the UK Biobank (up to  $N = 459\text{K}$ ); summary association statistics were computed using BOLT-LMM v2.3<sup>37,38</sup>. We obtained 41 independent traits (average  $N = 320\text{K}$ ) with genetic correlation less than 0.9 (computed using cross-trait LDSC<sup>39</sup>). For 6 traits, we analyzed two different sources (both publicly available and UK Biobank), resulting in total 47 summary statistics analyzed. For each trait, we report a trait identifier, trait description, reference, sample size, and heritability z-score.

**Table C.2: Informativeness for common disease of binary annotations derived from 11 Mendelian disease-derived missense scores and corresponding boosted scores.** We applied S-LDSC, conditional on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

**Table C.3: Correlations between published scores and genome-wide boosted scores** We report Pearson correlation ( $r$ ) between published scores and genome-wide boosted scores for (A) Mendelian disease-derived missense scores, (B) genome-wide Mendelian disease-derived scores, (C) other scores, and (D) all scores. We regard originally unscored SNPs as having a score of zero.

Table C.4: **Predictive accuracy of AnnotBoost in out-of-sample predictions of input published scores.** We evaluated the AnnotBoost model classification performance based on AUROC and AUPRC on a 20% held-out testing set (consisted of odd (resp. even) chromosome SNPs) through the 5-fold cross-validation. For each boosted score, we reported the average AUROCs of odd-chr and even-chr model.

Table C.5: **Expected heritability enrichments of binary annotations derived from boosted Mendelian disease-derived missense scores with significantly negative  $\tau^*$ .** We report the heritability enrichment that is expected based on the boosted annotation's overlap with the baseline-LD model and corresponding published annotations, by assuming that the  $\tau$  of the annotation is zero. In general, expected enrichments were significantly less than observed enrichments.

Table C.6: **Informativeness for common disease of binary annotations derived from boosted Mendelian disease-derived missense scores, restricted to non-coding regions.** We restricted boosted Mendelian missense annotations (whose top X% variants are already optimized) to non-coding SNPs and applied S-LDSC (conditioning on the baseline-LD model and corresponding published annotations). We report meta-analyzed enrichments,  $\tau$ , and  $\tau^*$  across 41 independent traits; for meta-analysis, we used the same weights from the meta-analysis of the unrestricted boosted annotations.

Table C.7: **Trait-specific informativeness for common disease of binary annotations.** Instead of meta-analyzing S-LDSC results across traits, we report the trait-specific results for 5,311 annotation-trait pairs (113 annotations x 47 traits; 113 annotations whose proportions of top SNPs were optimized). We indicated 152 annotation-trait pairs that were significant at  $FDR < 5\%$  on  $\tau^*$  p-value.

Table C.8: **Heterogeneity of heritability enrichment and  $\tau^*$  across traits.** We report the heterogeneity p-value of heritability enrichment and  $\tau^*$  across 41 traits, for each of annotations analyzed.

Table C.9: **Gene scores derived from published and boosted pathogenicity scores.** We report genes scores (maximum pathogenicity score of linked SNPs) for genes linked to each of 203 scores analyzed (consisting of 35 published, 82 boosted, 86 baseline-LD annotations).

Table C.10: **165 reference gene sets of biological importance.** We report (A) gene set’s short description, number of genes, and reference; (B) ENSG ID for each of 165 gene sets (limited to protein-coding genes).

Table C.11: **Excess overlap of 165 reference gene sets in each quintile bin of gene scores derived from input pathogenicity scores.** We report (A) the excess fold overlap of genes in each quintile bin of both published and boosted gene scores; (B) the standard error of the excess overlap; (C) the Pearson correlation ( $r$ ) among the gene scores; and (D) the Spearman correlation ( $\rho$ ) among the gene scores.

Table C.12: **Odds ratios for fold overlaps by gene quintiles from boosted scores and corresponding gene quintiles from published scores in 165 reference gene sets.** We computed the significance of the association between gene score quintiles from each of 82 published-boosted score pairs using Fisher’s exact test. We report the odds ratios and p-value for (A) fifth quintile (top 20%) gene sets, (B) fourth quintile gene sets, (C) third quintile gene sets, (D) second quintile gene sets, and (E) first quintile (bottom 20%) gene sets.

Table C.13: **Informativeness for common disease of binary annotations derived from 6 genome-wide Mendelian disease-derived scores and corresponding boosted scores.** We applied S-LDSC, conditional on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table C.14: **Informativeness for common disease of binary annotations derived from boosted genome-wide Mendelian pathogenicity scores, restricted to previously unscored variants.** We restricted boosted Mendelian annotations (whose coverage is  $< 100\%$ ; see Table 2.1) to previously unscored variants and applied S-LDSC (conditioning on the baseline-LD model and corresponding published annotations). We report meta-analyzed enrichments,  $\tau$ , and  $\tau^*$  across 41 independent traits; for meta-analysis, we used the same weights from the meta-analysis of the unrestricted boosted annotations.

**Table C.15: Summary of 18 additional genome-wide scores and 47 baseline-LD model annotations.** We provide a description of 18 additional genome-wide scores and 47 baseline-LD model annotations. These include 2 constraint-related annotations, 9 predicted regulatory annotations (7 epigenetic states and 2 predicted transcriptional/post-transcriptional regulatory effects), 7 gene-based annotations, and 47 annotations from the baseline-LD model. We report the coverage, out of 9,997,231 common SNPs we considered; we note that the baseline-LD model annotations consider genome-wide SNPs and do not have any missingness (see Methods). For binary annotations, the proportion of SNPs with value 1 is considered as coverage. Annotations are first ordered by the type and then by the year of publication. base (all SNPs), 10 MAF bins, and 6 LD-related annotations constitute baseline-LD-nofunct model (model with no functional annotations).

**Table C.16: Informativeness for common disease of binary annotations derived from 18 additional genome-wide scores + 47 baseline-LD model annotations and corresponding boosted scores.** We applied S-LDSC, conditional on the baseline-LD model and 8 Roadmap annotations (for annotations derived from published scores) or the baseline-LD model, 8 Roadmap annotations, and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

**Table C.17: Informativeness for common disease of binary annotations derived from boosted versions of additional genome-wide scores and baseline-LD model annotations, restricted to previously unscored variants.** We restricted boosted additional annotations (whose coverage is  $< 100\%$ ; see Table C.15) to previously unscored variants and applied S-LDSC (conditioning on the baseline-LD model and corresponding published annotations). We report meta-analyzed enrichments,  $\tau$ , and  $\tau^*$  across 41 independent traits; for meta-analysis, we used the same weights from the meta-analysis of the unrestricted boosted annotations.

Table C.18: **Informativeness for common disease of binary annotations derived from 18 additional genome-wide scores + 47 baseline-LD model annotations and corresponding boosted scores without conditioning on 8 Roadmap annotations.** We applied S-LDSC, conditional on the baseline-LD model (for annotations derived from published scores) or the baseline-LD model and corresponding published annotations (for annotations derived from boosted scores) and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table C.19: **Informativeness for common disease of marginally conditionally significant annotations derived from Mendelian disease-derived missense scores and genome-wide Mendelian disease-derived scores and corresponding boosted scores, when additionally conditioning on 8 Roadmap annotations.** We applied S-LDSC to marginally significant Mendelian missense and genome-wide Mendelian annotations, conditioning on the baseline-LD model, 8 Roadmap annotations, and published annotations (for boosted annotations), and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ .

Table C.20: **Genome-wide correlations among functional annotations.** We report Pearson correlation ( $r$ ) (A) among 11 jointly significant annotations, (B) between 11 jointly significant annotations and functional annotations from the baseline-LD model, and (C) among all 198 annotations (85 baseline-LD annotations (excluding ‘base’ annotation; all SNPs annotated with 1), 31 annotations derived from published scores, and 82 boosted annotations).

Table C.21: **Informativeness for common disease of 11 jointly significant binary annotations in a combined joint model.** We applied S-LDSC, conditioned on the baseline-LD model, 8 Roadmap annotations, and each other, and meta-analyzed results across 41 independent traits. We report proportion of SNPs, meta-analyzed enrichments, and  $\tau^*$ .

Table C.22: **Evaluation of heritability model fit using  $logl_{SS}$ .** We report  $logl_{SS}$ ,  $\Delta logl_{SS}$ , and AIC for GCTA (number of annotations = 1), baseline-LD-nofunct (17), baseline-LD (86), baseline-LD+joint (97), baseline-LD+marginal $\tau^* \geq 0.5$  (105), and baseline-LD+marginal model (150), after applying SumHer<sup>1</sup> across 30 UKBB summary statistics.  $\Delta logl_{SS}$  is computed as  $logl_{SS}$  of a given model -  $logl_{SS}$  of baseline-LD-nofunct model. AIC (to account for increase in model complexity) is computed as  $2 \times$  number of annotations -  $2 logl_{SS}$ , as described in ref. speed2020. We also report the p-value on the differences in  $\Delta logl_{SS}$ .

Table C.23: **Classification of fine-mapped disease SNPs: multi-score analysis.** We report AUROCs and AUPRCs of different combined scores in classifying 5 different independent SNP sets: 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup> (stringently defined by causal posterior probability  $\geq 0.95$ ; both functionally-informed and non functionally-informed), and 14,807 GWAS significant SNPs<sup>14,15</sup> (from 10 LD-, MAF-, and genomic element-matched control SNPs as well as the most matched low frequency and common SNPs in the reference panel<sup>40</sup>).

Table C.24: **Summary of fine-mapped SNPs analyzed and results of single-score classification analysis.** We report the reference SNP ID (RSID) of four fine-mapped SNP sets and GWAS significant SNPs and their corresponding control SNPs that are found in the reference panel ( $MAF \geq 0.5\%$ ). We also report their short description and the number of unique SNPs of 10 matched control SNPs and most matched control SNPs. We also report the summary of single-score classification analysis, including the number of scores, where boosted scores perform better than the corresponding published scores, and the average AUROC and AUPRC improvement. Numeric results for the single-score analysis can be found in Table C.25.

**Table C.25: Classification of fine-mapped disease SNPs: single-score analysis.** We report AUROCs and AUPRCs of published and boosted scores in classifying 5 different independent SNP sets: 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup> (stringently defined by causal posterior probability  $\geq 0.95$ ; both functionally-informed and non functionally-informed), and 14,807 GWAS significant SNPs<sup>14,15</sup> ((A) from 10 LD-, MAF-, and genomic element-matched control SNPs as well as (B) the most matched low frequency and common SNPs in the reference panel<sup>40</sup>). We report the S-LDSC results (including standardized enrichments) of derived annotations based on top prioritized SNPs side by side. The summary of single-score analysis can be found in Table C.24.

**Table C.26: Correlation between S-LDSC metrics and AUROCs from single-score analysis of classifying fine-mapped SNPs.** For (A) 75 published annotations (excluding small annotations with proportion of SNPs  $< 0.02\%$ ) and (B) 82 boosted annotations, we computed Pearson correlation ( $r$ ) between different S-LDSC metrics (standardized enrichment,  $\tau^*$ ) and AUROCs on classifying fine-mapped SNPs. We note AUROCs are computed on all SNPs, and S-LDSC results correspond to results on top prioritized SNPs. We further note it is 75 published annotations, not 82, because 7 scores were excluded in S-LDSC analysis due to small annotation size for 4 scores and not comparable enrichments for 3 continuous-valued baseline-LD annotations.

**Table C.27: Informativeness of the baseline-LD model before and after adding 11 jointly significant binary annotations** We applied S-LDSC to (A) the baseline-LD model + 8 Roadmap annotations and (B) the baseline-LD model + 8 Roadmap annotations + 11 jointly significant annotations and meta-analyzed results across 41 independent traits. We report proportion of SNPs, meta-analyzed enrichments, and  $\tau^*$ .



# Appendix D

## Supplementary figures for Chapter 2

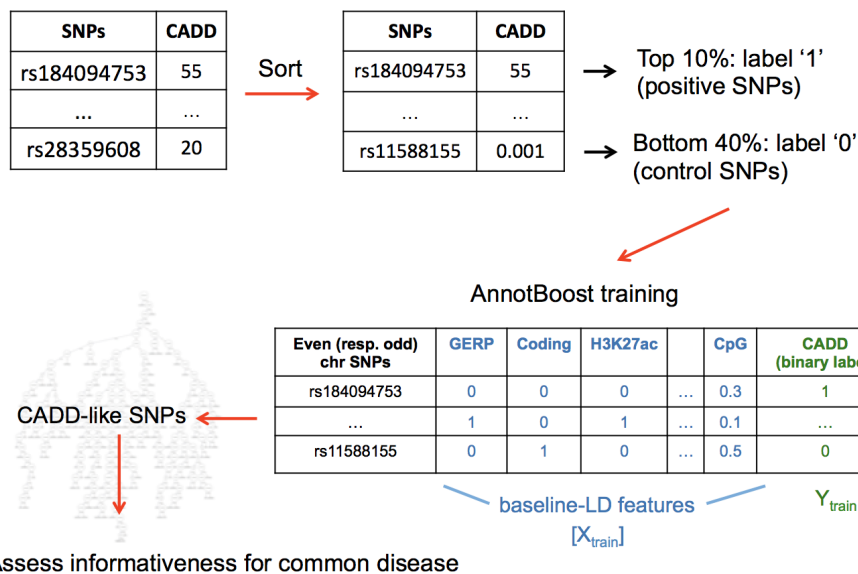


Figure D-1: **Overview of AnnotBoost framework.** We describe the AnnotBoost model training. AnnotBoost requires only one input, a pathogenicity score to boost, and generates a genome-wide (probabilistic) boosted pathogenicity score. From the input pathogenicity score (e.g. CADD as shown here), we built a classification model, each for even and odd chromosome SNPs using 75 baseline-LD annotations<sup>2,3</sup> as features. We assessed informativeness of annotations derived from published scores (input) and boosted scores (output) for common disease using S-LDSC<sup>4</sup>.



Figure D-2: **Feature importance of boosted Mendelian disease-derived missense pathogenicity scores.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD drives the prediction of 11 boosted missense scores by AnnotBoost. We report the signed impact of top 20 features for each of 11 predictive models: (A) PolyPhen-2, (B) PolyPhen-2-HVAR, (C) MetaLR, (D) MetaSVM, (E) PROVEAN, (F) SIFT 4G, (G) REVEL, (H) M-CAP, (I) Primate-AI, (J) MPC, and (K) MVP. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources).

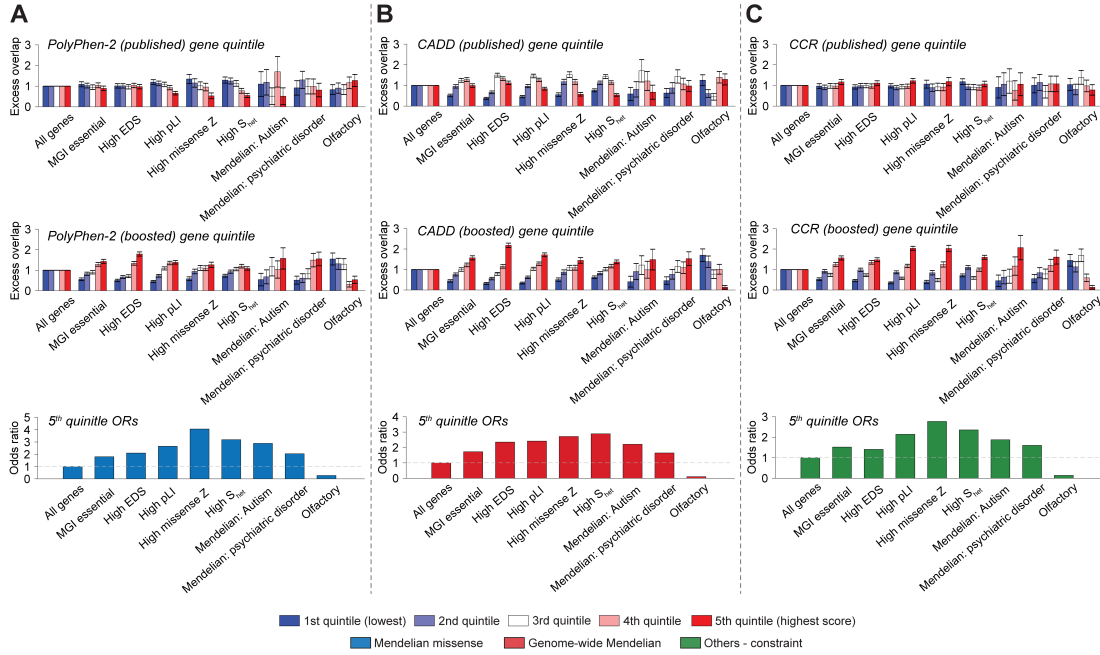
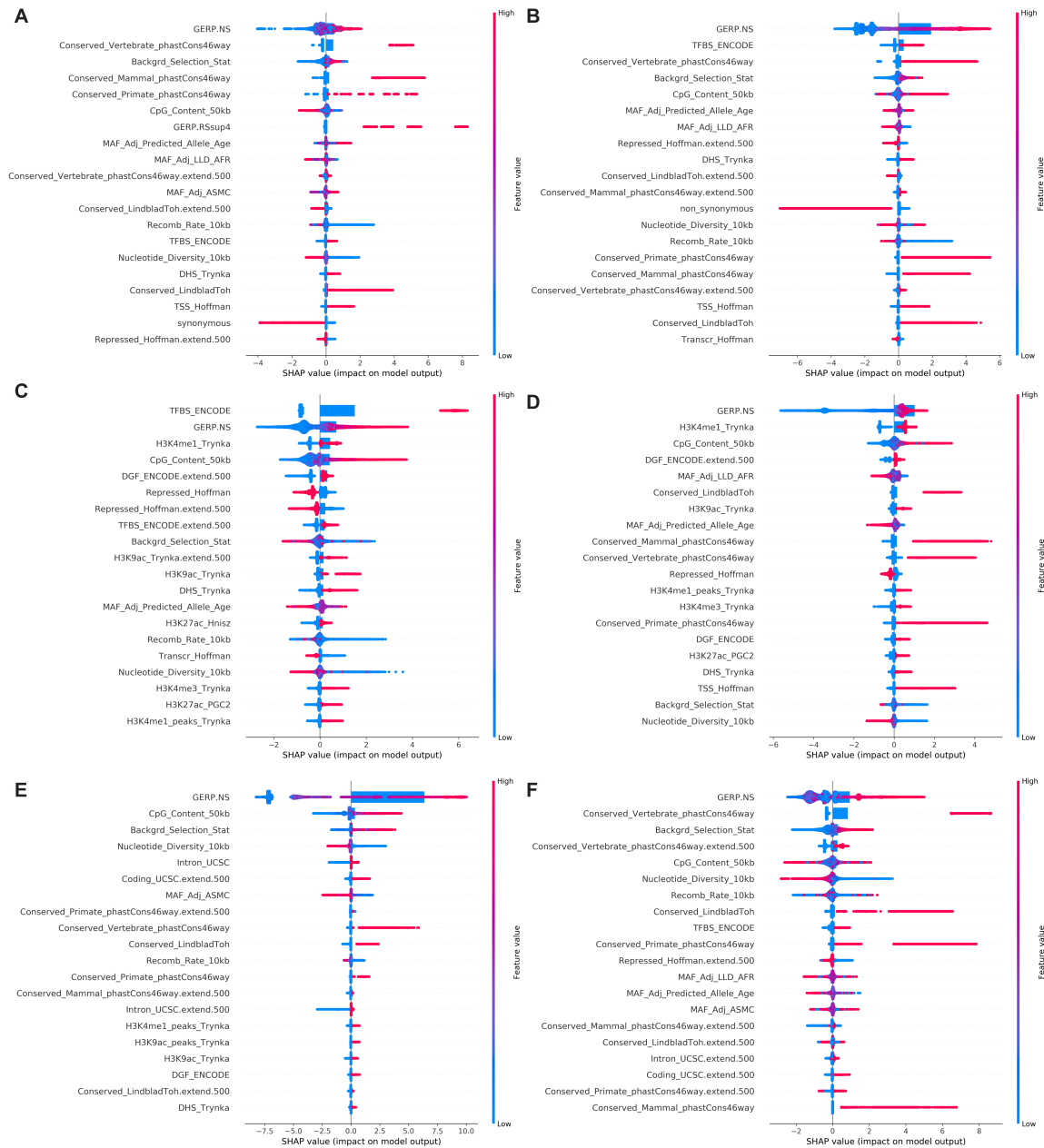


Figure D-3: **Excess overlap between gene scores derived from input pathogenicity scores and 165 reference gene sets of biological importance.** We report the excess overlap of genes linked to published and boosted scores in existing gene sets of biological importance (summarized in Table C.10): (A) PolyPhen-2<sup>6,7</sup> gene quintiles from published and boosted scores, (B) CADD<sup>8,9</sup> gene quintiles from published and boosted scores, and (C) CCR<sup>10</sup> gene quintiles from published and boosted scores. Error bars represent 95% confidence intervals. Numeric results for excess overlap and correlation among gene scores are shown in Table C.11. Numeric results for odds ratios and p-values from Fisher’s exact test between published gene quintiles and boosted gene quintiles are reported in Table C.12.



**Figure D-4: Feature importance of boosted genome-wide Mendelian disease-derived pathogenicity scores.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD drives the prediction of 11 boosted missense scores by AnnotBoost. We report the signed impact of top 20 features for each of 6 genome-wide Mendelian disease-derived pathogenicity scores: (A) CADD, (B) Eigen, (C) Eigen-PC, (D) ReMM, (E) NCBoost, and (F) ncER. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources).

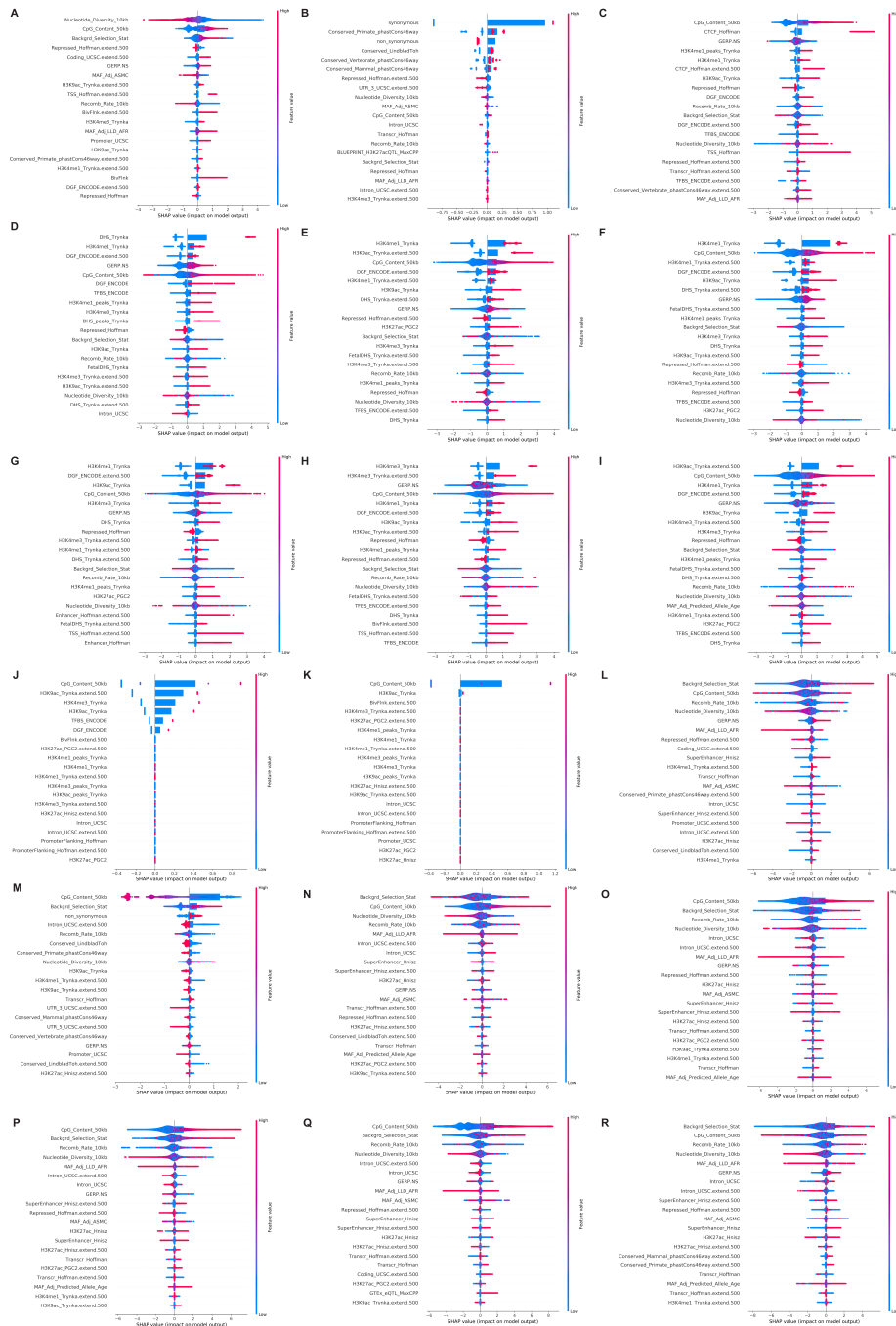


Figure D-5: Feature importance of boosted scores derived from 18 additional genome-wide scores and 47 baseline-LD model annotations. We applied SHAP<sup>5</sup> to assess which features from the baseline-LD drives the prediction of 18 boosted additional scores by AnnotBoost. We report the signed impact of top 20 features for each of 18 additional scores: (A) CDTs, (B) CCR, (C-I) DeepSEA-CTCF, -DNase, -H3K27ac, -H3K4me1, -H3K4me2, -H3K4me3, -H3K9ac, (J-K) DIS-DNA, (L) pLI, (M) LIMBR, (N-Q) Gene network connectivity-Saha, Greene, InWeb, Sonawane, (R) EDS. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see SHAP results of 47 boosted baseline-LD scores online; see Web resources).

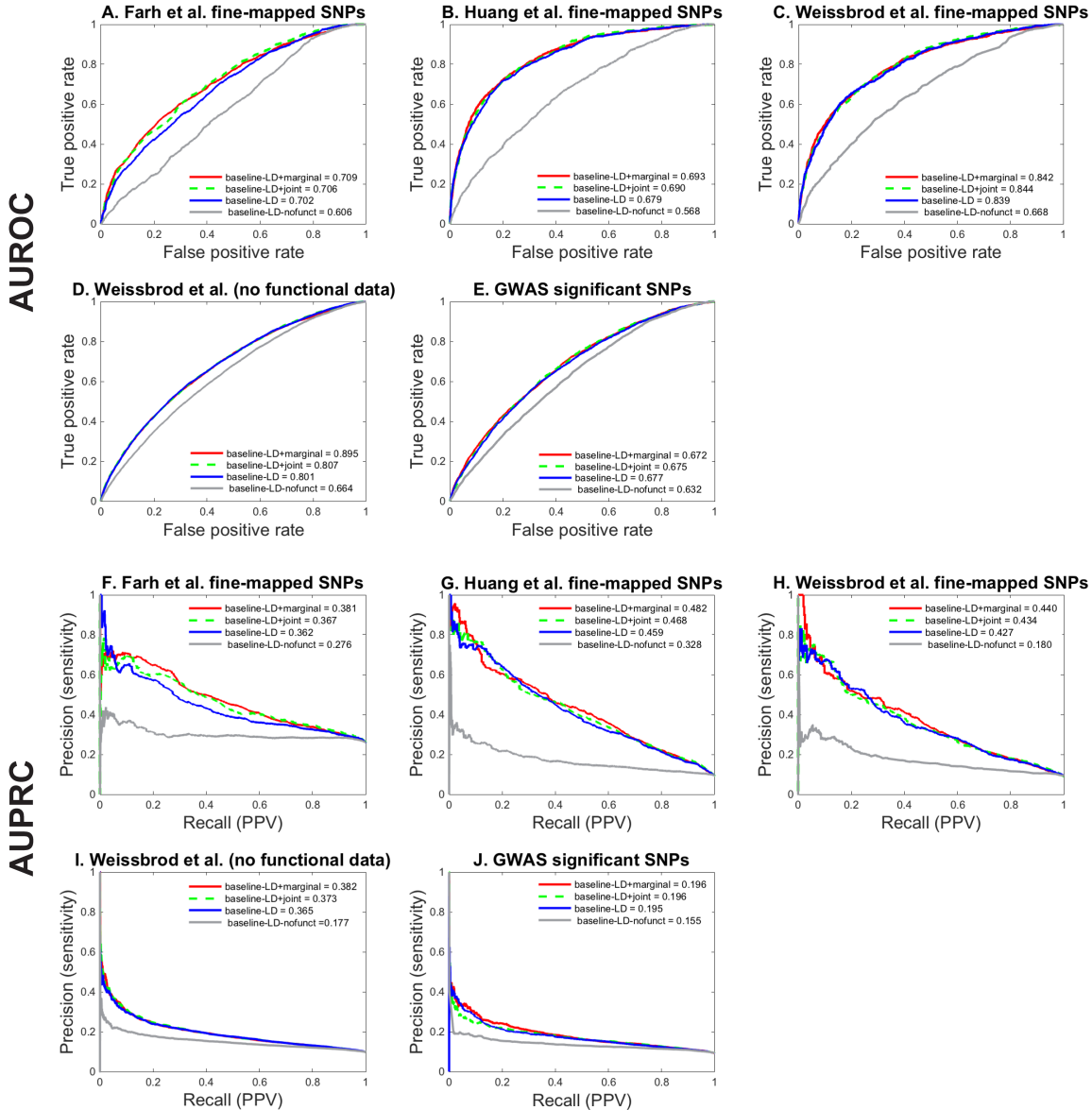


Figure D-6: **Classification of fine-mapped disease SNPs using aggregated scores.** We report the true positive rate, false positive rate, precision, and recall along with the classification accuracy (AUPRCs and AUROCs) of four aggregated scores on classifying 5 different independent SNP sets: (A, F) 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, (B, G) 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, (C, H) 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup>, (D, I) 1,379 fine-mapped SNPs without functional data for 47 traits from UK Biobank<sup>13</sup>, and (E, J) 14,807 GWAS significant SNPs<sup>14,15</sup>, from 10 LD-, MAF-, and genomic element-matched control SNPs. We report the average AUPRCs and AUROCs of even/odd-chromosome classifiers. Differences for AUROCs and AUPRCs attained between (1) baseline-LD and baseline-LD+joint model, (2) baseline-LD and baseline-LD+marginal model, and (3) baseline-LD+joint and baseline-LD+marginal were largely significant ( $p\text{-val} < 0.008$ ). Numerical results, including results using the most matched control SNPs (instead of 10), are reported in Table C.23.

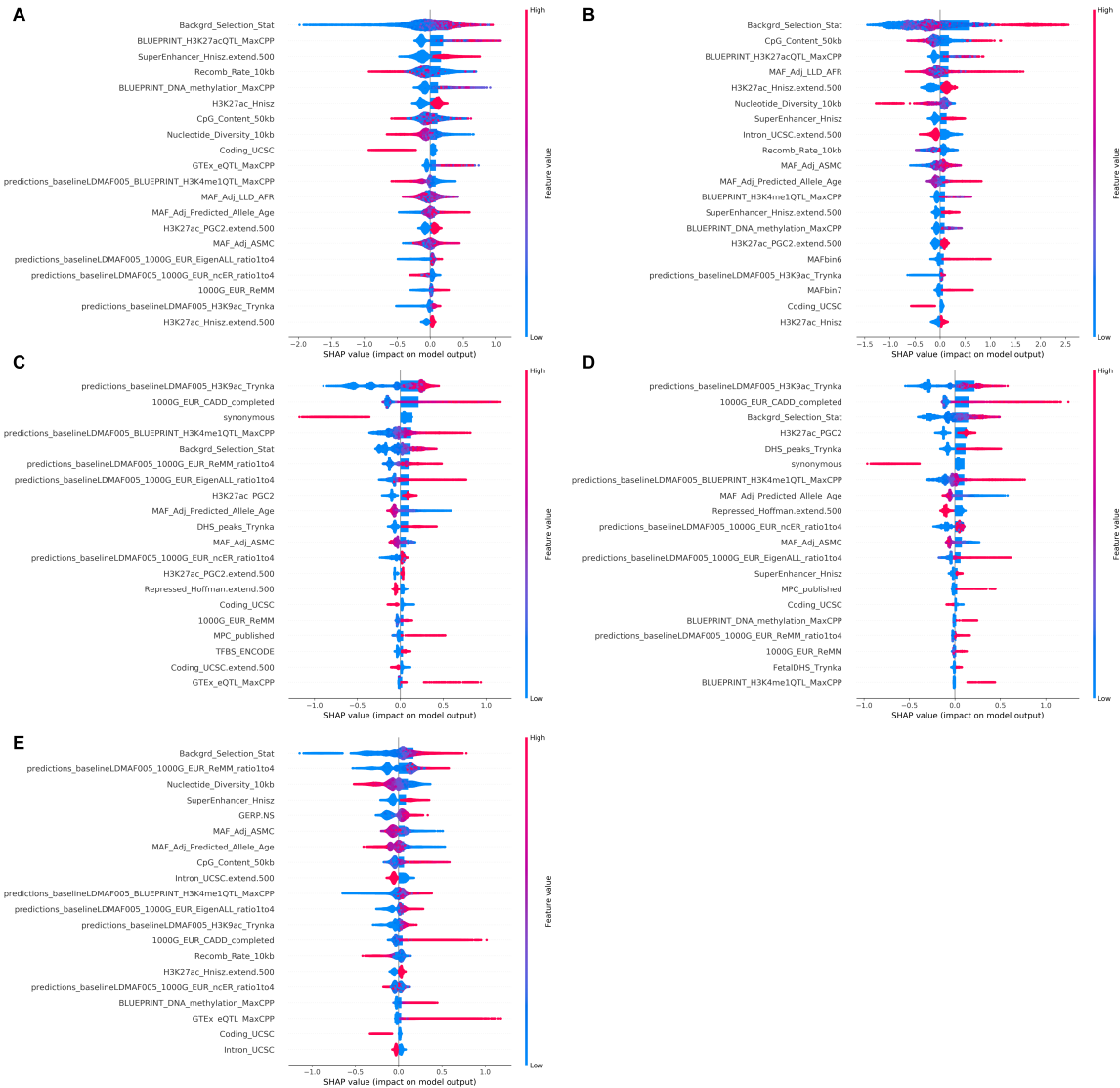


Figure D-7: **Feature importance of baseline-LD+joint model in predicting fine-mapped or GWAS significant SNPs.** We applied SHAP<sup>5</sup> to assess which features from the baseline-LD+joint drive the prediction of fine-mapped or GWAS significant SNPs from 10 matched control SNPs for each positive SNP. We report the signed impact of top 20 feature for each of 4 fine-mapped SNPs and GWAS significant SNPs: (A) Farh et al., (B) Huang et al., (C) Weissbrod et al., (D) Weissbrod et al. (fine-mapped without functional data), (E) GWAS significant SNPs. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources).

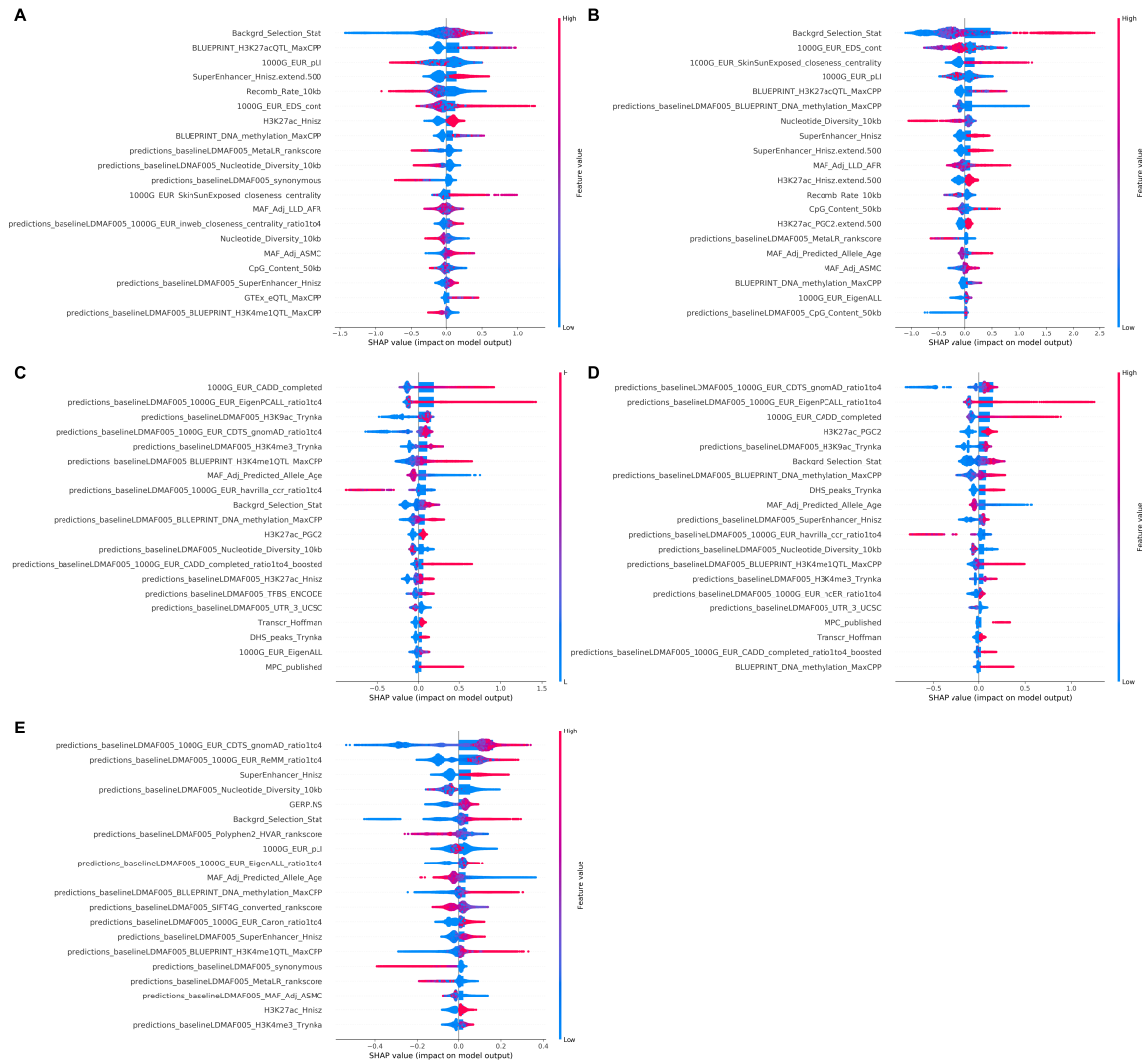


Figure D-8: Feature importance of baseline-LD+marginal model in predicting fine-mapped or GWAS significant SNPs. We applied SHAP<sup>5</sup> to assess which features from the baseline-LD+marginal drive the prediction of fine-mapped or GWAS significant SNPs from 10 matched control SNPs for each positive SNP. We report the signed impact of top 20 feature for each of 4 fine-mapped SNPs and GWAS significant SNPs: (A) Farh et al., (B) Huang et al., (C) Weissbrod et al., (D) Weissbrod et al. (fine-mapped without functional data), (E) GWAS significant SNPs. We obtained similar results for even/odd chromosome classifiers; we report odd chromosome results here (see full results online; see Web resources).



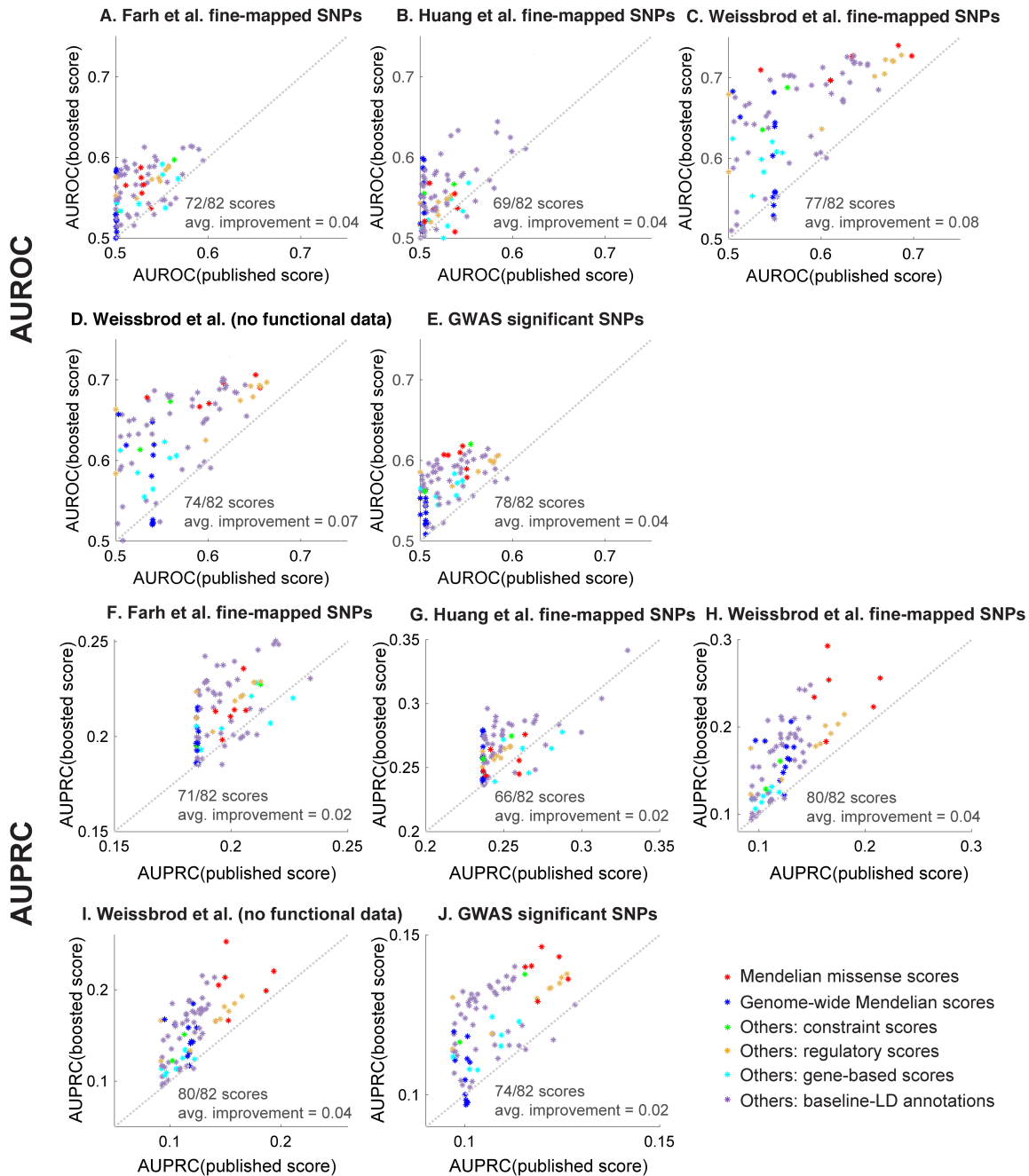


Figure D-9: **Classification of fine-mapped disease SNPs in single-score analysis using published and boosted scores.** We report the classification accuracy (AUROC and AUPRC) of each of the 82 boosted scores compared to the corresponding published score. We report AUROC (resp. AUPRC) on (A, F) 7,333 fine-mapped for 21 autoimmune diseases from Farh et al.<sup>11</sup>, (B, G) 3,768 fine-mapped SNPs for inflammatory bowel disease from Huang et al.<sup>12</sup>, (C, H) 1,851 fine-mapped SNPs for 47 traits from UK Biobank<sup>13</sup>, (D, I) 1,379 fine-mapped SNPs without functional data for 47 traits from UK Biobank<sup>13</sup>, and (E, J) 14,807 GWAS significant SNPs<sup>14,15</sup> from 10 LD-, MAF-, and genomic element-matched control SNPs. Numerical results, including results using the most matched control SNPs (instead of 10), are reported in Table C.25

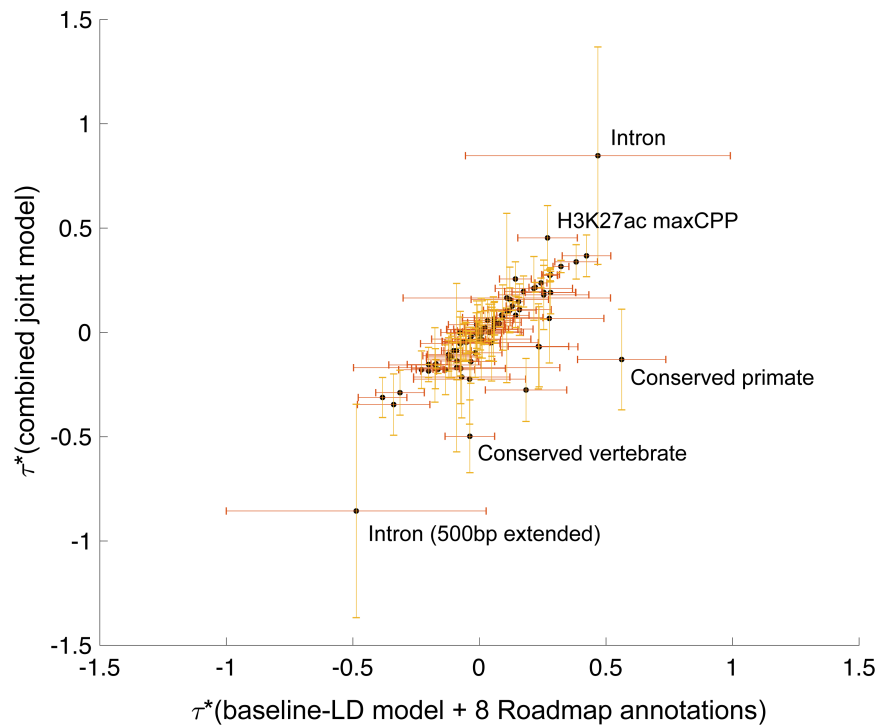


Figure D-10: **Informativeness of the baseline-LD model before and after adding 11 jointly significant binary annotations** We report meta-analyzed  $\tau^*$  of the baseline-LD model annotations, across 41 independent traits, from two different S-LDSC analyses: (1) the baseline-LD model + 8 Roadmap annotations and (2) the baseline-LD model + 8 Roadmap annotations + 11 jointly significant annotations. Error bars represent 95% confidence intervals. Numerical results are reported in Table C.27.

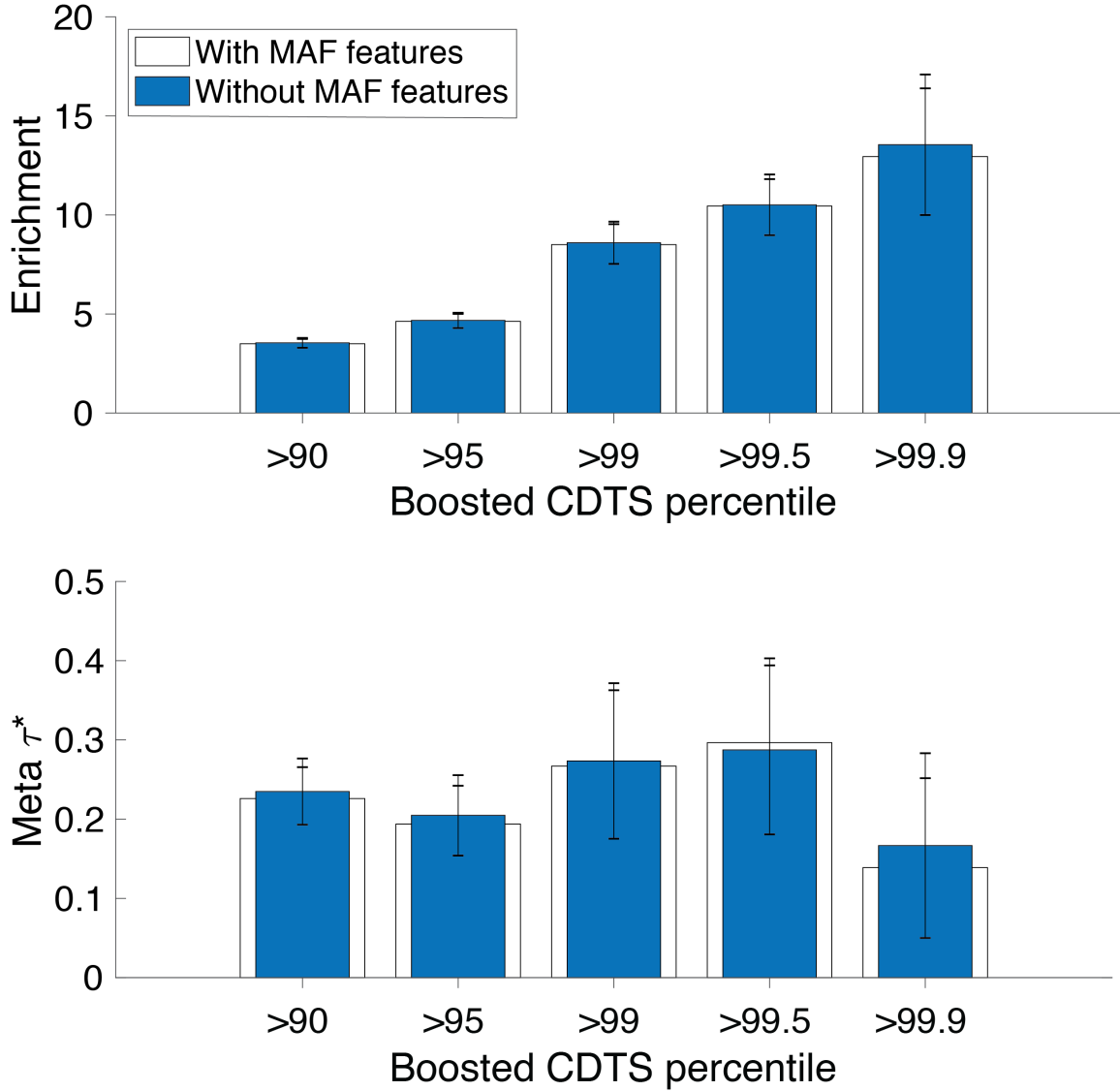


Figure D-11: **Informativeness for common disease of binary annotations derived from boosted CDTS scores with or without MAF features.** We applied AnnotBoost to CDTS annotation using all baseline-Ld features and all features excluding MAF bins. Then, we applied S-LDSC, conditioning on published binary CDTS annotations (five thresholds from 90th percentile to 99.9th percentile) and baseline-LD model annotations; and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . Error bars represent 95% confidence intervals.

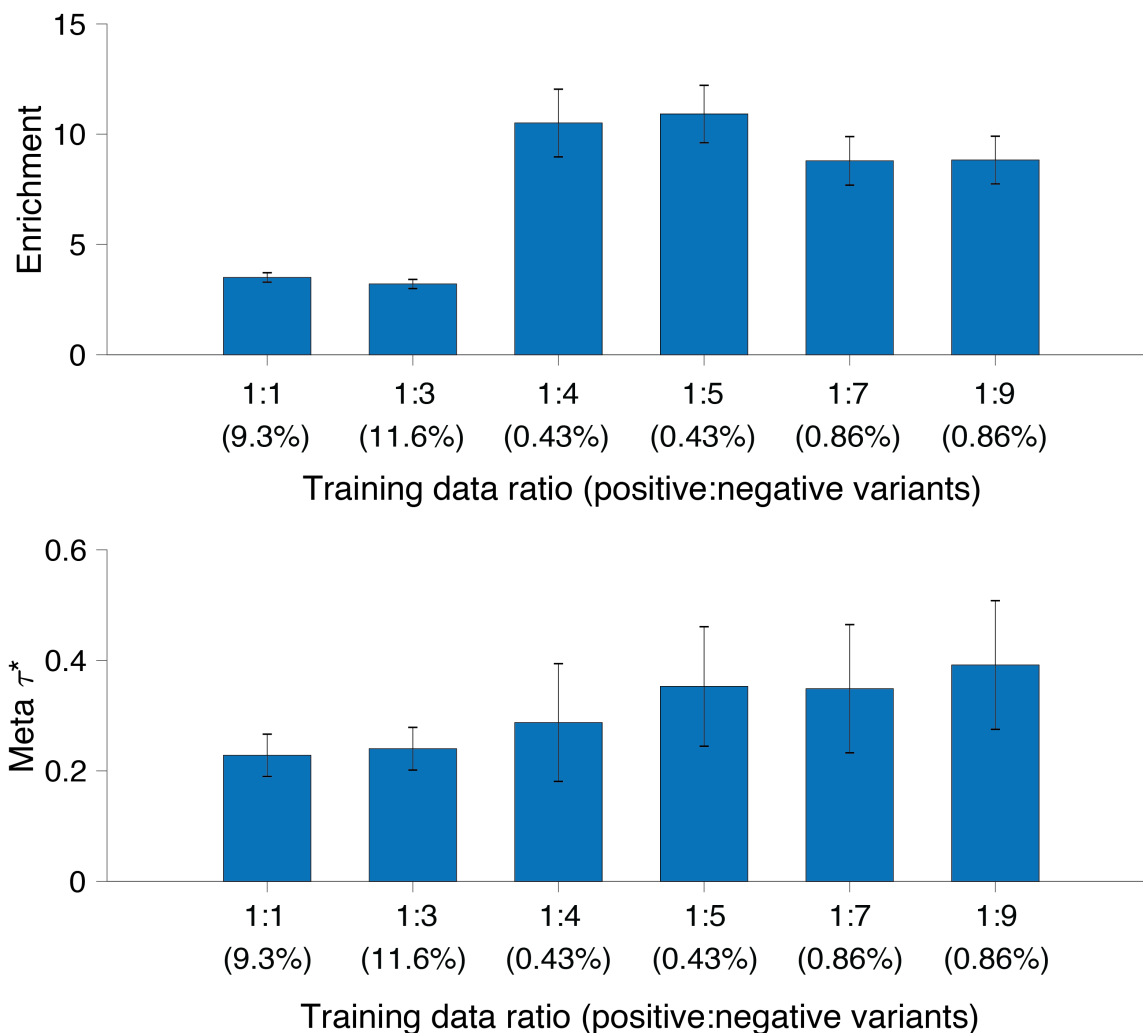


Figure D-12: **Informativeness for common disease of binary annotations derived from boosted CDTS scores using imbalanced data.** We applied AnnotBoost to CDTS annotation of varying training data. Then, we applied S-LDSC, conditioning on published binary CDTS annotations (five thresholds from 90th percentile to 99.9th percentile) and baseline-LD model annotations; and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . Error bars represent 95% confidence intervals.

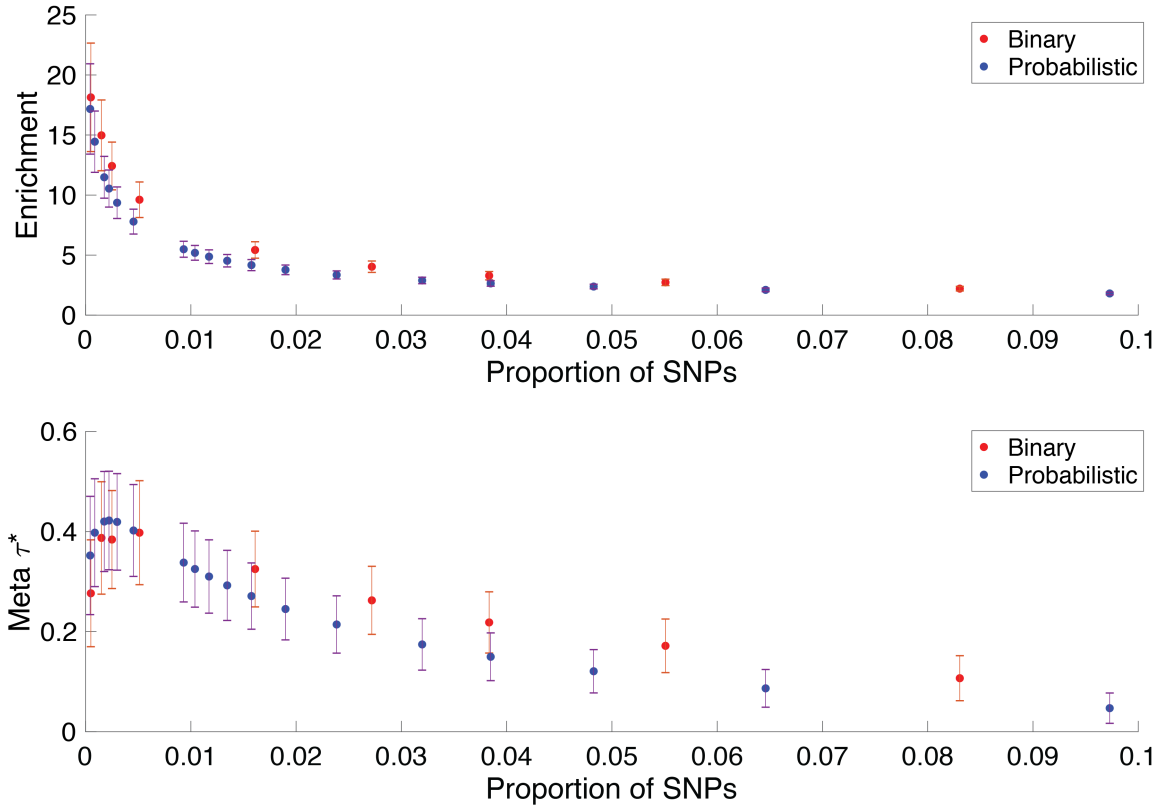


Figure D-13: **Informativeness for common disease of binary and probabilistic annotations derived from published CDTS scores.** We constructed binary and probabilistic annotations for published CDTS<sup>16</sup>. We applied S-LDSC, conditional on baseline-LD model annotations and meta-analyzed results across 41 independent traits. We report meta-analyzed enrichments and  $\tau^*$ . To construct probabilistic annotations of varying proportion of SNPs, we performed the following transformation to upweight the upper percentile and downweight the lower percentile SNPs:  $\frac{e^{\alpha \cdot annot}}{e^\alpha}$  with  $\alpha$  varied from 3 to 2000. Error bars represent 95% confidence intervals.

# Appendix E

## Supplementary tables for Chapter 3

See Kim et al.<sup>190</sup> for all supplementary tables. Titles and captions are provided below.

Table E.1: **List of 28 brain-related traits and 6 non-brain-related traits analyzed.** For each trait, we report a trait identifier, trait description, reference, sample size, and heritability z-score.

Table E.2: **Description of the baseline model.** For each of 53 functional annotations in the baseline model, we provide the description, proportion of SNPs annotated, and reference.

Table E.3: **Tissues and cell types analyzed.** We summarize 83 cell types (for 4 single-cell datasets) and 26 tissues (for bulk chromatin data) analyzed.

Table E.4: **Disease enrichments of cell-type annotations derived from fetal brain scATAC-seq data.** For each of 392 trait-cell type pairs (28 brain traits \* 14 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

Table E.5: **Matched cell types across datasets.** We report the matched cell types, defined as either exactly or closely matching cell types appearing in (A) fetal brain (scATAC-seq and scRNA-seq), (B) adult brain (scATAC-seq and scRNA-seq), (C) scATAC-seq (fetal and adult brain), and (D) scRNA-seq (fetal and adult brain).

Table E.6: **Disease enrichments of cell-type annotations derived from fetal brain scRNA-seq data.** For each of 952 trait-cell type pairs (28 brain traits \* 34 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

Table E.7: **Comparison between fetal brain scATAC-seq and fetal brain scRNA-seq cell-type annotations.** For 13 cell types appearing in both fetal brain scATAC-seq and scRNA-seq datasets, we report correlation between two corresponding annotations and  $\tau^*$  (conditioning on the baseline model, union of open chromatin regions, brain-specific enhancer-gene links of all genes, and to each other).

Table E.8: **Correlation among cell-type annotations.** We report the genome-wide Pearson correlation ( $r$ ) across cell-type annotations from 4 single-cell datasets.

Table E.9: **Disease enrichments of cell-type annotations derived from fetal brain scATAC-seq data in non-brain traits.** For 84 trait-cell type pairs (6 non-brain traits \* 14 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

Table E.10: **Disease enrichments of cell-type annotations derived from fetal brain scRNA-seq data in non-brain traits.** For 204 trait-cell type pairs (6 non-brain traits \* 34 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

Table E.11: **GREAT gene set enrichment of cell-type annotations derived from fetal brain datasets.** For each cell-type annotation (A) from scATAC-seq and (B) from scRNA-seq, we report the top enriched gene ontology (GO) terms, fold enrichment, binomial test p-value, and FDR-corrected q-value.

Table E.12: **Disease enrichments of cell-type annotations derived from fetal brain scRNA-seq data (using window-based approach).** We annotated SNPs +/- 100kb around genes, instead of brain-specific enhancer-gene links, and assessed the heritability enrichments. For each of 952 trait-cell type pairs (28 brain traits \* 34 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

**Table E.13: Disease-critical tissues using bulk chromatin data of fetal epigenomes.** For each of 980 disease-tissue-chromatin mark triplets (28 brain traits \* 5 fetal tissues \* 7 chromatin marks), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ . We considered 833 high-quality epigenomesboix2021 and selected 25 fetal biosamples related to brain tissues. We considered tier 1 assays (DNase-seq, H3K4me1, H3K4me3, H3K27ac, H3K36me3, H3K9me3 and H3K27me3), as most of observed data comes from tier 1 assays. Consistent with previous study<sup>2,4,17</sup>, we built tissue annotations by annotating chromatin peaks using the MACS2<sup>41</sup>.

**Table E.14: Disease enrichments of cell-type annotations derived from adult brain scATAC-seq data.** For each of 504 trait-cell type pairs (28 brain traits \* 18 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

**Table E.15: Disease enrichments of cell-type annotations derived from adult brain scRNA-seq data.** For each of 476 trait-cell type pairs (28 brain traits \* 17 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

**Table E.16: Comparison between adult brain scATAC-seq and adult brain scRNA-seq annotations.** For 6 cell types appearing in both adult brain scATAC-seq and scRNA-seq datasets, we report correlation between two corresponding annotations and  $\tau^*$  (conditioning on the baseline model, union of open chromatin regions, brain-specific enhancer-gene links of all genes, and to each other).

**Table E.17: Comparison between fetal brain scATAC-seq and adult brain scATAC-seq annotations.** For 3 cell types appearing in both fetal brain scATAC-seq and adult brain scATAC-seq datasets, we report correlation between two corresponding annotations and  $\tau^*$  (conditioning on the baseline model, union of open chromatin regions, brain-specific enhancer-gene links of all genes, and to each other).

**Table E.18: Comparison between fetal brain scRNA-seq and adult brain scRNA-seq annotations.** For 6 cell types appearing in both fetal brain scRNA-seq and adult brain scRNA-seq datasets, we report correlation between two corresponding annotations and  $\tau^*$  (conditioning on the baseline model, brain-specific enhancer-gene links of all genes, and to each other).



**Table E.19: Disease enrichments of cell-type annotations derived from adult brain scATAC-seq data in non-brain traits.** For 108 trait-cell type pairs (6 non-brain traits \* 18 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

**Table E.20: Disease enrichments of cell-type annotations derived from adult brain scATAC-seq data in non-brain traits.** For 102 trait-cell type pairs (6 non-brain traits \* 17 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

**Table E.21: Disease enrichments of cell-type annotations derived from scATAC-seq data, conditioning on only baseline model.** We assessed heritability enrichments of (A) adult scATAC-seq and (B) fetal scATAC-seq cell-type annotations, conditioning on only baseline model (instead of conditioning on the baseline model and the union of chromatin marks across cell types). We report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

**Table E.22: GREAT gene set enrichment of cell-type annotations derived from adult brain datasets.** For each cell-type annotation (A) from scATAC-seq and (B) from scRNA-seq, we report the top enriched gene ontology (GO) terms, fold enrichment, binomial test p-value, and FDR-corrected q-value.

**Table E.23: Disease enrichments of cell-type annotations derived from adult brain scRNA-seq data (using window-based approach).** We annotated SNPs +/- 100kb around genes, instead of brain-specific enhancer-gene links, and assessed the heritability enrichments. For each of 476 trait-cell type pairs (28 brain traits \* 17 cell types), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ .

**Table E.24: Disease-critical tissues using bulk chromatin data of adult epigenomes.** For each of 4,116 disease-tissue-chromatin mark triplets (28 brain traits \* 21 adult tissues \* 7 chromatin marks), we report the proportion of SNPs, heritability enrichments, and  $\tau^*$ . We considered 32 adult biosamples related to brain tissues and tier 1 assays (DNase-seq, H3K4me1, H3K4me3, H3K27ac, H3K36me3, H3K9me3 and H3K27me3). Donor information can be found in ref. boix2021.

Table E.25: **Genetic correlation among 28 brain-related traits.** We report the genetic correlation as estimated by S-LDSC across 28 GWAS summary association statistics.

# Appendix F

## Supplementary figures for Chapter 3

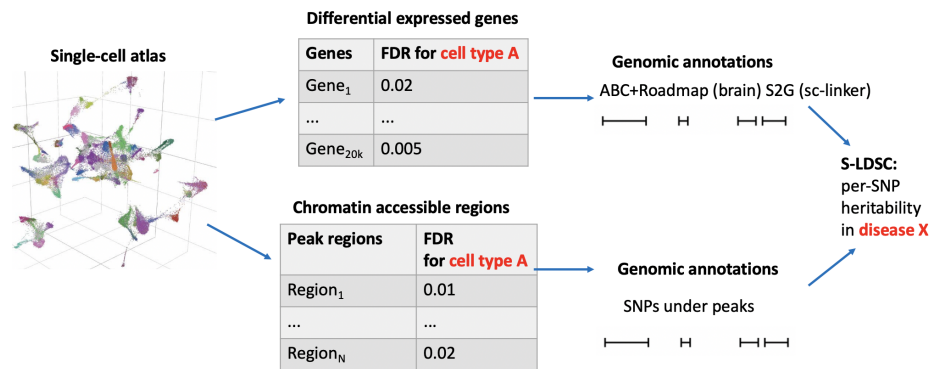


Figure F-1: **Overview of methods and analyses.** We describe the overview of methods building cell-type annotations from single-cell sequencing datasets and evaluating disease informativeness applying S-LDSC across GWAS summary statistics. ABC+Roadmap S2G refers to the brain-specific SNPs-to-Genes linking strategy using enhancer-gene links<sup>17-20</sup>. We separately analyzed fetal and adult brain data.

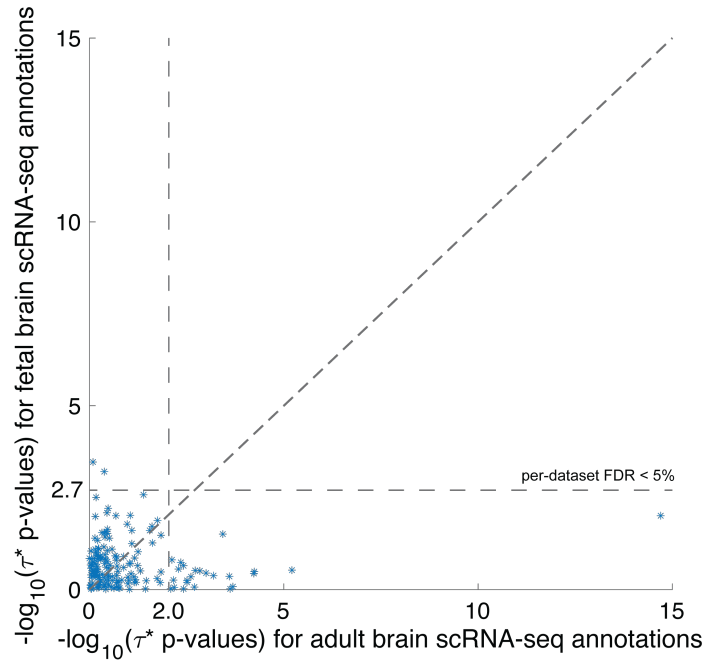


Figure F-2: **Comparison between fetal brain scRNA-seq and adult brain scRNA-seq cell-type annotations.** We report  $-\log_{10}(\tau^*$  p-values) of fetal brain scRNA-seq and adult brain scRNA-seq annotations for 6 matched cell types (astrocytes, endothelial cells, microglia, oligodendrocytes, excitatory neurons, inhibitory neurons), conditioning on the baseline model, union of open chromatin regions, and each other. Numeric results are found in Table E.18. Correlation among cell-type annotations is found in Table E.8.

# Bibliography

- [1] Speed, D., Holmes, J., and Balding, D. J. (2020). Evaluating and improving heritability models using summary statistics. *Nature Genetics* *52*, 458–462.
- [2] Finucane, H. K., Bulik-Sullivan, B., Gusev, A., Trynka, G., Reshef, Y., Loh, P.-R., Anttila, V., Xu, H., Zang, C., Farh, K., et al. (2015). Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nature Genetics* *47*, 1228.
- [3] Gazal, S., Finucane, H. K., Furlotte, N. A., Loh, P.-R., Palamara, P. F., Liu, X., Schoech, A., Bulik-Sullivan, B., Neale, B. M., Gusev, A., et al. (2017). Linkage disequilibrium-dependent architecture of human complex traits shows action of negative selection. *Nature Genetics* *49*, 1421.
- [4] Finucane, H. K., Reshef, Y. A., Anttila, V., Slowikowski, K., Gusev, A., Byrnes, A., Gazal, S., Loh, P.-R., Lareau, C., Shores, N., et al. (2018). Heritability enrichment of specifically expressed genes identifies disease-relevant tissues and cell types. *Nature Genetics* *50*, 621.
- [5] Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems* pp. 4765–4774.
- [6] Adzhubei, I. A., Schmidt, S., Peshkin, L., Ramensky, V. E., Gerasimova, A., Bork, P., Kondrashov, A. S., and Sunyaev, S. R. (2010). A method and server for predicting damaging missense mutations. *Nature methods* *7*, 248.
- [7] Adzhubei, I., Jordan, D. M., and Sunyaev, S. R. (2013). Predicting functional effect of human missense mutations using PolyPhen-2. *Current protocols in human genetics* *76*, 7–20.
- [8] Kircher, M., Witten, D. M., Jain, P., O’Roak, B. J., Cooper, G. M., and Shendure, J. (2014). A general framework for estimating the relative pathogenicity of human genetic variants. *Nature Genetics* *46*, 310.
- [9] Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J., and Kircher, M. (2018). CADD: predicting the deleteriousness of variants throughout the human genome. *Nucleic acids research* *47*, D886—D894.

- [10] Havrilla, J. M., Pedersen, B. S., Layer, R. M., and Quinlan, A. R. (2019). A map of constrained coding regions in the human genome. *Nature Genetics* *51*, 88–95.
- [11] Farh, K. K.-H., Marson, A., Zhu, J., Kleinewietfeld, M., Housley, W. J., Beik, S., Shoresh, N., Whitton, H., Ryan, R. J. H., Shishkin, A. A., et al. (2015). Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* *518*, 337.
- [12] Huang, H., Fang, M., Jostins, L., Mirkov, M. U., Boucher, G., Anderson, C. A., Andersen, V., Cleyneen, I., Cortes, A., Crins, F., et al. (2017). Fine-mapping inflammatory bowel disease loci to single-variant resolution. *Nature* *547*, 173.
- [13] Weissbrod, O., Hormozdiari, F., Benner, C., Cui, R., Ulirsch, J., Gazal, S., Schoech, A. P., Van De Geijn, B., Reshef, Y., Marquez-Luna, C., et al. (2019). Functionally-informed fine-mapping and polygenic localization of complex trait heritability. *BioRxiv* pp. 807792.
- [14] MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon, A., Milano, A., Morales, J., et al. (2016). The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Research* *45*, D896—D901.
- [15] Buniello, A., MacArthur, J. A. L., Cerezo, M., Harris, L. W., Hayhurst, J., Malangone, C., McMahon, A., Morales, J., Mountjoy, E., Sollis, E., et al. (2018). The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic acids research* *47*, D1005—D1012.
- [16] Di Iulio, J., Bartha, I., Wong, E. H. M., Yu, H.-C., Lavrenko, V., Yang, D., Jung, I., Hicks, M. A., Shah, N., Kirkness, E. F., et al. (2018). The human noncoding genome defined by genetic diversity. *Nature Genetics* *50*, 333.
- [17] Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., Ziller, M. J., et al. (2015). Integrative analysis of 111 reference human epigenomes. *Nature* *518*, 317.
- [18] Fulco, C. P., Nasser, J., Jones, T. R., Munson, G., Bergman, D. T., Subramanian, V., Grossman, S. R., Anyoha, R., Doughty, B. R., Patwardhan, T. A., et al. (2019). Activity-by-contact model of enhancer–promoter regulation from thousands of CRISPR perturbations. *Nature Genetics* *51*, 1664–1669.
- [19] Nasser, J., Bergman, D. T., Fulco, C. P., Guckelberger, P., Doughty, B. R., Patwardhan, T. A., Jones, T. R., Nguyen, T. H., Ulirsch, J. C., Natri, H. M., et al. (2020). Genome-wide maps of enhancer regulation connect risk variants to disease genes. *bioRxiv*.

- [20] Jagadeesh, K. A., Dey, K. K., Motoro, D. T., Gazal, S., Engreitz, J. M., Xavier, R. J., Price, A. L., and Regev, A. (2021). Identifying disease-critical cell types and cellular processes across the human body by integration of single-cell profiles and human genetics. *bioRxiv*.
- [21] Low, D. and Chen, K.-S. (2011). UBE3A regulates MC1R expression: a link to hypopigmentation in Angelman syndrome. *Pigment Cell & Melanoma Research* *24*, 944–952.
- [22] Simeonov, D. R., Wang, X., Wang, C., Sergeev, Y., Dolinska, M., Bower, M., Fischer, R., Winer, D., Dubrovsky, G., Balog, J. Z., et al. (2013). DNA variations in oculocutaneous albinism: an updated mutation list and current outstanding issues in molecular diagnostics. *Human Mutation* *34*, 827–835.
- [23] Okbay, A., Beauchamp, J. P., Fontana, M. A., Lee, J. J., Pers, T. H., Rietveld, C. A., Turley, P., Chen, G.-B., Emilsson, V., Meddens, S. F. W., et al. (2016). Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* *533*, 539.
- [24] Wood, A. R., Esko, T., Yang, J., Vedantam, S., Pers, T. H., Gustafsson, S., Chu, A. Y., Estrada, K., Luan, J., Kutalik, Z., et al. (2014). Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature Genetics* *46*, 1173.
- [25] Locke, A. E., Kahali, B., Berndt, S. I., Justice, A. E., Pers, T. H., Day, F. R., Powell, C., Vedantam, S., Buchkovich, M. L., Yang, J., et al. (2015). Genetic studies of body mass index yield new insights for obesity biology. *Nature* *518*, 197.
- [26] Jostins, L., Ripke, S., Weersma, R. K., Duerr, R. H., McGovern, D. P., Hui, K. Y., Lee, J. C., Schumm, L. P., Sharma, Y., Anderson, C. A., et al. (2012). Host–microbe interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature* *491*, 119.
- [27] Okada, Y., Wu, D., Trynka, G., Raj, T., Terao, C., Ikari, K., Kochi, Y., Ohmura, K., Suzuki, A., Yoshida, S., et al. (2014). Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature* *506*, 376.
- [28] Boyle, E. A., Li, Y. I., and Pritchard, J. K. (2017). An expanded view of complex traits: from polygenic to omnigenic. *Cell* *169*, 1177–1186.
- [29] Nagel, M., Jansen, P. R., Stringer, S., Watanabe, K., de Leeuw, C. A., Bryois, J., Savage, J. E., Hammerschlag, A. R., Skene, N. G., Muñoz-Manchado, A. B., et al. (2018). Meta-analysis of genome-wide association studies for neuroticism in 449,484 individuals identifies novel genetic loci and pathways. *Nature Genetics* pp. 1.

- [30] Pardiñas, A. F., Holmans, P., Pocklington, A. J., Escott-Price, V., Ripke, S., Carrera, N., Legge, S. E., Bishop, S., Cameron, D., Hamshere, M. L., et al. (2018). Common schizophrenia alleles are enriched in mutation-intolerant genes and in regions under strong background selection. *Nature Genetics* *50*, 381.
- [31] Savage, J. E., Jansen, P. R., Stringer, S., Watanabe, K., Bryois, J., de Leeuw, C. A., Nagel, M., Awasthi, S., Barr, P. B., Coleman, J. R. I., et al. (2018). Genome-wide association meta-analysis in 269,867 individuals identifies new genetic and functional links to intelligence. *Nature Genetics* pp. 1.
- [32] Wray, N. R., Ripke, S., Mattheisen, M., Trzaskowski, M., Byrne, E. M., Abdellaoui, A., Adams, M. J., Agerbo, E., Air, T. M., Andlauer, T. M. F., et al. (2018). Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. *Nature Genetics* *50*, 668.
- [33] Zhu, X. and Stephens, M. (2018). Large-scale genome-wide enrichment analyses identify new trait-associated genes and pathways across 31 human phenotypes. *Nature communications* *9*, 4361.
- [34] Cao, M., Zhang, H., Park, J., Daniels, N. M., Crovella, M. E., Cowen, L. J., and Hescott, B. (2013). Going the distance for protein function prediction: a new distance metric for protein interaction networks. *PLoS One* *8*, e76339.
- [35] Li, T., Kim, A., Mercer, J., Rosenbluh, J., Horn, H., Greenfeld, L., An, D., Zimmer, A., Liberzon, A., Bistline, J., et al. (2018). A unified web platform for network-based analyses of genomic data. *bioRxiv* pp. 196303.
- [36] Hormozdiari, F., Gazal, S., van de Geijn, B., Finucane, H. K., Ju, C. J.-T., Loh, P.-R., Schoech, A., Reshef, Y., Liu, X., O’Connor, L., et al. (2018). Leveraging molecular quantitative trait loci to understand the genetic architecture of diseases and complex traits. *Nature Genetics* *50*, 1041–1047.
- [37] Loh, P.-R., Bhatia, G., Gusev, A., Finucane, H. K., Bulik-Sullivan, B. K., Pollack, S. J., de Candia, T. R., Lee, S. H., Wray, N. R., Kendler, K. S., et al. (2015). Contrasting genetic architectures of schizophrenia and other complex diseases using fast variance-components analysis. *Nature Genetics* *47*, 1385.
- [38] Loh, P.-R., Kichaev, G., Gazal, S., Schoech, A. P., and Price, A. L. (2018). Mixed-model association for biobank-scale datasets. *Nature Genetics* *50*, 906–908.
- [39] Bulik-Sullivan, B., Finucane, H. K., Anttila, V., Gusev, A., Day, F. R., Loh, P.-R., Duncan, L., Perry, J. R. B., Patterson, N., Robinson, E. B., et al. (2015). An atlas of genetic correlations across human diseases and traits. *Nature Genetics* *47*, 1236.
- [40] 1000 Genomes Project Consortium, Others, Consortium, . G. P., and Others. (2015). A global reference for human genetic variation. *Nature* *526*, 68.



- [41] Feng, J., Liu, T., Qin, B., Zhang, Y., and Liu, X. S. (2012). Identifying ChIP-seq enrichment using MACS. *Nature protocols* 7, 1728–1740.
- [42] Kim, S. S., Dai, C., Hormozdiari, F., van de Geijn, B., Gazal, S., Park, Y., O’Connor, L., Amariuta, T., Loh, P.-R., and Finucane, H. (2019). Genes with high network connectivity are enriched for disease heritability. *The American Journal of Human Genetics* 104, 896–913.
- [43] Visscher, P. M., Wray, N. R., Zhang, Q., Sklar, P., McCarthy, M. I., Brown, M. A., and Yang, J. (2017). 10 years of GWAS discovery: biology, function, and translation. *The American Journal of Human Genetics* 101, 5–22.
- [44] Menche, J., Sharma, A., Kitsak, M., Ghiassian, S. D., Vidal, M., Loscalzo, J., and Barabási, A.-L. (2015). Uncovering disease-disease relationships through the incomplete interactome. *Science* 347, 1257601.
- [45] Barabási, A.-L., Gulbahce, N., and Loscalzo, J. (2011). Network medicine: a network-based approach to human disease. *Nature Reviews Genetics* 12, 56.
- [46] Chen, Y., Zhu, J., Lum, P. Y., Yang, X., Pinto, S., MacNeil, D. J., Zhang, C., Lamb, J., Edwards, S., Sieberts, S. K., et al. (2008). Variations in DNA elucidate molecular networks that cause disease. *Nature* 452, 429.
- [47] Szklarczyk, D., Franceschini, A., Wyder, S., Forslund, K., Heller, D., Huerta-Cepas, J., Simonovic, M., Roth, A., Santos, A., Tsafou, K. P., et al. (2014). STRING v10: protein–protein interaction networks, integrated over the tree of life. *Nucleic Acids Research* 43, D447–D452.
- [48] Chatr-Aryamontri, A., Oughtred, R., Boucher, L., Rust, J., Chang, C., Kolas, N. K., O’Donnell, L., Oster, S., Theesfeld, C., Sellam, A., et al. (2017). The BioGRID interaction database: 2017 update. *Nucleic Acids Research* 45, D369–D379.
- [49] Li, T., Wernersson, R., Hansen, R. B., Horn, H., Mercer, J., Slodkowitz, G., Workman, C. T., Rigina, O., Rapacki, K., Stærfeldt, H. H., et al. (2017). A scored human protein–protein interaction network to catalyze genomic interpretation. *Nature Methods* 14, 61.
- [50] Greene, C. S., Krishnan, A., Wong, A. K., Ricciotti, E., Zelaya, R. A., Himmelstein, D. S., Zhang, R., Hartmann, B. M., Zaslavsky, E., Sealfon, S. C., et al. (2015). Understanding multicellular function and disease with human tissue-specific networks. *Nature Genetics* 47, 569.
- [51] Saha, A., Kim, Y., Gewirtz, A. D. H., Jo, B., Gao, C., McDowell, I. C., Engelhardt, B. E., Battle, A., Aguet, F., Ardlie, K. G., et al. (2017). Co-expression networks reveal the tissue-specific regulation of transcription and splicing. *Genome Research* 27, 1843–1858.

- [52] Marbach, D., Lamparter, D., Quon, G., Kellis, M., Kutalik, Z., and Bergmann, S. (2016). Tissue-specific regulatory circuits reveal variable modular perturbations across complex diseases. *Nature Methods* *13*, 366.
- [53] Sonawane, A. R., Platig, J., Fagny, M., Chen, C.-Y., Paulson, J. N., Lopes-Ramos, C. M., DeMeo, D. L., Quackenbush, J., Glass, K., and Kuijjer, M. L. (2017). Understanding tissue-specific gene regulation. *Cell Reports* *21*, 1077–1088.
- [54] van der Wijst, M. G. P., de Vries, D. H., Brugge, H., Westra, H.-J., and Franke, L. (2018). An integrative approach for building personalized gene regulatory networks for precision medicine. *Genome Medicine* *10*, 96.
- [55] Yoon, S., Nguyen, H. C. T., Yoo, Y. J., Kim, J., Baik, B., Kim, S., Kim, J., Kim, S., and Nam, D. (2018). Efficient pathway enrichment and network analysis of GWAS summary data using GSA-SNP2. *Nucleic Acids Research* *46*, e60—e60.
- [56] Cowen, L., Ideker, T., Raphael, B. J., and Sharan, R. (2017). Network propagation: a universal amplifier of genetic associations. *Nature Reviews Genetics* *18*, 551.
- [57] Peters, L. A., Perrigoue, J., Mortha, A., Iuga, A., Song, W.-m., Neiman, E. M., Llewellyn, S. R., Di Narzo, A., Kidd, B. A., Telesco, S. E., et al. (2017). A functional genomics predictive network model identifies regulators of inflammatory bowel disease. *Nature Genetics* *49*, 1437.
- [58] Tacsan, M., Musso, G., Hao, T., Vidal, M., MacRae, C. A., and Roth, F. P. (2015). Selecting causal genes from genome-wide association studies via functionally coherent subnetworks. *Nature Methods* *12*, 154.
- [59] Gerstein, M. B., Kundaje, A., Hariharan, M., Landt, S. G., Yan, K.-K., Cheng, C., Mu, X. J., Khurana, E., Rozowsky, J., Alexander, R., et al. (2012). Architecture of the human regulatory network derived from ENCODE data. *Nature* *489*, 91.
- [60] Califano, A., Butte, A. J., Friend, S., Ideker, T., and Schadt, E. (2012). Leveraging models of cell regulation and GWAS data in integrative network-based association studies. *Nature Genetics* *44*, 841.
- [61] Lee, I., Blom, U. M., Wang, P. I., Shim, J. E., and Marcotte, E. M. (2011). Prioritizing candidate disease genes by network-based boosting of genome-wide association data. *Genome Research* pp. gr—118992.
- [62] Vanunu, O., Magger, O., Ruppin, E., Shlomi, T., and Sharan, R. (2010). Associating genes and protein complexes with disease via network propagation. *PLoS Computational Biology* *6*, e1000641.

- [63] Köhler, S., Bauer, S., Horn, D., and Robinson, P. N. (2008). Walking the interactome for prioritization of candidate disease genes. *The American Journal of Human Genetics* *82*, 949–958.
- [64] Trynka, G., Sandor, C., Han, B., Xu, H., Stranger, B. E., Liu, X. S., and Raychaudhuri, S. (2013). Chromatin marks identify critical cell types for fine mapping complex trait variants. *Nature Genetics* *45*, 124.
- [65] Maurano, M. T., Humbert, R., Rynes, E., Thurman, R. E., Haugen, E., Wang, H., Reynolds, A. P., Sandstrom, R., Qu, H., Brody, J., et al. (2012). Systematic localization of common disease-associated variation in regulatory DNA. *Science* pp. 1222794.
- [66] Consortium, E. P. and Others. (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature* *489*, 57.
- [67] Segrè, A. V., Groop, L., Mootha, V. K., Daly, M. J., Altshuler, D., Consortium, D., Investigators, M., and Others. (2010). Common inherited variation in mitochondrial genes is not enriched for associations with type 2 diabetes or related glycemic traits. *PLoS Genetics* *6*, e1001058.
- [68] de Leeuw, C. A., Mooij, J. M., Heskes, T., and Posthuma, D. (2015). MAGMA: generalized gene-set analysis of GWAS data. *PLoS Computational Biology* *11*, e1004219.
- [69] Pers, T. H., Karjalainen, J. M., Chan, Y., Westra, H.-J., Wood, A. R., Yang, J., Lui, J. C., Vedantam, S., Gustafsson, S., Esko, T., et al. (2015). Biological interpretation of genome-wide association studies using predicted gene functions. *Nature Communications* *6*, 5890.
- [70] de Leeuw, C. A., Neale, B. M., Heskes, T., and Posthuma, D. (2016). The statistical properties of gene-set analysis. *Nature Reviews Genetics* *17*, 353.
- [71] Boyle, E. A., Pritchard, J. K., and Greenleaf, W. J. (2018). High-resolution mapping of cancer cell networks using co-functional interactions. *bioRxiv* pp. 369751.
- [72] Consortium, I. H. ., Others, International HapMap 3 Consortium, and Others. (2010). Integrating common and rare genetic variation in diverse human populations. *Nature* *467*, 52.
- [73] Geer, L. Y., Marchler-Bauer, A., Geer, R. C., Han, L., He, J., He, S., Liu, C., Shi, W., and Bryant, S. H. (2010). The NCBI BioSystems database. *Nucleic Acids Research* *38*, D492–D496.
- [74] Caspi, R., Billington, R., Fulcher, C. A., Keseler, I. M., Kothari, A., Krummenacker, M., Latendresse, M., Midford, P. E., Ong, Q., Ong, W. K., et al. (2018). The MetaCyc database of metabolic pathways and enzymes. *Nucleic Acids Research* *46*, D633–D639.

- [75] Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y., and Morishima, K. (2017). KEGG: new perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Research* *45*, D353–D361.
- [76] Schaefer, C. F., Anthony, K., Krupa, S., Buchoff, J., Day, M., Hannay, T., and Buetow, K. H. (2009). PID: the Pathway Interaction Database. *Nucleic Acids Research* *37*, D674–D679.
- [77] Fabregat, A., Jupe, S., Matthews, L., Sidiropoulos, K., Gillespie, M., Garapati, P., Haw, R., Jassal, B., Korninger, F., May, B., et al. (2018). The Reactome Pathway Knowledgebase. *Nucleic Acids Research* *46*, D649–D655.
- [78] Pico, A. R., Kelder, T., Van Iersel, M. P., Hanspers, K., Conklin, B. R., and Evelo, C. (2008). WikiPathways: pathway editing for the people. *PLoS Biology* *6*, e184.
- [79] Cerami, E. G., Gross, B. E., Demir, E., Rodchenkov, I., Babur, Ö., Anwar, N., Schultz, N., Bader, G. D., and Sander, C. (2011). Pathway Commons, a web resource for biological pathway data. *Nucleic Acids Research* *39*, D685–D690.
- [80] Romero, P., Wagg, J., Green, M. L., Kaiser, D., Krummenacker, M., and Karp, P. D. (2004). Computational prediction of human metabolic pathways from the complete human genome. *Genome Biology* *6*, R2.
- [81] Yamamoto, S., Sakai, N., Nakamura, H., Fukagawa, H., Fukuda, K., and Takagi, T. (2011). INOH: ontology-based highly structured database of signal transduction pathways. *Database* *2011*.
- [82] Mi, H. and Thomas, P. (2009). PANTHER Pathway: An Ontology-Based Pathway Database Coupled with Data Analysis Tools PANTHER Pathway: An Ontology-Based Pathway Database Coupled with Data Analysis Tools. (Totowa, NJ: Humana Press).
- [83] Jewison, T., Su, Y., Disfany, F. M., Liang, Y., Knox, C., Maciejewski, A., Poelzer, J., Huynh, J., Zhou, Y., Arndt, D., et al. (2014). SMPDB 2.0: Big Improvements to the Small Molecule Pathway Database. *Nucleic Acids Research* *42*, D478–D484.
- [84] Kandasamy, K., Mohan, S. S., Raju, R., Keerthikumar, S., Kumar, G. S. S., Venugopal, A. K., Telikicherla, D., Navarro, J. D., Mathivanan, S., Pecquet, C., et al. (2010). NetPath: a public resource of curated signal transduction pathways. *Genome Biology* *11*, R3.
- [85] Eppig, J. T., Blake, J. A., Bult, C. J., Kadin, J. A., Richardson, J. E., and Group, M. G. D. (2014). The Mouse Genome Database (MGD): facilitating mouse as a model for human biology and disease. *Nucleic Acids Research* *43*, D726–D736.

- [86] Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D., Butler, H., Cherry, J. M., Davis, A. P., Dolinski, K., Dwight, S. S., Eppig, J. T., et al. (2000). Gene Ontology: tool for the unification of biology. *Nature Genetics* *25*, 25.
- [87] Liberzon, A., Subramanian, A., Pinchback, R., Thorvaldsdóttir, H., Tamayo, P., and Mesirov, J. P. (2011). Molecular signatures database (MSigDB) 3.0. *Bioinformatics* *27*, 1739–1740.
- [88] Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B. W. J. H., Jansen, R., De Geus, E. J. C., Boomsma, D. I., Wright, F. A., et al. (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nature Genetics* *48*, 245.
- [89] Özgür, A., Vu, T., Erkan, G., and Radev, D. R. (2008). Identifying gene-disease associations using centrality on a literature mined gene-interaction network. *Bioinformatics* *24*, i277—i285.
- [90] Bartha, I., di Iulio, J., Venter, J. C., and Telenti, A. (2018). Human gene essentiality. *Nature Reviews Genetics* *19*, 51.
- [91] Povey, S., Lovering, R., Bruford, E., Wright, M., Lush, M., and Wain, H. (2001). The HUGO gene nomenclature committee (HGNC). *Human genetics* *109*, 678–680.
- [92] Blake, J. A., Bult, C. J., Kadin, J. A., Richardson, J. E., Eppig, J. T., and Group, M. G. D. (2011). The Mouse Genome Database (MGD): premier model organism resource for mammalian genomics and genetics. *Nucleic Acids Research* *39*, D842—D848.
- [93] Georgi, B., Voight, B. F., and Bućan, M. (2013). From mouse to human: evolutionary genomics analysis of human orthologs of essential genes. *PLoS Genetics* *9*, e1003484.
- [94] Liu, X., Jian, X., and Boerwinkle, E. (2013). dbNSFP v2. 0: a database of human non-synonymous SNVs and their functional predictions and annotations. *Human Mutation* *34*, E2393—E2402.
- [95] Blekhman, R., Man, O., Herrmann, L., Boyko, A. R., Indap, A., Kosiol, C., Bustamante, C. D., Teshima, K. M., and Przeworski, M. (2008). Natural selection on genes that underlie human disease susceptibility. *Current Biology* *18*, 883–889.
- [96] Berg, J. S., Adams, M., Nassar, N., Bizon, C., Lee, K., Schmitt, C. P., Wilhelmsen, K. C., and Evans, J. P. (2013). An informatics approach to analyzing the incidentalome. *Genetics in Medicine* *15*, 36.
- [97] Lek, M., Karczewski, K. J., Minikel, E. V., Samocha, K. E., Banks, E., Fennell, T., O’Donnell-Luria, A. H., Ware, J. S., Hill, A. J., Cummings, B. B., et al.

- (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature* *536*, 285.
- [98] Cassa, C. A., Weghorn, D., Balick, D. J., Jordan, D. M., Nusinow, D., Samocha, K. E., O'Donnell-Luria, A., MacArthur, D. G., Daly, M. J., Beier, D. R., et al. (2017). Estimating the selective effects of heterozygous protein-truncating variants from human exome data. *Nature Genetics* *49*, 806.
- [99] Bartha, I., Rausell, A., McLaren, P. J., Mohammadi, P., Tardaguila, M., Chaturvedi, N., Fellay, J., and Telenti, A. (2015). The characteristics of heterozygous protein truncating variants in the human genome. *PLoS Computational Biology* *11*, e1004647.
- [100] Samocha, K. E., Robinson, E. B., Sanders, S. J., Stevens, C., Sabo, A., McGrath, L. M., Kosmicki, J. A., Rehnström, K., Mallick, S., Kirby, A., et al. (2014). A framework for the interpretation of de novo mutation in human disease. *Nature Genetics* *46*, 944.
- [101] Landrum, M. J., Lee, J. M., Benson, M., Brown, G. R., Chao, C., Chitipiralla, S., Gu, B., Hart, J., Hoffman, D., Jang, W., et al. (2017). ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Research* *46*, D1062—D1067.
- [102] Hamosh, A., Scott, A. F., Amberger, J. S., Bocchini, C. A., and McKusick, V. A. (2005). Online Mendelian Inheritance in Man (OMIM), a knowledgebase of human genes and genetic disorders. *Nucleic Acids Research* *33*, D514—D517.
- [103] Petrovski, S., Wang, Q., Heinzen, E. L., Allen, A. S., and Goldstein, D. B. (2013). Genic intolerance to functional variation and the interpretation of personal genomes. *PLoS Genetics* *9*, e1003709.
- [104] Lambert, S. A., Jolma, A., Campitelli, L. F., Das, P. K., Yin, Y., Albu, M., Chen, X., Taipale, J., Hughes, T. R., and Weirauch, M. T. (2018). The human transcription factors. *Cell* *172*, 650–665.
- [105] Wishart, D. S., Feunang, Y. D., Guo, A. C., Lo, E. J., Marcu, A., Grant, J. R., Sajed, T., Johnson, D., Li, C., Sayeeda, Z., et al. (2017). DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic Acids Research* *46*, D1074—D1082.
- [106] Wang, X. and Goldstein, D. B. (2020). Enhancer Domains Predict Gene Pathogenicity and Inform Gene Discovery in Complex Disease. *The American Journal of Human Genetics* *106*, 215–233.
- [107] Mainland, J. D., Li, Y. R., Zhou, T., Liu, W. L. L., and Matsunami, H. (2015). Human olfactory receptor responses to odorants. *Scientific Data* *2*, 150002.
- [108] Consortium, G. and Others. (2017). Genetic effects on gene expression across human tissues. *Nature* *550*, 204.

- [109] Auerbach, R. K., Chen, B., and Butte, A. J. (2013). Relating genes to function: identifying enriched transcription factors using the ENCODE ChIP-Seq significance tool. *Bioinformatics* *29*, 1922–1924.
- [110] Huang, D. W., Sherman, B. T., and Lempicki, R. A. (2008). Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nature Protocols* *4*, 44.
- [111] Yang, J., McAuley, J., and Leskovec, J. (2013). Community detection in networks with node attributes. In *IEEE 13th international conference on Data Mining (ICDM) IEEE* pp. 1151–1156.
- [112] Wang, B., Pourshafeie, A., Zitnik, M., Zhu, J., Bustamante, C. D., Batzoglou, S., and Leskovec, J. (2018). Network enhancement as a general method to denoise weighted biological networks. *Nature Communications* *9*, 3108.
- [113] Consortium, C.-D. G. o. t. P. G. and \textCross-Disorder Group of the Psychiatric Genomics Consortium and others. (2013). Identification of risk loci with shared effects on five major psychiatric disorders: a genome-wide analysis. *The Lancet* *381*, 1371–1379.
- [114] Ripke, S., Neale, B. M., Corvin, A., Walters, J. T. R., Farh, K.-H., Holmans, P. A., Lee, P., Bulik-Sullivan, B., Collier, D. A., Huang, H., et al. (2014). Biological insights from 108 schizophrenia-associated genetic loci. *Nature* *511*, 421.
- [115] Fatemi, S. H., Reutiman, T. J., Folsom, T. D., and Thuras, P. D. (2009). GABA A receptor downregulation in brains of subjects with autism. *Journal of Autism and Developmental Disorders* *39*, 223.
- [116] Cellot, G. and Cherubini, E. (2014). GABAergic signaling as therapeutic target for autism spectrum disorders. *Frontiers in Pediatrics* *2*, 70.
- [117] de Leeuw, C. A., Stringer, S., Dekkers, I. A., Heskes, T., and Posthuma, D. (2018). Conditional and interaction gene-set analysis reveals novel functional pathways for blood pressure. *Nature Communications* *9*, 3768.
- [118] West, N. R., Hegazy, A. N., Owens, B. M. J., Bullers, S. J., Linggi, B., Buonocore, S., Coccia, M., Görtz, D., This, S., Stockenhuber, K., et al. (2017). Oncostatin M drives intestinal inflammation and predicts response to tumor necrosis factor–neutralizing therapy in patients with inflammatory bowel disease. *Nature Medicine* *23*, 579.
- [119] Scarr, E., Udawela, M., Greenough, M. A., Neo, J., Seo, M. S., Money, T. T., Upadhyay, A., Bush, A. I., Everall, I. P., Thomas, E. A., et al. (2016). Increased cortical expression of the zinc transporter SLC39A12 suggests a breakdown in zinc cellular homeostasis as part of the pathophysiology of schizophrenia. *NPJ Schizophrenia* *2*, 16002.

- [120] Takagishi, T., Hara, T., and Fukada, T. (2017). Recent Advances in the Role of SLC39A/ZIP Zinc Transporters In Vivo. *International Journal of Molecular Sciences* *18*, 2708.
- [121] Charlesworth, B., Morgan, M. T., and Charlesworth, D. (1993). The effect of deleterious mutations on neutral molecular variation. *Genetics* *134*, 1289–1303.
- [122] O’Connor, L. J., Schoech, A. P., Hormozdiari, F., Gazal, S., Patterson, N., and Price, A. L. (2019). Polygenicity of complex traits is explained by negative selection. *The American Journal of Human Genetics* *105*, 456–476.
- [123] Parsana, P., Ruberman, C., Jaffe, A. E., Schatz, M. C., Battle, A., and Leek, J. T. (2017). Addressing confounding artifacts in reconstruction of gene co-expression networks. *bioRxiv* pp. 202903.
- [124] Speed, D., Hemani, G., Johnson, M. R., and Balding, D. J. (2012). Improved heritability estimation from genome-wide SNPs. *The American Journal of Human Genetics* *91*, 1011–1021.
- [125] Yang, J., Bakshi, A., Zhu, Z., Hemani, G., Vinkhuyzen, A. A. E., Lee, S. H., Robinson, M. R., Perry, J. R. B., Nolte, I. M., van Vliet-Ostaptchouk, J. V., et al. (2015). Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nature Genetics* *47*, 1114.
- [126] Speed, D., Cai, N., Johnson, M. R., Nejentsev, S., Balding, D. J., Consortium, U., and Others. (2017). Reevaluation of SNP heritability in complex human traits. *Nature Genetics* *49*, 986.
- [127] Gazal, S., Marquez-Luna, C., Finucane, H. K., and Price, A. L. (2019). Reconciling S-LDSC and LDAK functional enrichment estimates. *Nature Genetics* *51*, 1202–1204.
- [128] Kim, S. S., Dey, K. K., Weissbrod, O., Marquez-Luna, C., Gazal, S., and Price, A. L. (2020). Improving the informativeness of Mendelian disease-derived pathogenicity scores for common disease. *Nature communications* *11*.
- [129] Ionita-Laza, I., McCallum, K., Xu, B., and Buxbaum, J. D. (2016). A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nature Genetics* *48*, 214.
- [130] Smedley, D., Schubach, M., Jacobsen, J. O. B., Köhler, S., Zemojtel, T., Spielmann, M., Jäger, M., Hochheiser, H., Washington, N. L., McMurry, J. A., et al. (2016). A whole-genome analysis framework for effective identification of pathogenic regulatory variants in Mendelian disease. *The American Journal of Human Genetics* *99*, 595–606.
- [131] Vaser, R., Adusumalli, S., Leng, S. N., Sikic, M., and Ng, P. C. (2016). SIFT missense predictions for genomes. *Nature protocols* *11*, 1.



- [132] Ioannidis, N. M., Rothstein, J. H., Pejaver, V., Middha, S., McDonnell, S. K., Baheti, S., Musolf, A., Li, Q., Holzinger, E., Karyadi, D., et al. (2016). REVEL: an ensemble method for predicting the pathogenicity of rare missense variants. *The American Journal of Human Genetics* *99*, 877–885.
- [133] Jagadeesh, K. A., Wenger, A. M., Berger, M. J., Guturu, H., Stenson, P. D., Cooper, D. N., Bernstein, J. A., and Bejerano, G. (2016). M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity. *Nature Genetics* *48*, 1581.
- [134] Sundaram, L., Gao, H., Padigepati, S. R., McRae, J. F., Li, Y., Kosmicki, J. A., Fritzilas, N., Hakenberg, J., Dutta, A., Shon, J., et al. (2018). Predicting the clinical impact of human mutation with deep neural networks. *Nature Genetics* *50*, 1161.
- [135] Wells, A., Heckerman, D., Torkamani, A., Yin, L., Sebat, J., Ren, B., Telenti, A., and di Iulio, J. (2019). Ranking of non-coding pathogenic variants and putative essential regions of the human genome. *Nature communications* *10*, 1–9.
- [136] Caron, B., Luo, Y., and Rausell, A. (2019). NCBoost classifies pathogenic non-coding variants in Mendelian diseases through supervised learning on purifying selection signals in humans. *Genome biology* *20*, 32.
- [137] Eilbeck, K., Quinlan, A., and Yandell, M. (2017). Settling the score: variant prioritization and Mendelian disease. *Nature Reviews Genetics* *18*, 599.
- [138] Peltonen, L., Perola, M., Naukkarinen, J., and Palotie, A. (2006). Lessons from studying monogenic disease for common disease. *Human molecular genetics* *15*, R67—R74.
- [139] Blair, D. R., Lyttle, C. S., Mortensen, J. M., Bearden, C. F., Jensen, A. B., Khiabani, H., Melamed, R., Rabadan, R., Bernstam, E. V., Brunak, S., et al. (2013). A nondegenerate code of deleterious variants in Mendelian loci contributes to complex disease risk. *Cell* *155*, 70–80.
- [140] Teslovich, T. M., Musunuru, K., Smith, A. V., Edmondson, A. C., Stylianou, I. M., Koseki, M., Pirruccello, J. P., Ripatti, S., Chasman, D. I., Willer, C. J., et al. (2010). Biological, clinical and population relevance of 95 loci for blood lipids. *Nature* *466*, 707.
- [141] Kathiresan, S. and Srivastava, D. (2012). Genetics of human cardiovascular disease. *Cell* *148*, 1242–1257.
- [142] Chong, J. X., Buckingham, K. J., Jhangiani, S. N., Boehm, C., Sobreira, N., Smith, J. D., Harrell, T. M., McMillin, M. J., Wiszniewski, W., Gambin, T., et al. (2015). The genetic basis of Mendelian phenotypes: discoveries, challenges, and opportunities. *The American Journal of Human Genetics* *97*, 199–215.

- [143] Zhu, X., Need, A. C., Petrovski, S., and Goldstein, D. B. (2014). One gene, many neuropsychiatric disorders: lessons from Mendelian diseases. *Nature neuroscience* *17*, 773.
- [144] Katsanis, N. (2016). The continuum of causality in human genetic disorders. *Genome biology* *17*, 233.
- [145] Freund, M. K., Burch, K. S., Shi, H., Mancuso, N., Kichaev, G., Garske, K. M., Pan, D. Z., Miao, Z., Mohlke, K. L., Laakso, M., et al. (2018). Phenotype-specific enrichment of Mendelian disorder genes near GWAS regions across 62 complex traits. *The American Journal of Human Genetics* *103*, 535–552.
- [146] Zeng, J., De Vlaming, R., Wu, Y., Robinson, M. R., Lloyd-Jones, L. R., Yengo, L., Yap, C. X., Xue, A., Sidorenko, J., McRae, A. F., et al. (2018). Signatures of negative selection in the genetic architecture of human complex traits. *Nature genetics* *50*, 746.
- [147] Zhang, Y., Qi, G., Park, J.-H., and Chatterjee, N. (2018). Estimation of complex effect-size distributions using summary-level statistics from genome-wide association studies across 32 complex traits. *Nature genetics* *50*, 1318.
- [148] Schoech, A. P., Jordan, D. M., Loh, P.-R., Gazal, S., O’Connor, L. J., Balick, D. J., Palamara, P. F., Finucane, H. K., Sunyaev, S. R., and Price, A. L. (2019). Quantification of frequency-dependent genetic architectures in 25 UK Biobank traits reveals action of negative selection. *Nature communications* *10*, 790.
- [149] Liu, X., Jian, X., and Boerwinkle, E. (2011). dbNSFP: a lightweight database of human nonsynonymous SNPs and their functional predictions. *Human mutation* *32*, 894–899.
- [150] Liu, X., Wu, C., Li, C., and Boerwinkle, E. (2016). dbNSFP v3. 0: A one-stop database of functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Human mutation* *37*, 235–241.
- [151] Dong, C., Wei, P., Jian, X., Gibbs, R., Boerwinkle, E., Wang, K., and Liu, X. (2014). Comparison and integration of deleteriousness prediction methods for nonsynonymous SNVs in whole exome sequencing studies. *Human molecular genetics* *24*, 2125–2137.
- [152] Choi, Y., Sims, G. E., Murphy, S., Miller, J. R., and Chan, A. P. (2012). Predicting the functional effect of amino acid substitutions and indels. *PLoS one* *7*, e46688.
- [153] Choi, Y. and Chan, A. P. (2015). PROVEAN web server: a tool to predict the functional effect of amino acid substitutions and indels. *Bioinformatics* *31*, 2745–2747.

- [154] Stenson, P. D., Mort, M., Ball, E. V., Evans, K., Hayden, M., Heywood, S., Hussain, M., Phillips, A. D., and Cooper, D. N. (2017). The Human Gene Mutation Database: towards a comprehensive repository of inherited mutation data for medical research, genetic diagnosis and next-generation sequencing studies. *Human genetics* *136*, 665–677.
- [155] Samocha, K. E., Kosmicki, J. A., Karczewski, K. J., O’Donnell-Luria, A. H., Pierce-Hoffman, E., MacArthur, D. G., Neale, B. M., and Daly, M. J. (2017). Regional missense constraint improves variant deleteriousness prediction. *BioRxiv* pp. 148353.
- [156] Qi, H., Chen, C., Zhang, H., Long, J. J., Chung, W. K., Guan, Y., and Shen, Y. (2018). MVP: predicting pathogenicity of missense variants by deep learning. *bioRxiv* pp. 259390.
- [157] Landrum, M. J., Lee, J. M., Benson, M., Brown, G., Chao, C., Chitipiralla, S., Gu, B., Hart, J., Hoffman, D., Hoover, J., et al. (2015). ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic acids research* *44*, D862—D868.
- [158] UniProt Consortium. (2018). UniProt: a worldwide hub of protein knowledge. *Nucleic acids research* *47*, D506—D515.
- [159] Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining ACM* pp. 785–794.
- [160] Dey, K. K., de Geijn, B., Kim, S. S., Hormozdiari, F., Kelley, D. R., and Price, A. L. (2019). Evaluating the informativeness of deep learning annotations for human complex diseases. *bioRxiv (Nat. Comm. under revision)* *11*, 784439.
- [161] Wang, K., Li, M., and Hakonarson, H. (2010). ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic acids research* *38*, e164—e164.
- [162] Hart, T., Chandrashekhar, M., Aregger, M., Steinhart, Z., Brown, K. R., MacLeod, G., Mis, M., Zimmermann, M., Fradet-Turcotte, A., Sun, S., et al. (2015). High-resolution CRISPR screens reveal fitness genes and genotype-specific cancer liabilities. *Cell* *163*, 1515–1526.
- [163] Mohammadi, P., Castel, S. E., Cummings, B. B., Einson, J., Sousa, C., Hoffman, P., Donkervoort, S., Jiang, Z., Mohassel, P., Foley, A. R., et al. (2019). Genetic regulatory variation in populations informs transcriptome analysis in rare disease. *Science* *366*, 351–356.
- [164] Szklarczyk, D., Gable, A. L., Lyon, D., Junge, A., Wyder, S., Huerta-Cepas, J., Simonovic, M., Doncheva, N. T., Morris, J. H., Bork, P., et al. (2019). STRING v11: protein–protein association networks with increased coverage,

- supporting functional discovery in genome-wide experimental datasets. *Nucleic acids research* *47*, D607—D613.
- [165] Nelson, M. R., Wegmann, D., Ehm, M. G., Kessner, D., Jean, P. S., Verzilli, C., Shen, J., Tang, Z., Bacanu, S.-A., Fraser, D., et al. (2012). An abundance of rare functional variants in 202 drug target genes sequenced in 14,002 people. *Science* *337*, 100–104.
- [166] Consortium, G. and Others. (2020). The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* *369*, 1318–1330.
- [167] Eisenberg, E. and Levanon, E. Y. (2013). Human housekeeping genes, revisited. *TRENDS in Genetics* *29*, 569–574.
- [168] McCarthy, S., Das, S., Kretzschmar, W., Delaneau, O., Wood, A. R., Teumer, A., Kang, H. M., Fuchsberger, C., Danecek, P., Sharp, K., et al. (2016). A reference panel of 64,976 haplotypes for genotype imputation. *Nature genetics* *48*, 1279–1283.
- [169] Davis, J. and Goadrich, M. (2006). The relationship between Precision-Recall and ROC curves. In *Proceedings of the 23rd international conference on Machine learning* pp. 233–240.
- [170] Hormozdiari, F., van de Geijn, B., Nasser, J., Weissbrod, O., Gazal, S., Ju, C. J.-T., O’Connor, L., Hujoel, M. L. A., Engreitz, J., Hormozdiari, F., et al. (2019). Functional disease architectures reveal unique biological role of transposable elements. *Nature communications* *10*.
- [171] Davydov, E. V., Goode, D. L., Sirota, M., Cooper, G. M., Sidow, A., and Batzoglou, S. (2010). Identifying a high fraction of the human genome to be under selective constraint using GERP++. *PLoS computational biology* *6*, e1001025.
- [172] Zhou, J. and Troyanskaya, O. G. (2015). Predicting effects of noncoding variants with deep learning–based sequence model. *Nature methods* *12*, 931.
- [173] Zhou, J., Theesfeld, C. L., Yao, K., Chen, K. M., Wong, A. K., and Troyanskaya, O. G. (2018). Deep learning sequence-based ab initio prediction of variant effects on expression and disease risk. *Nature Genetics* *50*, 1171.
- [174] Zhou, J., Park, C. Y., Theesfeld, C. L., Wong, A. K., Yuan, Y., Scheckel, C., Fak, J. J., Funk, J., Yao, K., Tajima, Y., et al. (2019). Whole-genome deep-learning analysis identifies contribution of noncoding mutations to autism risk. *Nature Genetics* *51*, 973.
- [175] Hayeck, T. J., Stong, N., Wolock, C. J., Copeland, B., Kamalakaran, S., Goldstein, D. B., and Allen, A. S. (2019). Improved Pathogenic Variant Localization via a Hierarchical Model of Sub-regional Intolerance. *The American Journal of Human Genetics* *104*, 299–309.

- [176] Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., and O’Connell, J. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature* *562*, 203–209.
- [177] Kichaev, G., Yang, W.-Y., Lindstrom, S., Hormozdiari, F., Eskin, E., Price, A. L., Kraft, P., and Pasaniuc, B. (2014). Integrating functional data to prioritize causal variants in statistical fine-mapping studies. *PLoS genetics* *10*, e1004722.
- [178] Kichaev, G. and Pasaniuc, B. (2015). Leveraging functional-annotation data in trans-ethnic fine-mapping studies. *The American Journal of Human Genetics* *97*, 260–271.
- [179] Chen, W., McDonnell, S. K., Thibodeau, S. N., Tillmans, L. S., and Schaid, D. J. (2016). Incorporating functional annotations for fine-mapping causal variants in a Bayesian framework using summary statistics. *Genetics* *204*, 933–958.
- [180] Kichaev, G., Roytman, M., Johnson, R., Eskin, E., Lindstroem, S., Kraft, P., and Pasaniuc, B. (2017). Improved methods for multi-trait fine mapping of pleiotropic risk loci. *Bioinformatics* *33*, 248–255.
- [181] Hu, Y., Lu, Q., Powles, R., Yao, X., Yang, C., Fang, F., Xu, X., and Zhao, H. (2017). Leveraging functional annotations in genetic risk prediction for human complex diseases. *PLoS computational biology* *13*, e1005589.
- [182] Marquez-Luna, C., Gazal, S., Loh, P.-R., Kim, S. S., Furlotte, N., Auton, A., Price, A. L., Team, a. R., and Others. (2020). LDpred-funct: incorporating functional priors improves polygenic prediction accuracy in UK Biobank and 23andMe data sets. *bioRxiv*.
- [183] Kichaev, G., Bhatia, G., Loh, P.-R., Gazal, S., Burch, K., Freund, M. K., Schoech, A., Pasaniuc, B., and Price, A. L. (2019). Leveraging polygenic functional enrichment to improve GWAS power. *The American Journal of Human Genetics* *104*, 65–75.
- [184] Zeggini, E., Gloyn, A. L., Barton, A. C., and Wain, L. V. (2019). Translational genomics and precision medicine: Moving from the lab to the clinic. *Science* *365*, 1409–1413.
- [185] Jung, I., Schmitt, A., Diao, Y., Lee, A. J., Liu, T., Yang, D., Tan, C., Eom, J., Chan, M., Chee, S., et al. (2019). A compendium of promoter-centered long-range chromatin interactions in the human genome. *Nature genetics* *51*, 1442–1449.
- [186] Kanai, M., Akiyama, M., Takahashi, A., Matoba, N., Momozawa, Y., Ikeda, M., Iwata, N., Ikegawa, S., Hirata, M., Matsuda, K., et al. (2018). Genetic

- analysis of quantitative traits in the Japanese population links cell types to complex human diseases. *Nature genetics* *50*, 390–400.
- [187] Lam, M., Chen, C.-Y., Li, Z., Martin, A. R., Bryois, J., Ma, X., Gaspar, H., Ikeda, M., Benyamin, B., Brown, B. C., et al. (2019). Comparative genetic architectures of schizophrenia in East Asian and European populations. *Nature genetics* *51*, 1670–1678.
- [188] Seldin, M. F., Pasaniuc, B., and Price, A. L. (2011). New approaches to disease mapping in admixed populations. *Nature Reviews Genetics* *12*, 523–528.
- [189] Luo, Y., Li, X., Wang, X., Gazal, S., Mercader, J. M., Neale, B. M., Florez, J. C., Auton, A., Price, A. L., Finucane, H. K., et al. (2019). Estimating heritability and its enrichment in tissue-specific gene sets in admixed populations. *bioRxiv* pp. 503144.
- [190] Kim, S. S., Jagadeesh, K., Dey, K. K., Shen, A. Z., Raychaudhuri, S., Kellis, M., and Price, A. L. (2021). Leveraging single-cell ATAC-seq to identify disease-critical fetal and adult brain cell types. *bioRxiv*.
- [191] Price, A. L., Spencer, C. C. A., and Donnelly, P. (2015). Progress and promise in understanding the genetic basis of common diseases. *Proceedings of the Royal Society B: Biological Sciences* *282*, 20151684.
- [192] Hekselman, I. and Yeger-Lotem, E. (2020). Mechanisms of tissue and cell-type specificity in heritable traits and diseases. *Nature Reviews Genetics* *21*, 137–150.
- [193] Boix, C. A., James, B. T., Park, Y. P., Meuleman, W., and Kellis, M. (2021). Regulatory genomic circuitry of human disease loci by integrative epigenomics. *Nature* pp. 1–8.
- [194] Calderon, D., Bhaskar, A., Knowles, D. A., Golan, D., Raj, T., Fu, A. Q., and Pritchard, J. K. (2017). Inferring relevant cell types for complex traits by using single-cell gene expression. *The American Journal of Human Genetics* *101*, 686–699.
- [195] Ripke, S., Walters, J. T., O’Donovan, M. C., of the Psychiatric Genomics Consortium, S. W. G., et al. (2020). Mapping genomic loci prioritises genes and implicates synaptic biology in schizophrenia. *MedRxiv*.
- [196] Tanay, A. and Regev, A. (2017). Scaling single-cell genomics from phenomenology to mechanism. *Nature* *541*, 331–338.
- [197] Velmeshev, D., Schirmer, L., Jung, D., Haeussler, M., Perez, Y., Mayer, S., Bhaduri, A., Goyal, N., Rowitch, D. H., and Kriegstein, A. R. (2019). Single-cell genomics identifies cell type-specific molecular changes in autism. *Science* *364*, 685–689.

- [198] Cao, J., O’Day, D. R., Pliner, H. A., Kingsley, P. D., Deng, M., Daza, R. M., Zager, M. A., Aldinger, K. A., Blecher-Gonen, R., and Zhang, F. (2020). A human cell atlas of fetal gene expression. *Science* 370.
- [199] Domcke, S., Hill, A. J., Daza, R. M., Cao, J., O’Day, D. R., Pliner, H. A., Aldinger, K. A., Pokholok, D., Zhang, F., and Milbank, J. H. (2020). A human cell atlas of fetal chromatin accessibility. *Science* 370.
- [200] Corces, M. R., Shcherbina, A., Kundu, S., Gloudemans, M. J., Frésard, L., Granja, J. M., Louie, B. H., Eulalio, T., Shams, S., Bagdatli, S. T., et al. (2020). Single-cell epigenomic analyses implicate candidate causal variants at inherited risk loci for Alzheimer’s and Parkinson’s diseases. *Nature genetics* 52, 1158–1168.
- [201] Ulirsch, J. C., Lareau, C. A., Bao, E. L., Ludwig, L. S., Guo, M. H., Benner, C., Satpathy, A. T., Kartha, V. K., Salem, R. M., Hirschhorn, J. N., et al. (2019). Interrogation of human hematopoiesis at single-cell and single-variant resolution. *Nature Genetics* 51, 683.
- [202] Hook, P. W. and McCallion, A. S. (2020). Leveraging mouse chromatin data for heritability enrichment informs common disease architecture and reveals cortical layer contributions to schizophrenia. *Genome research* 30, 528–539.
- [203] Bryois, J., Skene, N. G., Hansen, T. F., Kogelman, L. J. A., Watson, H. J., Liu, Z., Brueggeman, L., Breen, G., Bulik, C. M., Arenas, E., et al. (2020). Genetic identification of cell types underlying brain complex traits yields insights into the etiology of Parkinson’s disease. *Nature genetics* 52, 482–493.
- [204] Kang, H. J., Kawasawa, Y. I., Cheng, F., Zhu, Y., Xu, X., Li, M., Sousa, A. M., Pletikos, M., Meyer, K. A., Sedmak, G., et al. (2011). Spatio-temporal transcriptome of the human brain. *Nature* 478, 483–489.
- [205] Pletikos, M., Sousa, A. M., Sedmak, G., Meyer, K. A., Zhu, Y., Cheng, F., Li, M., Kawasawa, Y. I., and Šestan, N. (2014). Temporal specification and bilaterality of human neocortical topographic gene expression. *Neuron* 81, 321–332.
- [206] Bakken, T. E., Miller, J. A., Ding, S.-L., Sunkin, S. M., Smith, K. A., Ng, L., Szafer, A., Dalley, R. A., Royall, J. J., Lemon, T., et al. (2016). A comprehensive transcriptional map of primate brain development. *Nature* 535, 367–375.
- [207] Li, T., Kim, A., Rosenbluh, J., Horn, H., Greenfeld, L., An, D., Zimmer, A., Liberzon, A., Bistline, J., and Natoli, T. (2018). GeNets: a unified web platform for network-based genomic analyses. *Nature methods* 15, 543–546.
- [208] Mallard, T. T., Linnér, R. K., Okbay, A., Grotzinger, A. D., de Vlaming, R., Meddens, S. F. W., Tucker-Drob, E. M., Kendler, K. S., Keller, M. C.,

- Koellinger, P. D., et al. (2019). Not just one p: Multivariate GWAS of psychiatric disorders and their cardinal symptoms reveal two dimensions of cross-cutting genetic liabilities. *BioRxiv* pp. 603134.
- [209] Korsunsky, I., Millard, N., Fan, J., Slowikowski, K., Zhang, F., Wei, K., Baglaenko, Y., Brenner, M., Loh, P.-r., and Raychaudhuri, S. (2019). Fast, sensitive and accurate integration of single-cell data with Harmony. *Nature methods* pp. 1–8.
- [210] Traag, V. A., Waltman, L., and Van Eck, N. J. (2019). From Louvain to Leiden: guaranteeing well-connected communities. *Scientific reports* *9*, 1–12.
- [211] van de Geijn, B., Finucane, H., Gazal, S., Hormozdiari, F., Amariuta, T., Liu, X., Gusev, A., Loh, P.-R., Reshef, Y., Kichaev, G., et al. (2020). Annotations capturing cell type-specific TF binding explain a large fraction of disease heritability. *Human molecular genetics* *29*, 1057–1067.
- [212] Ziffra, R. S., Kim, C. N., Wilfert, A., Turner, T. N., Haeussler, M., Casella, A. M., Przytycki, P. F., Kreimer, A., Pollard, K. S., Ament, S. A., et al. (2020). Single cell epigenomic atlas of the developing human brain and organoids. *bioRxiv* pp. 2012–2019.
- [213] McLean, C. Y., Bristor, D., Hiller, M., Clarke, S. L., Schaar, B. T., Lowe, C. B., Wenger, A. M., and Bejerano, G. (2010). GREAT improves functional interpretation of cis-regulatory regions. *Nature biotechnology* *28*, 495–501.
- [214] Moriguchi, S., Takamiya, A., Noda, Y., Horita, N., Wada, M., Tsugawa, S., Plitman, E., Sano, Y., Tarumi, R., ElSalhy, M., et al. (2019). Glutamatergic neurometabolite levels in major depressive disorder: a systematic review and meta-analysis of proton magnetic resonance spectroscopy studies. *Molecular psychiatry* *24*, 952–964.
- [215] Paul, K. N., Saafir, T. B., and Tosini, G. (2009). The role of retinal photoreceptors in the regulation of circadian rhythms. *Reviews in endocrine and metabolic disorders* *10*, 271–278.
- [216] Sabel, B. A., Wang, J., Cárdenas-Morales, L., Faiq, M., and Heim, C. (2018). Mental stress as consequence and cause of vision loss: the dawn of psychosomatic ophthalmology for preventive and personalized medicine. *EPMA journal* *9*, 133–160.
- [217] Benes, F. M. and Berretta, S. (2001). GABAergic interneurons: implications for understanding schizophrenia and bipolar disorder. *Neuropsychopharmacology* *25*, 1–27.
- [218] Dogan, B., Kazim Erol, M., Dogan, U., Habibi, M., Bulbuller, N., Turgut Coban, D., and Bulut, M. (2016). The retinal nerve fiber layer, choroidal thickness, and central macular thickness in morbid obesity: an evaluation using



- spectral-domain optical coherence tomography. *Eur Rev Med Pharmacol Sci* *20*, 886–891.
- [219] Canto, C. B., Onuki, Y., Bruinsma, B., van der Werf, Y. D., and De Zeeuw, C. I. (2017). The sleeping cerebellum. *Trends in Neurosciences* *40*, 309–323.
- [220] Batiuk, M. Y., Martirosyan, A., Wahis, J., de Vin, F., Marneffe, C., Kusserow, C., Koeppen, J., Viana, J. F., Oliveira, J. F., Voet, T., et al. (2020). Identification of region-specific astrocyte subtypes at single cell resolution. *Nature communications* *11*, 1–15.
- [221] Nagai, J., Rajbhandari, A. K., Gangwani, M. R., Hachisuka, A., Coppola, G., Masmanidis, S. C., Fanselow, M. S., and Khakh, B. S. (2019). Hyperactivity with disrupted attention by activation of an astrocyte synaptogenic cue. *Cell* *177*, 1280–1292.
- [222] Jansen, P. R., Watanabe, K., Stringer, S., Skene, N., Bryois, J., Hammerschlag, A. R., de Leeuw, C. A., Benjamins, J. S., Muñoz-Manchado, A. B., and Nagel, M. (2019). Genome-wide analysis of insomnia in 1,331,010 individuals identifies new risk loci and functional pathways. *Nature genetics* *51*, 394–403.
- [223] Dashti, H. S., Jones, S. E., Wood, A. R., Lane, J. M., Van Hees, V. T., Wang, H., Rhodes, J. A., Song, Y., Patel, K., and Anderson, S. G. (2019). Genome-wide association study identifies genetic loci for self-reported habitual sleep duration supported by accelerometer-derived estimates. *Nature communications* *10*, 1–12.
- [224] Demontis, D., Walters, R. K., Martin, J., Mattheisen, M., Als, T. D., Agerbo, E., Baldursson, G., Belliveau, R., Bybjerg-Grauholm, J., and Bækvad-Hansen, M. (2019). Discovery of the first genome-wide significant risk loci for attention deficit/hyperactivity disorder. *Nature genetics* *51*, 63–75.
- [225] Davies, G., Lam, M., Harris, S. E., Trampush, J. W., Luciano, M., Hill, W. D., Hagenaars, S. P., Ritchie, S. J., Marioni, R. E., Fawns-Ritchie, C., et al. (2018). Study of 300,486 individuals identifies 148 independent genetic loci influencing general cognitive function. *Nature communications* *9*, 1–16.
- [226] Howard, D. M., Adams, M. J., Clarke, T.-K., Hafferty, J. D., Gibson, J., Shih, M., Coleman, J. R. I., Hagenaars, S. P., Ward, J., and Wigmore, E. M. (2019). Genome-wide meta-analysis of depression identifies 102 independent variants and highlights the importance of the prefrontal brain regions. *Nature neuroscience* *22*, 343.
- [227] Stahl, E. A., Breen, G., Forstner, A. J., McQuillin, A., Ripke, S., Trubetskoy, V., Mattheisen, M., Wang, Y., Coleman, J. R. I., and Gaspar, H. A. (2019). Genome-wide association study identifies 30 loci associated with bipolar disorder. *Nature genetics* *51*, 793–803.

- [228] Nirenberg, S. and Meister, M. (1997). The light response of retinal ganglion cells is truncated by a displaced amacrine circuit. *Neuron* *18*, 637–650.
- [229] Li, M., Santpere, G., Kawasawa, Y. I., Evgrafov, O. V., Gulden, F. O., Pochareddy, S., Sunkin, S. M., Li, Z., Shin, Y., Zhu, Y., et al. (2018). Integrative functional genomic analysis of human brain development and neuropsychiatric risks. *Science* *362*.
- [230] Lee, B.-H. and Kim, Y.-K. (2010). The roles of BDNF in the pathophysiology of major depression and in antidepressant treatment. *Psychiatry investigation* *7*, 231.
- [231] Grande, I., Fries, G. R., Kunz, M., and Kapczinski, F. (2010). The role of BDNF as a mediator of neuroplasticity in bipolar disorder. *Psychiatry investigation* *7*, 243.
- [232] Favalli, G., Li, J., Belmonte-de Abreu, P., Wong, A. H. C., and Daskalakis, Z. J. (2012). The role of BDNF in the pathophysiology and treatment of schizophrenia. *Journal of psychiatric research* *46*, 1–11.
- [233] Toker, L., Mancarci, B. O., Tripathy, S., and Pavlidis, P. (2018). Transcriptomic evidence for alterations in astrocytes and parvalbumin interneurons in subjects with bipolar disorder and schizophrenia. *Biological psychiatry* *84*, 787–796.
- [234] Ferguson, B. R. and Gao, W.-J. (2018). PV interneurons: critical regulators of E/I balance for prefrontal cortex-dependent behavior and psychiatric disorders. *Frontiers in neural circuits* *12*, 37.
- [235] He, H., Mahnke, A. H., Doyle, S., Fan, N., Wang, C.-C., Hall, B. J., Tang, Y.-P., Inglis, F. M., Chen, C., and Erickson, J. D. (2012). Neurodevelopmental role for VGLUT2 in pyramidal neuron plasticity, dendritic refinement, and in spatial learning. *Journal of Neuroscience* *32*, 15886–15901.
- [236] Fernandez, F. and Garner, C. C. (2007). Over-inhibition: a model for developmental intellectual disability. *Trends in neurosciences* *30*, 497–503.
- [237] Zoghbi, H. Y. and Bear, M. F. (2012). Synaptic dysfunction in neurodevelopmental disorders associated with autism and intellectual disabilities. *Cold Spring Harbor perspectives in biology* *4*, a009886.
- [238] Runge, K., Mathieu, R., Bugeon, S., Lafi, S., Beurrier, C., Sahu, S., Schaller, F., Loubat, A., Herault, L., Gaillard, S., et al. (2020). Disruption of the transcription factor NEUROD2 causes an autism syndrome via cell-autonomous defects in cortical projection neurons. *BioRxiv* pp. 296889.
- [239] Ruzzo, E. K., Pérez-Cano, L., Jung, J.-Y., Wang, L.-k., Kashef-Haghighi, D., Hartl, C., Singh, C., Xu, J., Hoekstra, J. N., Leventhal, O., et al. (2019). Inherited and de novo genetic risk for autism impacts shared networks. *Cell* *178*, 850–866.

- [240] Chung, W.-S., Clarke, L. E., Wang, G. X., Stafford, B. K., Sher, A., Chakraborty, C., Joung, J., Foo, L. C., Thompson, A., Chen, C., et al. (2013). Astrocytes mediate synapse elimination through MEGF10 and MERTK pathways. *Nature* *504*, 394–400.
- [241] Sloan, S. A., Darmanis, S., Huber, N., Khan, T. A., Birey, F., Caneda, C., Reimer, R., Quake, S. R., Barres, B. A., and Pasca, S. P. (2017). Human astrocyte maturation captured in 3D cerebral cortical spheroids derived from pluripotent stem cells. *Neuron* *95*, 779–790.
- [242] Freimer, J. W., Shaked, O., Naqvi, S., Sinnott-Armstrong, N., Kathiria, A., Chen, A. F., Cortez, J., Greenleaf, W. J., Pritchard, J. K., and Marson, A. (2021). Systematic discovery and perturbation of regulatory genes in human t cells reveals the architecture of immune networks. *bioRxiv*.
- [243] Dixit, A., Parnas, O., Li, B., Chen, J., Fulco, C. P., Jerby-Arnon, L., Marjanovic, N. D., Dionne, D., Burks, T., Raychowdhury, R., et al. (2016). Perturb-Seq: dissecting molecular circuits with scalable single-cell RNA profiling of pooled genetic screens. *Cell* *167*, 1853–1866.
- [244] Kircher, M., Xiong, C., Martin, B., Schubach, M., Inoue, F., Bell, R. J. A., Costello, J. F., Shendure, J., and Ahituv, N. (2019). Saturation mutagenesis of twenty disease-associated regulatory elements at single base-pair resolution. *Nature communications* *10*, 1–15.
- [245] Hanna, R. E., Hegde, M., Fagre, C. R., DeWeirdt, P. C., Sangree, A. K., Szegetes, Z., Griffith, A., Feeley, M. N., Sanson, K. R., Baidi, Y., et al. (2021). Massively parallel assessment of human variants with base editor screens. *Cell* *184*, 1064–1080.
- [246] Ma, S., Zhang, B., LaFave, L. M., Earl, A. S., Chiang, Z., Hu, Y., Ding, J., Brack, A., Kartha, V. K., Tay, T., et al. (2020). Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin. *Cell* *183*, 1103–1116.
- [247] Schaid, D. J., Chen, W., and Larson, N. B. (2018). From genome-wide associations to candidate causal variants by statistical fine-mapping. *Nature Reviews Genetics* *19*, 491–504.
- [248] Mancuso, N., Freund, M. K., Johnson, R., Shi, H., Kichaev, G., Gusev, A., and Pasaniuc, B. (2019). Probabilistic fine-mapping of transcriptome-wide association studies. *Nature genetics* *51*, 675–682.