**SUPPLEMENTAL METHODS**

*Additional simulations comparing the pooled and traditional association methods*

We ran additional simulations that model different distributions of read depth across individuals and across the genome. First, we simulated pooled reads by sampling the number of reads from a Poisson distribution with $\lambda$ = total number of reads (the Poisson distribution assumes that the mean and the variance of the number of reads are equal). For the traditional non-ASM method, we randomly distributed the number of reads that we sampled across individuals as we did in the previous simulations and then sampled allele/methylation status pairs with replacement for each read from the read's individual.

Since sequencing data are known to be over-dispersed, the number of reads is often modeled as being sampled from a negative binomial instead of Poisson distribution (Anders and Huber 2010). We therefore also sampled the number of reads from a negative binomial distribution. We fit negative binomial parameters to the numbers of reads covering CpGs by finding maximum likelihood values. After sampling a number of reads from the negative binomial distribution, we scaled the number of reads to $\frac{(Num.Reads\ at\ Position)*(Num.Reads\ in\ Simulation)}{Mean\ Num.Reads\ Across\ Positions}$ before sampling. For the traditional non-ASM method, we randomly distributed the post-scaling number of reads that we sampled across individuals as we did in the previous simulations and then sampled allele, methylation status pairs with replacement for each read from the read's individual.

In addition to evaluating the power of each method, we also compared its false positive rate to its true positive rate. We simulated false positives by creating a distribution of reads with the same MAF and MMF as our real distribution but 0.0 correlation between allele and methylation status and sampling reads from that distribution. We ran our true positive and false positive simulations at p-value

cutoffs 0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001, 0.00005, and 0.00001.  We did this for all of our

Fisher's exact test negative binomial simulations with 100 individuals and all combinations of numbers

of reads, effect sizes, MAFs, and MMFs.  We generated ROC curves to illustrate our results.


*Details of data processing*

We trimmed reads using Trim Galore! version 0.2.8

(http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/).  We used the default parameters

with the exception of --stringency 4, --quality 35, and --paired.  These parameter adjustments prevented

us from removing ends of reads where only a few bases overlap with the adapter, removed read ends

that were not especially high-quality, and forced Trim Galore! to account for our reads being paired-end.

We chose these parameters after trying multiple parameter settings on a subset of our data because

they enabled the most reads to map uniquely to the genome.  Next, we converted all HapMap Phase II

(Frazer et al. 2007) and 1000 Genomes Phase I Integrated Version 3 (The 1000 Genomes Project

Consortium 2010) single nucleotide polymorphisms (SNPs) with MAF > 0.04 in human genome version

hg19 (The International Human Genome Sequencing Consortium 2001) into N's in order to eliminate

sources of reference bias when mapping (Degner et al. 2009).  We mapped reads to the autosomes in

hg19 (The International Human Genome Sequencing Consortium 2001) using Bismark version 0.12.3

(Krueger and Andrews 2011) with Bowtie 2 version 2.2.3 (Langmead and Salzberg 2012).  Bismark

converts all Cs to Ts and all Gs to As before mapping, maps them to both a C-to-T and a G-to-A-

converted genome, and then converts the Ts and As back to their original bases (Krueger and Andrews

2011). We used the default parameters with the exception of, for mapping, --bowtie2 and, for extracting

methylation, --ignore_r2 7, -p, and --no_overlap so that we could remove incorrect methylation calls at

the 5' end of read 2 due to DNA repair (Supplemental Figure 27), account for our reads being paired-

end, and not double-count cytosines on both ends of the reads. In addition, when calling methylation statuses, we removed the 11 bases at the 3' end of read one and the 31 bases at the 3' end of read two because we noticed substantial methylation degradation towards the 3' end that seemed independent of read sequence (Supplemental Figure 27). We also used Bismark (Krueger and Andrews 2011) to map reads to the lambda phage genome (Leinonen et al. 2011) and to evaluate the observed methylation; since lambda phage is completely unmethylated, any observed methylation is due to failure in bisulfite treatment or sequencing errors.

We also filtered the reads in multiple ways. Bismark divides all cytosines into four categories: cytosines followed by guanines (CpGs), cytosines followed by non-guanines followed by guanines (CHGs), cytosines followed by at least two non-guanines (CHHs), and cytosines followed by Ns (CNs). For our analysis, we focused on CpGs. The CpGs with called methylation statuses should not contain most SNPs because we masked SNPs with Ns, so such CpGs would have become CNs. After running Bismark, we removed duplicates from the mapped reads using rmdup from Samtools version 0.1.19 (Li et al. 2009). We also removed reads that overlapped with regions in the ENCODE black list (The ENCODE Project Consortium 2012). Thus, we were left with reliable mapped reads and methylation calls.

After removing duplicates, we determined the allele of each SNP, insertion, and deletion from 1000 Genomes AFR in each read (The 1000 Genomes Project Consortium, 2010). We used only the 1000 Genomes AFR variants because calling variants from pooled data is a challenging problem (Nielsen et al. 2011; Li 2014), and this panel should contain most of the variants in these individuals (The 1000 Genomes Project Consortium, 2010); however, when genomic variant positions are not available, they can also be inferred from the reads using established methods like GATK (McKenna et al. 2010) or Bis-SNP (Liu et al. 2012). When identifying alleles from reads, we did not include any cytosine/thymine (C/T) SNPs or adenine/guanine (A/G) SNPs. During the bisulfite treatment, unmethylated cytosines are

converted into uracils (that become thymines during PCR) and, as a result, the guanines that complement them become adenines during PCR. Therefore, for C/T SNPs, we cannot distinguish between unmethylated cytosines and thymines from the original reads, and for A/G SNPs on the reverse strand, we cannot distinguish between adenines that complement unmethylated cytosines (that have become thymines) and adenines from the original reads.

*Estimating the fraction of each individual in the pool*

Although unnecessary for our method, we evaluated how well our pool represented each individual by using our reads to estimate the frequency of each individual in our pool. In order to do this, we first computed the number of reads covering each allele of each SNP from HapMap Phase II (Frazer et al. 2007). We then solved the constrained optimization problem

$$\mathrm{argmin}_f \frac{1}{2}(Xf - y)^2$$

subject to

$$f \geq 0$$

$$\sum_{j}^{60} f_j = 1,$$

where $y$ is the weighted vector of alternate allele frequencies in the pool, $f$ is the vector of individual frequencies in the pool, $f_j$ is the entry in the vector $f$ for individual $j$, and $X$ is a weighted (number of variants) x (number of individuals) matrix of genotypes (on a 0 to 1 scale, where 0 is homozygous reference allele, 0.5 is heterozygous, and 1 is homozygous alternate allele) of each individual for each SNP. We weighted $y$ and $X$ by multiplying them by the number of reads at the current SNP and then

4

dividing them by the sum of the numbers of reads across all SNPs; this allows SNPs with more reads to contribute more to the optimization. We solved the optimization problem using Matlab's lsqlin (Coleman and Li 1996) with initial individual frequencies of 1/60. We should note that this does require genotype information, which may not always be available for pooled samples, but the results from this are not necessary for our pooling method.

*Obtaining data for overlaps between mQTLs, eQTLs, dsQTLs, CTCF-binding-QTLs, and GWAS hits*

For eQTLs from Pickrell *et al*., we downloaded the final_eqtl_list, and final_sqtl_list files from eqtl.uchicago.edu/RNA_Seq_data/results (Pickrell et al. 2010). For eQTLs from the Geuvadis Consortium, a larger, more recent study, we downloaded the YRI89 files from ftp://ftp.ebi.ac.uk/pub/databases/microarray/data/experiment/GEUV/E-GEUV-1/analysis_results; for SNPs that were tested for eQTLs in Geuvadis YRI, we downloaded data from http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502 (Lappalainen et al. 2013). For dsQTLs, we downloaded the files from eqtl.uchicago.edu/dsQTL_data/QTLs; for SNPs that were tested for dsQTLs, we downloaded data from http://eqtl.uchicago.edu/dsQTL_data/GENOTYPES (Degner et al. 2012). We then used liftOver (Kent et al. 2002) to convert the SNP coordinates from hg18 to hg19 and finally combined short-range and long-range dsQTLs. For CTCF-binding-QTLs and SNPs that were tested for CTCF-binding-QTLs, we downloaded data from http://www.ebi.ac.uk/birney-srv/CTCF-QTL (Ding et al. 2014). For the GWAS data, we downloaded the GWAS Catalog (Welter et al. 2014) on January 14, 2014. For all QTL and GWAS datasets except for the CTCF-binding-QTLs, we used SNAP Proxy Search (Johnson et al. 2008) with the YRI population panel and the default distance limit to identify 1000 Genomes Pilot 1 YRI SNPs in perfect LD and $r^2 \geq 0.8$ LD with the SNPs in the dataset and used liftOver (Kent et al. 2002) to convert SNP coordinates from hg18 to hg19. For the CTCF-binding-QTLs, we used the same

procedure for finding SNPs in LD as was used in the other studies except that we used SNAP's CEU

population panel (Johnson et al. 2008) because this study was done in individuals with European

ancestry; we used this study even though it came from a different population because it is the only

existing CTCF-binding-QTL study.

*Obtaining mQTLs from the Zhang et al. data*

We compared our mQTLs to those in the Zhang *et al.* study because it is one of the two largest

CpG microarray studies of Yoruban LCLs (Zhang et al. 2014).  We obtained a list of filtered CpGs from the

authors, where the filtering included all of the metrics described in their paper.  We downloaded their

Supplemental Table 2, which has their YRI mQTLs.  We computed the p-value for the overlap between

CpGs with mQTLs in both studies using a hypergeometric test, where the background was all CpGs

tested in both studies.

*Obtaining mQTLs from the Banovich et al. data*

We also compared our mQTLs to those in the Banovich *et al.* study because it is the other of the

two largest CpG microarray studies of Yoruban LCLs (Banovich et al. 2014).  We obtained a list of filtered

CpGs from the authors, where the filtering included all of the metrics described in their paper.  We

downloaded their mQTLs from http://giladlab.uchicago.edu/data/meQTL_summary_table.txt (Banovich

et al. 2014).  We computed the p-value for the overlap between CpGs with mQTLs in both studies using

a hypergeometric test, where the background was all CpGs tested in both studies.  We also intersected

their CpGs with mQTLs with those found by Zhang et al. (Zhang et al. 2014) to compare the results of

these two earlier studies.

*Overlapping CpGs in our study with CpG islands and surrounding regions*

We determined the fraction of CpGs with mQTLs in our data-set that are also in CpG islands or shores. We did this by downloading CpG islands from http://rafalab.jhsph.edu/CGI/model-based-cpg-islands-hg19.txt (Irizarry et al. 2009; Wu et al. 2010). Because CpG shores are generally defined as the 2 kb surrounding CpG islands in each direction (Price et al. 2013), we extended each CpG island by 2 kb in each direction using BEDTools slopBed version 2.16.1 (Quinlan and Hall 2010). We then identified the number of CpGs with mQTLs that are in both the CpG islands and the extended CpG islands.

**SUPPLEMENTARY REFERENCES**

Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol* **11**: R106.

Coleman TF, Li Y. 1996. A Reflective Newton Method for Minimizing a Quadratic Function Subject to Bounds on Some of the Variables. *SIAM J Optim* **6**: 1040–1058.

Degner JF, Marioni JC, Pai AA, Pickrell JK, Nkadori E, Gilad Y, Pritchard JK. 2009. Effect of read-mapping biases on detecting allele-specific expression from RNA-sequencing data. *Bioinformatics* **25**: 3207–3212.

Degner JF, Pai AA, Pique-Regi R, Veyrieras J-B, Gaffney DJ, Pickrell JK, De Leon S, Michelini K, Lewellen N, Crawford GE, et al. 2012. DNase I sensitivity QTLs are a major determinant of human expression variation. *Nature* **482**: 390–394.

Ding Z, Ni Y, Timmer SW, Lee B-K, Battenhouse A, Louzada S, Yang F, Dunham I, Crawford GE, Lieb JD, et al. 2014. Quantitative Genetics of CTCF Binding Reveal Local Sequence Effects and Different Modes of X-Chromosome Association. *PLoS Genet* **10**: e1004798.

Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve LL, Gibbs RA, Belmont JW, Boudreau A, Hardenbol P, Leal SM, et al. 2007. A second generation human haplotype map of over 3.1 million SNPs. *Nature* **449**: 851–861.

Johnson AD, Handsaker RE, Pulit SL, Nizzari MM, O'Donnell CJ, de Bakker PIW. 2008. SNAP: a web-based tool for identification and annotation of proxy SNPs using HapMap. *Bioinformatics* **24**: 2938–2939.

Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler D. 2002. The Human Genome Browser at UCSC. *Genome Res* **12**: 996–1006.

Krueger F, Andrews SR. 2011. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics* **27**: 1571–1572.

Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**: 357–359.

Lappalainen T, Sammeth M, Friedländer MR, 't Hoen P a C, Monlong J, Rivas MA, Gonzàlez-Porta M, Kurbatova N, Griebel T, Ferreira PG, et al. 2013. Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* **501**: 506–511.

Leinonen R, Akhtar R, Birney E, Bower L, Cerdeno-Tárraga A, Cheng Y, Cleland I, Faruque N, Goodgame N, Gibson R, et al. 2011. The European Nucleotide Archive. *Nucleic Acids Res* **39**: D28–D31.

Li H. 2014. Toward better understanding of artifacts in variant calling from high-coverage samples. *Bioinformatics* **30**: 2843-2851.

Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079.

Liu Y, Siegmund KD, Laird PW, Berman BP. 2012. Bis-SNP: Combined DNA methylation and SNP calling for Bisulfite-seq data. *Genome Biol* **13**: R61.

Nielsen R, Paul JS, Albrechtsen A, Song YS. 2011. Genotype and SNP calling from next-generation sequencing data. *Nat Rev Genet* **12**: 443–451.

Pickrell JK, Marioni JC, Pai AA, Degner JF, Engelhardt BE, Nkadori E, Veyrieras J-B, Stephens M, Gilad Y, Pritchard JK. 2010. Understanding mechanisms underlying human gene expression variation with RNA sequencing. *Nature* **464**: 768–772.

The 1000 Genomes Project Consortium. 2010. A map of human genome variation from population-scale sequencing. Nature **467**: 1061-1073.

The ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**: 57–74.

The International Human Genome Sequencing Consortium. 2001. Initial sequencing and analysis of the human genome. *Nature* **409**: 860–921.

Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P, Manolio T, Hindorff L, et al. 2014. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res* **42**: D1001–D1006.

**Supplemental Figure 1: Pooled ASM vs. traditional non-ASM method, zero-variance simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads. In the simulations for this figure, the reads are randomly sampled with replacement from the pool/individuals, meaning that, for each read, we select an allele-methylation status combination from the distribution of combinations in our pool/for the individual. There are 100 individuals, and there are 0.1 minor allele and minor methylation status frequencies. The reads from the pool are randomly distributed across the individuals, so that each individual has approximately (but not exactly) the same number of reads.

**Supplemental Figure 2: Pooled ASM vs. traditional non-ASM method, Poisson distribution simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads. In the simulations for this figure, reads are sampled from a Poisson distribution, there are 100 individuals, and there are 0.1 minor allele and minor methylation status frequencies.

**Supplemental Figure 3: Pooled ASM vs. traditional non-ASM method, F-test p-value simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads. In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.1 minor allele and minor methylation status frequencies. p-values for mQTLs were calculated using the p-value from the F-test for the regression that predicts methylation status as a function of allele/genotype instead of Fisher's Exact Test.

**Supplemental Figure 4: Pooled ASM vs. traditional non-ASM method, correlation p-value simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads. In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.1 minor allele and minor methylation status frequencies. p-values for mQTLs were calculated using the asymptotic p-value from the Pearson correlation instead of Fisher's Exact Test.

**Supplemental Figure 5: Pooled ASM vs. traditional non-ASM method, MAF = 0.3 simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads. In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.3 minor allele and minor methylation status frequencies.

**Supplemental Figure 6: Pooled ASM vs. traditional non-ASM method, MAF = 0.5 simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads.  In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.5 minor allele and minor methylation status frequencies.

**Supplemental Figure 7: Pooled ASM vs. traditional non-ASM method, 25 individuals simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads.  In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 25 individuals, and there are 0.1 minor allele and minor methylation status frequencies.

**Supplemental Figure 8: Pooled ASM vs. traditional non-ASM method, 400 individuals simulations**

Plots showing the correlation of the allele and methylation status versus the fraction of simulations that identify the mQTL with p-value < 0.001 for 40, 160, and 640 reads.  In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 400 individuals, and there are 0.1 minor allele and minor methylation status frequencies.

**Supplemental Figure 9: ROCs for pooled ASM vs. traditional non-ASM method, MAF = 0.1 simulations**

ROC curves comparing the false positive versus true positive rates for simulations for 40, 160, and 640 reads and effect sizes 1.0 and 0.5.  In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.1 minor allele and minor methylation status frequencies.

**Supplemental Figure 10: ROCs for pooled ASM vs. traditional non-ASM method, F-test simulations**

ROC curves comparing the false positive versus true positive rates for simulations for 40, 160, and 640 reads and effect sizes 1.0 and 0.5.  In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.1 minor allele and minor methylation status frequencies.  p-values for mQTLs were calculated using the p-value from the F-test for the regression that predicts methylation status as a function of allele/genotype instead of Fisher's Exact Test.

**Supplemental Figure 11: ROCs for pooled ASM vs. traditional non-ASM method, MAF = 0.3 simulations**

ROC curves comparing the false positive versus true positive rates for simulations for 40, 160, and 640 reads and effect sizes 1.0 and 0.5. In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.3 minor allele and minor methylation status frequencies.

**Supplemental Figure 12: ROCs for pooled ASM vs. traditional non-ASM method, MAF = 0.5 simulations**

ROC curves comparing the false positive versus true positive rates for simulations for 40, 160, and 640 reads and effect sizes 1.0 and 0.5. In the simulations for this figure, reads are sampled from a negative binomial distribution, there are 100 individuals, and there are 0.5 minor allele and minor methylation status frequencies.

**Supplemental Figure 13: p-Values from pooled ASM vs. traditional non-ASM method simulations**

Histograms of p-values for variant-CpG pairs from 10,000 simulations. These simulations were done for 40 reads, perfect correlation between allele and methylation status, 0.1 minor allele and methylation status frequencies, and 100 individuals. The number of reads in each simulation was sampled from a negative binomial distribution, and p-values were computed using Fisher's Exact Test. **a)** $-\log_{10}$p-values for pooling method. **b)** $-\log_{10}$p-values for traditional method.

**Supplemental Figure 14: Estimated fraction of each individual's DNA in the pool**

**Supplemental Figure 15: Distances between variants and corresponding CpGs for mQTLs**

**Supplemental Figure 16: Pyrosequencing validation of an mQTL that is in strong LD with a GWAS hit**

Shown is an mQTL involving a SNP in 0.80 LD with a SNP previously associated with basal cell carcinoma that was validated in a different set of 30 LCLs. **a)** Pooled bisulfite sequencing for the mQTL showing strong association. **b)** Pyrosequencing validation of the mQTL in 30 additional YRI individuals did not confirm our findings. Light blue points are the methylation percentages from individuals, and crosses are the mean methylation percentages for individuals of each genotype.

**Supplemental Figure 17: Pyrosequencing validation of an mQTL that is a dsQTL**

Shown is an mQTL involving a SNP previously associated with open chromatin that was validated in a different set of 30 LCLs.  **a)** Pooled bisulfite sequencing for the mQTL, showing strong association.  **b)** Pyrosequencing validation of the mQTL in 30 additional YRI individuals shows that the mQTL is not limited to the individuals in our study.  Light blue points are the methylation percentages from individuals, and crosses are the mean methylation percentages for individuals of each genotype.

**Supplemental Figure 18: Pyrosequencing validation of an mQTL that is in strong LD with a GWAS hit**

Shown is an mQTL involving a SNP in 0.84 LD with a SNP previously associated with hypertension risk in short sleep duration that was validated in a different set of 30 LCLs.  **a)** Pooled bisulfite sequencing for the mQTL, showing strong association.  **b)** Pyrosequencing validation of the mQTL in 30 additional YRI individuals shows that the mQTL is not limited to the individuals in our study.  Light blue points are the methylation percentages from individuals, and crosses are the mean methylation percentages for individuals of each genotype.

**Supplemental Figure 19: Pyrosequencing validation of an mQTL that is in perfect LD with a GWAS hit**

Shown is an mQTL involving a SNP in perfect LD with a SNP previously associated with venous thromboembolism that was validated in a different set of 30 LCLs.  **a)** Pooled bisulfite sequencing for the mQTL, showing strong association.  **b)** Pyrosequencing validation of the mQTL in 30 additional YRI individuals shows that the mQTL is not limited to the individuals in our study.  Light blue points are the methylation percentages from individuals, and crosses are the mean methylation percentages for individuals of each genotype.

**Supplemental Figure 20: Pyrosequencing validation of an mQTL that is in strong LD with a GWAS hit**

Shown is an mQTL involving a SNP in 0.86 LD with a SNP previously associated with prostate cancer that was validated in a different set of 30 LCLs.  **a)** Pooled bisulfite sequencing for the mQTL, showing strong association.  **b)** Pyrosequencing validation of the mQTL in 30 additional YRI individuals shows that the mQTL is not limited to the individuals in our study.  Light blue points are the methylation percentages from individuals, and crosses are the mean methylation percentages for individuals of each genotype.

**Supplemental Figure 21: Pyrosequencing validation of an mQTL that is an exon-level eQTL**

Shown is an mQTL involving a SNP previously associated with exon-level expression that was validated in a different set of 30 LCLs.  **a)** Pooled bisulfite sequencing for the mQTL, showing strong association.  **b)** Pyrosequencing validation of the mQTL in 30 additional YRI individuals shows that the mQTL is not limited to the individuals in our study.  Light blue points are the methylation percentages from individuals, and crosses are the mean methylation percentages for individuals of each genotype.

**Supplemental Figure 22: Numbers of tested CpGs and mQTLs in pooled vs. Zhang *et al*. dataset**

**a)** Illustration of the number of CpGs tested for mQTLs in our pooled dataset and in the Zhang *et al*. array data-set.  **b)** Illustration of the number of CpGs with mQTLs in our pooled dataset and in the Zhang *et al*. array data-set.

**Supplemental Figure 23: Numbers of tested CpGs and mQTLs in pooled vs. Banovich *et al*. dataset**

**a)** Illustration of the number of CpGs tested for mQTLs in our pooled dataset and in the Banovich *et al*. array data-set.  **b)** Illustration of the number of CpGs with mQTLs in our pooled data-set and in the Banovich *et al*. array dataset.

**Supplemental Figure 24: Fold-enrichments of CpGs with mQTLs in each chromatin state**

Numbers of chromatin states correspond to the numbers listed in Table 1.

**Supplemental Figure 25: p-Values for mQTL enrichment in open chromatin from LCLs vs. others**

Histograms of $-\log_{10}$p-values for mQTL enrichment in open chromatin regions from different cell types. The histogram for LCLs is in red, and the histogram for all other cell types is in light blue.

**Supplemental Figure 26: Fold-enrichments of mQTLs in TF-binding sites**

This bar graph contains the 12 TF-binding sites that are enriched for mQTLs in Supplementary Table 1.

**Supplemental Figure 27: Methylation bias that depends on position in read**

M-bias plots generated by Bismark for one library in one sequencing lane of one flowcell.  The dark blue line in the plot is the percentage of CpG methylation.  The part of the read underlined by the dark blue line was removed during methylation status calling using Bismark's bismark_methylation_extractor's ignore_r2 option.  The parts of the reads underlined by the dark red lines were removed during methylation status calling by in-house scripts.  Other libraries and sequencing lanes and flowcells for this library have M-bias plots that look similar to this one.

**Supplemental Figure 28: Using LD to combine data from reads with the same CpG**

Reads for each variant, CpG pair can be used to generate a contingency table, even though the individual from which each read was generated is not known.  Reads from the same CpG that have different variants that are in perfect LD can be combined, and the alleles of the first variant for reads with the second variant can be imputed when making the contingency table.

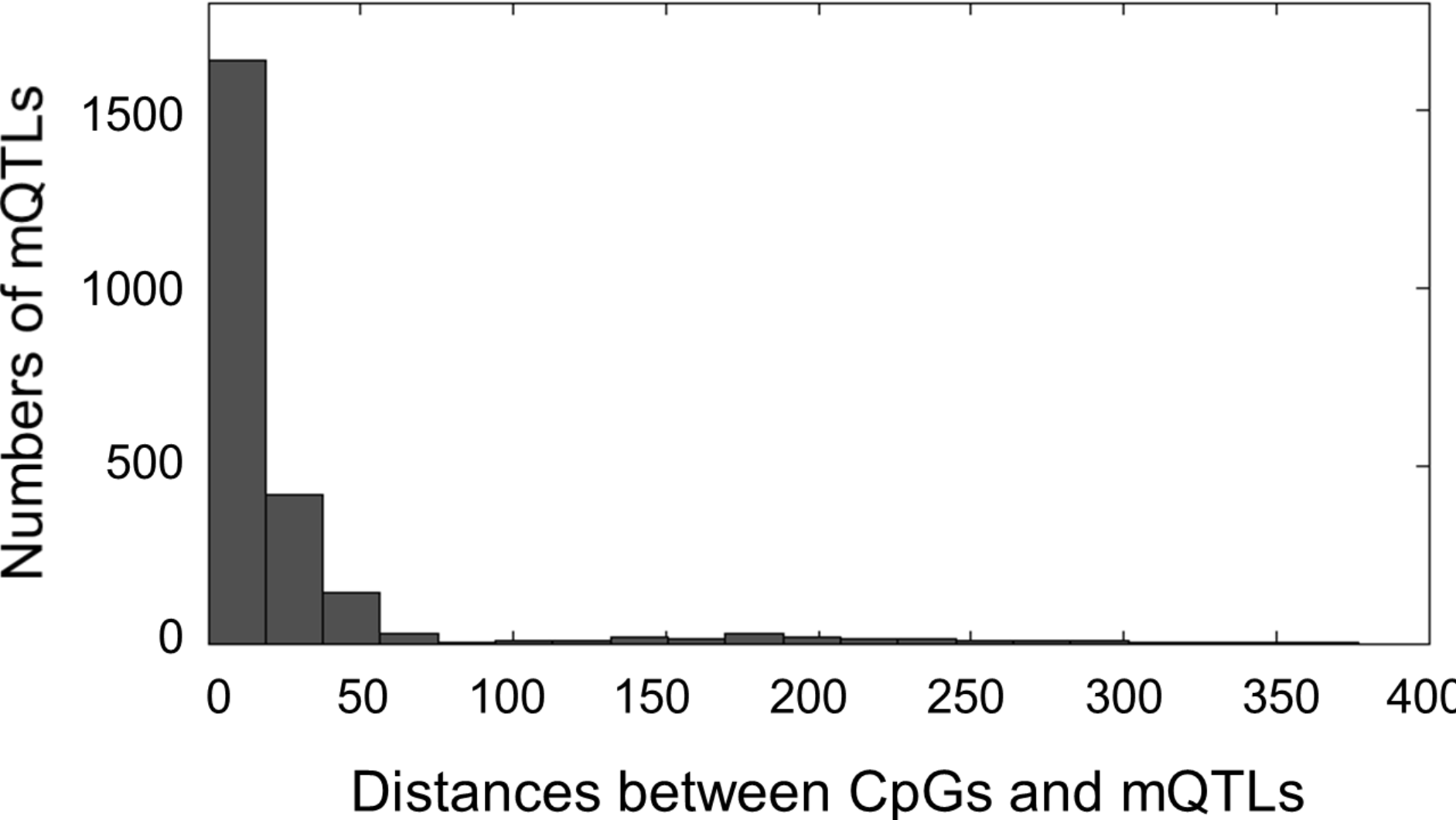**Supplemental Figure 1**

**Supplemental Figure 2**

**Supplemental Figure 3**

# Supplemental Figure 4

**Supplemental Figure 6**

**Supplemental Figure 7**

**Supplemental Figure 8**

# Supplemental Figure 9

**Supplemental Figure 10**

**Supplemental Figure 11**

**Supplemental Figure 12**

# Supplemental Figure 13

**Supplemental Figure 14**

Numbers of mQTLs vs Distances between CpGs and mQTLs

# Supplemental Figure 16



a
Legend: Methylated, Unmethylated

$p = 8 \times 10^{-4}$

b

$p > 1$

**Supplemental Figure 17**

# Supplemental Figure 18

**a**

Legend: Methylated (light blue), Unmethylated (dark blue)
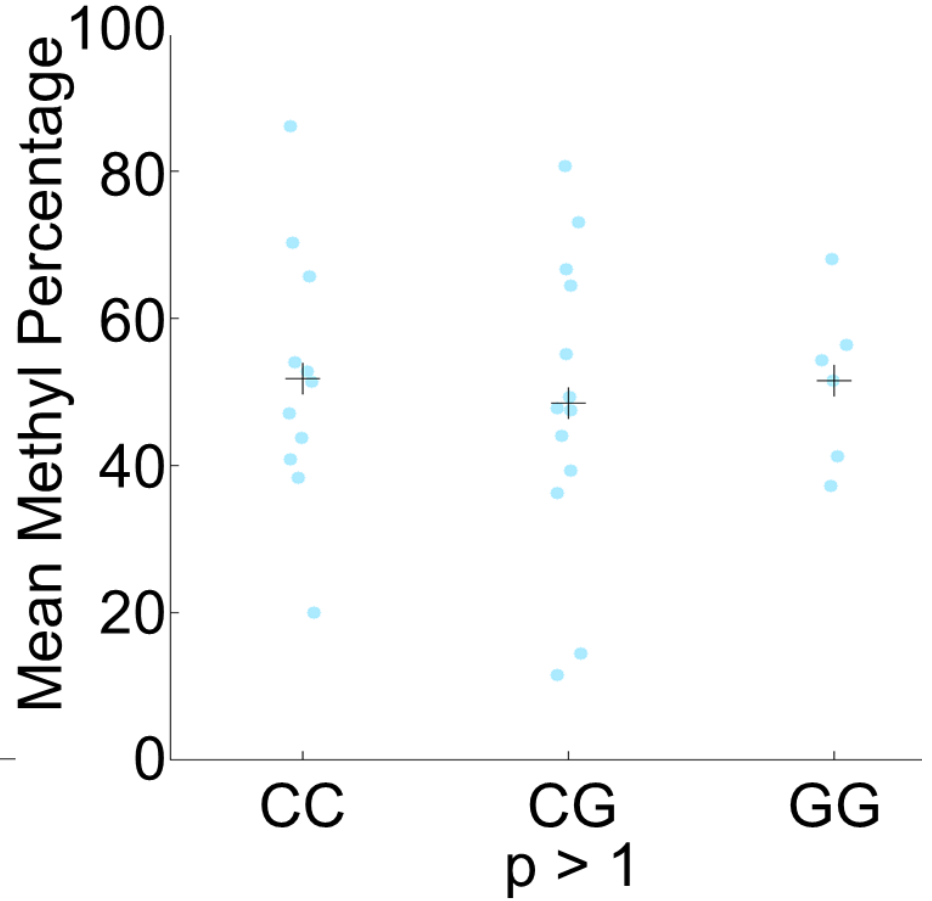


$p = 3 \times 10^{-5}$

**b**



$p = 6 \times 10^{-9}$

# Supplemental Figure 19

**a**

Methylated
Unmethylated



$p = 1 \times 10^{-5}$

**b**



$p = 4 \times 10^{-13}$

# Supplemental Figure 20



**a** Methylated / Unmethylated

Number of Reads — C Allele, A Allele — p = 1 x 10$^{-4}$

**b** Mean Methyl Percentage — CC, CA, AA — p = 2 x 10$^{-5}$

# Supplemental Figure 21
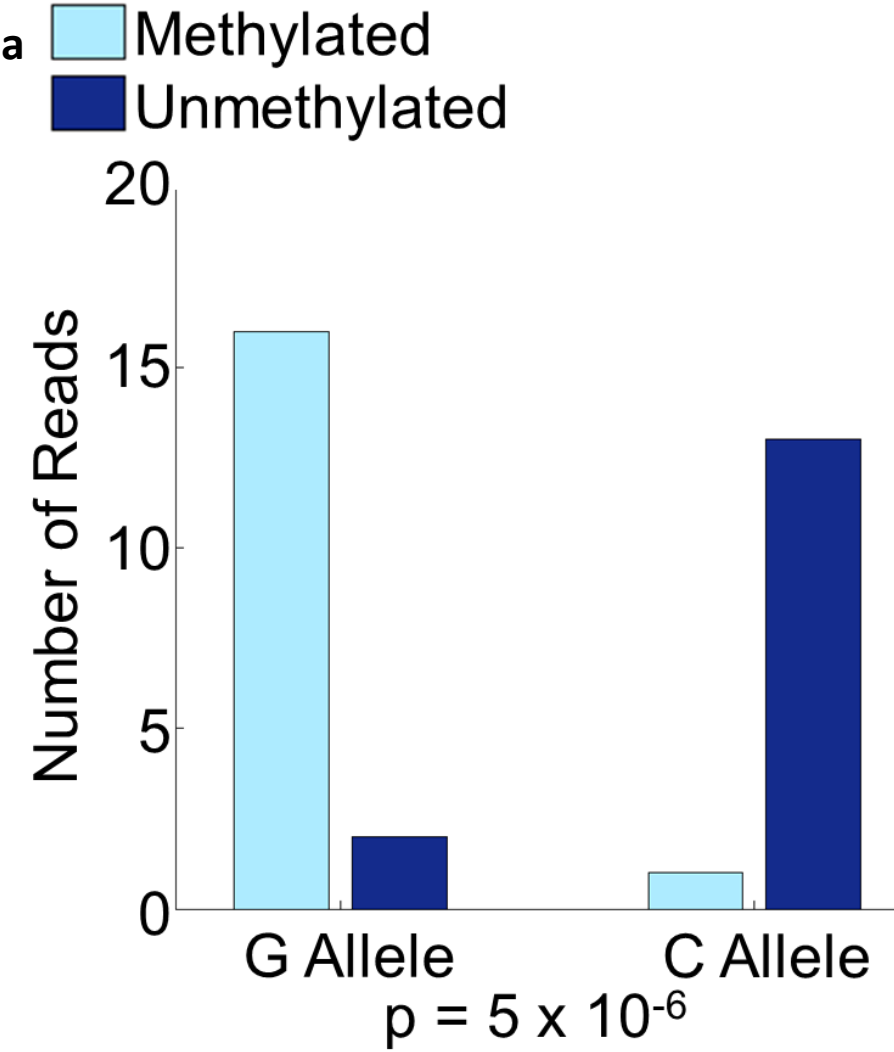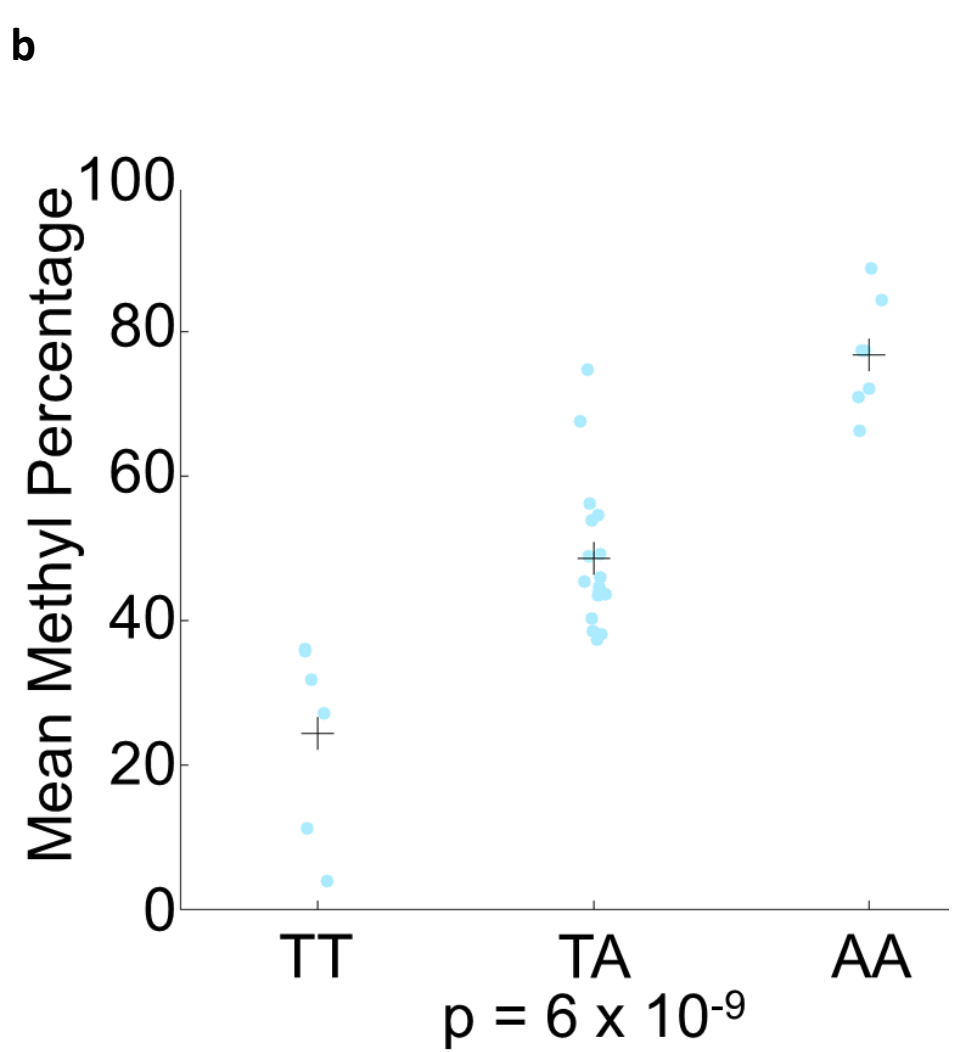
**Supplemental Figure 22**

**Supplemental Figure 23**

**Supplemental Figure 24**

**Supplemental Figure 25**



Histogram with x-axis labeled "$-\log_{10}$(p-Values for DNase Overlap mQTL Enrichment)" ranging from 0 to 25, and y-axis labeled "Numbers of Cell Types" ranging from 0 to 200.

**Supplemental Figure 26**

**Supplemental Figure 27**

Allele 1   Methyl Status   Allele 2

|  | Unmethyl | Methyl |
|---|---|---|
| Ref Allele (T) | 5 | 0 |
| Alt Allele (G) | 1 | 4 |

Perfect LD

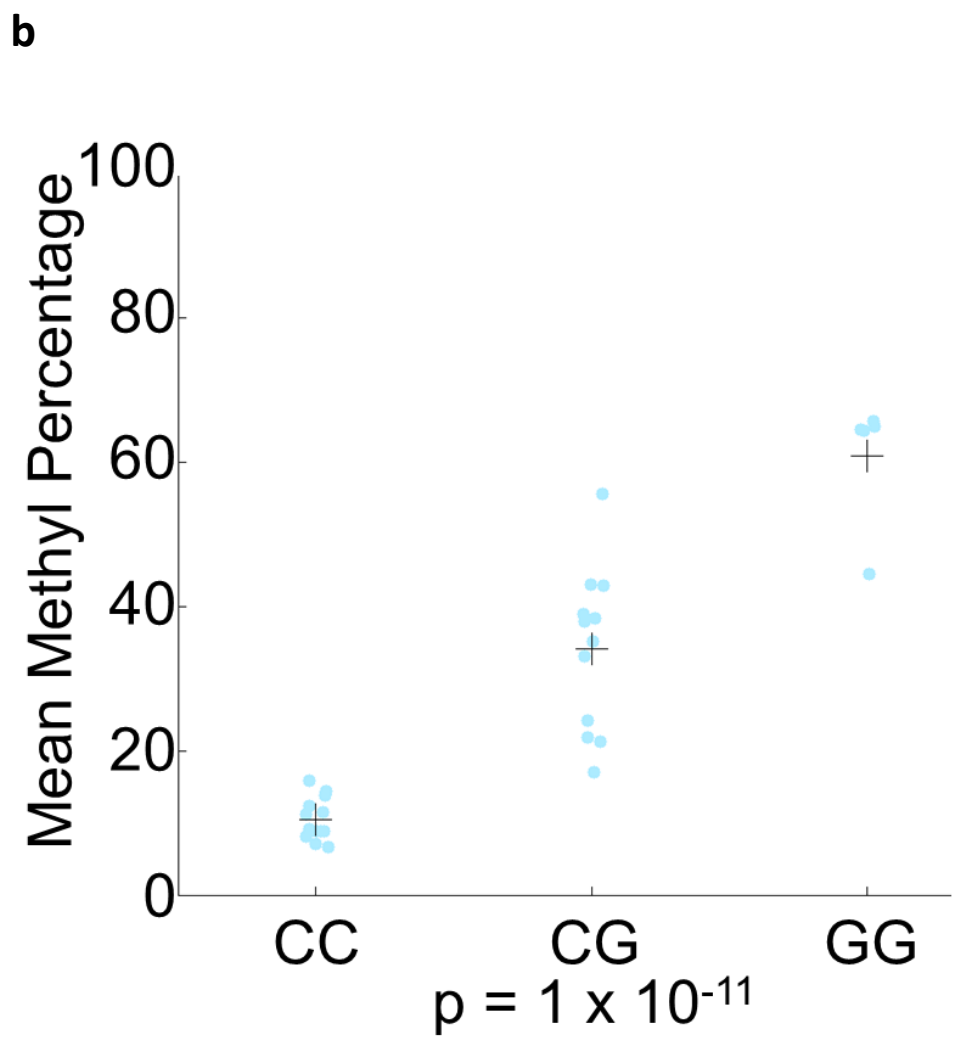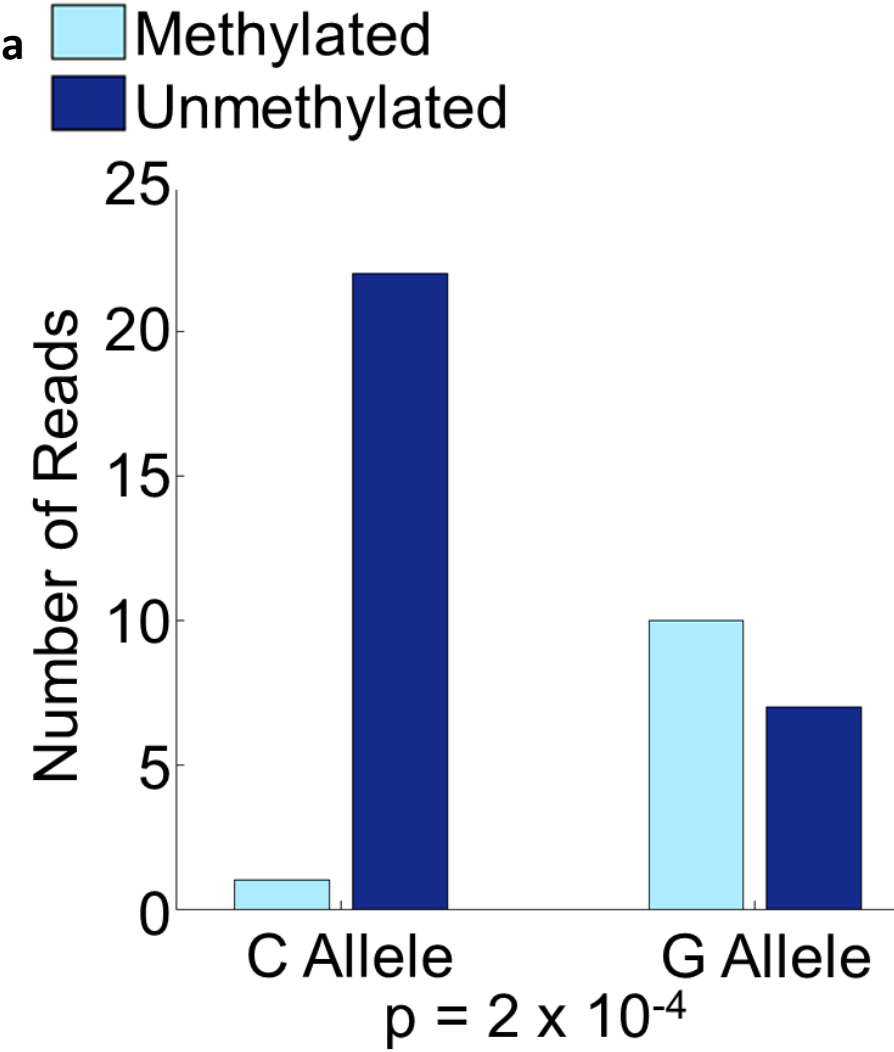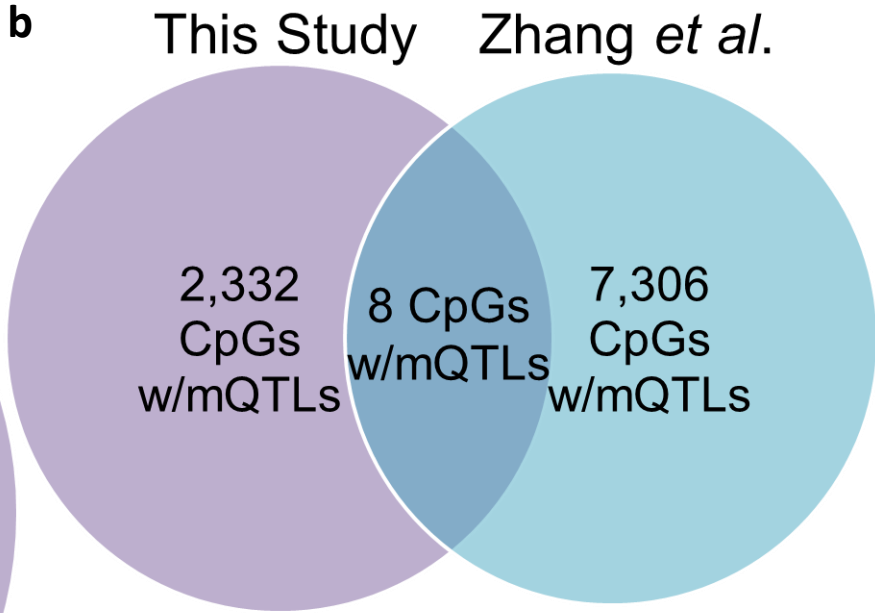| Histone modification/ transcription factor | Number of mQTLs in GM12878 region/peak | p-value for mQTL enrichment |
| --- | --- | --- |
| H2Az | 51 | $1.13 \times 10^{-1}$ |
| H3K4me1 | 92 | > 1 |
| H3K4me2 | 90 | > 1 |
| H3K4me3 | 41 | > 1 |
| H3K27ac | 71 | $4.04 \times 10^{-1}$ |
| H3K27me3 | 3 | > 1 |
| H3K36me3 | 19 | > 1 |
| H3K79me2 | 44 | > 1 |
| H3K9ac | 34 | > 1 |
| H3K9me3 | 13 | > 1 |
| H4K20me1 | 2 | > 1 |
| Atf2 | 15 | $2.71 \times 10^{-1}$ |
| Atf3 | 0 | > 1 |
| **Batf** | **16** | **$2.55 \times 10^{-2}$** |
| Bcl11a | 8 | $7.32 \times 10^{-1}$ |
| Bcl3 | 5 | > 1 |
| Bclaf1 | 0 | > 1 |
| Bhlhe40 | 6 | > 1 |
| Brca1 | 0 | > 1 |
| Cebpb | 3 | > 1 |
| c-Fos | 0 | > 1 |
| Chd1 | 1 | > 1 |
| Chd2 | 5 | > 1 |
| c-Myc | 0 | > 1 |
| CoREST | 0 | > 1 |
| **CTCF** | **31** | **$2.93 \times 10^{-5}$** |
| E2f4 | 0 | > 1 |
| Ebf1 | 10 | > 1 |
| Egr1 | 3 | > 1 |
| Elf1 | 6 | > 1 |
| Elk1 | 1 | > 1 |
| Ets1 | 1 | > 1 |
| Ezh2 | 0 | > 1 |
| **Foxm1** | **16** | **$4.15 \times 10^{-2}$** |
| Gabp | 1 | > 1 |
| Gcn5 | 0 | > 1 |
| Ikzf1 | 3 | > 1 |
| Irf3 | 0 | > 1 |
| Irf4 | 9 | $4.46 \times 10^{-1}$ |
| JunD | 0 | > 1 |
| Max | 3 | > 1 |
| Maz | 4 | > 1 |
| **Mef2a** | **10** | **$1.06 \times 10^{-2}$** |
| Mef2c | 6 | $1.23 \times 10^{-1}$ |

| | | |
|---|---|---|
| Mta3 | 5 | > 1 |
| Mxi1 | 1 | > 1 |
| Nfat | 8 | > 1 |
| Nfe2 | 0 | > 1 |
| **Nfic** | 27 | **$1.82 \times 10^{-4}$** |
| Nfkb | 7 | > 1 |
| Nfya | 0 | > 1 |
| Nfyb | 3 | > 1 |
| Nrf1 | 0 | > 1 |
| NRSF | 2 | > 1 |
| p300 | 7 | > 1 |
| Pax5C | 10 | > 1 |
| Pax5N | 11 | $1.90 \times 10^{-1}$ |
| Pbx3 | 0 | > 1 |
| Pml | 6 | > 1 |
| Pol2, 4h | 9 | > 1 |
| Pol2 | 9 | > 1 |
| Pol2-S2 | 1 | > 1 |
| Pol3 | 0 | > 1 |
| **Pou2f2** | 14 | **$7.56 \times 10^{-4}$** |
| **PU.1** | 21 | **$1.09 \times 10^{-3}$** |
| **Rad21** | 19 | **$5.87 \times 10^{-4}$** |
| Rfx5 | 0 | > 1 |
| **Runx3** | 37 | **$6.65 \times 10^{-5}$** |
| Rxra | 0 | > 1 |
| Sin3a | 3 | > 1 |
| Six5 | 2 | > 1 |
| **Smc3** | 18 | **$7.11 \times 10^{-3}$** |
| Sp1 | 6 | > 1 |
| Spt20 | 71 | > 1 |
| Srf | 4 | > 1 |
| Stat1 | 1 | > 1 |
| Stat3 | 1 | > 1 |
| Stat5 | 7 | > 1 |
| Taf1 | 2 | > 1 |
| Tblr1 | 4 | > 1 |
| Tbp | 4 | > 1 |
| **Tcf12** | 11 | **$2.30 \times 10^{-2}$** |
| Tcf3 | 9 | $2.81 \times 10^{-1}$ |
| Tr4 | 1 | > 1 |
| Usf1 | 2 | > 1 |
| Usf2 | 6 | > 1 |
| Whip | 8 | > 1 |
| Yy1 | 9 | > 1 |
| Zbtb33 | 1 | > 1 |
| Zeb1 | 1 | > 1 |

| Znf143 | 12 | 2.57 x 10$^{-3}$ |
|--------|----|------------------|
| Znf274 | 0 | > 1 |
| Zzz3 | 0 | > 1 |

**Supplemental Table 1: Enrichment of mQTLs in GM12878 regions/peaks**

All p-values are Bonferroni-corrected. GM12878 peaks for twelve TFs are enriched for mQTLs; these TFs are shown in bold. The sites for all of these TFs contain at least ten mQTLs.

| Molecular QTL | Number of mQTL intersections | Number of mQTL intersections, including variants in perfect LD | Number of mQTL intersections, including variants in $r^2 \geq 0.8$ LD |
|---|---|---|---|
| eQTL from Pickrell *et al*. | 0 | 3 | 5 |
| sQTL from Pickrell *et al*. | 0 | 0 | 0 |
| Exon-level eQTL from GEUVADIS Consortium | 28 | 28 | 32 |
| Gene-level eQTL from GEUVADIS Consortium | 5 | 5 | 9 |
| dsQTL from Degner *et al*. | 34 | 48 | 74 |
| CTCF-binding-QTL from Ding *et al*. | 15 | 18 | 28 |

**Supplemental Table 2: Numbers of mQTL intersections with other molecular QTL datasets**

Other QTL data-sets were expanded to incorporate variants at different levels of LD in 1000 Genomes. dsQTLs were from the combined list of dsQTLs with both distance cutoffs used in Degner *et al*.

| Molecular QTL | Fold enrichment of mQTL intersections | Fold enrichment of mQTL intersections, including variants in perfect LD | Fold enrichment of mQTL intersections, including variants in $r^2 \geq 0.8$ LD |
|---|---|---|---|
| eQTL from Pickrell *et al*. | N/A | 3.16 | 1.95 |
| sQTL from Pickrell *et al*. | N/A | 0.00 | 0.00 |
| Exon-level eQTL from GEUVADIS Consortium | 1.91 | 1.78 | 1.68 |
| Gene-level eQTL from GEUVADIS Consortium | 1.58 | 1.46 | 2.22 |
| dsQTL from Degner *et al*. | 10.75 | 7.27 | 4.98 |
| CTCF-binding-QTL from Ding *et al*. | 2.70 | 2.48 | 2.41 |

**Supplemental Table 3: Fold enrichments for mQTL intersections with other molecular QTL datasets**

Fold enrichment is (number of mQTL intersections)/(expected number of mQTL intersections). Other QTL datasets were expanded to incorporate variants at different levels of LD in 1000 Genomes. dsQTLs were from the combined list of dsQTLs with both distance cutoffs used in Degner *et al*. N/A indicates that no molecular QTLs were tested for having mQTLs.

| mQTL rsID | CpG position | GWAS SNP rsID | LD ($r^2$) between mQTL and GWAS SNP | GWAS trait | GWAS paper PMID(s) |
|---|---|---|---|---|---|
| rs10888935 | chr1:56060953 | rs10888935 | mQTL is GWAS SNP | Inflammatory biomarkers | 22228203 |
| rs10797916 | chr1:183838064 | rs4651156 | 1.00 | Response to antidepressants | 20360315 |
| rs10737680 | chr1:196679457 | rs10737680 | mQTL is GWAS SNP | Age-related macular degeneration | 23455636, 20385819 |
| rs801736 | chr11:65928440 | rs564343 | 0.81 | Obesity (early onset extreme) | 23563609 |
| rs9517668 | chr13:99923789 | rs7335046 | 0.80 | Basal cell carcinoma | 21700618 |
| rs9806806 | chr16:9916205 | rs8058295 | 0.87 | Autism spectrum disorder, attention deficit-hyperactivity disorder, bipolar disorder, major depressive disorder, and schizophrenia (combined) | 23453885 |
| rs617201 | chr17:30893145 | rs225212 | 0.84 | Hypertension risk in short sleep duration | 22322875 |
| rs1113144 | chr18:71774787 | rs9945428 | 1.00 | Venous thromboembolism | 23509962 |
| rs885252 | chr2:39850639 | rs7587205 | 0.96 | Response to angiotensin II receptor blocker therapy (opposite direction w/ diuretic therapy) | 22566498 |
| rs16826873 | chr2:198898223 | rs1016883 | 0.80 | Ulcerative colitis | 23128233 |
| rs4809455 | chr20:61660749 | rs6089829 | 0.90 | Prostate cancer | 22219177 |
| rs4809456 | chr20:61660870 | rs6089829 | 0.86 | Prostate cancer | 22219177 |
| rs2837821 | chr21:42161839 | rs2837828 | 1.00 | Neutrophil count | 21507922 |
| rs1904394 | chr3:2653208 | rs4370013 | 0.86 | Blood pressure | 17903302 |
| rs2673051 | chr3:45732458 | rs2742417 | 0.93 | Response to antidepressant treatment | 22041458 |
| rs4859682 | chr4:77410304 | rs4859682 | mQTL is GWAS SNP | Glomerular filtration rate | 23535967 |
| rs830885 | chr5:52021828 | rs830884 | 1.00 | Response to platinum-based agents | 22020760 |
| rs2400797 | chr5:101781191 | rs1502844 | 0.92 | Schizophrenia | 19571808 |
| rs7705033 | chr5:122774795 | rs7705033 | mQTL is GWAS SNP | Visceral adipose tissue/subcutaneous adipose tissue ratio | 22589738 |
| rs2504567 | chr6:26662913 | rs1056667 | 0.90 | Educational attainment | 23722424 |
| rs1405069 | chr6:36922682 | rs1405069 | mQTL is GWAS SNP | Chemerin levels | 20237162 |
| rs71572559 | chr6:66905309 | rs3857536 | 0.88 | Blood trace element | 23720494 |

**Supplemental Table 4: mQTLs in strong LD with GWAS SNPs**

Five mQTLs are GWAS SNPs, and seventeen others are in LD ($r^2$ ≥ 0.8 in YRI) with GWAS SNPs.

| Assay | CpG position | Primer names | Primer sequences |
|---|---|---|---|
| Assay 1 | chr1:196679457 | HF CpG Assay_1 R1 | /5BiodT/AT TTT CTA ACC CTT CAC CCT CCA TAA |
| | | HF CpG Assay_1 F1 | AGT GAG TAA AGG ATT TTA TGA TAT TGG |
| | | HF CpG Assay_1 S1 | AGG TTT ATA TGT TTA TTG TTT AGT |
| Assay 2 | chr5:122774795 | HF CpG Assay_2 R1 | /5BiodT/AT TAC TAC TAC TTA CCA AAA ACT CTT AAA C |
| | | HF CpG Assay_2 F1 | GGA AAG GAA GTG AGG TAG TAA AA |
| | | HF CpG Assay_2 S1 | GAA GTG AGG TAG TAA AAA TAA TA |
| Assay 3 | chr18:71774787 | HF CpG Assay 3 R1* | /5BiodT/CA AAA CAA ACC ACT ATC CCA AAA T |
| | | HF CpG Assay 3 F1 | TGT TAT GGA GGT TTT GGT TTA ATA G |
| | | HF CpG Assay 3 S1 | GGA GGT TTT GGT TTA ATA GA |
| Assay 4 | chr17:3089145 | HF CpG Assay 4 F1* | /5BiodT/AG TGT TGG GAT TAT AGA TGT GAG TT |
| | | HF CpG Assay 4 R1 | ACC CTC TCC TCA AAC AAA TCT AAA TC |
| | | HF CpG Assay 4 S1 | AAA CTA TAT CTA CCT CCC |
| Assay 5 | chr13:99923889 | HF CpG Assay_5 R1* | /5BiodT/CA CTA TCC TAT CAA ACC ATT ATA CTA A |
| | | HF CpG Assay_5 F1 | TTG AGG GAG AAT TTG ATA ATT TGA GA |
| | | HF CpG Assay_5 S1 | AAA AGA ATG GGA AAT AAT GAA |
| Assay 6 | chr21:30158046 | HF CpG Assay 6 R1* | /5BiodT/TT CCC ACT TTA ACT CTT ACT TCA ATA CTA |
| | | HF CpG Assay 6 F1 | GAG TTA GAA ATT TAG GTT GGG TTT AGG |
| | | HF CpG Assay 6 S1 | TGT TTA TGT GTA GGG AAT |
| Assay 7 | chr15:65164853 | HF CpG assay_7 R1* | /5BiodT/CC TAC CAC CAC CCC TAA CTA ATT TTA TAT |
| | | HF CpG assay_7 F1 | AGA GGA AAT AGT AGG ATG TAA GTA GA |
| | | HF CpG assay_7 S1 | GGA GTT AGA GAT TAG TTT GGT TAA |
| Assay 8 | chr20:61660870 | HF CpG Assay 8 R1* | /5BiodT/TT ATC TTT CCT CAT TAA ACC TCT ACT |
| | | HF CpG Assay 8 F1 | GTT TTA GTT TTT TAG TTT GGA TTG GAT AA |
| | | HF CpG Assay 8 S1 | GTT ATA AGT TTT TTT TGG ATT TAG GG |

**Supplemental Table 5: Pyrosequencing primers**

We used these primers for the pyrosequencing validation of mQTLs.