

Noname manuscript No.
(will be inserted by the editor)

Annotation Regression for Genome-Wide Association Studies with an Application to Psychiatric Genomic Consortium Data

Sunyoung Shin · Sündüz Keleş

Received: December, 2015 / Accepted: date

Abstract Although genome-wide association studies (GWAS) have been successful at finding thousands of disease-associated genetic variants (GVs), identifying causal variants and elucidating the mechanisms by which genotypes influence phenotypes are critical open questions. A key challenge is that a large percentage of disease-associated GV are potential regulatory GV located in noncoding regions, making them difficult to interpret. Recent research efforts focus on going beyond annotating GV by integrating functional annotation data with GWAS to prioritize GV. However, applicability of these approaches are challenged by high dimensionality and heterogeneity of functional annotation data. Furthermore, existing methods often assume global associations of GV with annotation data. This strong assumption is susceptible to violations for GV involved in many complex diseases. To address these issues, we develop a general regression framework, named **Annotation Regression for GWAS (ARoG)**. ARoG is based on finite mixture of linear regression models where GWAS association measures are viewed as responses and functional annotations as predictors. This mixture framework addresses heterogeneity of impacts of GV by grouping them into clusters and high dimensionality of the functional annotations by enabling annotation selection within each cluster. ARoG employs permutation testing to evaluate the significance of selected annotations. Computational experiments indicate that ARoG can discover distinct associations between disease risk and functional annotations. Application of ARoG to autism and schizophrenia data from Psychiatric Genomics Con-

Sunyoung Shin
Department of Statistics, Department of Biostatistics and Medical Informatics, University of Wisconsin Madison
E-mail: shin@stat.wisc.edu

Sündüz Keleş
Department of Statistics, Department of Biostatistics and Medical Informatics, University of Wisconsin Madison
E-mail: keles@stat.wisc.edu

sortium led to identification of GVs that significantly affect interactions of several transcription factors with DNA as potential mechanisms contributing to these disorders.

Keywords Finite mixture of regressions, Functional genomic data, Genome-wide association studies, Integrative analysis, Regularized variable selection

1 Introduction

Although genome-wide association studies (GWAS) have successfully identified thousands of genetic loci associated with human diseases, design and analysis of these studies are challenged in two critical aspects. First, existing GWAS have revealed that for many common disorders, the typical genetic architecture encompasses many genetic variants (GV) with individually small effects on the phenotype [28], indicating the need for larger sample sizes to reliably identify them. Second, roles of a large proportion of identified GVs remain elusive since they reside in non-coding regions. Today, with the availability of affordable whole genome sequencing, our ability to elucidate the mechanisms by which genotypes influence phenotypes is far behind our ability to identify phenotype-associated variants.

In parallel to the rapid developments in the design and analysis of GWAS, large consortia projects such as the Encyclopedia of DNA Elements (ENCODE [14, 34]), Roadmap Epigenome Mapping Consortium (REMC [24]), the Genotype-Tissue Expression Consortium (GTEx [29]), and the International Human Epigenome Consortium (IHEC [3]) as well as many investigator-driven projects are generating diverse data types of RNA transcription (RNA-seq), DNA accessibility (DNase-seq), DNA methylation (Methyl-seq), protein-DNA interactions (ChIP-seq/exo), protein-RNA interactions (CLIP-seq), and chromatin state (Histone ChIP-seq) across diverse cell/tissue types. Although there is a growing literature on methods for utilizing one or more classes of these functional annotation data to support GWAS results, these data have rarely played more than an indirect role in assessing evidence for association in these approaches. They are commonly used to follow up identified significant GVs or prioritize them for causality [10, 13, 17, 21]. Specifically, [15, 21, 30] used annotations to model prior probabilities of association of GVs under Bayesian framework or hierarchical modeling. [10, 17] developed models to integrate binarized functional annotation data and GWAS summary statistics such as p-values and z-scores of GVs. [13] correlated annotations with GWAS association status, and used the estimated odds of association to estimate the posterior odds of association of GVs for prioritization. A significant shortcoming of these methods is that they aim to globally relate associations of GVs to functional annotation data despite the fact that the same disease mechanism might be governed by distinct functional annotations. For example, disruption of an important pathway may arise by GVs in coding regions of the genes and/or in one or more of their regulatory mechanisms. Regulatory GVs may

have a variety of mechanisms such as transcription factor (TF) binding, histone modifications, enhancer activity through chromatin architecture, DNA methylation, and alternative splicing [20]. Another significant shortcoming is that these approaches use functional annotation in an agnostic manner by ignoring the relevance of the tissue/cell type that the annotation is drawn from to the tissue/cell type that is most relevant to disease etiology. In other words, they are not equipped to automatically select important annotations. Furthermore, several of them can only use annotation data in specific formats (e.g., most recent genetic analysis incorporating pleiotropy and annotation (GPA) [10] requires binary annotation variables) due to computational impediments. A useful common key feature of many existing methods is that they use population level GWAS data in the form of summary statistics of GVs from GWAS [10, 13, 17, 30] as opposed to individual subject-level data which is not immediately available publicly [32].

To overcome these challenges, we develop a regression framework named **Annotation Regression for GWAS (ARoG)** and integrate GWAS and functional annotation data. ARoG models GWAS association measures, e.g., z-scores from univariate analysis of GWAS, as a linear function of functional annotations. It employs a mixture of linear regressions framework to accommodate the heterogeneity of associations between GWAS association measures and functional annotations. It aims to capture locally distinct associations that would not be revealed with an analysis that assumes homogeneity of these associations. A critical aspect of ARoG is that it can automatically select relevant annotations among a large number of annotations with penalization techniques. ARoG works with all the commonly used functional annotation types captured by categorical or continuous variables. The rest of the paper is organized as follows. Section 2 presents empirical observations regarding GWAS association measures and functional annotations using Psychiatric Genomics Consortium (PGC) data. Section 3 develops ARoG and discusses implementation details. In Section 4, we analyze PGC autism and schizophrenia data and identify GVs that are associated with these diseases and have the potential to modulate TF-DNA interactions. Section 5 presents computational experiments with a wide variety of settings including a PGC analysis-driven one. In Section 6, we provide concluding remarks and discuss extensions.

2 Exploring Psychiatric Genomics Consortium Data with Functional Annotations

PGC has conducted analysis of combined GWAS data from separate studies. Specifically, they examined five psychiatric disorders: attention deficit/hyperactivity disorder (ADHD), autism spectrum disorder (ASD), bipolar disorder (BIP), major depressive disorder (MDD), and schizophrenia (SCZ) [4, 11]. As a result, they identified 4 genome-wide significant loci for BIP [22], and more than 100 genome-wide significant loci for SCZ [25, 26]. However, this analysis did not lead to any reproducible genome-wide significant loci for ADHD and

MDD, and the analysis on ASD is in progress. Their GWAS summary datasets are publicly available at <http://www.med.unc.edu/pgc/downloads>. In what follows, we focus on AUT and SCZ data.

2.1 Autism GWAS

We generated a set of candidate SNPs by starting with the intersection of SNPs genotyped in all five disorder datasets from the PGC cross-disorder study [11]. After lifting the original genomic coordinates from hg18 to hg19, we obtained 1,219,561 SNPs common to all five disorders. We next selected the subset of the SNPs with a Benjamini-Hochberg (BH) adjusted association p-value smaller than or equal to 0.1 in any of the five disorders [6]. This led to a total of 1,430 SNPs. Next, we included 761 linkage disequilibrium (LD) partners of these SNPs as identified by the SNAP tool [16] with an $r^2 \geq 0.8$ to one or more of the 1,430 SNPs. As part of pre-processing, we discarded SNPs with more than one nucleotide on the reference genome, SNPs with nucleotide mismatches between the PGC dataset and the SNP Database dbSNP [2], and SNPs not listed in dbSNP [2]. This resulted in a total of 2,191 SNPs for ARoG analysis. Supplementary Figure 1(a) displays the histogram of the autism z-scores for these sets of SNPs and illustrates that, as expected, LD partners tend to contribute z-scores around zero to the overall distribution since they had BH adjusted p-values larger than 0.1 in the initial selection step. The manhattan plot of the p-values in Figure 1(a) indicates that the SNPs with the strongest association are on chr 5 (6 of them) and chr 6 (3 of them) with p-values ranging from 5.37×10^{-7} to 1.93×10^{-7} . All of these have raw p-values less than 10^{-6} ; however, they make neither the conventional GWAS p-value cutoff of 5×10^{-8} nor the Bonferroni cutoff of 4.1×10^{-8} specific for this study, indicating that common practice for GWAS analysis would not confidently identify significant SNPs from this study.

Currently, most integrative analysis methods consider functional annotations enabled by the large scale analysis results of consortia projects such as ENCODE (e.g., [10, 17, 21]). In our exposition, we consider a specific class of functional annotation, namely, impact of single nucleotide polymorphisms (SNPs) on TF binding. Specifically, we used atSNP [36] to create an annotation score matrix for the 2,191 SNPs. atSNP quantifies the impact of SNPs, i.e., the likelihood that a given SNP disrupts or enhances the binding sites from a given set of position weight matrices (PWMs) characterizing the class of sequences TFs recognize. atSNP operates by scanning through subsequences overlapping with the SNP position with reference and SNP alleles for the best matches of both to a given PWM. It quantifies the significance of the best matches with the reference and SNP alleles by p-values. Then, the log ratio of the two p-values are defined as the atSNP annotation score, which empirically reflects the change in the ranks of the PWM matches of the alleles. SNPs likely to enhance or disrupt binding of given TF have large absolute atSNP scores for the corresponding PWM while SNPs with little potential impact on binding

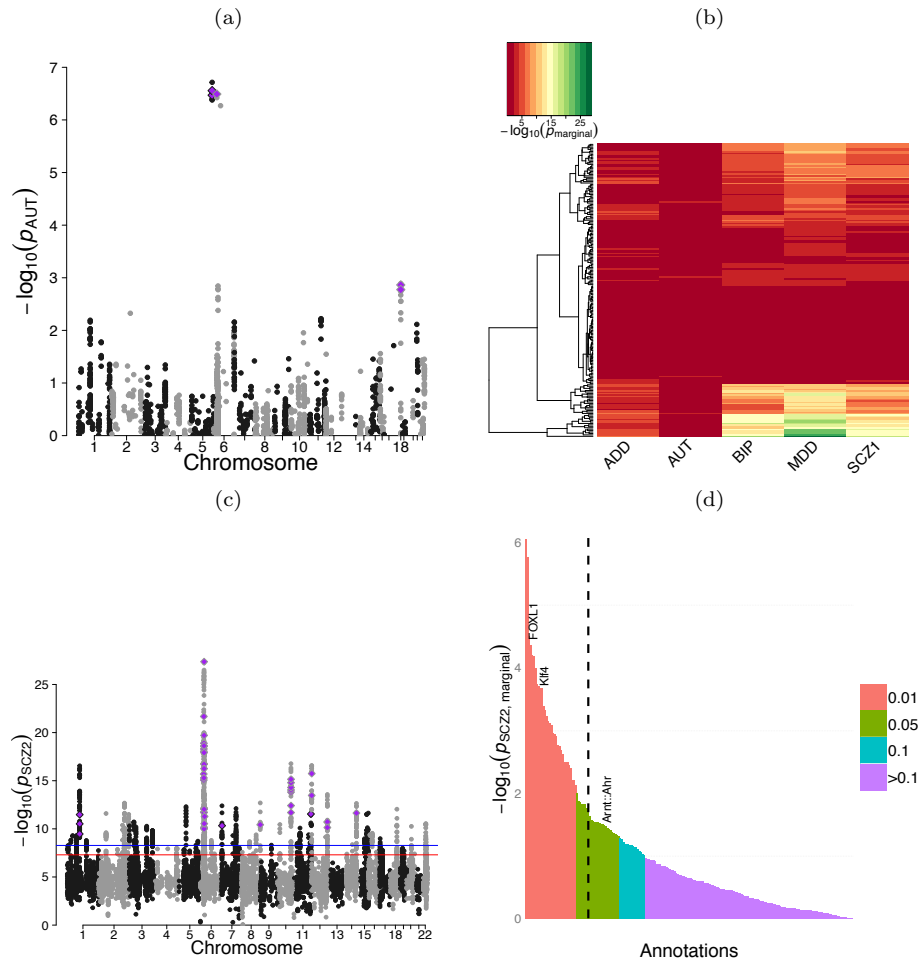


Fig. 1: (a) Manhattan plot for autism association p-values across all 2,191 SNPs. ARoG SNPs identified in Section 4.1 are marked with purple diamonds. (b) Heatmap of $-\log_{10}$ p-values from marginal regressions of z-scores on annotation scores. (c) Manhattan plot for SCZ2 association p-values across all 11,386 SNPs. ARoG SNPs identified in Section 4.3 are marked with purple diamonds. Blue and red horizontal lines depict the Bonferroni cut-off at significance level of 0.05 and the conventional p-value cutoff of 5×10^{-8} . (d) Ranking of transcription factors based on marginal regressions of the SCZ2 z-scores on annotations scores for each TF. Transcription factors FOXL1, Klf4, and Arnt::Ahr that are identified as associated with the z-scores in Section 4.3 are labeled. The dashed vertical line depicts the BH cut-off at significance level 0.1.

have scores close to zero. We refer to [36] for further computational details. We considered the JASPAR CORE database [18] for vertebrates with 205 PWMs as our motif library and scored the SNP set. atSNP evaluates the impact of SNPs on the binding affinities of the TFs by evaluating and comparing the best sub-sequence matches overlapping the SNP position to the PWM with both SNP and the reference alleles. Large positive and negative atSNP scores indicate enhancement and disruption of TF binding, respectively. Supplementary Figure 1(b) displays the heatmap of the resulting annotation score matrix along with the z-scores. Here, only SNPs colored as dark green or red are likely to lead to significant changes (as assessed by atSNP p-values) in TF binding. As part of our exploratory analysis, we first regressed z-scores from each of the five disorders on each annotation score separately. Figure 1(b) displays the $-\log_{10}$ transformed p-values from these marginal regressions. We note that the overall association of the z-scores and functional annotations for some diseases are apparent (e.g., MDD). However, for autism, none of the annotations can be deemed as contributing to the variation in the autism z-scores based on this global marginal analysis as all BH adjusted p-values are greater than 0.1 (Supplementary Figure 1(c)).

2.2 Schizophrenia GWAS

PGC provides two mega analyses of schizophrenia GWAS. We will refer to the first study, results of which are included in the above five disorder exploratory analysis as SCZ1 [25] and the second as SCZ2 [26]. SCZ1 and SCZ2 have genotypes for 1,252,901 and 9,444,230 SNPs, respectively. Intersection of these two leads to 1,179,262 SNPs. We filtered out SNPs with BH adjusted p-values larger than 0.01 for both studies, and retained 8,029 SNPs. Similar to the autism analysis, we also excluded SNPs with multiple reference alleles, with allele mismatches between PGC dataset and dbSNP, and SNPs that are not in dbSNP. We next extended this set by including their LD partners with $r^2 \geq 0.8$. Our final set of SNPs for this analysis included 11,386 SNPs. z-scores of the 11,386 SNPs from both SCZ1 and SCZ2 have a bimodal distribution (Supplementary Figure 1(d)). However, SCZ2 has many more statistically significant SNPs as evidenced by long tails of the z-score distribution. This can be viewed as increased precision since the sample size of SCZ2 is four times as large as that of SCZ1. In what follows, we perform ARoG analysis on SCZ2 dataset, and use SCZ1 dataset for validation. The manhattan plot in Figure 1(c) indicates that genome-wide significant SNPs from SCZ2 spread throughout the genome.

We next generated an $11,386 \times 205$ annotation score matrix using atSNP with the JASPAR PWM library. Supplementary Figure 1(e) displays the heatmap of the resulting annotation score matrix along with the SCZ1 and SCZ2 z-scores and indicates that only a small proportion of SNPs change the binding affinities of TFs, and only a few TFs have noticeable binding affinity changes due to the SNPs. Marginal regressions of z-scores on annotation scores

identify 40 annotations which are significant when adjusted for multiple testing by the BH procedure at level 0.1 (Figure 1(d)). However, given the large sample size, i.e., the number of SNPs, in this marginal analysis, we view these associations as suggestive and turn our attention to developing a framework that can identify subgroups of SNPs with association measures explained by a subgroup of functional annotations.

3 A Mixture of Linear Regressions Framework for Incorporating Functional Annotations into GWAS Analysis

In the ARoG framework, both the predictors and the response measure effects of SNPs: the predictors capture effects on transcription factor binding affinities and the response on disease/trait. ARoG associates them in a regression framework, and aims to increase the detection power for association SNPs with this additional functional information.

3.1 Basic Annotation Regression for GWAS (ARoG(I))

Let $z_i \in \mathbb{R}, i = 1, \dots, n$ denote z-scores for n SNPs from a GWAS and, $\mathbf{x}_i = (x_{i0}, \dots, x_{ip}) \in \mathbb{R}^{p+1}$ denote a vector for p functional annotations of the i -th SNP with the first element of 1 as the intercept term. ARoG assumes that n SNPs can be partitioned into K clusters with finite mixture of linear regression models. Following the notation of FMRLasso developed by [27], we denote the prior probability of the k -th cluster as π_k , its regression parameters as $\boldsymbol{\beta}_k = (\beta_{k0}, \dots, \beta_{kp})^T$, and its variance as σ_k^2 . The z-score of a SNP from the k th cluster with a functional annotation vector, \mathbf{x} , is assumed to be normally distributed with mean $\mathbf{x}^T \boldsymbol{\beta}_k$ and variance σ_k^2 . The conditional density function of z given \mathbf{x} is then

$$f_{\boldsymbol{\xi}}(z|\mathbf{x}) = \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(z - \mathbf{x}^T \boldsymbol{\beta}_k)^2}{2\sigma_k^2}\right), \quad (1)$$

where $\boldsymbol{\xi} = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K, \sigma_1, \dots, \sigma_K, \pi_1, \dots, \pi_{K-1}) \in \mathbb{R}^{K \cdot (p+1)} \times \mathbb{R}_{>0}^K \times \Pi$, and $\Pi = \{\boldsymbol{\pi} : \pi_k > 0 \text{ for } k = 1, \dots, K-1, \text{ and } \sum_{k=1}^{K-1} \pi_k < 1\}$ with $\pi_K = 1 - \sum_{k=1}^{K-1} \pi_k$. [27] considered a reparametrized form of this density for scale-invariant estimation and efficient computation. Specifically, they reparametrized the regression parameters and the variances as follows:

$$\boldsymbol{\phi}_k = \boldsymbol{\beta}_k / \sigma_k, \quad \rho_k = \sigma_k^{-1}, \quad k = 1, \dots, K.$$

We can rewrite equation (1) with the new parameters as

$$f_{\boldsymbol{\theta}}(z|\mathbf{x}) = \sum_{k=1}^K \pi_k \frac{\rho_k}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (\rho_k z - \mathbf{x}^T \boldsymbol{\phi}_k)^2\right),$$

where $\boldsymbol{\theta} = (\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_K, \rho_1, \dots, \rho_K, \pi_1, \dots, \pi_{K-1}) \in \mathbb{R}^{k \cdot (p+1)} \times \mathbb{R}_{>0}^k \times \Pi$ and Π is the same set as above with $\pi_K = 1 - \sum_{k=1}^{K-1} \pi_k$.

FMRLasso penalizes the negative log-likelihood with an l_1 norm penalty [31]:

$$\begin{aligned} -\frac{1}{n} l_{\text{pen}, \lambda}(\boldsymbol{\theta}) &= -\frac{1}{n} l(\boldsymbol{\theta}) + \lambda \sum_{k=1}^K \pi_k \|\boldsymbol{\phi}_k\|_1 \\ &= -\frac{1}{n} \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k \frac{\rho_k}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} (\rho_k z_i - \mathbf{x}_i^T \boldsymbol{\phi}_k)^2 \right) \right) + \lambda \sum_{k=1}^K \pi_k \|\boldsymbol{\phi}_k\|_1, \end{aligned}$$

where λ is a tuning parameter and $\|\boldsymbol{\phi}_k\|_1 = \sum_{j=1}^p |\phi_{kj}|$. The penalty term weighs the contributions from each cluster by the corresponding prior probabilities. For a given number of clusters, K , and a given tuning parameter, λ , we define the FMRLasso estimator as $\tilde{\boldsymbol{\theta}}_{\lambda, K} = (\tilde{\boldsymbol{\phi}}_1, \dots, \tilde{\boldsymbol{\phi}}_K, \tilde{\rho}_1, \dots, \tilde{\rho}_K, \tilde{\pi}_1, \dots, \tilde{\pi}_{K-1})$. [27] also suggested an unweighted penalty term $\sum_k \|\boldsymbol{\phi}_k\|_1$ along with another weighted penalty term of the form $\sum_k \pi_k^{0.5} \|\boldsymbol{\phi}_k\|_1$. The unweighted penalty term tends to perform poorly in unbalanced cases, where the numbers of observations across clusters differ significantly [27]. Therefore, ARoG utilizes the weighted penalty term with π_k , which performs well in both the balanced and unbalanced cases.

A key issue in the mixture linear regression model is the selection of the optimal number of clusters and the optimal tuning parameter. We use a modified Bayesian Information Criteria (BIC), defined by [27] as

$$BIC = -2l(\hat{\boldsymbol{\theta}}_{\lambda, K}) + \log(n) d_e,$$

where $d_e = K + (K-1) + \sum_{r=1, \dots, K; j=1, \dots, p} 1_{\{\hat{\pi}_{r,j} \neq 0\}}$ is the effective number of parameters. We perform a grid search over a set of (λ, K) and find the optimal combinations, $(\hat{\lambda}, \hat{K})$, achieving the smallest modified BIC. [27] showed that single cluster model with $\lambda_{\max} = \max_{\{j=1, \dots, p\}} \frac{\langle y, \mathbf{x}_j \rangle}{\sqrt{n} \|\mathbf{V}\|}$ selects no variables and suggested λ_{\max} as the upper bound for the value of the tuning parameter. ARoG increases this upper bound three to six times since multiple clusters may require a larger tuning parameter to avoid selecting false positive annotations. With a slight abuse of notation, we denote the ARoG(I) parameter estimates as $\tilde{\boldsymbol{\theta}} \equiv \hat{\boldsymbol{\theta}}_{\hat{\lambda}, \hat{K}}$. The annotation coefficient for the k th cluster is obtained with $\tilde{\boldsymbol{\beta}}_k = \tilde{\boldsymbol{\phi}}_k / \tilde{\rho}_k$ and the variance for the k th cluster is estimated with $\tilde{\sigma}_k = 1 / \tilde{\rho}_k$. We obtain the posterior probability that SNP i belongs to the k th cluster as

$$\tilde{\gamma}_{ik} \equiv P(k | \mathbf{x}_i, z_i, \tilde{\boldsymbol{\theta}}) = \frac{\tilde{\pi}_k \frac{\tilde{\rho}_k}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} (\tilde{\rho}_k z_i - \mathbf{x}_i^T \tilde{\boldsymbol{\phi}}_k)^2 \right)}{\sum_{k=1}^K \tilde{\pi}_k \frac{\tilde{\rho}_k}{\sqrt{2\pi}} \exp \left(-\frac{1}{2} (\tilde{\rho}_k z_i - \mathbf{x}_i^T \tilde{\boldsymbol{\phi}}_k)^2 \right)}.$$

ARoG assigns the SNPs to the clusters for which they have the largest posterior probabilities for. This generates K SNP sets with members $\mathcal{C}_k = \{i : k =$

$\operatorname{argmax}_m \tilde{\gamma}_{im}$ and the set sizes are $|\mathcal{C}_k| = n_k$, $k = 1, \dots, K$. We denote the annotation set of each cluster as $S_k = \{j \in \{1, \dots, p\} : \tilde{\beta}_{kj} \neq 0\}$ with numbers of annotations $|S_k| = p_k \leq p$ and define the annotations selected by the model as the union of the selected annotations across all the clusters, $S = \cup_k S_k \subseteq \{1, \dots, p\}$.

3.2 Permutation Testing for ARoG

The ARoG framework follows up the penalized likelihood-based selection with a permutation testing to evaluate the significance of the selected annotations, which typically have small effect sizes based on our data analysis results in Section 4. We specifically test whether the maximum absolute value of the coefficient of each functional annotation across all the clusters can arise by chance association. We randomly permute the z -scores of the SNPs a large number of times (at least 1000 times), and refit ARoG to each permuted dataset. At each fit, we record the maximum absolute value of the estimated coefficients for each annotation across all the clusters. This collection generates functional annotation specific null distributions. Then the p -value for the j -th annotation is defined as the proportion of datasets with the maximum absolute values of the estimated coefficients larger than $\max_k |\tilde{\beta}_{kj}|$. We utilize the Benjamini-Hochberg false discovery rate (FDR) procedure [7] at level 0.1 to account for the multiplicity of the annotations.

3.3 Two-step ARoG (ARoG(II))

Basic ARoG filters false positive annotations with a global penalization across all clusters; however, it is still prone to selecting a nonignorable number of false positive variables as both the simulations of [27] and our computational experiments in Section 5 illustrate. To reduce this effect and thereby increase specificity, we propose and study two-step ARoG. The two-step ARoG implements a cluster level penalization and a refit estimation after the initial global penalization by ARoG(I). The additional penalization further screens out false positives. The refit step aims to improve parameter estimation by alleviating the shrinkage effect towards small coefficients and large standard deviations led by FMRLasso. This step is similar to relaxed Lasso of [19] which employs another level of Lasso in the context of standard multivariate linear regression model. Both two-step ARoG and relaxed Lasso aim to filter out false positive variables resulting from the initial penalization and thereby lead to better or comparable prediction with more accurate variable selection. Refitting has been widely used as a simple but practical tool to overcome the biased estimation of Lasso [9].

In two-step ARoG, after the initial FMRLasso, we consider a standard multivariate linear regression model for each cluster with its corresponding

selected annotation set. First, the two-step ARoG adds an l_1 penalty term to the residual sum of squares within each cluster as follows:

$$\hat{\boldsymbol{\beta}}_{\lambda_k, k}^{S_k} = \operatorname{argmin}_{\boldsymbol{\beta}_k^{S_k}} \frac{1}{n_k} \sum_{i \in \mathcal{C}_k} \left(z_i - \mathbf{x}_i^{S_k T} \boldsymbol{\beta}_k^{S_k} \right)^2 + \lambda_k \|\boldsymbol{\beta}_k^{S_k}\|_1, \quad (2)$$

where $\boldsymbol{\beta}_{\lambda_k, k}^{S_k} \in \mathbb{R}^{p_k+1}$, $\mathbf{x}_i^{S_k} = [1; \mathbf{x}_{ij}; j \in S_k] \in \mathbb{R}^{p_k+1}$. The tuning parameter selection for each cluster is through BIC and the coefficients estimated with the optimal tuning parameter are simply denoted as $\hat{\boldsymbol{\beta}}_k^{S_k} = \hat{\boldsymbol{\beta}}_{\lambda_k, k}^{S_k}$. Next, based on the clusterwise Lasso, we obtain a smaller annotation set, $M_k = \{j \in S_k : \hat{\beta}_{kj}^{S_k} \neq 0\}$, with size $|M_k| = d_k < p_k$, and refit a least squares regression with this annotation set:

$$\hat{\boldsymbol{\beta}}_k^{M_k} = \operatorname{argmin}_{\boldsymbol{\beta}_k^{M_k}} \frac{1}{n_k} \sum_{i \in \mathcal{C}_k} \left(z_i - \mathbf{x}_i^{M_k T} \boldsymbol{\beta}_k^{M_k} \right)^2, \quad (3)$$

where $\mathbf{x}_i^{M_k} = [1; \mathbf{x}_{ij}; j \in M_k] \in \mathbb{R}^{d_k+1}$. We then have the two-step ARoG annotation score coefficients for the k -th cluster as

$$\hat{\beta}_{kj} = \begin{cases} \hat{\beta}_{kj}^{M_k} & j \in M_k, \\ 0 & j \in \{1, \dots, p\} / M_k. \end{cases} \quad (4)$$

Similar to ARoG(I), we define ARoG(II) annotations as the union of selected annotations over the clusters, $M = \cup_k M_k \subset S$. There is a trade-off between basic ARoG and two-step ARoG since cluster-level Lasso tends to gain specificity and lose sensitivity with more aggressive annotation screening. The level of the trade-off varies on a case by case basis. We further discuss this issue with computational experiments in Section 5.1. The permutation testing described in Section 3.2 is also part of two-step ARoG. We denote basic ARoG as ARoG(I) and two-step ARoG as ARoG(II) in the remainder of this paper.

3.4 Numerical Implementation

We implement ARoG with publicly available R packages `fmrlasso` and `glmnet`. The `fmrlasso` package fits FMRLasso with a block coordinate descent generalized expectation-maximization algorithm (BCD-GEM) proposed by [27]. It alternates between an expectation step (E-step) and a generalized maximization step (generalized M-step), which updates the prior probabilities, $\boldsymbol{\phi}$ at once, then updates the reparametrized regression coefficients, $\boldsymbol{\pi}$ and standard deviations, $\boldsymbol{\rho}$. Each cluster is decoupled into K distinct optimization problems, and the BCD-GEM applies coordinate updates to each optimization problem separately. For the initialization of the E-step, we first summarize the annotation score vector for each SNP by its l_2 norm and then use hierarchical clustering with the distance between two SNPs as the l_1 norm between their z-scores and the summarized annotation scores. This distance criterion ensures that SNPs with similar z-scores and similar variability in the functional

annotations are more similar to each other. Using this hierarchical clustering, we then assign each SNP to a single cluster for any given K . We set the posterior probabilities of the SNPs for their assigned clusters as $9/(K+8)$, and the posterior probabilities for the other clusters as $1/(K+8)$. For the M-step, we initialize $\phi_{kj}^{(0)} = 0$, $\rho_k^{(0)} = 2$, $\pi_k^{(0)} = 1/K$, $k = 1, \dots, K$, and $j = 1 \dots, p$. Finally, we implement the cluster level Lasso of ARoG(II) with a coordinate descent algorithm using `glmnet`.

4 ARoG Analysis of PGC Data

4.1 PGC Autism GWAS

We next fit ARoG(I) and ARoG(II) to the autism dataset described in Section 2 and varied the number of clusters as $K = 1, \dots, 10$. ARoG resulted in $\hat{K} = 3$ as the optimal number of clusters. Table 1 presents parameter estimates from both ARoGs. Refitting for ARoG(II) is performed after each SNP is assigned to the cluster for which it has the highest posterior probability for based on ARoG(I) and leads to reestimation of both the regression parameters and the cluster-specific variances. Both ARoGs have the first and the second clusters as intercept-only models and select FOXL1 and Nkx2-5 for the third cluster. We kept the ARoG(I) intercept estimate for the first cluster since no other SNPs were assigned to this cluster. Estimated coefficients for both TFs of the cluster 3 indicate that the SNP-driven increases in binding affinities for FOXL1 and Nkx2-5 associate with the increased autism risk in cluster 3. We further support the significance of these associations with a permutation test described in Section 3.2 (Supplementary Figure 2(a)). The third cluster has a total of thirteen SNPs, nine of which constitute the most genome-wide significant SNPs depicted in the Manhattan plot Figure 1(a). As an alternative multivariate approach to ARoG, we also used ordinary least squares (OLS) and Lasso regression to select the most relevant annotations from the set of 205. OLS did not select any annotations with a BH adjustment on the OLS p-values at level 0.1 and had unadjusted p-values of 0.332 and 0.014 for FOXL1 and Nkx2-5, respectively. We utilized 5-fold cross-validation to tune the l_1 penalty parameter for Lasso and obtained 11 Lasso-selected annotations including Nkx2-5 and FOXL1. However, neither of these survived the permutation testing implemented in a way similar to that of ARoG's (Section 3.2) (Supplementary Figure 2(a)). This analysis suggests that ARoG is indeed exploiting associations detectable only when appropriate subgroups of SNPs are considered.

Next, we investigated the effects of the selected annotations, FOXL1 and Nkx2-5, on autism z-scores of cluster 3. Figure 2(a) highlights significant positive associations of the z-scores with FOXL1 and Nkx2-5 annotation scores, respectively, for SNPs of cluster 3. Considering all the SNPs lead to weak positive associations without statistical support from marginal regressions. Figure 2(b) displays the heatmap of the z-scores and FOXL1 and Nkx2-5

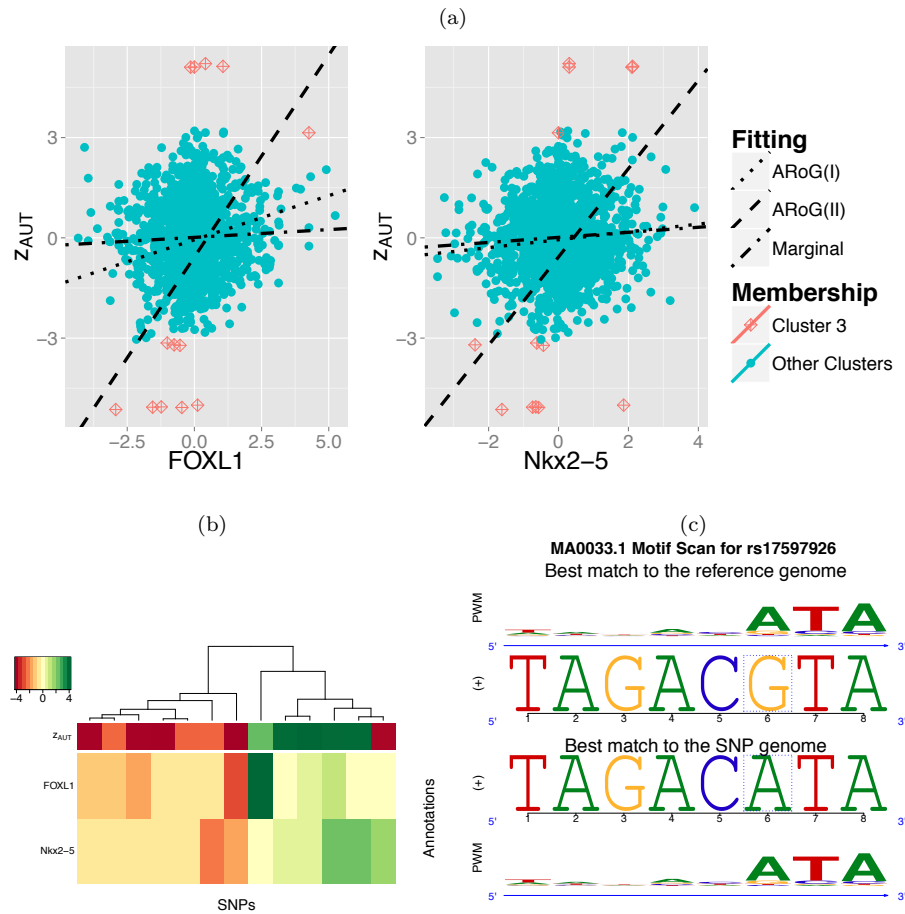


Fig. 2: (a) Autism z-scores vs. annotation scores for TFs FOXL1 and Nkx2-5 selected for cluster 3 along with the marginal linear regression line fit and ARoG estimates (intercept and slopes for FOXL1 and Nkx2-5). (b) Hierarchical clustering of ARoG selected annotations and SNPs in Cluster 3 with their AUT z-scores. (c) Composite sequence logo of SNP rs17597926 with the FOXL1 position weight matrix. The middle two rows represent best matching genomic subsequences to the FOXL1 PWM with the reference (G) and SNP rs17597926 alleles (A), respectively. The dashed box marks the SNP location. Top and bottom rows display FOXL1 PWM sequence logos aligned to the best reference and SNP allele matches.

Table 1: ARoG parameter estimates with PGC AUT data.

	Cluster	1	2	3
	Estimated prior prob. ($\tilde{\pi}_k$)	0.0929	0.8724	0.0347
ARoG(I)	SSD ($\tilde{\sigma}_k$)	0.1739	1.0153	2.3258
	(Intercept)	-0.0686	0.0253	-0.0580
	FOXL1	0	0	0.2639
	Nkx2-5	0	0	0.1171
ARoG(II)	SSD ($\tilde{\sigma}_k$)	0.1739	0.9926	3.7191
	(Intercept)	-0.0686	0.0183	-0.5724
	FOXL1	0	0	1.2075
	Nkx2-5	0	0	1.3215

annotation scores of cluster 3 SNPs organized by hierarchical clustering and supports that the variation in z-scores is well explained by these two annotation scores. We note that atSNP further quantifies the annotation scores by testing whether the observed change in TF binding affinity due to SNP is significant and reports corresponding p-values. We used these p-values to along with raw GWAS association p-values to further refine the SNPs in cluster 3 and define *ARoG SNPs* as SNPs leading to significant TF binding affinity changes and having marginal association with the disorder. Specifically, we considered the subset of cluster 3 SNPs with raw GWAS p-value of at most 0.005 and atSNP p-value of at most 0.01 for FOXL1 or Nkx2-5, resulting in a single ARoG SNP, rs17597926. The unadjusted p-value of this SNP from autism GWAS is 0.0017 and the resulting atSNP FOXL1 p-value is 0.0008. The composite logo plot in Figure 2(c) indicates that rs17597926 is creating a potential FOXL1 binding site. rs17597926 is located within the 5th intron of the TCF4 gene, known to interact with helix-loop-helix proteins and regulate neurodevelopment [12]. Furthermore, this SNP has been identified as a *cis*-eQTL for TCF4 in a recent brain expression GWAS [35]. This provides additional support for potential regulatory role of rs17597926 as a mediator of TCF4 gene in psychiatric disorders.

4.2 Comparison with GPA

In addition to the ARoG analysis, we also applied the GPA approach of [10] to the autism data. We would like to emphasize that GPA and ARoG approaches utilize functional annotations from different angles: GPA goes after global signals utilizing all the SNPs genotyped whereas ARoG aims to identify local signals by focusing on a smaller set of signals with potential significance. GPA is based on a joint generative model of association p-values of the SNPs and annotation data and aims to identify annotations that the disease-associated SNPs are enriched for. It aims to simultaneously identify null (SNPs not associated with the phenotype) and non-null (SNPs associated with the phenotype) and quantify the enrichment of a given annotation within these SNP sets. It

specifically tests whether equal proportions of non-null and null SNPs carry the annotation. Although it can handle multiple annotations simultaneously, our results from two application schemes of “one annotation at a time” versus “all annotations simultaneously” showed extreme differences which could potentially be attributable to the violation of the GPA independent assumption of the annotations conditional on the SNPs null versus non-null status. As a result, we focus on applying GPA one annotation at a time. GPA works with binary annotations; therefore, we first created a binary annotation score matrix by running atSNP [36] on all the 1,210,570 SNPs and thresholding atSNP p-values at 0.05. These SNPs are a subset of the 1,219,561 genotyped SNPs in the PGC study that were also in the dbSNP database. We applied GPA to each of the 205 annotations separately and estimated the proportions of null and non-null SNPs associated with each annotation. Results from GPA hypothesis testing for annotation enrichment did not identify any annotation as significantly enriched for autism-associated SNPs (Supplementary Figure 4). This is consistent with our marginal analysis in Section 2 where none of the annotations exhibited significant marginal associations with the autism z-scores. The estimated fold enrichments of FOXL1 and Nkx2-5 in the GPA analysis were 1.003 (s.e. 0.113) and 0.912 (s.e. 0.146), respectively. Both of these levels were too small to be detected with this analysis that considered only two global classes of SNPs (null and non-null).

4.3 PGC Schizophrenia GWAS

ARoG application to the SCZ2 dataset described in Section 2 with numbers of clusters $K = 1, \dots, 10$ led to best BIC values with $K = 6$ and $K = 7$, and with only a 0.01% difference between the two. We carried out the rest of the analysis with $K = 6$. ARoG did not select any annotations for clusters 1-4.

ARoG(I) selects FOXL1, Klf4, Prrx2, and NKX3-1 annotations for cluster 5, and Arnt::Ahr, E2F1, FOXL1, Klf4, Foxq1, Prrx2, ARID3A, and E2F4 annotations for cluster 6. Among the selected annotations, the pair of E2F1 and E2F4 and the pair of Prrx2 and ARID3A share similar sequence logos, respectively, thus, they have relatively high correlations of 0.8096 and 0.6234 within each other in the annotation score matrix. ARoG(II) retains FOXL1 and Klf4 for cluster 5 and Arnt::Ahr and FOXL1 for Cluster 6 (Table 2). ARoG(II) tends to have regression coefficients increased and standard deviation estimates decreased compared to ARoG(I). Overall, both cluster 5 and 6 are populated with the most genome-wide significant SNPs depicted in the Manhattan plot of Figure 1(c). Permutation testing results for ARoG(I) and ARoG(II) indicate significance of the selected annotations with BH adjustment at level 0.1 (Supplementary Figure 2(b)). OLS analysis of this dataset selected 15 annotations with unadjusted permutation p-values smaller than 0.05; however, none of these survived the multiple testing correction with the BH adjustment at level 0.1. In contrast, Lasso with cross validation tuning se-

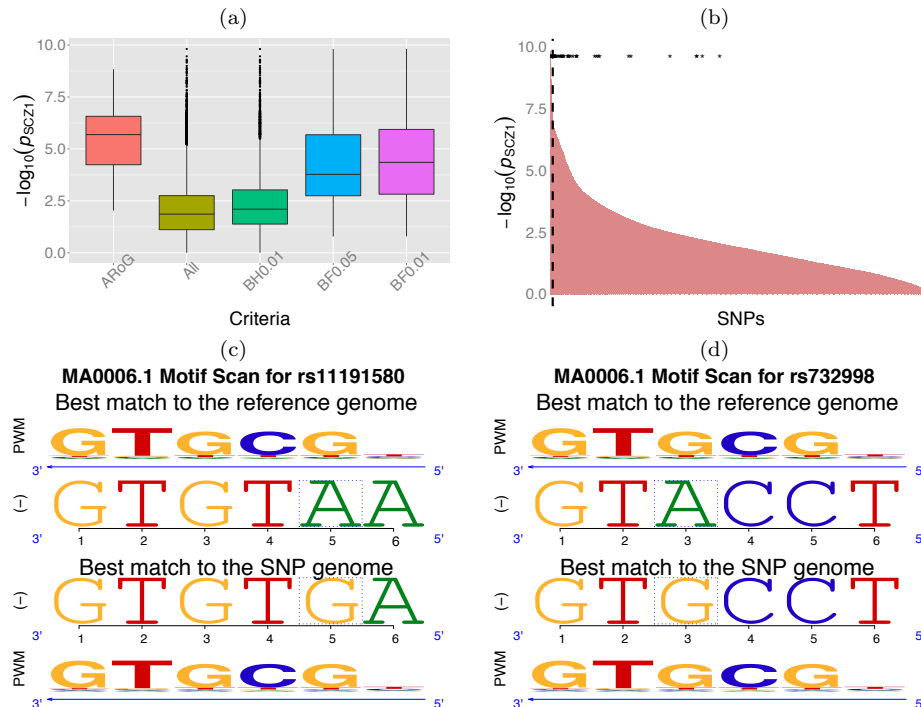


Fig. 3: (a) SCZ1 p-values for multiple SNP sets generated based on the SCZ2 data: The SNP sets include ARoG: ARoG SNPs; All: All the 11,386 SNPs; BH0.01: SNPs with SCZ2 BH-adjusted p-values less than or equal to 0.01; BF0.05/BF0.01: SNPs with SCZ2 Bonferroni adjusted p-values less than or equal to 0.05/0.01. (b) SNPs ranked based on their SCZ1 significance levels. ARoG SNPs are marked with asterisks. The vertical dashed line depicts Bonferroni cut-off of SCZ1 analysis under significance level of 0.05. (c) Composite sequence logo of rs11191580 with the Arnt::Ahr PWM: the SNP enhances the binding of Arnt::Ahr. (d) Composite sequence logo of SNP rs732998 with the Arnt::Ahr PWM.

lected 17 annotations, of which only two (Foxq and Zfx annotations) survived the same multiple testing adjustment (Supplementary Figure 2(b)).

The scatter plots of the SCZ2 z-scores against the selected annotations exhibit associations in clusters 5 and 6 with a similar global trend across the whole SNP set (Supplementary Figure 3). Next, we defined a set of ARoG SNPs as the SNPs with Bonferroni corrected p-values less than 0.05 and at-SNP p-values less than 0.01. As a result, ARoG SNPs included 14 SNPs from cluster 5 and 30 SNPs from cluster 6. Supplementary Table 1 presents genomic locations, GWAS p-values, and RegulomeDB scores [8] of the ARoG SNPs. RegulomeDB scores range from 1 to 7. SNPs with score 1 are likely to

Table 2: ARoG parameter estimates with PGC SCZ2 data

	Cluster	1	2	3	4	5	6
	Estimated prior prob. ($\tilde{\pi}_k$)	0.2335	0.2173	0.2202	0.2084	0.0615	0.0591
	Membership dist.	3533	3307	2066	1859	309	312
ARoG(I)	SSD ($\tilde{\sigma}_k$)	0.4122	0.3825	0.8805	0.9001	1.7423	1.8453
	(Intercept)	4.1206	-4.1117	-4.5537	4.6279	-5.7453	5.8446
	Arnt::Ahr	0	0	0	0	0	0.0859
	FOXL1	0	0	0	0	-0.2335	-0.1394
	Klf4	0	0	0	0	0.1027	0.0653
ARoG(II)	SSD ($\tilde{\sigma}_k$)	0.3336	0.2941	0.9605	0.9427	1.5425	1.6721
	(Intercept)	4.0791	-4.0776	-4.8039	4.9482	-7.0816	7.1786
	Arnt::Ahr	0	0	0	0	0	0.1814
	FOXL1	0	0	0	0	-0.2879	-0.4295
	Klf4	0	0	0	0	0.2654	0

affect TF binding and linked to expression of a gene target. SNPs with score 2 are likely to affect binding, those with score 3 are less likely to affect binding, and those with score 4 to 6 have minimal binding evidence. We refer Table 2 of [8] for further details. RegulomeDB scores of 15 ARoG SNPs indicate a high likelihood of impact on TF binding and gene expression, further providing evidence for potential importance of these SNPs to schizophrenia. We next compared the SCZ1 association measures (p-values) of ARoG SNPs to other SNP sets one could have identified from the initial set of 11,386 SNPs to evaluate which SNP sets are more supported by the SCZ1 study (Figure 3(a)). The other SNP sets one could define without using additional functional annotation are BH0.01 (SNPs defined by BH correction at level 0.01 on the SCZ2 p-values), BF0.05 (SNPs defined by Bonferroni correction at level 0.05 on the SCZ2 p-values), and BF0.01 (SNPs defined by Bonferroni correction at level 0.01 on the SCZ2 p-values). ARoG SNPs are on average more significant (more reproducible in the SCZ1) than the SNP sets. Comparison of ARoG SNPs with a randomly selected SNP set of the same size from BF0.05 also indicated that ARoG SNPs are on average more significant, illustrating that the use of the functional annotation information is biasing the selection towards SNPs with reproducible associations. Figure 3(b) displays ranking of SCZ1 p-values of all the 11,386 SNPs and illustrates that most ARoG SNPs are located near the most significant SNPs with respect to SCZ1. Four of these SNPs reach genome-wide significance with Bonferroni adjustment at level 0.05 in the SCZ1 study.

Next, we assessed whether any of the ARoG SNPs were among the schizophrenia associated SNPs from dbGaP [1]. dbGaP harbors 249 SNPs associated with schizophrenia and 42 of these are among the 11,386 SNPs we utilized. Two of the ARoG SNPs (rs11191580 and rs10224497, located at chr10:104,906,211 and chr7:2,149,967) are among the dbGaP SNPs. SNP rs11191580 leads to enhancement of Arnt::Ahr binding while rs10224497 seems to disrupt Arnt::Ahr binding. Since, overall, we observed that enhanced Arnt::Ahr binding associated with increased schizophrenia risk in clusters 5 and 6 of the ARoG results,

we further investigated rs11191580. rs11191580 is located within the 3rd intron of Nt5C2 and has rs732998, located within the 4th intron of Nt5C2, as a perfect LD partner. Their composite logo plots support that these SNPs might indeed enhance the binding of Arnt::Ahr (Figure 3(c), (d)). Furthermore, the association of rs11191580 is also validated in SCZ1 with p-value of 2.23×10^{-8} . Although rs732998 does not quite make the genome-wide significance cut-off, it also exhibits significant association in SCZ1 with p-value of 9.50×10^{-8} . In summary, these two SNPs that are in perfect LD and both lead to sequence changes that are likely to improve the binding of the Arnt/Ahr complex. This complex regulates genes in response to the carcinogenic environmental contaminant 2,3,7,8-tetrachlorodibenzo-p-dioxin (TCDD). [5] showed that transcription factor Ahr is frequently detected in brain and over-expression of AhR causes neural differentiation of Neuro2a cells. Furthermore, recent studies support that dioxins and related chemicals influence neural development, and the AhR-signaling pathway might mediate the impact of dioxins on the nervous system [33].

Finally, as we have done for the autism dataset in Section 4.2, we also applied GPA to the schizophrenia dataset using all 1,175,307 SNPs that were in the dbSNP database out of the 1,179,262 genotyped SNPs. This analysis identified non-null SNPs as significantly depleted for Zfp423 annotation (Supplementary Figure 5) under Bonferroni adjusted significance level of 0.1. However, the probability that a SNP is non-null is estimated as 0.3157, therefore the estimated non-null set is likely to include a large number of SNPs unassociated with schizophrenia. This implies that the Zfp423 depletion is likely to be a false positive finding.

5 Simulation Studies

We evaluated ARoG(I) and ARoG(II) with synthetic datasets and PGC data-driven simulated datasets. As alternative multivariate methods, we included Lasso regression, and OLS with BH correction on the regression coefficient p-values at level of 0.05. Both ARoGs enable clustering of SNPs and heterogeneous annotation coefficients across the clusters while both Lasso regression and OLS assume the homogeneity of the annotation effects. Our simulation scheme focuses on detection of relevant annotations among the many weak annotation scores. We may implement GPA also by taking sparse binary annotations after binarizing the scores. However, from our experiments, GPA fittings were unstable or sometimes failed as the severe sparsity of annotations lets their estimators located near or in the boundary of the parameter space. We generated 100 simulated datasets under each scenario, where each simulated dataset consisted of training data, validation data, and test data. The validation dataset was used for selection of the optimal tuning parameter, and its sample size was increased 100 times compared to that of the training dataset. The test error was calculated as the negative log-likelihood on the test dataset with the same sample size as the training dataset.

Table 3: Simulation settings for Section 5.1.

Model	n	Cluster	β	σ	π
Sparse (I)	100	1	$(0, \underbrace{3, \dots, 3}_{5 \text{ repetitions}}, \underbrace{0, \dots, 0}_{100 \text{ repetitions}})$	0.5	0.5
		2	$(0, \underbrace{-1, \dots, -1}_{5 \text{ repetitions}}, \underbrace{0, \dots, 0}_{100 \text{ repetitions}})$	0.5	0.5
Intermediate (I)	1000	1	$(0, \underbrace{1.5, \dots, 1.5}_{55 \text{ repetitions}}, \underbrace{0, \dots, 0}_{50 \text{ repetitions}})$	1	0.5
		2	$(0, \underbrace{-0.5, \dots, -0.5}_{55 \text{ repetitions}}, \underbrace{0, \dots, 0}_{50 \text{ repetitions}})$	1	0.5
Dense (I)	1000	1	$(0, \underbrace{-1, \dots, -1}_{30 \text{ repetitions}}, \underbrace{1, \dots, 1}_{30 \text{ repetitions}}, \underbrace{0, \dots, 0}_{45 \text{ repetitions}})$	0.5	0.5
		2	$(0, \underbrace{-1, \dots, -1}_{30 \text{ repetitions}}, \underbrace{0, \dots, 0}_{30 \text{ repetitions}}, \underbrace{1, \dots, 1}_{30 \text{ repetitions}}, \underbrace{0, \dots, 0}_{15 \text{ repetitions}})$	0.5	0.5

We report for each method the test error, numbers of true (TPs) and false positives (FPs), adjusted rand index (ARI) [23], receiver operating characteristic (ROC) curves, and precision recall curves. TPs and FPs for ARoG are defined by pooling selected annotations across the identified clusters. We use adjusted rand index (ARI) to measure the similarity between the true SNP clusters and estimated SNP clusters. ROC curves and precision-recall curves present the performance of annotation selection in thresholds (cut-off for p-values for OLS, tuning parameters for Lasso and both ARoGs) agnostic manner. In these curves, we plot the average values of true positive rate (TPR), false positive rate (FPR), and precision across 100 simulation replications.

5.1 Synthetic Data

We generated data from several Gaussian finite mixture regression models varying the sparsity level of annotation signals as sparse, intermediate, and dense (Table 3). The columns of the predictor matrix X are generated from an independent standard normal distribution. Supplementary Figures 6 and 7 and Figure 4 present the results of these simulations.

The Sparse (I) setting is a small n , large p setting; hence, our comparisons only include ARoG(I), ARoG(II), and Lasso. Supplementary Figure 6 shows that ARoG(I) has the smallest test error and both ARoGs have a median ARI of about 0.7. ARoG(I) and ARoG(II) have the same ARI by design since they share the same clustering assignment and Lasso has ARI of 0 since it does not perform clustering. Both ARoGs have high true positive rates; however, ARoG(I) has an inflated false positive rate with an average of 10 more false positives compared to ARoG(II). Lasso tends to underselect annotations and on average has four false negatives. ROC and precision-recall curves indicate that ARoGs outperform Lasso significantly. We note that Lasso, ARoG(I),

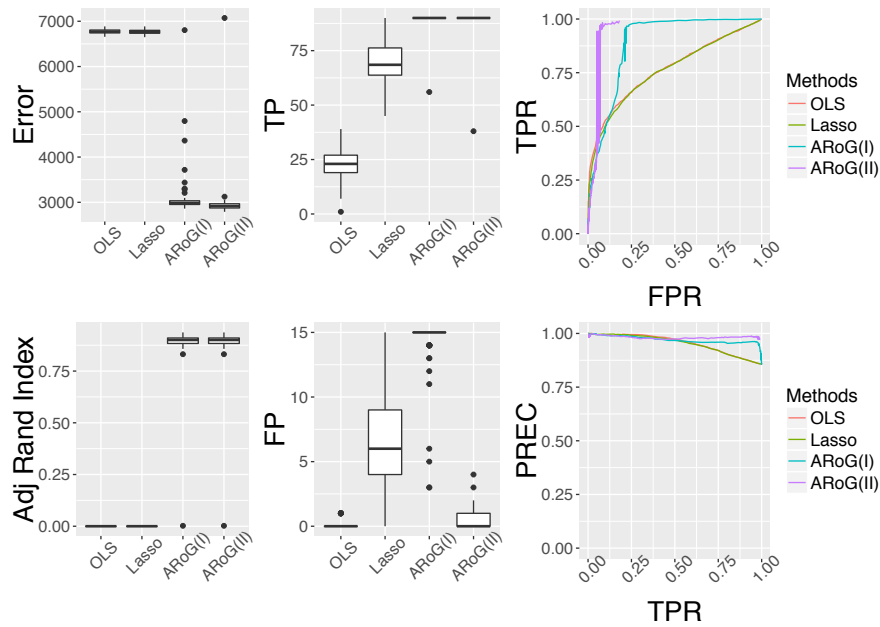


Fig. 4: Simulation results for Dense(I) setting of Table 3.

and ARoG(II) do not select annotations in the sequentially augmenting order as tuning parameters decreasing; thus, the ROC curves are not monotonically increasing. We also investigated this sparse setting by increasing the sample size to 1000 and observed almost perfect performance by all methods with an area under the curve of 1.

Supplementary Figure 7 presents the simulation results from the intermediate annotation setting of Table 3, where almost half of the regression parameters are set as zero. The test error evaluation, ROC curves, and precision recall curves clearly indicate that ARoGs outperform OLS and Lasso. Both ARoG(I) and ARoG(II) have perfect TP rate with the optimal tuning parameters; however their numbers of FPs are substantially different. ARoG(II), on average, selects 2 false positives whereas ARoG(I) selects more than 45 false positives. This emphasizes the significance of the cluster-level Lasso application of ARoG(II) for reducing the numbers of false positives. In the ROC curves, ARoG(II) has competitive performance with ARoG(I), and performs marginally better in the top left corner, where both accurately identify all the true positives. This region is also where the optimal tuning parameters based on the validation set resides in and where ARoG(II) filters out many false positives.

Figure 4 presents the results for the Dense (I) setting. These results also highlight superior performance of ARoGs in terms of prediction error. This setting also highlights the contrast between two implementations of ARoG.

ARoG(I) tends to select all annotations by essentially failing to do variable selection whereas ARoG(II) is able to filter out false positives. OLS tends to underselect annotations and Lasso tends to select more compared to OLS; however, can only partially recover the true positives. ARoG(II) has the best ROC curve and precision-recall curve performances. This setting includes 60 annotations that are not shared between clusters in contrast to the previous settings where the annotations were shared by multiple clusters. The selection performances on these cluster-specific variables are less stable; as a result, the TPRs of both ARoGs heavily fluctuate between 0.4 and 1 in the top left corner of the ROC curves where both have small false positive rates. Overall, we conclude that the differences between ARoG(I) and ARoG(II) become more pronounced as the sparsity level decreases and ARoG(II) outperforms ARoG(I) in dense settings where ARoG(I) tends to have much higher FPR than ARoG(II).

These computational experiments involved completely simulated datasets where the predictor matrix had independent columns and was not designed to be sparse. Supplementary Figure 8 presents results from a sparse setting where the actual annotation predictor matrix from the PGC autism GWAS is used to simulate data (Supplementary Table 2). This predictor matrix is sparse compared to the randomly generated predictor matrix in the above simulations. The overall conclusions from this setting agree well with the Sparse (I) setting.

5.2 PGC Analysis-Driven Data

We next evaluated the performance of ARoG in two simulation settings based on the AUT and SCZ2 data analyses of Sections 4.1 and 4.3. Each setting had the annotation score matrix from each application as the predictor matrix and the ARoG(II) estimates of prior probabilities, standard deviations, and regression slopes as the parameters. These data-driven simulation studies aim to capture the typical signal to noise levels observed in these type of studies. Figure 5 displays the results for the AUT simulation setting. Both ARoGs reduce the prediction error by about 8%. ARoG(I) and ARoG(II) tend to select at least one annotation 86 and 78 times out of 100 repetitions, respectively. Specifically, FOXL1 is selected 71 and 61 times and Nkx2-5 is selected 49 and 42 times by ARoG(I) and ARoG(II), respectively. ARoG(II) on average filters out 2 false positives more compared to ARoG(I). Both OLS and Lasso tend to miss true positives and thus fail to recover the underlying associations. Based on the ROC and precision-recall curves, ARoG(I) has the best tuning parameter-free performance followed by ARoG(II). OLS performs almost the same as random guess with an ROC curve on the 45 degree line.

Supplementary Figure 9 presents the results for the SCZ2 simulation setting. Both ARoGs perform very well in terms of prediction error. Similar to the AUT simulations, both OLS and Lasso fail to select any annotations. ARoG(I) and ARoG(II) tend to recover the true positives to some extent by selecting at least one correct annotation except in one simulated dataset with ARoG(II).

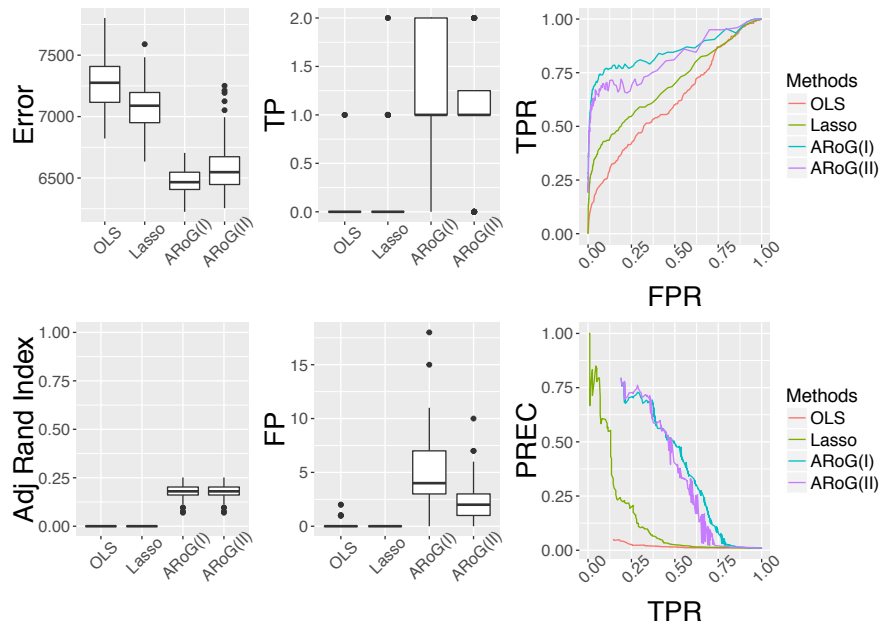


Fig. 5: Simulation results for autism data analysis driven setting.

In this setting, a trade-off between two ARoGs is clear since ARoG(I) seems better at identifying all the true positives, namely, *Arnt::Ahr*, *FOXL1*, and *Klf4*, whereas ARoG(II) more aggressively eliminates false positives. ARoG(I) has a median of 11 FPs, with more than 25 FPs in 13 of the simulated datasets. In contrast, ARoG(II) has a median of 2 FPs, with less than 10 FP annotations in almost all simulated datasets. In terms of the ROC and precision-recall curve comparisons, ARoG(I) exhibits a better tuning parameter-free performance compared to ARoG(II). Both OLS and Lasso perform similar to random guesses.

6 Discussion

We presented an integrative framework, named ARoG, for incorporating functional annotation data into GWAS analysis. The key idea behind ARoG is that even when a set of SNPs disrupt a global mechanism, e.g., pathway, that lead to disease, they might be achieving this by disrupting various sub-mechanisms. Some might be disrupting coding sequences, some transcription factor binding sites, some methylation profile or chromatin accessibility. ARoG capitalizes on this idea and aims to identify clusters of SNPs for which GWAS association measures can be explained by a subset of functional annotations. ARoG utilizes FMRLasso [27] which enables selection among large numbers of func-

tional annotations. We illustrated ARoG with an application to PGC data focusing on autism and schizophrenia disorders and by utilizing the impact of SNPs on transcription factor binding affinities as functional annotations. Our analysis led to identification of SNPs which do not quite make the genome-wide significance cut-offs; however, potentially worthy of following up since their GWAS associations are supplemented by their significant effects on TF binding affinities. This versatile framework provides many directions for useful extensions. First, the fact that it can select among annotations makes it applicable with larger sets of functional annotations including TF ChIP-seq, DNase I-accessibility, Histone ChIP-seq, and DNA methylation. Second, we have currently focused our analysis on one disorder at a time; however, ARoG framework can be easily extended to consider multiple disorders-related GWAS simultaneously.

Acknowledgements

This research was supported by National Institutes of Health grants HG007019, HG003747, and U54AI117924. The authors thank the editor and two referees for their helpful comments.

References

1. dbGaP: The Database of Genotypes and Phenotypes. <http://www.ncbi.nlm.nih.gov/gap>.
2. dbSNP: Short Genetic Variations. <http://www.ncbi.nlm.nih.gov/SNP/>.
3. International human epigenome consortium. <http://ihec-epigenomes.org/research/projects/>
4. Psychiatric Genomics Consortium. <http://www.med.unc.edu/pgc>.
5. Akahoshi, E., Yoshimura, S., Ishihara-Sugano, M.: Over-expression of AhR (aryl hydrocarbon receptor) induces neural differentiation of Neuro2a cells: neurotoxicology study. *Environmental health : a global access science source* **5**(1), 24 (2006)
6. Benjamini, Y., Hochberg, Y.: Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B* **57**, 289–300 (1995)
7. Benjamini, Y., Hochberg, Y.: Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistics Society, Series B* **57**, 289–300 (1995)
8. Boyle, A.P., Hong, E.L., Hariharan, M., Cheng, Y., Schaub, M.A., Kasowski, M., Karczewski, K.J., Park, J., Hitz, B.C., Weng, S., Cherry, J.M., Snyder, M.: Annotation of functional variation in personal genomes using regulomedb. *Genome Research* **22**(9), 1790–7 (2012)
9. Candès, E., Tao, T.: The Dantzig selector: Statistical estimation when p is much larger than n . *Annals of Statistics* **35**(6), 2313–2351 (2007)

10. Chung, D., Yang, C., Li, C., Gelernter, J., Zhao, H.: GPA: A Statistical Approach to Prioritizing GWAS Results by Integrating Pleiotropy and Annotation. *PLoS Genetics* **10**(11), e1004787 (2014)
11. Cross-disorder Working Group of the Psychiatric Genomics Consortium: Identification of risk loci with shared effects on five major psychiatric disorders: a genome-wide analysis. *Lancet* **381**(9875), 1371–9 (2013)
12. Forrest, M.P., Hill, M.J., Quantock, A.J., Martin-Rendon, E., Blake, D.J.: The emerging roles of TCF4 in disease and development. *Trends in Molecular Medicine* **20**(6), 322–331 (2014)
13. Gagliano, S.A., Barnes, M.R., Weale, M.E., Knight, J.: A Bayesian Method to Incorporate Hundreds of Functional Characteristics with Association Evidence to Improve Variant Prioritization. *PLoS ONE* **9**(5), e98122 (2014). DOI 10.1371/journal.pone.0098122
14. Gerstein, M.B., Kundaje, A., Hariharan, M., Landt, S.G., Yan, K.K., Cheng, C., Mu, X.J., Khurana, E., Rozowsky, J., Alexander, R., Min, R., Alves, P., Abyzov, A., Addleman, N., Bhardwaj, N., Boyle, A.P., Cayting, P., Charos, A., Chen, D.Z., Cheng, Y., Clarke, D., Eastman, C., Euskirchen, G., Fietze, S., Fu, Y., Gertz, J., Grubert, F., Harmanci, A., Jain, P., Kasowski, M., Lacroute, P., Leng, J., Lian, J., Monahan, H., O’Geen, H., Ouyang, Z., Partridge, E.C., Patacsil, D., Pauli, F., Raha, D., Ramirez, L., Reddy, T.E., Reed, B., Shi, M., Slifer, T., Wang, J., Wu, L., Yang, X., Yip, K.Y., Zilberman-Schapira, G., Batzoglou, S., Sidow, A., Farnham, P.J., Myers, R.M., Weissman, S.M., Snyder, M.: Architecture of the human regulatory network derived from ENCODE data. *Nature* **489**(7414), 91–100 (2012)
15. Iversen, E.S., Lipton, G., Clyde, M.A., Monteiro, A.N.: Functional annotation signatures of disease susceptibility loci improve SNP association analysis. *BMC Genomics* **15**, 398 (2014)
16. Johnson, A.D., Handsaker, R.E., Pulit, S.L., Nizzari, M.M., O’Donnell, C.J., de Bakker, P.I.W.: SNAP: a web-based tool for identification and annotation of proxy SNPs using HapMap. *Bioinformatics* **24**(24), 2938–9 (2008)
17. Kichaev, G., Yang, W.Y., Lindstrom, S., Hormozdiari, F., Eskin, E., Price, A.L., Kraft, P., Pasaniuc, B.: Integrating functional data to prioritize causal variants in statistical fine-mapping studies. *PLoS Genetics* **10**(10), e1004722 (2014). DOI 10.1371/journal.pgen.1004722
18. Mathelier, A., Zhao, X., Zhang, A.W., Parcy, F., Worsley-Hunt, R., Arenillas, D.J., Buchman, S., Chen, C., Chou, A., Ienasescu, H., Lim, J., Shyr, C., Tan, G., Zhou, M., Lenhard, B., Sandelin, A., Wasserman, W.W.: JASPAR 2014: an extensively expanded and updated open-access database of transcription factor binding profiles. *Nucleic Acids Research* **42**(D1), D142–D147 (2014)
19. Meinshausen, N.: Relaxed Lasso. *Computational Statistics and Data Analysis* **52**(1), 374–393 (2007)
20. Pai, A.A., Pritchard, J.K., Gilad, Y.: The genetic and mechanistic basis for variation in gene regulation. *PLoS Genetics* **11**(1), e1004857 (2015).

- DOI 10.1371/journal.pgen.1004857
21. Pickrell, J.K.: Joint Analysis of Functional Genomic Data and Genome-wide Association Studies of 18 Human Traits. *The American Journal of Human Genetics* **94**(4), 559 – 573 (2014)
 22. Psychiatric GWAS Consortium Bipolar Disorder Working Group: Large-scale genome-wide association analysis of bipolar disorder identifies a new susceptibility locus near ODZ4. *Nature Genetics* **43**(10), 977–983 (2011)
 23. Rand, W.M.: Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association* **66**(336), 846–850 (1971)
 24. Roadmap Epigenomics Consortium: Integrative analysis of 111 reference human epigenomes. *Nature* **518**(7539), 317–330 (2015). URL <http://view.ncbi.nlm.nih.gov/pubmed/25693563>
 25. Schizophrenia Working Group of the Psychiatric Genomics Consortium: Genome-wide association study identifies five new schizophrenia loci. *Nature genetics* **43**(10), 969–76 (2011)
 26. Schizophrenia Working Group of the Psychiatric Genomics Consortium: Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421–427 (2014)
 27. Städler, N., Bühlmann, P., van de Geer, S.: l_1 -penalization for mixture regression models. *TEST* **19**(2), 209–256 (2010)
 28. Stranger, B.E., Stahl, E.A., Raj, T.: Progress and promise of genome-wide association studies for human complex trait genetics. *Genetics* **187**(2), 367–83 (2011)
 29. The GTEx Consortium: The Genotype-Tissue Expression (GTEx) pilot analysis: Multi-tissue gene regulation in humans. *Science* **348**(6235), 648–660 (2015)
 30. Thompson, J.R., Gögele, M., Weichenberger, C.X., Modenese, M., Attia, J., Barrett, J.H., Boehnke, M., De Grandi, A., Domingues, F.S., Hicks, A.A., Marroni, F., Pattaro, C., Ruggeri, F., Borsani, G., Casari, G., Parmigiani, G., Pastore, A., Pfeufer, A., Schwenbacher, C., Taliun, D., Consortium, C., Fox, C.S., Pramstaller, P.P., Minelli, C.: SNP Prioritization Using a Bayesian Probability of Association. *Genetic Epidemiology* **37**(2), 214–221 (2013)
 31. Tibshirani, R.: Regression Selection and Shrinkage via the Lasso. *Journal of the Royal Statistical Society B* **58**, 267–288 (1994). URL <http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.35.7574>
 32. Wasserman, W.W., Long, N., Dickson, S.P., Maia, J.M., Kim, H.S., Zhu, Q., Allen, A.S.: Leveraging prior information to detect causal variants via multi-variant regression. *PLoS Computational Biology* **9**(6), e1003093 (2013)
 33. Xie, H.Q., Xu, H.M., Fu, H.L., Hu, Q., Tian, W.J., Pei, X.H., Zhao, B.: AhR-mediated effects of dioxin on neuronal acetylcholinesterase expression in vitro. *Environmental health perspectives* **121**(5), 613–8 (2013)
 34. Yue, F., Cheng, Y., Breschi, A., Vierstra, J., Wu, W., Ryba, T., Sandstrom, R., Ma, Z., Davis, C., Pope, B.D., Shen, Y., Pervouchine, D.D., Djebali, S., Thurman, R.E., Kaul, R., Rynes, E., Kirilusha, A., Marinov, G.K.,

- Williams, B.A., Trout, D., Amrhein, H., Fisher-Aylor, K., Antoshechkin, I., DeSalvo, G., See, L.H., Fastuca, M., Drenkow, J., Zaleski, C., Dobin, A., Prieto, P., Lagarde, J., Bussotti, G., Tanzer, A., Denas, O., Li, K., Bender, M.A., Zhang, M., Byron, R., Groudine, M.T., McCleary, D., Pham, L., Ye, Z., Kuan, S., Edsall, L., Wu, Y.C., Rasmussen, M.D., Bansal, M.S., Kellis, M., Keller, C.A., Morrissey, C.S., Mishra, T., Jain, D., Dogan, N., Harris, R.S., Cayting, P., Kawli, T., Boyle, A.P., Euskirchen, G., Kundaje, A., Lin, S., Lin, Y., Jansen, C., Malladi, V.S., Cline, M.S., Erickson, D.T., Kirkup, V.M., Learned, K., Sloan, C.A., Rosenbloom, K.R., Lacerda de Sousa, B., Beal, K., Pignatelli, M., Flicek, P., Lian, J., Kahveci, T., Lee, D., Kent, J.W., Ramalho Santos, M., Herrero, J., Notredame, C., Johnson, A., Vong, S., Lee, K., Bates, D., Neri, F., Diegel, M., Canfield, T., Sabo, P.J., Wilken, M.S., Reh, T.A., Giste, E., Shafer, A., Kutayavin, T., Haugen, E., Dunn, D., Reynolds, A.P., Neph, S., Humbert, R., Hansen, R.S., De Bruijn, M., Selleri, L., Rudensky, A., Josefowicz, S., Samstein, R., Eichler, E.E., Orkin, S.H., Levasseur, D., Papayannopoulou, T., Chang, K.H., Skoultchi, A., Gosh, S., Disteche, C., Treuting, P., Wang, Y., Weiss, M.J., Blobel, G.A., Cao, X., Zhong, S., Wang, T., Good, P.J., Lowdon, R.F., Adams, L.B., Zhou, X.Q., Pazin, M.J., Feingold, E.A., Wold, B., Taylor, J., Mortazavi, A., Weissman, S.M., Stamatoyannopoulos, J.A., Snyder, M.P., Guigo, R., Gingeras, T.R., Gilbert, D.M., Hardison, R.C., Beer, M.A., Ren, B., The Mouse ENCODE Consortium: A comparative encyclopedia of DNA elements in the mouse genome. *Nature* **515**(7527), 355–364 (2014). URL <http://dx.doi.org/10.1038/nature13992>
35. Zou, F., Chai, H.S., Younkin, C.S., Allen, M., Crook, J., Pankratz, V.S., Carrasquillo, M.M., Rowley, C.N., Nair, A.A., Middha, S., Maharjan, S., Nguyen, T., Ma, L., Malphrus, K.G., Palusak, R., Lincoln, S., Bisceglia, G., Georgescu, C., Kouri, N., Kolbert, C.P., Jen, J., Haines, J.L., Mayeux, R., Pericak-Vance, M.A., Farrer, L.A., Schellenberg, G.D., Petersen, R.C., Graff-Radford, N.R., Dickson, D.W., Younkin, S.G., Ertekin-Taner, N.: Brain expression genome-wide association study (eGWAS) identifies human disease-associated variants. *PLoS genetics* **8**(6), e1002707 (2012)
36. Zuo, C., Shin, S., Keleş, S.: atSNP: transcription factor binding affinity testing for regulatory SNP detection. *Bioinformatics* **31**(20), 3353–3355 (2015)