

## Matrix Ash

Sarah Urbut <sup>1,2</sup>, Gao Wang <sup>1</sup>, Matthew Stephens <sup>1,3,†</sup>, with the GTEX Consortium<sup>¶</sup>

**1** Department of Human Genetics/ University of Chicago, Chicago, IL USA

**2** Pritzker School of Medicine/Growth and Development Training Program/University of Chicago, Chicago, IL USA

**3** Department of Statistics/ University of Chicago, Chicago, IL USA

## Abstract

### Author Summary

Variation in gene expression is an important mechanism underlying susceptibility to complex disease. The simultaneous genome-wide assay of gene expression and genetic variation allows the mapping of the genetic factors that underpin individual differences in quantitative levels of expression (expression QTLs; eQTLs). By analyzing these effects across multiple tissues, we exploit the information that the effect of the gene-snp pair in one tissue can provide about its effect in alternative tissues. Furthermore, quantifying the effect size as opposed to simply calling QTLs present or absent reveal many patterns of sharing of effects among tissues which differ in both sign and magnitude. We provide a novel framework for estimating effect sizes across multiple subgroups, considering the evidence contained in all subgroups jointly, which provides a powerful and detailed insight into quantitative heterogeneity present in the genome.

# Introduction

Variation in gene expression is an important mechanism underlying susceptibility to complex disease.

The simultaneous genome-wide assay of gene expression and genetic variation allows the mapping of the genetic factors that underpin individual differences in quantitative levels of expression (expression QTLs; eQTLs). The availability of this information provides immediate insight into a biological basis for disease associations identified through genome-wide association (GWA) studies, and can help to identify networks of genes involved in disease pathogenesis ([1,2]). Available methods are limited not only in their ability to *jointly analyze data on all tissues* to maximize power, but also in simultaneously *allowing for both qualitative and quantitative differences among eQTLs* present in each tissue.

Initial approaches to quantify the effect of a particular SNP on gene expression considered only one tissue at a time, and ignored the effect of the SNP on gene expression in other tissues. This fails to exploit the power of shared genetic variation in effects on expression - i.e. the information that the effect of the gene-snp pair in one tissue can provide about the effect in another- and limits our understanding of multiple-tissue phenotypes. Furthermore, even past attempts at quantifying heterogeneity of eQTLs using the data across tissues jointly were limited in both the number of tissues considered, and also the level of heterogeneity considered. Qualitative heterogeneity refers to calling a snp ‘active’ or ‘inactive’ in a given tissue. For example, previous work has referred to the setting in which the gene-snp pair is active in all tissues as ‘shared’ and active in only one as ‘tissue-specific’ ([3,4]). However, a QTL may be ‘active’ in all or many tissues and with varying magnitude or sign; we refer to this as quantitative heterogeneity. Indeed, our initial motivation came from our analysis of GTEx pilot data, in which we see evidence that many (50%) QTLs are shared across all nine tissues. In this context, a QTL is called based on whether it demonstrates significant posterior probability of being active in a particular tissue. Applying our hierarchical model to the dataset from Dimas *et al* [5] with 3 tissues, we found just 8% of eQTLs are specific to a single tissue, with an estimated 88% of eQTLs being common to fibroblasts, LCL cells and T-Cells [3]. Not all eQTLs are shared by all tissues; some tissues may share eQTLs more than others. To allow for this, our previous hierarchical model attempted to infer the extent of such sharing by estimating the proportion of eQTL which were shared in various ‘configurations’ or patterns of binary activity in which a QTL was ‘called’ or ‘absent’ in each tissue. However, these binary configurations are still *limited in their ability to capture continuous variation* in levels of activity among tissues in which the SNP is considered active. In fact, as the number of tissues considered increases, perhaps the more interesting and biologically relevant question becomes one of quantitative heterogeneity - that is, how do the patterns of effect vary across tissues in which the SNP is called ‘active’. Quantifying the effect sizes of the gene-snp pair across tissues considering the evidence contained in all tissues jointly thus reveals new patterns of activity across tissues, which differ in their relationship in sign and magnitude within and between tissues. Thus effects can be ‘shared’ but not ‘consistent’ across tissues. The structure of this paper is as follows: we will describe our approach in brief for modeling and estimating these effect sizes across tissues, demonstrate the utility of such an approach on simulated data, and apply our method to data from the GTEx dataset, version 6.0, where we analyze effects size for 16,069 genes across 44 tissues.

## 0.1 Methods Overview: Approach

We aim to learn about patterns of sharing across tissues within a SNP and among SNPs, which join to help us better understand the global and snp-specific patterns of effects of genetics on gene expression. This allows us to make comparisons among tissues in which the QTL is called active, and among gene-snp pairs with a similar degree of activity in a given tissue. Thus as an additional level of combining information, we assume that each eQTL may follow a particular pattern of activity characterized by its effects across tissues. Within these groups, the tissues exhibit characteristic patterns of sharing, which can be captured by considering the covariance structure of the genetic effects among tissues.

This lends itself to a mixture model, in which we assume all the gene-snp pairs arise from a mixture of a finite number of multivariate normal (MVN) distribution, each characterised by the covariance matrix from which the vector of effects is thought to arise. For each of  $J$  gene-snp pairs, we observe an  $R$  dimensional vector of standardized effect sizes  $\hat{b}$  and their standard error and assume that these effects descend from some true effect size  $b$ .

$$b_j | \pi, U \sim \sum_{k,l} \pi_{k,l} N_R(0, \omega_l U_k) \quad (1)$$

Here, the covariance matrix  $U_k$  captures the particular patterns of sharing - variation in effect sizes within and between tissues, while  $\omega_l$  determines the scale of each pattern - the magnitude of the effect size. Thus we recognize that while two eQTL may obey a similar pattern or shape, the absolute scale may vary. For example, two eQTL may both have strong correlation between tissues 1 and 2 with consistently larger effects in 2, but the absolute size of the effects may vary between SNPs.

**Figure: Colorful patterns, caption indicates different patterns of activity across tissues in groups 1-3, same pattern or shape for snps of group 4 but tendency towards small effects**

## 0.2 Previous Approach

Previous work from our lab considered the idea of configuration - i.e., that a tissue was simply 'active' or 'inactive' in a particular tissues - and thus for  $R$  tissues, there were  $2^R$  possible configurations, which becomes computationally infeasible as  $R$  grows.

Furthermore, this considered only the idea that the variance in effect sizes between two tissues was the same across tissues thought to be active and the covariances were also the same among tissues thought to be active in a given 'configuration', and thus failed to incorporate the much richer covariance structure between tissues. For example, many gene SNP pairs might follow a pattern in which it is common to be 'active' across all tissues, but some QTL may have consistently larger effects in liver, lung and thyroid while other QTL may possess consistently larger effects in brain tissues and still another class of gene-snp pairs may show consistently quantitatively specific activity in Whole Blood but non-trivial effects in other tissues.

As a critical innovation on our previous method (1001[3, 4]) the covariance matrices used in our method contain distinct diagonal and off-diagonal elements which reflect data-specific patterns of variation within and covariance between subgroups (tissues). This captures the variation in effect sizes within and between subgroups better than restricting effects to simply 'shared' or 'unshared' between subgroups.

Because we can't know the 'true covariance matrix' for each gene-snp pair, we aim to assemble a list which sufficiently captures the various patterns, and then 'learn' the relative proportions of each pattern of sharing from the data. One can now model each

vector of effect sizes  $\mathbf{b}$  each as arising from a mixture that captures all the covariance patterns.

The primary novelty of this approach is *to estimate this multivariate posterior distribution on the effect size in a data-sensitive way* - i.e., using the mixture model to capture information about the covariance structure among subgroups (here, tissues). We deem this model hierarchical because these prevailing patterns of activity are learned from the larger dataset, and influence our inference about a given gene-snp pair. Thus we might identify a situation in which it is common to have large effects in certain tissues and not others, and thus if a given observed gene-snp pair demonstrates a small effect in one of the 'off issues', we might be inclined to conclude that it is indeed a member of this particular class and shrink the small effect in this tissue accordingly without adjusting the more 'active tissues' similarly. However, if we observe the same small effect in a setting in which 'similar tissues' have large effects, we might 'shrink' this effect size less, due to our high prior belief in the SNP's effectiveness garnered from adjacent tissues. Thus we deem this method 'Adaptive Shrinkage' because the appropriate amount of shrinkage is learned from the data. Critically, our method is dually adaptive, in the sense that we learn the relative abundance of effect sizes and directions from the overarching data set: thus observed effects are nudged towards prevailing patterns and sizes, according to the learned proportions of each.

Because our prior belief in consistency is strong in this particular dataset, we identify many more 'significant associations' in settings where perhaps the observed univariate statistic in one tissue is small but otherwise large in additional tissues, nudging these effects towards something more consistent. This is in contrast to a univariate shrinkage approach, in which all effects of the same size would be 'shrunk' equivalently, due to lack of information garnered from adjacent tissues.

In fact, shrinkage towards  $\mathbf{0}$  of small effects is a result, not a necessity - since the majority of the prior weight is on small  $\omega$  components which emphasize components with small prior variance of the effect size  $b$ , many of the modest observed effects will be smoothed or shrunk towards the prior mean,  $\mathbf{0}$ . An additional novelty is that in learning something about the effect size in each tissue for a given gene-snp pair, we can make statements about the degree of heterogeneity - that is the proportion of the time we expect a SNP to have effects of different sign. We will be confident in our ability to identify the direction of the effect for A SNP with a large effect and relatively high precision, and thus we can use an estimate of the posterior mean in each component and the proportion to quantify the distribution of gene SNP pairs who have effects of opposite direction (or lack convincing evidence of effects in a consistent direction across tissues).

## Materials and Methods

Let  $\mathbf{b}_j$  represents the genetic effect of SNP-gene pair  $j$  across  $R = 44$  tissues.

We assume the following mixture prior for the  $R$  dimensional vector of true effects,

$$\mathbf{b}_j | \pi, \mathbf{U}, \omega \sim \sum_{\mathbf{k}, l} \pi_{\mathbf{k}, l} N_{\mathbf{R}}(\mathbf{x}; \mathbf{0}, \omega_l \mathbf{U}_{\mathbf{k}}) \quad (2)$$

Where  $N_{\mathbf{R}}(\cdot; \mathbf{0}, \omega_l \mathbf{U}_{\mathbf{k}})$  denotes the density of a normal distribution with mean  $\mathbf{0}$  and variance  $\omega_l \mathbf{U}_{\mathbf{k}}$ .

Each component of the mixture distribution is characterized by these prior covariance matrices,  $\mathbf{U}_{\mathbf{k}}$  which capture the pattern of effects across tissues. Critically, this prior distribution is the same for all  $J$  - hence the hierarchical incorporation of shared information.

### 0.3 Covariance Matrices

For a given  $\omega_l$ , we specify 4 ‘types’ of  $R \times R$  prior covariance matrices  $\mathbf{U}_{k,l}$ .

1.  $\mathbf{U}_{k=1,l} = \omega_l \mathbf{I}_R$
2.  $\mathbf{U}_{k=2,l} = \omega_l \mathbf{X}_z$  The (naively) estimated tissue covariance matrix as estimated from the column-centered  $J \times R$  matrix of  $Z$  statistics,  $Z_{center}$ :  $\frac{1}{J} Z_{center}^t Z_{center}$
3.  $\mathbf{U}_{k=3,l} = \omega_l \frac{1}{J} \mathbf{V}_{1..p} d_{1..p}^2 \mathbf{V}_{1..p}^t$  is the rank  $p$  eigenvector approximation of the tissue covariance matrices, i.e., the sum of the first  $p$  eigenvector approximations, where  $\mathbf{V}_{1..p}$  represent the eigenvectors of the covariance matrix of tissues and  $d_{1..p}$  are the first  $p$  eigenvalues.
4.  $\mathbf{U}_{k=4:4+Q-1,l} = \frac{1}{J} ((\Lambda \mathbf{F})^t \Lambda \mathbf{F})_q$  corresponding to the  $q_{th}$  sparse factor representation of the tissue covariance matrix
5.  $\mathbf{U}_{k=4+Q,l} = \frac{1}{J} (\Lambda \mathbf{F})^t \Lambda \mathbf{F}$  is the sparse factor representation of the tissue covariance matrix, estimated using all  $q$  factors.
6.  $\mathbf{U}_{k=5+Q:R+4+Q,l} = \frac{1}{J} ([100..]'^t [100...])$
7.  $\mathbf{U}_{k=R+5+Q,l} = \frac{1}{J} ([111...]'^t [111...])$
8.  $[1000...]$  or  $[111...]$  represent configurations such that given membership,  $\mathbf{b}_j$  arise from the same prior variance.

### 0.4 Deconvolution

To retrieve a ‘denoised’ or ‘deconvoluted’ estimate of the non-single rank dimensional reduction matrices, we then perform deconvolution after initializing the EM algorithm with the matrices specified in (2), (3) and (5). The final results of this iterative procedure preserves the rank of the initialization matrix, and allows us to use the ‘true’ effect at each component component  $\mathbf{b}_j$  as missing data in deconvoluting the prior covariance matrices. In brief, this algorithm works by treating not only the component identity but also the true effect  $\mathbf{b}_j$  as unobserved data, and maximizing the likelihood over the expectation of the complete data likelihood, considering the values  $\mathbf{b}_j$  as extra missing data (in addition to the indicator variables  $q_{ij}$ ) (Bovy et al, 2014).

## 0.5 Likelihood

By maximum likelihood in each tissue separately, we can easily obtain the observed estimates of the standardized genotype effect sizes,  $\hat{\mathbf{b}}_j$ , and their observed squared standard errors recorded on the diagonal of an  $R \times R$  matrix noted  $\hat{V}_j = (\hat{\mathbf{b}}_j)$ . We assume that the matrix of standard errors of  $\hat{\mathbf{b}}_j$ ,  $V_j$  as approximated by  $\hat{V}_j$  is diagonal and that  $\hat{V}_j$  is an accurate point estimate for the standard error and that these standard errors are independent between tissues.

If we now view  $\hat{\mathbf{b}}_j$  and  $\hat{V}_j$  as *observed data*, we can write a new “likelihood” using only the sufficient statistics,  $\hat{\mathbf{b}}_j$  and  $\hat{V}_j$ :

$$\hat{\mathbf{b}}_j | \mathbf{b}_j \sim_R (\mathbf{x}; \mathbf{b}_j, \hat{V}_j) \quad (3)$$

## 0.6 Posterior Quantities

We aim to report posterior quantities for a given gene-snp pair  $\mathbf{j}$ . We know that for a single multivariate *Normal* the posterior on  $\mathbf{b} | U$  is simply:

$$\mathbf{b} | \hat{\mathbf{b}} \sim N_R(\tilde{\boldsymbol{\mu}}, \tilde{U})$$

where:

- $\tilde{\boldsymbol{\mu}} = \tilde{U}(\hat{V}^{-1}\hat{\mathbf{b}})$
- $\tilde{U} = (U^{-1} + \hat{V}^{-1})^{-1}$ .

Let us concatenate the list of all  $K \times L$  combinations of prior covariance matrices  $U_k$  and their scaling parameters  $\omega_l$  into a  $K \times L$  list and assign this length  $K$  for simplicity of notation.

Now each  $U_k$  imparts information about both *scale* and *direction*. Furthermore, a mixture-multivariate normal prior and a normal likelihood yields a mixture multivariate posterior, where the final posterior distribution is simply a weighted combination of multivariate normal distributions, where for each gene-snp pair  $\mathbf{j}$  is now characterized by it's posterior mean  $\tilde{\boldsymbol{\mu}}_{jk}$  and covariance  $\tilde{U}_{jk} = (U_k^{-1} + \hat{V}_j^{-1})^{-1}$ .

$$\mathbf{b}_j | \hat{\mathbf{b}}_j, \hat{V}_j, \hat{\boldsymbol{\pi}} \sim \sum_k^K \tilde{\pi}_{jk} N_R(\mathbf{x}; \tilde{\boldsymbol{\mu}}_{jk}, \tilde{U}_{jk}) \quad (5)$$

Where again,  $N_R(\cdot; \mathbf{0}, \omega_l U_k)$  denotes the density of a normal distribution with mean  $\tilde{\boldsymbol{\mu}}_k$  and variance  $\tilde{U}_k$  and the posterior mixture weight  $\tilde{\pi}_k$  is simply

$$\tilde{\pi}_{jk} = \frac{p(\hat{\mathbf{b}}_j | \hat{V}_j, z_j = k) \hat{\pi}_k}{\sum_{k=1}^K p(\hat{\mathbf{b}}_j | \hat{V}_j, z_j = k) \hat{\pi}_k} \quad (6)$$

Where  $z_j = k$  is the latent variable indicator of the component identity and each  $\hat{\pi}_k$  represents the Maximum Likelihood Estimate of the prior mixture weights assigned to each component.

## 0.7 Reported Quantities

For every gene-snp pair 'j', we aim to report the effect size as the posterior mean, defined as:

$$E(\mathbf{b}_j | \hat{\mathbf{b}}_j, \hat{V}_j, \hat{\boldsymbol{\pi}}) = \sum_k^K \tilde{\pi}_k \tilde{\boldsymbol{\mu}}_k \quad (7)$$

And the local false sign rate, or posterior probability of incorrectly identifying the sign of the effect for a given tissues 'r' as :

$$P(b_{jr}) = 1 - \max_k [\sum_k p(b_{j,r} > 0 | \hat{\mathbf{b}}_j, \hat{V}_j, z_j = k) \tilde{\pi}_k, \sum_k p(b_{j,r} < 0 | \hat{\mathbf{b}}_j, \hat{V}_j, z_j = k) \tilde{\pi}_{jk}] \quad (8)$$

# 1 Results

## 1.1 Demonstrating Features of the Method

To get a sense of the accuracy of our novel approach to estimating multivariate effects, we simulated two types of data.

In the first set, in which we expect our method to be superior to both univariate methods and methods in which the configuration approach is utilized, we simulate 50,000 gene-snp pairs, with only 400 representing true signal. This represents roughly 500 genes with 100 snps in cis, 80% of which contain one active QTL. Thus naturally, if the gene contains such a QTL, it is the same QTL among all tissues in which the tissue is active. This puts a dual burden on both features of the method: The small number of true associations present in these simulations tests whether the method accurately encourages small observed effects toward zero while preserving the true signal when it exists. Furthermore, the multivariate nature of these simulated effects when they exist tests the ability of the method to accurately infer patterns of sharing from the dataset. These true effects are thus simulated from the ‘learned’ covariance matrices representing  $U_k$  2-9, and thus aim to emulate the patterns of sharing present in real biological data. We compare with univariate ‘shrinkage’ method Ash (Stephens et al, unpublished) as well as the eqtlBMA-lite (Flutre et al, 2013) which uses the singleton and fully consistent (i.e., active in only one tissue, or active with the same effect size in all tissues) configurations to estimate these effects jointly. We call this the ‘sharing’ (S) scenario.

One might expect that our method would prove superior only in the setting in which true effects are shared among all tissues, and thus fail in the setting of tissue specificity. Thus, building on the situation above, we add a simulation in which 35% of the true effects are active in only one tissue, according to 5 different patterns of tissue specificity. We call this the ‘tissue-specific’ (TS) scenario.

We see that in terms of both power and accuracy, Matrix Ash is superior in both setting to both univariate methods, and to joint analysis under the setting of orthogonal configurations.

$$RMSE \sqrt{(\sum_{jr} (b_{jr} - E(b_{jr}|Data))^2)} \quad (9)$$

**Table 1.** Accuracy Comparison: RMSE

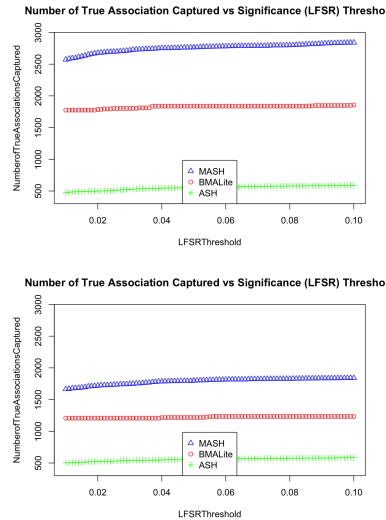
Inference Method	MASH	ASH	eqtlBMA-lite
RMSE <sub>S</sub>	0.010	0.030	0.047
RMSE <sub>TS</sub>	0.008	0.025	0.043
cor.with.truth <sub>S</sub>	0.99	0.94	0.84
cor.with.truth <sub>TS</sub>	.99	0.94	0.82

**Table 2. Accuracy Analysis** Here, we compare the ability of matrix ash (‘MASH’) to capture the true effect size estimates. We compare with univariate-shrinkage method ‘ASH’ and configuration-specific joint approach ‘eqtl-BMA-lite’. We report the Root Mean Squared Error (RMSE) and the correlation with the truth.

To demonstrate the ability of Matrix Ash to powerfully capture these accurately estimated effect sizes, we compare the proportion of true associations called significant at a given significance threshold among the three methods. Indeed, Matrix Ash proves superior to both methods under each condition (i.e., sharing or tissue-specific).

**Figure: Power vs Accuracy Figure:Shared simulation**





In introducing a method to quantify the heterogeneity of effect sizes, we have developed a ‘heterogeneity index’ which attempts to capture the heterogeneity in magnitude among tissues in which the gene-SNP pair is active. For each gene-snp pair  $\mathbf{j}$ , we normalize its vector of effects

$$RMSE_{HI} \sqrt{(\sum (true_{HI} - estimated_{HI})^2)} \quad (10)$$

**Table 3.** Accuracy Comparison: RMSE

Inference Method	MASH	ASH	eqtlBMA-lite
$HI_S$	39.38	40.87	39.78
$HI_{TS}$	39.98	40.77	39.51

## 1.2 Adaptive Shrinkage: The Multivariate Approach

To demonstrate the utility of shrinking effect size estimates jointly, we consider the estimated effect sizes vs their observed input summary statistics using our joint (Matrix Ash) and comparing to a univariate shrinkage method (Ash). On simulated data, we can also then plot the estimated effect sizes against the true values, again comparing among methods. Here, we show the results under the setting of tissue specificity, to analyze the behavior of eQTL of each class.

**Figure: Insert Simulations from TSpecific Conditions here**

In both Matrix Ash and univariate methods, values with large standard errors will be shrunk more harshly (Ash Paper, Stephens et al). Comparing estimated Z statistics (i.e.,  $E(Z_{jr}|Data)$ ) vs the observed ‘raw’ input values (i.e.,  $\hat{b}_{jr}$ ) allows us to understand the behavior of multivariate vs univariate methods once the standard error has been considered. In this simulated data, where there is an abundance of small effects, both univariate and multivariate methods tend to shrink small observed values towards prior mean at  $\mathbf{0}$  as their likelihoods will be maximized by component with small  $\omega$ . Critically, considering observed Z statistics of the same size, Matrix Ash does not shrink all small values are shrunk to the same extent, due to the power of joint analysis to consider the

effects across tissues in inferring the final vector of effect sizes. Thus the method is dually ‘adaptive’ by considering the abundance of both effect sizes and shapes in the overall data-set. Here, acknowledging consistency, small effects in one tissue will be augmented in the presence of larger effects other tissues, resulting in dramatic power increases.

Furthermore, when we plot the estimated effects against truth and segregate these effects by class (e.g., active and shared, active and tissue-specific, or null), we see that the correlation among the true and estimated effect ( $E(Z|D)$ ) sizes is much tighter using our multivariate approach. Similarly, truly null effects are shrunk more tightly, due to the fact that in the presence of consistency, small effects across subgroups will lead us to have a high prior belief that an additional small observed effect in that eQTL is also likely to be close to 0. Importantly, tissue-specific QTLs are still captured using our joint approach, demonstrating that if tissue-specific patterns exist in the data, our prior belief will capture this phenomenon and accordingly our posterior estimates will reflect the underlying tissue-specific nature at a given tissue-specific SNP.

Together, these results demonstrate the tremendous power increase of using a multivariate method and the accuracy of estimating patterns of sharing from the data rather than imposing forced configurations which fail to capture the heterogeneity of effect sizes among tissues.

## 2 Real Data

Here, we report the analysis on the top SNP for each of 16,069 genes, where the ‘top’ snp is defined as the SNP with the largest observed univariate Z-statistic in absolute value across tissues. As described above and demonstrated in simulations, in the setting of an abundance of small effects in data set, Matrix Ash tends to shrink small observed values towards the prior mean at **0**. It should be noted that this is a result specific to a particular data set, and in that sense ‘adaptive’ - indeed, if small effects were rare and large effects abundant, such shrinkage would not occur.

But perhaps more importantly, the striking increase in power when compared to univariate methods is noted. There are a total of 44 tissues x 16,069 gene-snp pair associations considered, or 707,036 total tissue-level effect size coefficients. At an  $lfsr$  threshold of 0.05, we identify 393,414 significant snp-gene-tissue effects ( $b_{jr}$ ). Using estimates shrunk according to a univariate approach (again, Ash), we identify only 91,755, meaning that using univariate methods we would be confident our ability to identify the sign in only 13% of cases, while using our joint procedure for estimating effects, we would confidently argue the SNP has a non-zero effect for a gene in a particular tissue over half (55%) of the time. As described, this tremendous increase in power arises from the fact that in the presence of a data set possessing consistency, as learned by the hierarchical model, small effects in the presence of a gene containing large effects in alternative tissues will be augmented to reflect such consistency, thus increasing our confidence in its size and direction. While the number of associations capture is slightly greater using the BMA lite approach, we note that the likelihood of the data set under this model is much much worse ( $-1298672$  vs  $-1267997.5$ , see supplementary data ‘Testing and Training’ procedure). Indeed, eQTLBMA would put the vast majority of the prior weight on the fully ‘consistent’ configuration ( $\hat{\pi}$  figure, as SNPs demonstrating activity across all tissues, regardless of how heterogeneous among subgroups, are forced into this configuration. Simulations above demonstrate the lack of accuracy arising from such an approach.

### Figure: Real Data Scatterplot

**Table 4.** Power Comparison

Metric	LFSR <sub>MatrixAsh</sub>	LFSR <sub>ASH</sub>	eQTL-BMALite
Significant $\mathbf{b}_{jr} \leq 0.05$	393414	91,755	401552

**Table 5. Power** Restricting our analysis to thresholding by local false sign rates, we can quantify the number of associations we identify at a given local false sign rate threshold using the original summary statistics and posterior means computed using multivariate Matrix Ash and Univariate Ash. We can see that Matrix Ash calls nearly twice and 4 times as many associations significant when compared to univariate approach, and is comparable to less-accurate joint approach

**Figure:  $\hat{\pi}$  barplot in MASH vs BMALite**

To further contrast our approach with existing joint methods on this data-set, consider a two-tissue example, in which a configuration type approach recognizes only patterns constrained to lie along the x and y axis or along the x-y line. Matrix ash allows for patterns which show consistently larger effects in one tissue over another, with varying amounts of correlation among tissues. In these example from real data, we can see that while eQTL the green, blue and yellow class appear consistently active in both tissues plotted (Brain and Muscle), the blue and yellow effects have consistently larger effects in Skeletal Muscle than brain, while SNPs of the green class show the reverse. Similarly, eQTL of the yellow close show only loose correlation among the pair of tissues, while eQTL of the green class show strong prediction of activity in one from activity in the other. Similarly, we can see that eQTL of the pink class tend to show tissue-specificity in Whole Blood relative to testis, while eQTL of the black class show tissue specificity in testis relative to whole blood.

**Figure: Colored Scatterplots of 2 tissues, contrasting effects between tissues, color coded by responsibility**

## 2.1 A qualitative description of heterogeneity in the GTEX data

Indeed, from the prior weight assigned to the ‘learned matrices’ above coupled with the simulation results in the previous sections, we can see that Matrix Ash is able to accurately parse shared configurations, thus resolving the relationship among tissues in which the QTL is active. *See Supplement for Heatmap of all 8 learned matrices.* For instance, learned matrix  $U_k = 3$  captures gene-snp pairs with large, correlated effects in brain, and is the most prevalent pattern of sharing in the larger data set, as reflected by it’s prior weight summed across effect size (see pi barplot). Matrix  $U_k = 2$  captures SNPs with small effects in brain and larger effects in thyroid and transformed cell-types (e.g., fibroblasts, lymphocytes). Several of the lower rank matrices whose patterns receive high prior weighting (e.g.,  $U_k = 4, 5, 8$  and  $9$ ) show somewhat tissue specific (i.e., high prior variance in only one tissue-type) effects in testes and whole blood, consistent with our conclusions that whole blood and testes indeed demonstrate an abundance of tissue-specific gene-snp pairs. Here we examine several example gene-snp pairs with a high posterior probability of arising from the covariance patterns captured by our model. We deem this posterior probability of arising from a particular pattern as a high ‘loading’ or ‘responsibility.’

## 2.2 Examples of strong loading on UK3, UK9 and eQTL-BMA lite

**Figure: High loading on UK3: Captures correlation in sign, quantitative heterogeneity in magnitude along diagonal emphasizes the utility of continuous approach**

In this particular example, strong, shared effects in brain tissues match an underlying pattern of shared effects present in the larger data set and thus allows this gene-snp pair to find its true match. Brain effect sizes thus borrow strength from one another, and accordingly, the posterior estimates tend to nudge the brains towards a consistent, shared effect. Similarly, an overall tendency towards consistency in sign in the larger data set, as captured by the hierarchical model and reflected in the positive correlation in sign among all tissues, tends to ‘flip’ erratic off directions towards the prevailing positive direction. Heterogeneity in magnitude among the other tissues is reflected in the variety of banding intensity along the diagonal.

**Figure: UK 9: Quantitative specificity in magnitude in Testes/Whole Blood**

In this example, though the particular pattern featured ( $U_k = 9$ ) captures correlation in sign among all tissues, significant quantitative heterogeneity is again reflected in the intensity of the banding along the diagonal, in this dramatically dichotomous between testes/whole blood and all other tissues. Here, we introduce the idea of quantitative specificity - e.g., that a SNP can be modestly ‘active’ in all tissues though to dramatically different degrees. here, though this matrix was learned (and not forced, as in eQTLBMA-lite) from the data, the pattern of quantitative tissue specificity in testes and whole blood is evident. Again, erratic, off-directions are flipped in sign. We refer to this as quantitative specificity, because the effects are quantitatively unique to particular tissues - e.g., significantly larger in magnitude in testis than all other tissues - and yet considered non-zero in all tissues. This is in contrast to qualitative specificity, described below, in which we would conclude that the QTL is active in only one tissue.

**Figure: Qualitative specificity in Testis Example: High loading on eqtlbmalite config mat**

Lastly, the inclusion of the eQTLBMA lite configurations (in which the SNP has a non-zero effect in only one tissue) coupled with the learned patterns of tissue specificity evident in matrices  $U_k : 5 - 9$  serve to allow the preservation of qualitatively specific effects. Here, we show a gene-snp pair demonstrating high loading on one of the eqtlbma-lite configuration matrices - indeed, we reject the significance of the effect size estimates in all tissues but testes, a pattern consistent with the presence of tissue-specificity described below. Together, these results cement the resolution afforded by methods which can distinguish among tissues in which a QTL is called active, beyond reducing genetic effects to binary ‘on’ or ‘off’ conclusions.

## 2.3 Tissue Specificity

One of the criticisms of a joint approach might be its loss of tissue-specificity. That is, by considering effects across subgroups in estimating the effect size, one might lose sight of tissue-specific activity when it exists. Here, we demonstrate our ability to recognize such specificity both quantitatively, as described above through learned patterns of sharing which specify consistently larger effects in one tissue over others, and

qualitatively through forced prior effect size mass on 0. For each tissue, we can ask how many gene-snp pairs meet a given significance threshold in that tissue alone.

### Figure: Number of QTL per Tissue Plot

Furthermore, tissue specific eQTL demonstrate the smoothing feature of this joint shrinkage approach: for gene SNP pairs which demonstrate strong effects in only one tissue, the weaker erratic tissue are shrunk towards the prior mean at 0, resulting in a tissue specific smoothing.

### Figure: Tissue Specific Smoothing Plot

## 2.4 Quantifying Heterogeneity

Armed with a vector of effect size estimates across 44 tissues,  $b_j$ , we can move beyond asking in how many tissues is a given gene-snp pair significant, and ask about the relationship in effect size and direction among tissues in which the gene-snp pair is active. From a biological standpoint, we might produce that effects of a different sign are rare. Similar to the heterogeneity index described in the simulation framework above which attempted to describe heterogeneity in magnitude, we can plot the number of tissues in which the sign is differ than the effect with maximum absolute value. Considering this results with and without the inclusion of the brain tissues, which appear to behave as a strongly correlated group, we observe several phenomenon. The majority of gene-snp pairs are consistent in sign (indeed, only about 20% of genes show two significant effects of a different sign when including brain, and even fewer (14.8%) when excluding brains) and removing brains from our analysis tends to push us towards consistency, suggesting that brain appears to behave as a large tissue-specific entity. After normalizing each gene-snp effect size coefficient  $b_{jr}$  by the effect size with the maximum value for the gene, we can also ask what proportion of these are positive. We come to similar conclusions with 83% and 87% with and without brain demonstrating positive normalized effects, respectively.

### Figure: Sign Heterogeneity Distribution Figure

Furthermore, we can now quantify the heterogeneity index in magnitude described in the simulation framework above, and ask, for each gene, in how many tissues is the effect greater than equal to a significant fraction, here 50% of the maximum effect. Again, homogenous genes will tend to be featured towards the right of the distribution with maximal value at 44. Intuitively, the majority of their effects across tissues are similar in magnitude, while heterogeneous genes will be featured towards the left of the distribution, with tissue-specific genes possessing a heterogeneity index of 1. Again, excluding brain form the analysis tends to nudge us towards a belief in consistency.

Taken together, these results suggest the presence of consistency in sign in our data set, and a bimodal distribution of heterogeneity in magnitude.

### Figure: Magnitude Heterogeneity Index Distribution

Attempting to understand which genes tend to behave the most homogeneously or heterogeneously, we can plot the value used to normalize each gene, e.g., the ‘maximum’ effect size across tissue of the gene, against the normalized values. We can see that if a large effect is present, it tends to be in the presence of homogenous effects across the board, while small normalizing effects tend to be in the presence of effects that are more variable in sign and magnitude. Furthermore, aggregating the gene-snp pairs at a given

**Table 6.** Heterogeneity Comparison

Data	All Tissues	No Brains
Consistent in Sign $E(b_{jrnorm} D) > 0$	0.833	0.880
$E(\text{Consistent Sign} \text{PosteriorMean} \text{LFSR}\leq 0.05)$	0.802	0.852
$E(\text{At least 50\% max value})$	0.354	0.449

**Table 7. Heterogeneity Analysis** After normalizing each gene-snp-effect size coefficient by the effect size with maximal value at that gene, we can ask how many of these gene-snp effect coefficients are positive. Similarly, at a given significance threshold, we can ask how many gene-pairs contain effects of different signs across tissues. At an arbitrary LFSR threshold of 0.05 for instance, we note that 80% of genes are homogenous in sign when all tissues are considered. Excluding brains from our analysis, this rises to 85%. To evaluate consistency in magnitude, we can ask how many gene-snp-tissue effects are greater than 50% of the maximal effect across tissues for the pair. Again, we see that excluding brains from our analysis tends to push this towards consistency.

heterogeneity index and classifying them by the effect used to normalize (e.g., the ‘max effect’) we can see that gene-snp pairs with greater Heterogeneity indices tend to have larger effects.

**Figure: Insert Plot of Normalized Effect vs Max value (e.g., biplot)**  
**Figure: Insert Median Max Effect by HI Index** **Figure: QTL chart by HI index: homogenous or heterogenous**

### 3 Testing and Training

In order to determine the optimal number and rank of the covariance matrices, we divide our data set into a training and test data set, each containing 8000 genes.

In the training set, we proceed as above: choosing the top SNP for each of the 8000 genes, creating a list of covariance matrices through deconvolution and grid selection of these top 'training gene-snp' pairs.

Then, within the training data, we similarly choose a random set of gene-snp pairs (restricting our analysis to genes contained in the training set. Again, we choose 20,000 random-gene snp pairs and use the EM algorithm to learn the mixture proportions  $\pi$  from this data set.

We then use the KxL vector of  $\pi$  from the training set to estimate the log likelihood of each data point in the test data set. If our model is 'overfit' to the training data set, than a larger number of covariance matrices may actually decrease the test log-likelihood.

I found that the K=1188 set of covariance matrices containing the Identity, the denoised empirical covariance matrix, rank 5 SFA approximation and rank 3 SVD approximation as well as 5 single-rank SFA factors and the 45 *eqtl.bma.lite* configurations maximized this likelihood.

### 4 Training and Testing Procedure: Estimating Hierarchical Weights

We wish to choose the model which best maximizes the probability of observing the data set.

Incomplete Data likelihood:

$$L(\pi; \hat{\mathbf{b}}) = \prod_{j=1}^J \sum_k^K \pi_k P(\hat{\mathbf{b}}_j | z_j = k) \quad (11)$$

- To estimate the hierarchical prior weights  $\pi_k$ : compute the likelihood at each each gene snp pair  $j$  by evaluating the probability of observing  $\hat{\mathbf{b}}_j$  given that we know the true  $\mathbf{b}_j$  arises from component  $k$
- Use the EM algorithm to estimate the optimal combination of weights: How often does this particular covariance matrix occur in the data?

We then use these weights to estimate the test set log likelihood.

### Discussion

### Supporting Information

### Acknowledgments

## References

1. Nicolae DL, Gamazon E, Zhang W, Duan S, Dolan ME, Cox NJ. Trait-Associated SNPs Are More Likely to Be eQTLs: Annotation to Enhance Discovery from GWAS. *PLoS Genetics*. 2010 Apr;6(4). Available from: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2848547/>.
2. Veyrieras JB, Kudaravalli S, Kim SY, Dermitzakis ET, Gilad Y, Stephens M, et al. High-Resolution Mapping of Expression-QTLs Yields Insight into Human Gene Regulation. *PLoS Genet*. 2008 Oct;4(10):e1000214. Available from: <http://dx.doi.org/10.1371/journal.pgen.1000214>.
3. Flutre T, Wen X, Pritchard J, Stephens M. A Statistical Framework for Joint eQTL Analysis in Multiple Tissues. *PLoS Genet*. 2013 May;9(5):e1003486. Available from: <http://dx.doi.org/10.1371/journal.pgen.1003486>.
4. Wen X, Stephens M. Bayesian methods for genetic association analysis with heterogeneous subgroups: From meta-analyses to gene-environment interactions. *The Annals of Applied Statistics*. 2014 Mar;8(1):176–203. ArXiv:1111.1210 [stat]. Available from: <http://arxiv.org/abs/1111.1210>.
5. Dimas AS, Deutsch S, Stranger BE, Montgomery SB, Borel C, Attar-Cohen H, et al. Common regulatory variation impacts gene expression in a cell type-dependent manner. *Science (New York, NY)*. 2009 Sep;325(5945):1246–1250.