

Joint Fine Mapping of GWAS and eQTL Detects Target Gene and Relevant Tissue

Farhad Hormozdiari¹, Ayellet V. Segrè², Martijn van de Bunt^{3,4}, Xiao Li²,
Jong Wha J Joo¹, Michael Bilow¹, Jae Hoon Sul^{5,6}, Sriram Sankararaman¹,
Bogdan Pasaniuc^{7,8}, and Eleazar Eskin^{*1,8}

¹Department of Computer Science, University of California, Los Angeles, CA

²Cancer Program, The Broad Institute of Massachusetts Institute of Technology and
Harvard University, Cambridge, MA

³Oxford Centre for Diabetes, Endocrinology & Metabolism, University of Oxford, Oxford,
United Kingdom

⁴Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, United
Kingdom

⁵Department of Psychiatry and Biobehavioral Sciences, University of California, Los
Angeles, CA

⁶Semel Center for Informatics and Personalized Genomics, University of California, Los
Angeles, CA

⁷Department of Pathology and Laboratory Medicine, University of California, Los Angeles,
CA

⁸Department of Human Genetics, University of California, Los Angeles, CA

*Corresponding author: eeskin@cs.ucla.edu

Abstract

The vast majority of genome-wide association studies (GWAS) risk loci fall in non-coding regions of the genome. One possible hypothesis is that these GWAS risk loci alter the disease risk through their effect on gene expression in different tissues. In order to understand the mechanisms driving a GWAS risk locus, it is helpful to determine which gene is affected in specific tissue types. If the same variant responsible for a GWAS locus also affects gene expression, the relevant gene and tissue may play a role in the disease mechanism. Identifying whether or not the same variant is causal in both GWAS and eQTL studies is challenging due to the uncertainty induced by linkage disequilibrium (LD) and the fact that some loci harbor multiple causal variants. However, current methods that address this problem assume that each locus contains a single causal variant. In this paper, we present a new method, eCAVIAR, that is capable of accounting for LD while computing the quantity we refer to as the colocalization posterior probability (CLPP). The CLPP is the probability that the same variant is responsible for both the GWAS and eQTL signal. eCAVIAR has several key advantages. First, our method can account for more than one causal variant in any loci. Second, it can leverage summary statistics without accessing the individual genotype data. We use both simulated and real datasets to demonstrate the utility of our method. Utilizing data from the Genotype-Tissue Expression (GTEx) project, we demonstrate that computing CLPP can prioritize likely relevant tissues and target genes for a set of Glucose and Insulin-related traits loci. eCAVIAR is available at <http://genetics.cs.ucla.edu/caviar/>

1 Introduction

Genome-wide association studies (GWAS) have successfully detected thousands of genetic variants associated with various traits and diseases [32, 33, 42, 44]. The vast majority of genetic variants detected by GWAS fall in non-coding regions of the genome, and it is unclear how these non-coding variants affect traits and diseases [27]. One potential approach to identify the mechanism of these non-coding variants on diseases is through integrating expression quantitative trait loci (eQTL) studies and GWAS[27]. This approach is based on the concept that a GWAS variant, in some tissue, affects expression at a nearby gene, and that both the gene and tissue may play a role in disease mechanism [12, 18].

Unfortunately, integrating GWAS and eQTL studies is challenging for two reasons. First, the correlation structure of the genome or linkage disequilibrium (LD) [31] produces an inherent ambiguity in interpreting results of genetic studies. Second, some loci harbor more than one causal variant for any given disease. We know that associate statistics of a variant can be affected by other variants in LD [3, 9, 22, 31]. For example, two variants in LD, their associate statistics capture a fraction of the effect of each other. Although GWAS have benefited from LD in the human genome by tagging only a subset of common variants to capture a majority of common variants, a fine mapping process, which attempts to detect true causal variants that are responsible for association signal at the locus, becomes more challenging. Colocalization determines whether a single variant is responsible for both GWAS and eQTL signals in a locus. Thus, colocalization requires correctly identifying the causal variant in both studies.

Recently, researchers proposed a series of methods [13, 15, 18, 26, 30, 43] to integrate GWAS and eQTL studies. One such method is PrediXscan [12], which imputes gene expression followed by association of the imputed expression with trait. However, this method does not provide a basis for determining colocalization of GWAS causal variants and eQTL causal variants. Another class of methods integrates GWAS and eQTL studies to provide insight about the colocalization. For example, regulatory trait concordance (RTC) [26] detects variants that are causal in both studies while accounting for the LD. RTC is based on the assumption that removing the effect of causal variants from eQTL studies reduces or eliminates any significant association signal at that locus. Thus, when the GWAS causal variant is colocalized with the eQTL causal variant, re-computing

the marginal statistics for the eQTL variant conditional on the GWAS causal variant will remove any significant association signal observed in the locus. Sherlock [15], another method, is based on a Bayesian statistical framework that matches the association signal of GWAS with those of eQTL for a specific gene in order to detect if the same variant is causal in both studies. Similar to RTC, Sherlock accounts for the uncertainty of LD. QTLMatch [30] is another proposed method to detect cases where the most significant GWAS and eQTL variants are colocalized due to causal relationship or coincidence. COLOC [13, 43], a method expanded from QTLMatch, is the state of the art method that colocalizes GWAS and eQTL signals. COLOC utilizes approximate Bayes factor to estimate the posterior probabilities for a variant is causal in both GWAS and eQTL studies. Unfortunately, all existing methods assume presence of only one causal variant in any given locus for both GWAS and eQTL studies. As we show below, this assumption reduces the accuracy of results when the locus contains multiple causal variants.

In this paper, we present a novel probabilistic model for integrating GWAS and eQTL data. For each study, we use only the reported summary statistics and simultaneously perform statistical fine mapping to optimize integration. Our approach, eCAVIAR (eQTL and GWAS CAusal Variants Identification in Associated Regions), extends the CAVIAR [16] framework to explicitly estimate the posterior probability of the same variant being causal in both GWAS and eQTL studies while accounting for the uncertainty of LD. We apply eCAVIAR to colocalize variants that pass the genome-wide significance threshold in GWAS. For any given peak variant identified in GWAS, eCAVIAR considers a collection of variants around that peak variant as one single locus. For example, this collection includes the peak variant itself, 50 variants that are upstream of this peak variant, and 50 variants that are downstream of this peak variant. Then, for all the variants in a locus, we consider their marginal statistics obtained from the eQTL study in all tissues and all genes. We only consider genes and tissues in which at least one of the genes is an eGene [5, 40]. eGenes are genes that have at least one significant variant (corrected p-value for multiple hypothesis of at least 10^{-5}) associated with the gene expression of that gene. We assume that the posterior probability of the same variant being causal in both GWAS and eQTL studies are independent. Thus, this posterior probability is equal to the product of posterior probabilities for a given variant is causal in GWAS and eQTL. We refer to the amount of support for a variant responsible for the associated signals in both studies as the quantity of colocalization posterior probability (CLPP).

Our framework allows for multiple variants to be causal in a single locus, a phenomenon that is widespread in eQTL data and referred to as allelic heterogeneity. We utilize data from the Genotype-Tissue Expression (GTEx) project (Release v6, dbGaP Accession phs000424.v6.p1 available at: <http://www.gtexportal.org>) [5] to identify likely relevant tissues. Our approach can accurately quantify the amount of support for a variant responsible for the associated signals in both studies and identify scenarios where there is support for an eQTL mediated mechanism. Moreover, we can identify scenarios where the variants underlying both studies are clearly different. Utilizing simulation datasets, we show that eCAVIAR has high accuracy in detecting target genes and relevant tissues. Furthermore, we observe that the amount of CLPP depends on the complexity of the LD.

We apply our method to colocalize Meta-Analyses of Glucose and Insulin-related traits Consortium (MAGIC) [10, 34, 36, 37] GWAS dataset and GTEx eQTL dataset (Release v6, dbGaP Accession phs000424.v6.p1 available at: <http://www.gtexportal.org>). Our results provide insight into disease mechanisms by identifying specific GWAS loci that share a causal variant with eQTL studies in a tissue. In addition, we identify several loci where GWAS and eQTL causal variants appear to be different. eCAVIAR is available at <http://genetics.cs.ucla.edu/caviar/index.html>

2 Results

2.1 Overview of eCAVIAR

The goal of our method is to identify target genes and the most relevant tissues for a given GWAS risk locus while accounting for the uncertainty of LD. Target genes are genes that their expression levels may affect the phenotype (e.g. disease status) of interest. Our method detects the target gene and the most relevant tissue by utilizing our proposed quantity of colocalization posterior probability (CLPP). eCAVIAR estimates CLPP, which is the probability that the same variant is causal in both eQTL and GWAS studies. eCAVIAR computes CLPP by utilizing the marginal statistics (e.g., z-score) obtained from GWAS and eQTL analyses, as well as the LD structure of genetic variants in each locus. LD can be computed from genotype data or approximated from existing datasets such as 1000 Genomes data [1, 2] or HapMap [4]. We show in the Methods section that the marginal statistics of both GWAS and eQTL follow a multivariate normal distribution

(MVN) given the causal variants and effect sizes for both studies. We use the MVN to estimate the CLPP. We show CLPP is equal to the product of the posterior probability of the variant is causal in GWAS and the posterior probability of the variant is causal in eQTL. Computing the posterior probability of a causal variant is computationally intractable. Therefore, we assume a presence of at most six causal variants in a locus.

The estimated CLPP for a GWAS risk locus and a gene, which is obtained from eQTL studies, can be used to infer specific disease mechanisms. First, we identify genes that have expression levels affected by a GWAS variant. These genes are referred to as target genes. Second, we identify in which tissues the eQTL variant has an effect. To identify target genes, we compute CLPP for all genes in the GWAS risk locus. Genes that have a significantly higher CLPP in comparison are selected as target genes (Figure 1a). Similarly, we compute CLPP for all tissues and identify relevant tissues as those with comparatively high values of CLPP (Figure 1b). Examining this figure, it appears that the GWAS risk locus affects the Gene4 and the relevant tissues are liver and blood. However, pancreas is not a relevant tissue for this GWAS risk locus. Another application of CLPP is to identify loci where the causal variants between GWAS and eQTL studies are different. We can identify these loci if CLPP is low for all variants in the loci, and if there are statistically significant variants in both GWAS and eQTL studies.

To better motivate the behavior of CLPP, we consider the following four scenarios in Figure 2. In the first scenario, the same variant has effects in both GWAS and eQTL studies. Thus, its CLPP is high (Figure 2a). In the second scenario, we consider that the variant is associated with a phenotype in GWAS and not associated with gene expression. In this case, the quantity of CLPP is low (Figure 2b). In the third scenario, we consider that the variant is not associated with a phenotype in GWAS. However, it is associated with expression of a gene. In this case, CLPP is not computed for this variant. Rather, we compute CLPP for GWAS risk loci that are considered significant. In the fourth scenario, we have a variant that appears significant in both GWAS and eQTL. However, other variants in GWAS or eQTL are also significant due to high LD with the causal variant. The complex LD (see Figure 2c) of these variants results in a low CLPP. Here, we remain uncertain about which variants are actual causal variants. Finally, Figure 2d illustrates an example in which there is more than one causal variant. This demonstrates that underestimation of CLPP can result from assuming presence of a single causal variant. In this example, we have a

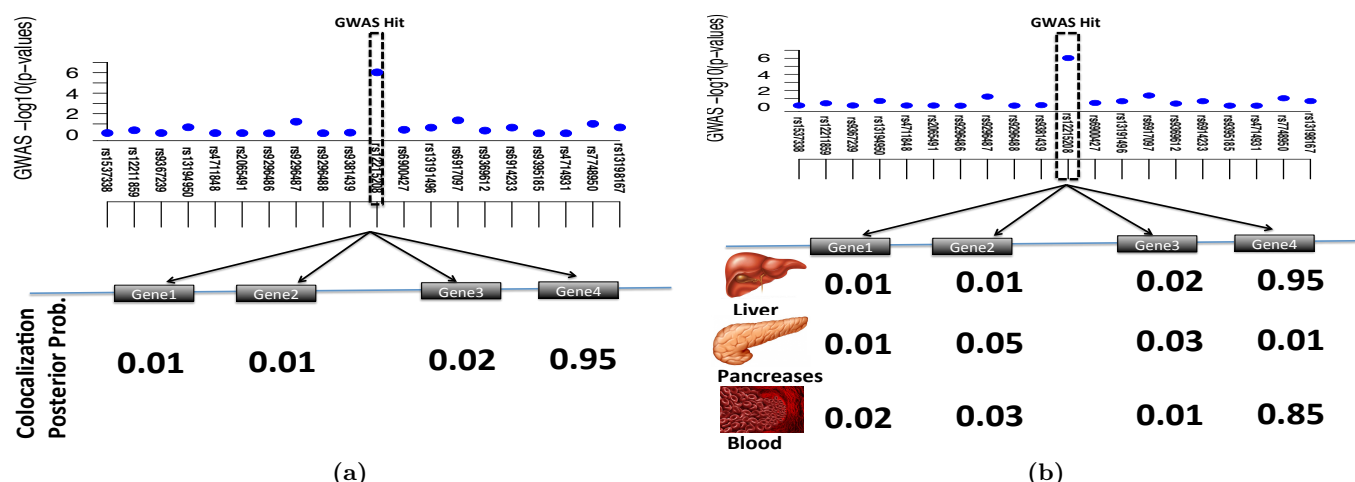


Figure 1. Overview of our method for detecting the target gene and most relevant tissue. We compute the CLPP for all genes and all tissues. Panel (a) illustrates a simple case where we only have one tissue and we want to find the target gene. We consider all the genes for this GWAS risk locus, and we observe Gene4 has the highest CLPP. Thus, in panel (a) the target gene is Gene4. In Panel (b), we have 3 tissues and we utilize the quantity of CLPP. Thus, the target gene is Gene4 again. Moreover, in this example, we consider that liver and blood are relevant tissues for this GWAS risk locus, while pancreas is not relevant to this GWAS risk locus.

locus with 35 variants (SNPs) and we have two causal variants (SNP6 and SNP26) that are not in high LD with each other. If we assume we have only one causal variant there are 35 possible causal variants for this locus and most of the causal variants have very low likelihood. The likelihood of two variants in which the SNP6 or SNP26 are selected as causal have similar likelihood and their likelihood is much higher than other variants. In this example, the estimated posterior probability of SNP6 or SNP26 being causal is equal to 50%. Thus, the estimated CLPP for SNP6 or SNP26 is 25%. However, if we allow more than one causal variant in the locus, all sets of causal variants have very low likelihood values except the set with both SNP6 and SNP26 selected as causal. In this case, the posterior probability of SNP6 or SNP26 being causal is close to 1. Thus, in this case we assume that we have more than one causal variant in this locus, the CLPP of SNP6 and SNP26 are close to 1.

2.2 eCAVIAR Accurately Computes the CLPP

In this section, we use simulated datasets in order to assess the accuracy of our method. We simulated summary statistics utilizing the multivariate normal distribution (MVN) that is utilized in previous studies [14, 16, 17, 21, 46]. More details on simulated data are provided in Section

3.3.2. In one set of simulations, we fix the effect size of a genetic variant so that the statistical power for the causal variant is 50%. In another set, we fix the effect size so that the power is 80%. We consider two cases. In the first case, we only have one causal variant in both studies. In the second case, we have more than one causal variant in these studies. For both cases, we simulated two datasets. In the first dataset, we implanted a shared causal variant. We generated 1000 simulated studies, which we then use to compute the true positive rate (TP). In the second dataset, we implanted a different causal variant in eQTL and GWAS. We filter out cases where the most significant variant is different between the two studies. Similarly to the previous case, we generated 1000 simulated studies.

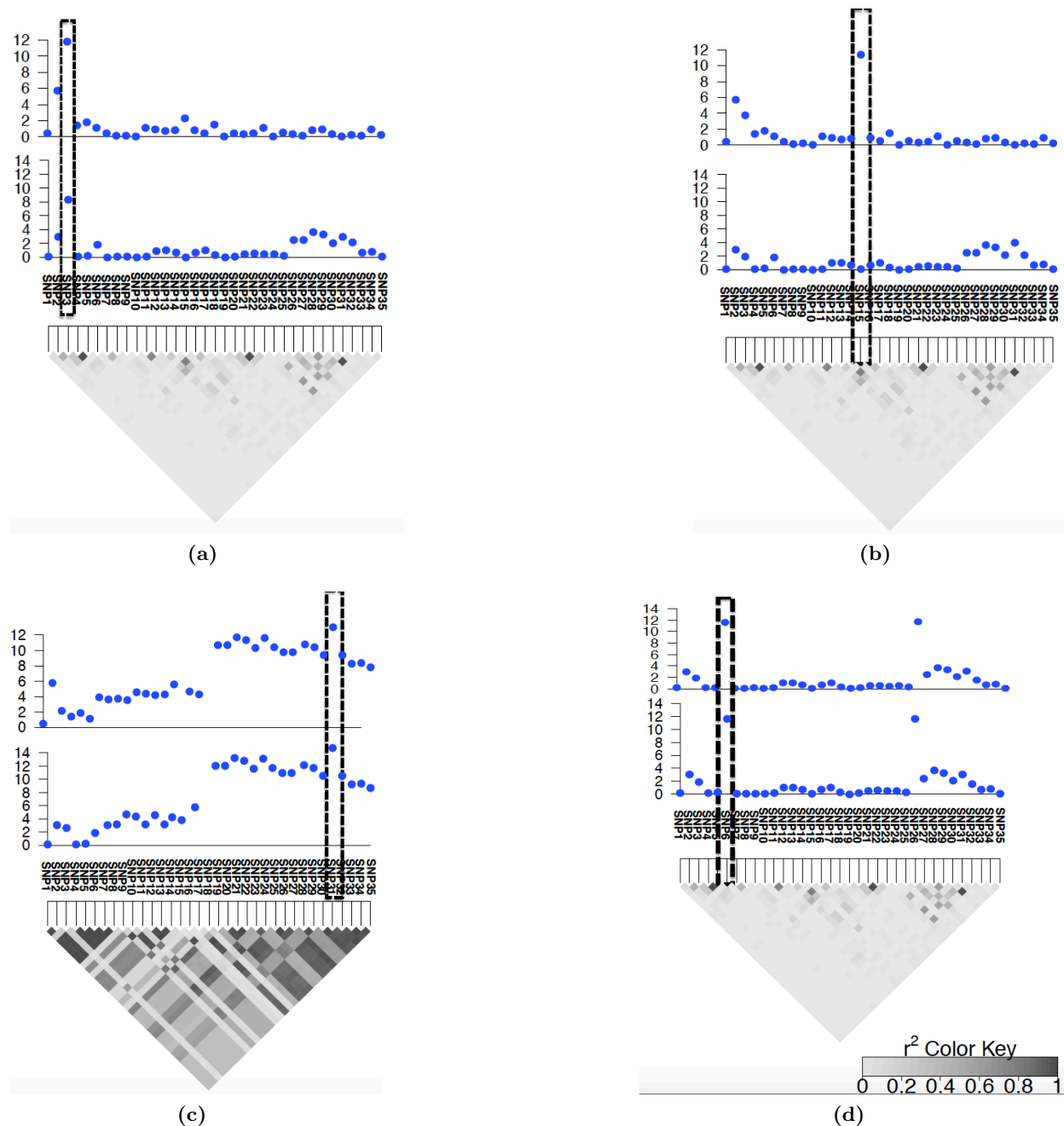


Figure 2. Overview of eCAVIAR. At high level eCAVIAR aligns the causal variants in eQTL and GWAS. The x-axis is the variant (SNP) location. The y-axis is the significant score (-log of p-value) for each variant. The grey triangle indicates the LD structure where every diamond in this triangle indicates the Pearson's correlation. The darker the diamond the higher the correlation and the lighter the diamond the lower the correlation between the variants. The case where the causal variants are aligned the colocalization posterior probability (CLPP) is high for the variant that is embedded in the dashed black rectangle as shown in panel (a). However, the case where the causal variants are not aligned (the causal variants are not the same variants) then the quantity of CLPP is low for the variant that is embedded in the dashed black rectangle as shown in panel (b). In the case, the LD is high, which implies the uncertainty is high due to LD, the CLPP value is low for the variant that is embedded in the dashed black rectangle as shown in panel (c). Panel (d) illustrates a case where in a locus we have two independent causal variants. If we consider that we only have one causal variant in a locus, then the CLPP of the causal variants are estimated to be 0.25. However, if we allow to have more than one causal variant in the locus, eCAVIAR estimates the CLPP to be 1.

2.2.1 eCAVIAR is Accurate in the Case of One Causal Variant

We apply eCAVIAR to the simulated datasets and compute the CLPP for each case. We use different cut-offs to determine whether or not a variant is shared between two studies. For each cut-off, we compute the false positive rate (FP) and true positive rate (TP). The baseline method checks if the most significant variant in GWAS is the most significant variant in eQTL study. We refer to this method as the Shared Peak SNP (SPS) method. The results are shown in Figures 3a and 3d. Moreover, we plot the same results in receiver operating characteristic (ROC) curve (Figure S1). We observe our method has higher TP and lower FP compared to SPS. However, eCAVIAR has low TP when the cut-off for CLPP is high. Furthermore, eCAVIAR has an extremely low FP. Our results imply that eCAVIAR has high confidence for selecting loci to be colocalized between the GWAS and eQTL. eCAVIAR is conservative in selecting a locus to be colocalized. Given the high cut-off of CLPP, eCAVIAR can miss some true colocalized loci. However, loci that are selected by eCAVIAR to be colocalized are likely to be predicted correctly.

The computed CLPP depends on the complexity of the LD at the locus. We apply eCAVIAR to the simulated datasets and compute the CLPP (Figure S2). Here, the average quantity of CLPP decreases as we increase the Pearson's correlation (r) between paired variants. This effect increases complexity of LD between the two variants. Furthermore, the 95% confidence intervals for the computed quantity increases as we increase the Pearson's correlation. This result implies that the computed CLPP can be small for a locus with complex LD, even when a variant is colocalized in both GWAS and eQTL studies.

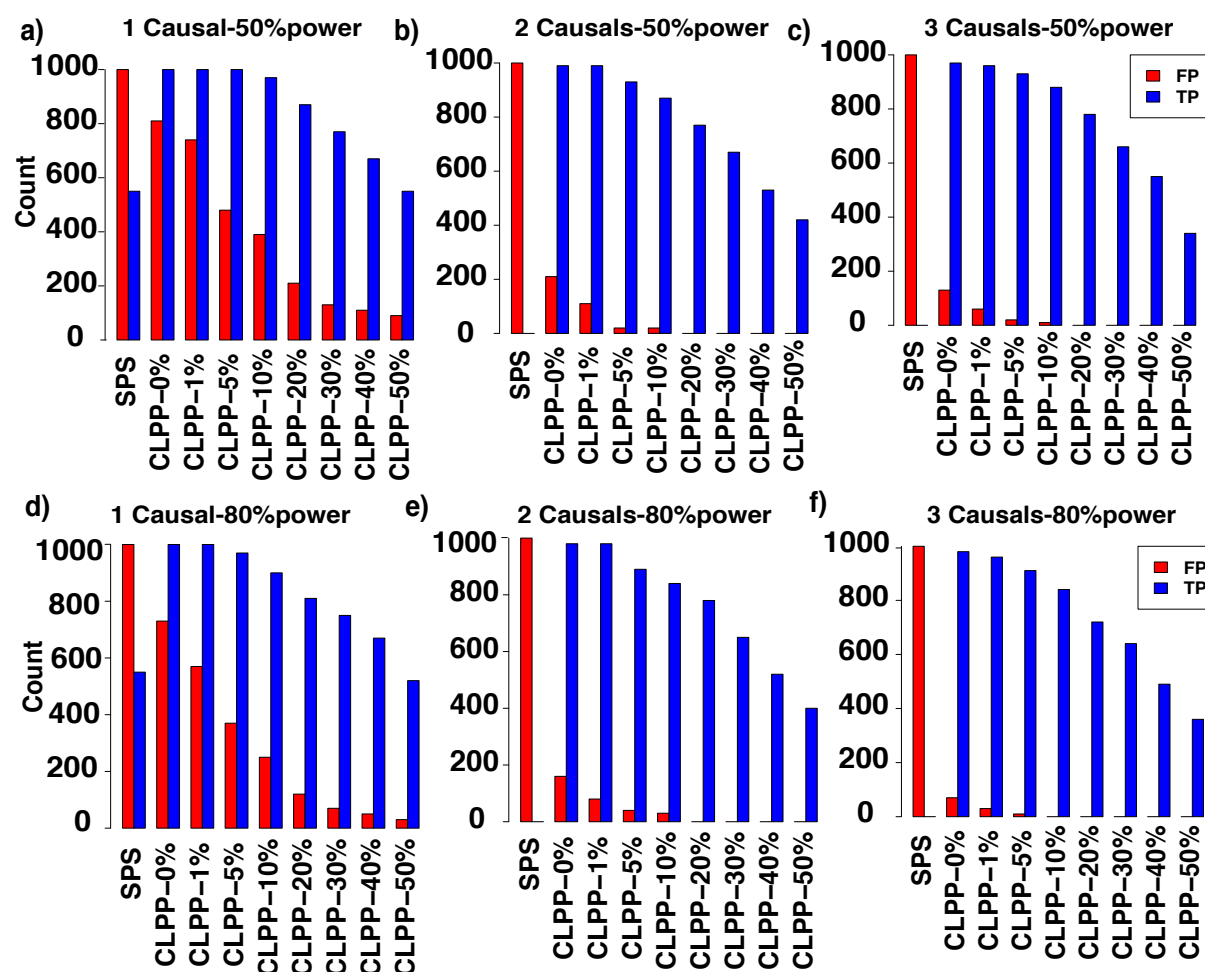


Figure 3. eCAVIAR is robust to the presence of allelic heterogeneity. We simulate marginal statistics directly from the LD structure for eQTL and GWAS. We implant one, two or three causal variants in both studies. Panels (a), (b), and (c) indicate the result for one, two, and three causal variants respectively where the statistical power on the causal variants is 50%. Panels (d), (e), and (f) indicate the result for one, two, and three causal variants respectively where the statistical power on the causal variants is 80%. eCAVIAR has a low TP for high cut-off, and eCAVIAR has low FP. This indicates that eCAVIAR has high confidence in detecting a locus to be colocalized between GWAS and eQTL, even in the presence of allelic heterogeneity.

2.2.2 eCAVIAR is Robust to the Presence of Allelic Heterogeneity

The presence of more than one causal variant in a locus is a phenomenon referred to as allelic heterogeneity (AH). AH may confound the association statistics in a locus, and colocalization for a locus harboring AH is challenging. In order to investigate the effect of AH, we perform the following simulations. We implanted two or three causal variants in both GWAS and eQTL, and we then generated the marginal statistics using MVN as mentioned in the previous section. Next, we compute TP and FP for eCAVIAR and SPS (see Figure 3). Figures 3a, 3b and 3c illustrate results of one, two, and three causal variants, respectively, when the statistical power is 50%. In a similar way, Figures 3d, 3e, and 3f illustrate results of one, two, and three causal variants, respectively, when the statistical power is 80%. Interestingly, SPS has a very low TP when there are two or three causal variants (see Figure 3). This implies that SPS is not accurate when AH is present. Similar to cases with one single casual variant (see Figures 3a and 3d), eCAVIAR has a very low FP when there are two or three causal variants (see Figures 3b, 3c, 3e, and 3f). This implies that eCAVIAR has high confidence in detecting a locus to be colocalized between GWAS and eQTL.

2.3 eCAVIAR is More Accurate than Existing Methods

We compare the results of eCAVIAR with RTC [26] and COLOC [13], two well known methods for eQTL and GWAS colocalization. We can use the previous section to generated simulated datasets; however, RTC is not designed to work with summary statistics. In order to provide a dataset compatible with RTC, we simulated eQTL and GWAS phenotypes under a linear additive model where we use simulated genotypes obtained from HAPGEN2 [38]. More details on the simulated datasets are provided in Section 3.3.3.

We compare the accuracy, precision, and recall rate of all three methods. Each method computes a probability for a variant to be causal in both eQTL and GWAS. In order to determine this probability for our comparison, we need to select two cut-off thresholds. We devised one threshold for detecting variants that are colocalized in both studies and another threshold to detect variants that are not colocalized. Here, we consider a variant to be causal in both studies if the probability of colocalization is greater than the colocalization cut-off threshold. The second cut-off threshold

is used to detect variants that are not causal in both studies. We consider a variant is non-causal in both studies if the probability of colocalization is less than the non-colocalization cut-off threshold. In our experiment, we set the non-colocalization cut-off threshold to be 0.001, and for the colocalization cut-off threshold, we vary this value from 0.0001 to 0.9.

eCAVIAR outperform existing methods when the locus has one causal variant. We observe that all three methods have a similarly high recall rate (see Figure S4). eCAVIAR has much higher accuracy and precision in comparison to RTC (see Figure 4). Next, we consider the performance of the three methods when the locus has allelic heterogeneity. We use the same simulation described in this section, but in this case we implant two causal variants instead of one causal variant. In this setting, eCAVIAR has higher accuracy and precision when compared to COLOC and RTC. However, RTC has a slightly higher recall rate in comparison to eCAVAIR. Moreover, RTC tends to perform better than COLOC in the presence of allelic heterogeneity (see Figure 5). This result indicates eCAVIAR is more accurate than existing methods—even in the presence of allelic heterogeneity. However, if there exists only one causal variant in a locus, COLOC has better performance than RTC. In cases with more than one causal variant, RTC has better performance. These results are obtained when we set the non-colocalization cut-off threshold to be 0.001. We change this value to 0.0001 to check the robustness of eCAVIAR. We observe for different values of non-colocalization still eCAVIAR outperforms existing methods (see Figures S5 and S6) .

eCAVIAR has better performance compared to COLOC and RTC, the pioneering methods for eQTL and GWAS colocalization. COLOC and RTC require different input data to perform the colocalization. COLOC only requires the marginal statistics from GWAS and eQTL studies. Unlike eCAVIAR, COLOC and RTC do not require the LD structure of genetic variants in a locus. However, RTC requires individual level data (genotypes and phenotypes) and is not applicable to datasets for which we have access to the summary statistics.

2.4 Integrating GTEx and MAGIC Datasets Using eCAVIAR

We utilize the Meta-Analyses of Glucose and Insulin-related traits Consortium (MAGIC) dataset and GTEx dataset [5] to detect the target gene and most relevant tissue for each GWAS risk locus. MAGIC datasets consist of 8 phenotypes [10]. These phenotypes are as follows: FastingGlucose, FastingInsulin, FastingProinsulin, HOMA-B (β -cell function), HOMA-IR (insulin resistance), and

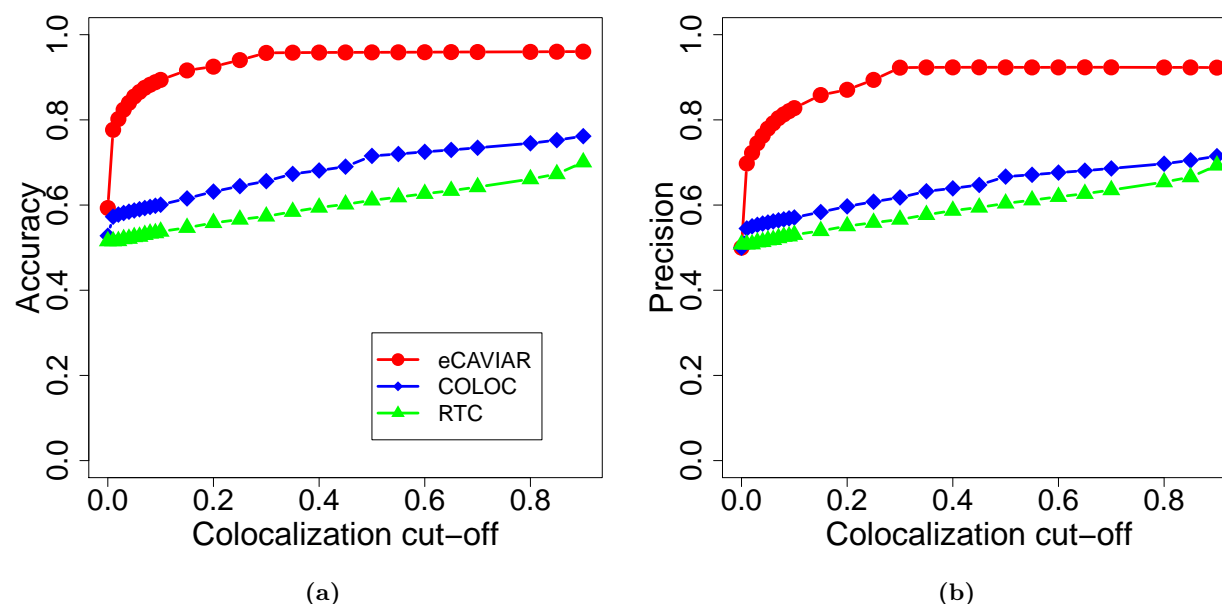


Figure 4. eCAVIAR is more accurate than existing methods for regions with one causal variant. We compare the accuracy and precision of eCAVIAR with the two existing methods (RTC and COLOC). The x-axis is the colocalization cut-off threshold. In these datasets we implant one causal variant, and we utilize simulated genotypes. We simulate the genotypes using HAPGEN2 [38] software. We use the European population from the 1000 Genomes data [1, 2] as the starting point to simulate the genotypes. Panels (a) and (b) illustrate the accuracy and precision respectively for all the three methods. We compute TP (true positive), TN (true negative), FN (false negative), and FP (false positive) for the set of simulated datasets where we generate the marginal statistics utilizing the linear model. Accuracy is the ratio of (TP+TN) and (TP+FP+FN+TN), $Accuracy = \frac{TP+TN}{TP+FP+FN+TN}$, and precision is the ratio of TP and (TP+FP), $Precision = \frac{TP}{TP+FP}$. We set the non-colocalization cut-off threshold to 0.001. We observe eCAVIAR and COLOC have higher accuracy and precision compared to RTC.

Hb1Ac (Hemoglobin A1c test for Diabetes), 2-hour glucose and 2-hour insulin after an oral glucose tolerance test. In our analysis, we use FastingGlucose (FG) and FastingProinsulin(FP) phenotypes which have the strongest and most association signals. FG phenotypes have 15 variants and FP phenotypes have 10 variants that are reported significantly associated with phenotypes from previous studies [10, 37]. We consider 44 tissues provided by GTEx consortium (Release v6, dbGaP Accession phs000424.v6.p1 available at: <http://www.gtexportal.org>) [5]. Table S2 lists tissues and the number of individuals for each tissue.

We want to detect the most relevant tissue and a target gene for each of previously reported significant variants in GWAS. eCACIAR utilizes the marginal statistics of all variants in a locus obtained from GWAS and eQTL. We obtain each locus by considering 50 variants upstream and

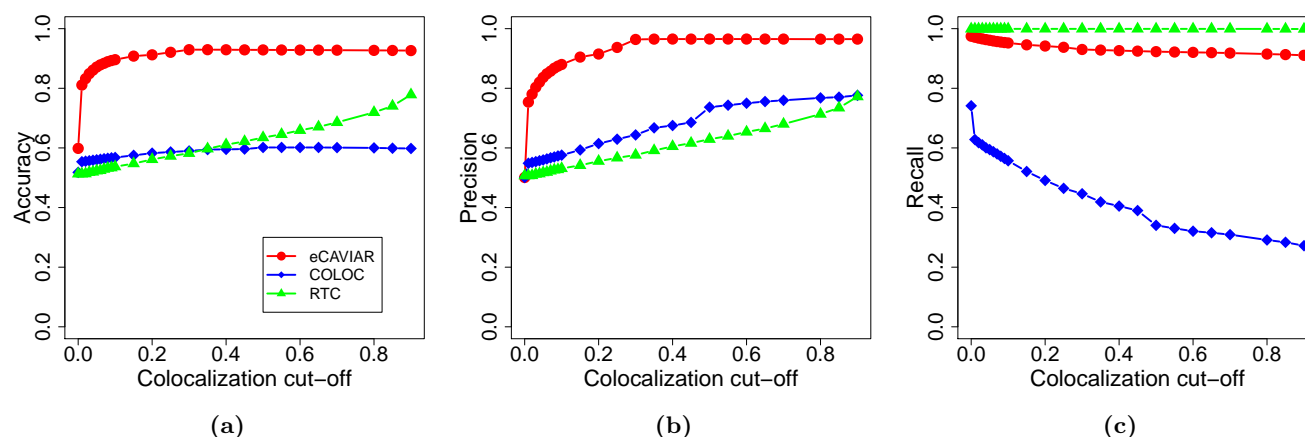


Figure 5. eCAVIAR is more accurate than existing methods in presence of allelic heterogeneity. We use similar process to generate the datasets as shown in Figure 4. However, in this case, we implant two causal variants. We simulate the genotypes using HAPGEN2 [38] software. We use the European population from the 1000 Genomes data [1, 2] as the starting point to simulate the genotypes. We compare the accuracy, precision, and recall rate. In these results, eCAVIAR tends to have higher accuracy and precision compared to the RTC and COLOC. However, RTC have slightly higher recall rate.

downstream of the reported variant. Then, we consider genes where at least one of the variants in the locus is significantly associated with the gene expression of that gene. Thus, for one GWAS variant, there may exist multiple genes in one tissue that satisfy these requirements, and we consider these pairs of variants and genes as potential colocalization loci. Table S3 and Table S4 list the potential colocalization loci for FG and FP phenotypes, respectively. For any given variant, we use CLPP to detect the most relevant tissue and a target gene. We then select the gene and tissue that have highest CLPP as the target gene and the most relevant tissue, respectively.

Table 1 and Table 2 indicate the result of eCAVIAR for FastingGlucose and FastingProinsulin, respectively. This result shows genetic variants that are causal in both eQTL and GWAS. We only considered variants that are reported to be significant with FG [10] and FP [37] phenotypes. We use the cut-off threshold of 0.01 (1%) to conclude two causal variants are shared.

Many of the significant variants have CLPP values which are in the range where it is difficult to make any conclusions about whether the causal variants are shared or not. However, we detect a large number of loci where the GWAS causal variants are clearly distinct from the causal variants in the eQTL data (Table 3). There are several genes that can be excluded in all tissues tested (e.g., *SEC22A* at the rs11708067 FG-locus where there is non-colocalization). Similarly, there

are instances where a gene that has been implicated previously as the likely gene mediating the GWAS association signal shows non-overlap with eQTLs for that gene in several GTEx tissues. An example of this can be found in *ADCY5*, also at the rs11708067 FG-locus. In pancreatic islet eQTL data the GWAS variant itself was also the primary eQTL signal for *ADCY5*. This could suggest that the phenotype acts through a tissue-specific regulatory element active in islets but not in GTEx tissues.

chrom	pos.	rsID	Relevant Tissue	Target Gene
7	44235668	rs4607517	Colon Sigmoid ($N=124$) Thyroid ($N=278$)	<i>GCK</i> <i>GCK</i>
11	47336320	rs7944584	Artery Tibial ($N=285$) Artery Tibial ($N=285$) Nerve Tibial ($N=256$) Nerve Tibial ($N=256$) Nerve Tibial ($N=256$) Pituitary ($N=87$)	<i>MDK</i> <i>MADD</i> <i>NR1H3</i> <i>CELF1</i> <i>RAPSN</i> <i>MADD</i>
11	45873091	rs11605924	Whole Blood($N=338$)	<i>MAPK8IP1</i>

Table 1. eCAVIAR joint analysis of FastingGlucose and GTEx dataset. We use N to indicate the number of individuals in each tissue which we have access to summary statistics from GTEx data.

chrom	pos.	rsID	Relevant Tissue	Target Gene
11	47293799	rs10501320	Artery Tibial ($N=285$) Artery Tibial ($N=285$) Esophagus Mucosa ($N=241$) Esophagus Muscularis ($N=218$) Pituitary ($N=87$)	<i>MDK</i> <i>MADD</i> <i>MADD</i> <i>C1QTNF4</i> <i>ARHGAP1</i>
11	72432985	rs11603334	Pituitary ($N=87$) Skin Sun Exposed Lower leg ($N=302$)	<i>PDE2A</i> <i>ARAP1</i>
15	71109147	rs1549318	Adipose Visceral Omentum ($N=185$) Cells Transformed Fibroblasts ($N=272$) Ovary ($N=85$)	<i>LARP6</i> <i>LARP6</i> <i>LARP6</i>

Table 2. eCAVIAR joint analysis of FastingProinsulin and GTEx dataset. We use N to indicate the number of individuals in each tissue which we have access to summary statistics from GTEx data

For a majority of loci in which we identify a single variant causal for both GWAS and eQTL, our results show that more than one gene is implicated. At rs7944584 (FG) / rs10501320 (FP) and rs1549318 (FP) there is support from other eQTL studies for causal roles for *MADD* in human pancreatic islets of Langerhans [41] and *LARP6* in adipose tissue [37], respectively. However, assessing potential candidacy of the roles of the other implicated genes will require additional sources of information, such as Capture-C experiments [19] to demonstrate chromatin interactions

between causal variant and gene promotor and/or in vitro function validation in relevant model systems. One other potential reason for the limited number of variants causal for both GWAS and eQTL signal identified for FG is that this phenotype is thought to mainly act through the pancreatic islet, with many of the identified loci containing compelling biological candidates such as transcription factors involved in pancreas development (e.g., *FOXA2* and *PDX1*) [10, 35]. Since pancreatic islets are not part of the GTEx dataset, applying eCAVIAR to data from this cell-type might provide further mechanistic insights.

Phenotype	chrom	pos.	rsID	GWAS p-value	eQTL p-value ¹	# Gene	# Tissue
FG	2	27741237	rs780094	2.49E-12	2.95E-55	17	30
	2	169763148	rs560887	4.61E-75	1.36E-14	5	20
	3	123065778	rs11708067	8.72E-09	4.28E-42	5	34
	9	4289050	rs7034200	0.0001204	9.95E-12	8	7
	10	113042093	rs10885122	8.41E-11	7.73E-11	2	3
	11	61571478	rs174550	1.48E-08	1.03E-125	24	29
	11	92708710	rs10830963	1.26E-68	7.49E-06	7	6
FP	1	99177253	rs9727115	5.285e-06	7.04e-16	3	12
	10	114758349	rs7903146	3.48e-18	7.92e-33	7	26
	15	62383155	rs4502156	3.80e-11	8.48e-14	7	15
	17	2262703	rs4790333	2.15e-08	5.39e-75	21	33

Table 3. The loci where the causal variants between the eQTL and GWAS are different. We utilize FastingGlucose (FG) and FastingProinsulin (FP) phenotypes. Number of genes and tissues indicate the genes and tissues, respectively which we apply eCAVIAR for a GWAS risk variant. The complete list of genes and tissue are provided in Tables S3 and S4 for FG and FP phenotypes, respectively. eCAVIAR utilizes the marginal statistics of all variants in a locus obtained from GWAS and eQTL. We obtain each locus by considering 50 variants upstream and downstream of the reported variant. Then, we consider genes where at least one of the variants in the locus is significantly associated with the gene expression of that gene. Thus, for one GWAS variant, we can have multiple genes in one tissue that satisfy our condition.

3 Material and Methods

3.1 CAVIAR model for Fine-mapping

Standard GWAS and Indirect Association. We collect quantitative traits for N individuals and genotypes for all the individuals at M SNPs (variants). In this case, we collect data for one phenotype and gene expression of each gene. We assume that both the phenotype and the gene expression have at least one significant variant. To simplify the description of our method, we assume that the number of individuals and the pairwise Pearson’s correlations of genotype (LD) in both GWAS and eQTL are the same. In the supplementary material, we use a more general

¹The p-value of the most significant variant in eQTL among all genes and all tissues

model where the number of individuals and LD in both GWAS and eQTL are not the same. Let $Y^{(p)}$ indicate a $(N \times 1)$ vector of the phenotypic values where $y_j^{(p)}$ denotes the phenotypic value for the j -th individuals. We use $Y^{(e)}$ to indicate a $(N \times 1)$ vector of gene expression collected for one gene of interest, for which there exists one significant variant associated with the gene expression of that gene. Let G indicate a $(N \times M)$ matrix of genotype information where G_i is a $(N \times 1)$ vector of minor allele counts for all the N individuals at the i -th variant. In this setting g_{ji} indicates the j -th element from vector G_i that indicates the minor allele count for the j -th individual. In diploid genomes such as human, we can have 3 possible minor allele counts, $g_{ji} = \{0, 1, 2\}$. We normalize both the phenotypes and the genotypes to mean zero and variance one, where the X is the normalized matrix of the G . Let X_i denote a $(N \times 1)$ vector of normalized minor allele counts for the i -th variant. We assume “additive” Fisher’s polygenic model, which is widely used by GWAS community. In the Fisher’s polygenic model, the phenotypes and genotypes follows a normal distribution. The additive assumption implies each variant contribute linearly to the phenotype. Thus, we consider the following linear model:

$$Y^{(p)} = \mu^{(p)} \mathbf{1} + \sum_{i=1}^M \beta_i^{(p)} X_i + \mathbf{e}^{(p)},$$

where $\mu^{(p)}$ is the phenotypic mean, $\beta_i^{(p)}$ is the effect size of the i -th variant towards the phenotype of interest, and $\mathbf{e}^{(p)}$ is the environment and measurement error toward the collected phenotype. In this model, we assume $\mathbf{e}^{(p)}$ is a vector of i.i.d and normally distributed. Let $\mathbf{e}^{(p)} \sim N(0, \sigma_e^{(p)2} \mathbf{I})$ where $\sigma_e^{(p)}$ is a covariance scalar and \mathbf{I} is a $(N \times N)$ identity matrix. Similarly, for the gene of interest that we perform eQTL, we assume the additive Fisher’s polygenic model. Thus, we use the following model:

$$Y^{(e)} = \mu^{(e)} \mathbf{1} + \sum_{i=1}^M \beta_i^{(e)} X_i + \mathbf{e}^{(e)},$$

In standard GWAS when we test the significance of a variant, we test each variant independently. Let’s consider the c -th variant is causal. In the testing process, we use the following model:

$$Y^{(p)} = \mu^{(p)} \mathbf{1} + \beta_c^{(p)} X_c + \mathbf{e}^{(p)}.$$

We use the same testing model for the eQTL study as well :

$$Y^{(e)} = \mu^{(e)}\mathbf{1} + \beta_c^{(e)}X_c + \mathbf{e}^{(e)},$$

utilizing the maximum likelihood, we compute the optimal estimate of $\beta_c^{(p)}$. We use $\hat{\beta}_c^{(p)}$ to indicate the optimal estimate of $\beta_c^{(p)}$ that is computed as $\hat{\beta}_c^{(p)} = \frac{X_c^T Y^{(p)}}{X_c^T X_c}$, $\hat{\beta}_c^{(p)} \sim N(\beta_c^{(p)}, \frac{\sigma_e^{(p)2}}{X_c^T X_c})$ and the marginal statistics for this variant is computed as $S_c^{(p)} = \frac{\hat{\beta}_c^{(p)}}{\hat{\sigma}_e^{(p)}} \sqrt{X_c^T X_c} \sim N(\lambda_c^{(p)}, 1)$ where $\lambda_c^{(p)}$ is non-centrality parameter (NCP). Next, we consider the indirect association. We test variant i , which is not causal yet is in LD with a causal variant c . It is shown in previous works [14, 16, 17, 21] the statistic is computed as $S_i^{(p)} = \frac{\hat{\beta}_i^{(p)}}{\hat{\sigma}_e^{(p)}} \sqrt{X_i^T X_i} \sim N(r_{ic}\lambda_c^{(p)}, 1)$, where r_{ic} is the genotype correlation between the variants i and c . We refer to $r_{ic}\lambda_c^{(p)}$ as the LD-induced NCP. Moreover, it is known the covariance between the computed statistics $S_i^{(p)}$ and $S_c^{(p)}$ is equal to r_{ic} [3, 9, 22, 29].

In our setting, we have the marginal statistics of M variants for phenotype of interest and the gene expression. Let $S^{(p)} = \{s_1^{(p)}, s_2^{(p)}, \dots, s_M^{(p)}\}$ and $S^{(e)} = \{s_1^{(e)}, s_2^{(e)}, \dots, s_M^{(e)}\}$ indicate the marginal statistics for the phenotype of interest and the gene expression, respectively. The joint distribution of the marginal statistics given the true NCPs follows a multivariate normal (MVN) distribution, which is known from previous studies [14, 16, 17, 21]. Thus, we have:

$$(S^{(p)}|\Lambda^{(p)}) \sim \mathcal{N}(\Sigma\Lambda^{(p)}, \Sigma), \quad (1)$$

Using the fact that eQTL on the gene of interest can be viewed as a GWAS on the expression level of that gene, we have the following distribution for the marginal statistics:

$$(S^{(e)}|\Lambda^{(e)}) \sim \mathcal{N}(\Sigma\Lambda^{(e)}, \Sigma). \quad (2)$$

where Σ is the pairwise Pearson's correlations of genotypes. Let $\Lambda^{(p)} = \{\lambda_1^{(p)}, \lambda_2^{(p)}, \dots, \lambda_M^{(p)}\}$ and $\Lambda^{(e)} = \{\lambda_1^{(e)}, \lambda_2^{(e)}, \dots, \lambda_M^{(e)}\}$ be the true NCPs for all the variants of desired phenotype and gene expression, respectively. The true NCP for a non-causal variant is zero and non-zero for causal variant. Let $\Sigma\Lambda^{(e)}$ and $\Sigma\Lambda^{(p)}$ be the LD-induced NCPs for desired phenotype and gene expression, respectively.

CAVIAR Generative Model for Single phenotype. We introduce a new variable $C^{(p)}$ which

is an $(M \times 1)$ binary vector. We refer to this binary vector as causal status. The causal status indicates which variants are causal and which are not. We set $c_i^{(p)}$ to be one if the i -th variant is causal and zero otherwise. In CAVIAR [16, 17], we introduce a prior on the vector of NCPs utilizing the MVN distribution. This prior on the vector of NCPs given the causal status vector is defined as follow:

$$(\Lambda^{(p)}|C^{(p)}) \sim \mathcal{N}\left(0, \sigma^{(p)}\Sigma_c^{(p)}\right), \quad (3)$$

where $\Sigma_c^{(p)}$ is a diagonal matrix and $\sigma^{(p)}$ is a constant which indicates the variance of our prior over the GWAS NCPs. We set $\sigma^{(p)}$ to 5.2 [16, 17]. The diagonal elements of $\Sigma_c^{(p)}$ are set to one or zero where variants that are selected causal in $C^{(p)}$ their corresponding diagonal elements in $\Sigma_c^{(p)}$ are set to one; otherwise, we set them to zero. Utilizing this prior as a conjugate prior, in CAVIAR, we compute the likelihood of each possible causal status. The joint distribution of the marginal statistics given the causal status is as follows:

$$(S^{(p)}|C^{(p)}) \sim \mathcal{N}\left(0, \Sigma + \sigma^{(p)}\Sigma_c^{(p)}\Sigma\right), \quad (4)$$

In a similar way, for the gene of interest which we perform eQTL, we have:

$$(\Lambda^{(e)}|C^{(e)}) \sim \mathcal{N}\left(0, \sigma^{(e)}\Sigma_c^{(e)}\right), \quad (5)$$

where $\Sigma_c^{(e)}$ is a diagonal matrix and $\sigma^{(e)}$ is set to 5.2 [16, 17]. The diagonal elements of $\Sigma_c^{(e)}$ are set to one or zero where variants that are selected causal in $C^{(e)}$ their corresponding diagonal elements in $\Sigma_c^{(e)}$ are set to one; otherwise, we set them to zero

3.2 eCAVIAR Computes the Colocalization Posterior Probability for GWAS and eQTL

Given the marginal statistics for GWAS and eQTL, which are denoted by $S^{(p)}$ and $S^{(e)}$, respectively, we want to compute the colocalization posterior probability (CLPP). CLPP is the probability that the same variant is causal in both studies. For simplicity, we compute CLPP for the i -th variant. We define CLPP for the i -th variant as $P\left(c_i^{(p)} = 1, c_i^{(e)} = 1 | S^{(p)}, S^{(e)}\right)$, and we use ϕ_i to indicate

the CLPP for the i -th variant. We utilize the law of total probability to compute the summation probability of all causal status where the i -th variant is causal in both GWAS and eQTL and other variants can be causal or non-causal. Thus, the above equation can be extend as follows:

$$\begin{aligned}\phi_i &= P\left(c_i^{(p)} = 1, c_i^{(e)} = 1 | S^{(p)}, S^{(e)}\right) \\ &= \sum_{C_{/i}^{*(p)} \in \{0,1\}^{M-1}} \sum_{C_{/i}^{*(e)} \in \{0,1\}^{M-1}} P\left(C_{/i}^{(p)} = C_{/i}^{*(p)}, C_{/i}^{(e)} = C_{/i}^{*(e)}, c_i^{(p)} = 1, c_i^{(e)} = 1 | S^{(p)}, S^{(e)}\right) \\ &= \sum_{C^{*(p)} \in \{0,1\}^M} \sum_{C^{*(e)} \in \{0,1\}^M} P\left(C^{(p)} = C^{*(p)}, C^{(e)} = C^{*(e)} | S^{(p)}, S^{(e)}\right) \mathbb{I}(c_i^{*(p)} = 1, c_i^{*(e)} = 1) \quad (6)\end{aligned}$$

where $C_{/i}^{(p)}$ and $C_{/i}^{(e)}$ are causal status vectors for all the variants excluding the i -th variant for the phenotype of interest and gene expression, respectively. Let $\mathbb{I}()$ be an indicator function which is defined as follows:

$$\mathbb{I}(c_i^{*(p)} = 1, c_i^{*(e)} = 1) = \begin{cases} 1 & c_i^{*(p)} \text{ and } c_i^{*(e)} \text{ are causal} \\ 0 & o/w \end{cases} \quad (7)$$

Utilizing the Bayes' rule, we compute the CLPP as follows:

$$\phi_i = \frac{\sum_{C^{*(p)}} \sum_{C^{*(e)}} P\left(S^{(p)}, S^{(e)} | C^{(p)} = C^{*(p)}, C^{(e)} = C^{*(e)}\right) P\left(C^{*(p)}, C^{*(e)}\right) \mathbb{I}(c_i^{*(p)} = 1, c_i^{*(e)} = 1)}{\sum_{C^{*(p)}} \sum_{C^{*(e)}} P\left(S^{(p)}, S^{(e)} | C^{(p)} = C^{*(p)}, C^{(e)} = C^{*(e)}\right) P\left(C^{*(p)}, C^{*(e)}\right)} \quad (8)$$

where $P\left(C^{*(p)}, C^{*(e)}\right)$ is the prior probability of the causal status of $C^{*(p)}$ and $C^{*(e)}$ for the GWAS and eQTL respectively. We assume the prior probability over the causal status for the GWAS and eQTL are independent, $P\left(C^{*(p)}, C^{*(e)}\right) = P\left(C^{*(p)}\right) P\left(C^{*(e)}\right)$. To compute the prior of causal status, we use the same assumptions that is widely used in the fine mapping methods [6, 16, 17], where the probability of causal status follows a Binomial distribution with the probability of variant being causal is equal to γ . Thus, this prior is equal to $P\left(C^{*(p)}\right) = \prod_{i=1}^M \gamma^{c_i^{*(p)}} (1 - \gamma)^{1 - c_i^{*(p)}}$ and γ is set to 0.01 [8, 11, 16, 39].

GWAS and eQTL studies are usually performed on independent sets of individuals. Further-

more, given the causal status of both the GWAS and eQTL, the marginal statistics for these two studies are independent. We have $P(S^{(p)}, S^{(e)} | C^{*(p)}, C^{*(e)}) = P(S^{(p)} | C^{*(p)}) P(S^{(e)} | C^{*(e)})$. Thus, we simplify the Equation (8), and the CLPP is computed as follows:

$$\phi_i = \frac{\sum_{C^{*(p)}} P(S^{(p)} | C^{(p)} = C^{*(p)}) P(C^{*(p)}) \mathbb{I}(c_i^{*(p)} = 1)}{\sum_{C^{*(p)}} P(S^{(p)} | C^{(p)} = C^{*(p)}) P(C^{*(p)})} \times \frac{\sum_{C^{*(e)}} P(S^{(e)} | C^{(e)} = C^{*(e)}) P(C^{*(e)}) \mathbb{I}(c_i^{*(e)} = 1)}{\sum_{C^{*(e)}} P(S^{(e)} | C^{(e)} = C^{*(e)}) P(C^{*(e)})} \quad (9)$$

The above equation indicates the probability that the same variant is causal in both GWAS and eQTL is independent. This probability is equal to the multiplication of two probabilities: probability that the variant is causal in GWAS and the probability of the same variant is causal in the eQTL study. Thus, we compute the CLPP as, $P(c_i^{(p)} = 1, c_i^{(e)} = 1 | S^{(p)}, S^{(e)}) = P(c_i^{(p)} = 1 | S^{(p)}) \times P(c_i^{(e)} = 1 | S^{(e)})$ where $P(c_i^{(p)} = 1 | S^{(p)})$ and $P(c_i^{(e)} = 1 | S^{(e)})$ are computed from the first part and second part of Equation (9), respectively.

3.3 Generating simulated datasets

3.3.1 Simulating Genotypes

We first simulated genotype data starting from the real genotypes obtained from European population in the 1000 Genomes data [1, 2]. In order to simulate the genotypes we utilize HAPGEN2 [38] software which is widely used to generate genotypes. We focus on the chromosome 1 and the GWAS variants that are obtained from the NHGRI catalog [45]. We consider 200-kb windows around the lead SNP to generate a locus. Then, we filter out monomorphic SNPs and SNPs with low minor allele frequency ($MAF \leq 0.01$) inside a locus.

3.3.2 Simulating Summary Statistics Directly from LD Structure

We generate LD matrix for a locus by computing the genotype Pearson's correlations between each pair of variants. Then, we generate marginal summary statistics for each locus, assuming the marginal summary statistics follows MVN that is utilized in previous studies [14, 16, 17, 21, 46]. We measure the strength of a causal variant based on NCP. We set the NCP of the causal variant in order to obtain a certain statistical power. The NCP of the non-causal variants are set to zero.

The statistical power is the probability of detecting a variant to be causal under the assumption that the causal variant is present. The statistical power is computed as follows:

$$\text{Power} = 1 - \frac{1}{\sqrt{2\pi}} \int_{\Phi^{-1}(\alpha/2)+\lambda}^{\Phi^{-1}(1-\alpha/2)+\lambda} e^{-\frac{1}{2}x^2} dx = \Phi(\Phi^{-1}(\alpha/2) + \lambda) + 1 - \Phi(\Phi^{-1}(\alpha/2) + \lambda)$$

where α is the significant threshold. Moreover, Φ and Φ^{-1} denote the cumulative density function (CDF) and inverse of CDF for the standard normal distribution. In our experiment, the NCP is computed for the genome-wide significant level ($\alpha = 10^{-8}$). We use binary search to compute the value of NCP for a desired statistical power.

3.3.3 Simulating Summary Statistics Utilizing Linear Additive Model

We utilize 100 variants in a locus to generate the simulated phenotypes from the simulated genotypes. We simulate the phenotypes assuming the linear additive model, which is as follows:

$$Y = \sum_{i=1}^M \beta_i X_i + e \quad (10)$$

where $e \sim N(0, \sigma_e^2)$. We generate the effect size of the causal variant from a normal distribution with mean zero and variance equal to σ_g^2/M_c where M_c indicates the number of causal variants in a locus. Furthermore, we set the effect size to zero for variants that are not causal. Thus, the effect size for each variant is simulated as follows:

$$\begin{cases} \beta_i = 0 & \text{if } i\text{-th variant is non-causal} \\ \beta_i \sim N(0, \sigma_g^2/M_c) & \text{if } i\text{-th variant is causal} \end{cases}$$

After simulating phenotype for all the individuals, we utilize linear regression to estimate the effect sizes and the marginal statistics for all the M variants in a locus. In our simulations M is equal to 100.

4 Discussion

Integrating GWAS and eQTL provides insights into the underlying mechanism for genetic variants detected in GWAS. In this paper, we propose a quantity that can measure CLPP, the probability that the same variant is causal in both GWAS and eQTL studies, while accounting for the LD. Utilizing CLPP, we can identify target genes and relevant tissues. Moreover, we can detect loci where the causal variants are different between the two studies with high confidence. We observe from our analysis that, in most cases, GWAS risk loci and eQTL are different.

As most GWAS loci were discovered to lie outside of coding regions, it is implicitly assumed that these implicated loci will affect the regulation of genes. However, our results produce a lower than expected number of variants colocalized between both GWAS and eQTL studies. This points to a more complicated relationship between gene regulation and disease. It is likely that future studies will shed some light to explain this observation.

One conjecture is that the GWAS loci in fact do affect expression, but are secondary signals compared to the stronger associations found in current eQTL studies. As eQTL studies include an increasing number of individuals, we will be able to prove or disprove this conjecture. Furthermore, the heterogeneity of tissues may render it hard to detect eQTLs specific to a disease-relevant cell type that comprises only a fraction of the tissue. A second possibility is that GWAS variants affect other aspects of gene regulation such as splicing, or regulation at a level other than transcription regulation. Several studies have shown that alternative splicing may explain the causal mechanism of complex disease associations (e.g., a variant associated with multiple sclerosis that leads to exon skipping in *SP140* [24]). Methods that identify variants associated with differences in relative expression of alternative transcript isoforms or exon junction abundances are being applied to the latest version of GTEx data [25, 28]. As we obtain more functional genomics information and are able to measure quantities such as protein abundance, we will be able to systematically catalogue variants which affect regulation at levels other than transcription. A third possibility is that GWAS loci are eQTL loci but only in certain conditions such as development which are not the conditions where expression levels are measured. Regardless, our study demonstrates strong evidence in support of the idea that most GWAS loci are not strong eQTL loci and that the mechanism of GWAS loci affecting gene regulation is more complicated than we expected.

Broadly, we identify an analogy between colocalization and fine-mapping methods. Fine-mapping methods can be categorized into three main classes. One class relies on just the computed marginal statistics that are obtained from GWAS or eQTL. In this class of methods, the probability that a variant is causal depends on the rank of a variant, which is obtained from the marginal statistics. Recently, Maller et. al [23] have proposed a new fine-mapping method that utilizes the Bayes factor. This method provides results similar to those produced by approaches that rank variants based solely on their marginal statistics. Maller et. al [23] method for fine-mapping is similar in nature to COLOC [13], which is used for colocalization. The second class of methods is based on a conditional model where we re-compute the marginal statistics of all variants by conditioning on variants selected as causal. The conditional method for fine-mapping and RTC [26] have some similarity in nature. The third class of methods is CAVIAR [16, 17], CAVIARBF [7] and FINEMAP [6], which assumes a presence of more than one causal variant in a region. These probabilistic based methods use the MVN distribution. In these methods, we detect a set of variants that can capture all the causal variants with a predefined probability. Thus, eCAVIAR is analogous in process to CAVIAR, CAVIARBF, and FINEMAP.

eCAVIAR is a probabilistic method that integrates GWAS and eQTL signals to detect biological mechanisms. eCAVIAR has several advantages over prior approaches. First, it can account for multiple causal variants in any given locus. Second, it leverages summary statistics without accessing the raw individual data. In addition, eCAVIAR can provide confidence levels for the colocalization of a GWAS risk variant. Utilizing the confidence level, we can categorize a variant to three categories: variants which colocalize, variants which do not colocalize, and variants which are ambiguous to detect their colocalization status for the current data. High-throughput technologies have made it possible to obtain multi-tissue eQTL studies. Leveraging multi-tissue eQTL studies such as GTEx and eCAVIAR can advance discovery of new biological mechanisms for GWAS risk loci.

5 Acknowledgments

FH, JWJJ, MB and EE are supported by National Science Foundation grants 0513612, 0731455, 0729049, 0916676, 1065276, 1302448, 1320589 and 1331176, and National Institutes of Health grants

K25-HL080079, U01-DA024417, P01-HL30568, P01-HL28481, R01-GM083198, R01-ES021801, R01-MH101782 and R01-ES022282. EE is supported in part by the NIH BD2K award, U54EB020403. MvdB is supported by a Novo Nordisk postdoctoral fellowship run in partnership with the University of Oxford. AVS and XL are supported by contract HHSN268201000029C (Broad Institute). We acknowledge the support of the NINDS Informatics Center for Neurogenetics and Neurogenomics (P30 NS062691).

6 Web Resources

eCAVIAR is available <http://genetics.cs.ucla.edu/caviar/>

References

1. 1000 Genomes Project Consortium (2010). A map of human genome variation from population-scale sequencing. *Nature*, **467**(7319), 1061–1073.
2. 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from 1,092 human genomes. *Nature*, **491**(7422), 56–65.
3. Abecasis, G. R., Noguchi, E., Heinzmann, A., *et al.* (2001). Extent and distribution of linkage disequilibrium in three genomic regions. *The American Journal of Human Genetics*, **68**(1), 191–197.
4. Altshuler, D. M., Gibbs, R. A., Peltonen, L., *et al.* (2010). Integrating common and rare genetic variation in diverse human populations. *Nature*, **467**(7311), 52–8.
5. Ardlie, K. G., Deluca, D. S., Segrè, A. V., *et al.* (2015). The genotype-tissue expression (gtex) pilot analysis: Multitissue gene regulation in humans. *Science*, **348**(6235), 648–660.
6. Benner, C., Spencer, C. C., Ripatti, S., and Pirinen, M. (2015). Finemap: Efficient variable selection using summary data from genome-wide association studies. *bioRxiv*, page 027342.
7. Chen, W., Larrabee, B. R., Ovsyannikova, I. G., *et al.* (2015). Fine mapping causal variants with an approximate bayesian method using marginal test statistics. *Genetics*, **200**(3), 719–736.
8. Darnell, G., Duong, D., Han, B., and Eskin, E. (2012). Incorporating prior information into association studies. *Bioinformatics*, **28**(12), i147–i153.
9. Dunning, A. M., Durocher, F., Healey, C. S., *et al.* (2000). The extent of linkage disequilibrium in four populations with distinct demographic histories. *The American Journal of Human Genetics*, **67**(6), 1544–1554.
10. Dupuis, J., Langenberg, C., Prokopenko, I., *et al.* (2010). New genetic loci implicated in fasting glucose homeostasis and their impact on type 2 diabetes risk. *Nature Genetics*, **42**(2), 105–116.

11. Eskin, E. (2008). Increasing power in association studies by using linkage disequilibrium structure and molecular function as prior information. *Genome research*, **18**(4), 653–660.
12. Gamazon, E. R., Wheeler, H. E., Shah, K. P., *et al.* (2015). A gene-based association method for mapping traits using reference transcriptome data. *Nature Genetics*, **47**(9), 1091–1098.
13. Giambartolomei, C., Vukcevic, D., Schadt, E. E., *et al.* (2014). Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS Genetics*, **10**(5), e1004383.
14. Han, B., Kang, H. M., and Eskin, E. (2009). Rapid and accurate multiple testing correction and power estimation for millions of correlated markers. *PLoS Genetics*, **5**(4), e1000456.
15. He, X., Fuller, C. K., Song, Y., *et al.* (2013). Sherlock: detecting gene-disease associations by matching patterns of expression qtl and gwas. *The American Journal of Human Genetics*, **92**(5), 667–680.
16. Hormozdiari, F., Kostem, E., Kang, E. Y., Pasaniuc, B., and Eskin, E. (2014). Identifying causal variants at loci with multiple signals of association. *Genetics*, **198**(2), 497–508.
17. Hormozdiari, F., Kichaev, G., Yang, W.-Y., Pasaniuc, B., and Eskin, E. (2015). Identification of causal genes for complex traits. *Bioinformatics*, **31**(12), i206–i213.
18. Huang, Y.-T., Liang, L., Moffatt, M. F., Cookson, W. O., and Lin, X. (2015). igwas: Integrative genome-wide association studies of genetic and genomic data for disease susceptibility using mediation analysis. *Genetic epidemiology*, **39**(5), 347–356.
19. Hughes, J. R., Roberts, N., McGowan, S., *et al.* (2014). Analysis of hundreds of cis-regulatory landscapes at high resolution in a single, high-throughput experiment. *Nature Genetics*, **46**(2), 205–212.
20. Joo, J. W. J., Hormozdiari, F., Eskin, E., and Han, B. (2016). Multiple testing correction in linear mixed models. *Genome Biology*.
21. Kostem, E., Lozano, J. A., and Eskin, E. (2011). Increasing power of genome-wide association studies by collecting additional single-nucleotide polymorphisms. *Genetics*, **188**(2), 449–460.

22. Kruglyak, L. (1999). Prospects for whole-genome linkage disequilibrium mapping of common disease genes. *Nature Genetics*, **22**(2), 139–144.
23. Maller, J. B., McVean, G., Byrnes, J., *et al.* (2012). Bayesian refinement of association signals for 14 loci in 3 common diseases. *Nature Genetics*, **44**(12), 1294–1301.
24. Matesanz, F., Potenciano, V., Fedetz, M., *et al.* (2015). A functional variant that affects exon-skipping and protein expression of SP140 as genetic mechanism predisposing to multiple sclerosis. *Human Molecular Genetics*, **24**(19), 5619–5627.
25. Monlong, J., Calvo, M., Ferreira, P. G., and Guigó, R. (2014). Identification of genetic variants associated with alternative splicing using sQTLseeker. *Nature Communications*, **5**, 4698.
26. Nica, A. C., Montgomery, S. B., Dimas, A. S., *et al.* (2010). Candidate causal regulatory effects by integration of expression qtls with complex trait genetic associations. *PLoS Genetics*, **6**(4), e1000895.
27. Nicolae, D. L., Gamazon, E., Zhang, W., *et al.* (2010). Trait-associated snps are more likely to be eqtls: annotation to enhance discovery from gwas. *PLoS Genetics*, **6**(4), e1000888.
28. Ongen, H. and Dermitzakis, E. T. (2015). Alternative splicing QTLs in european and african populations. *The American Journal of Human Genetics*, **97**(4), 567–575.
29. Pickrell, J., Clerget-Darpoux, F., and Bourgain, C. (2007). Power of genome-wide association studies in the presence of interacting loci. *Genet Epidemiol*, **31**(7), 748–762.
30. Plagnol, V., Smyth, D. J., Todd, J. A., and Clayton, D. G. (2009). Statistical independence of the colocalized association signals for type 1 diabetes and rps26 gene expression on chromosome 12q13. *Biostatistics*, **10**(2), 327–334.
31. Pritchard, J. K. and Przeworski, M. (2001). Linkage disequilibrium in humans: models and data. *The American Journal of Human Genetics*, **69**(1), 1–14.
32. Rietveld, C. A., Medland, S. E., Derringer, J., *et al.* (2013). Gwas of 126,559 individuals

- identifies genetic variants associated with educational attainment. *Science*, **340**(6139), 1467–1471.
33. Ripke, S., O’Dushlaine, C., Chambert, K., *et al.* (2013). Genome-wide association analysis identifies 13 new risk loci for schizophrenia. *Nature Genetics*, **45**(10), 1150–1159.
34. Saxena, R., Hivert, M.-F., Langenberg, C., *et al.* (2010). Genetic variation in gipr influences the glucose and insulin responses to an oral glucose challenge. *Nature genetics*, **42**(2), 142–148.
35. Scott, R. A., Lagou, V., Welch, R. P., *et al.* (2012). Large-scale association analyses identify new loci influencing glycemic traits and provide insight into the underlying biological pathways. *Nature Genetics*, **44**(9), 991–1005.
36. Soranzo, N., Sanna, S., Wheeler, E., *et al.* (2010). Common variants at 10 genomic loci influence hemoglobin a1c levels via glycemic and nonglycemic pathways. *Diabetes*, **59**(12), 3229–3239.
37. Strawbridge, R. J., Dupuis, J., Prokopenko, I., *et al.* (2011). Genome-wide association identifies nine common variants associated with fasting proinsulin levels and provides new insights into the pathophysiology of type 2 diabetes. *Diabetes*, **60**(10), 2624–2634.
38. Su, Z., Marchini, J., and Donnelly, P. (2011). Hapgen2: simulation of multiple disease snps. *Bioinformatics*, **27**(16), 2304–2305.
39. Sul, J. H., Han, B., He, D., and Eskin, E. (2011). An optimal weighted aggregated association test for identification of rare variants involved in common diseases. *Genetics*, **188**(1), 181–188.
40. Sul, J. H., Raj, T., de Jong, S., *et al.* (2015). Accurate and fast multiple-testing correction in eqtl studies. *The American Journal of Human Genetics*, **96**(6), 857–68.
41. van de Bunt, M., Fox, J. E. M., Dai, X., *et al.* (2015). Transcript expression data from human islets links regulatory signals from genome-wide association studies for type 2 diabetes and glycemic traits to their downstream effectors. *PLoS Genetics*, **11**(12), e1005694.
42. Visscher, P. M., Brown, M. A., McCarthy, M. I., and Yang, J. (2012). Five years of gwas discovery. *The American Journal of Human Genetics*, **90**(1), 7–24.

43. Wallace, C., Rotival, M., Cooper, J. D., *et al.* (2012). Statistical colocalization of monocyte gene expression and genetic risk variants for type 1 diabetes. *Human molecular genetics*, page dds098.
44. Wellcome Trust Case Control Consortium (2007). Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*, **447**(7145), 661–78.
45. Welter, D., MacArthur, J., Morales, J., *et al.* (2014). The nhgri gwas catalog, a curated resource of snp-trait associations. *Nucleic acids research*, **42**(D1), D1001–D1006.
46. Zaitlen, N., Paaniuc, B., Gur, T., Ziv, E., and Halperin, E. (2010). Leveraging genetic variability across populations for the identification of causal variants. *The American Journal of Human Genetics*, **86**(1), 23–33.