

This PDF file includes:

Captions for Supplementary Tables 1 to 6

Supplementary Notes

Supplementary Figures 1 to 6

Other supplementary material for this manuscript includes:

Supplementary Tables 1 to 6 (Excel format)

Reporting Summary (PDF)

Editorial Policy checklist (PDF)

Supplementary Tables

Supplementary Table 1. GTEx v8 cohort information. Tissue-related labelling, tissue and expressed genes sample size information.

Supplementary Table 2. eQTM findings. Summary statistics for eQTM discovery sets, replication of eQTMs in independent eQTM datasets, eQTM tissue specificity and eQTM contributing factors.

Supplementary Table 3. e/mQTL enrichment in functional annotations. Summary statistics for eQTL and mQTL enrichment in open chromatin regions, transcription factor binding sites, chromatin states and gene body and regulatory regions.

Supplementary Table 4. Colocalization of ovarian cancer GWAS with HOXD-locus QTLs. Summary statistics for colocalization of ovarian breast cancer GWAS results with eQTLs and mQTLs corresponding to eGenes and mCpGs of the HOXD locus.

Supplementary Table 5. QTL-GWAS colocalization. GWASs metadata, significant QTL-GWAS colocalization associations and utilized priors.

Supplementary Table 6. Trait-linked mCpG gene candidates. Integration of mQTL-specific trait-linked mCpGs with curated promoter- and enhancer-gene target predictions.

Supplementary Notes

Data and Samples

The GTEx v8 data consists of 17,382 RNA-seq samples from 948 post-mortem donors, with genotype data for 838 donors from whole genome sequencing (WGS) available in a phased analysis freeze VCF. The GTEx biospecimen collection, molecular phenotype data production and quality control are described in detail in ⁴. The eGTEx project ¹¹ seeks to complement the gene expression traits determined in the GTEx project with other molecular traits across the same tissues and individuals, including methylation.

Here, we have generated and analyzed DNA methylation (DNAm), and analyzed existing ⁴ gene expression data from 9 tissue types: colon transverse, kidney cortex, lung, muscle skeletal, ovary, prostate, testis, whole blood and breast mammary tissues (Supplementary Table 1). Altogether, we analyzed a total of 987 DNAm and 3,872 gene expression samples - depending on the tissue and analysis - corresponding to 424 and 938 individuals, respectively, as well as genotype data from a total of 830 individuals.

Reduction of methylomes' dimensionality

Considering the 987 profiled methylomes, β values were logit-transformed to M-values. The dimensionality of the methylome set was reduced to two dimensions with the t-Distributed Stochastic Neighbor Embedding (t-SNE) approach ⁶³ implemented in the *Rtsne* R package (parameters: perplexity = 20, theta = 0.5, max_iter = 5000, pca=TRUE). We observed a set of samples that did not optimally cluster with their corresponding tissue when projected to the t-SNE dimensions. Hence, we estimated cluster dissimilarity for each sample: for each dimension, t-SNE values were transformed to z-scores, and samples with t-SNE values greater than 2.5 standard deviations in either t-SNE dimension were flagged as potential tissue type mismatches but retained. In total, 12 samples were flagged - 7 ovary, 4 prostate and 1 colon samples. The untransformed t-SNE estimates of the 975 unflagged samples are visualized in Fig. 1a.

Evaluation of tissue similarity based on methylation and gene expression

Hierarchical clustering of the tissues (see Methods) resulted in the obtention of one clustering tree for the DNAm-based distance and another for the transcription-based distance (Fig. 1b), as well as bootstrap probabilities (BP, obtained by normal bootstrap resampling), which is a

measure of clustering support with respect to the data. BP values range from 0 to 100; higher value indicates stronger support for the clusters. All nodes of the DNAm-based tree exhibited maximum support (BP = 100). Most (7/9) nodes of the transcription-based tree exhibited maximum support, except the ones corresponding to muscle with ovary (BP = 53) and kidney (BP = 97) clusters.

We observe, both for gene expression and DNAm profiles, that testis and blood exhibit a lower degree of similarity relative to other tissue types. Testis is characterized by higher gene expression compared to other tissues, as it has been previously shown^{4,9}. In contrast, blood appears to be the most divergent tissue for DNAm; CpGs highly methylated in whole blood and lowly methylated in ovary are prominent features of the tissue-specific DNAm signatures (Fig. 1b).

Definition of expression quantitative trait methylation (eQTM) mapping sets

The number of CpG-gene tests performed per tissue varied as a function of the number of genes expressed per tissue; we analyzed a total of 5,350,829 CpG-gene pairs. For CpGs with at least one significant (FDR < 0.05) eQTM (see Methods), we defined as significant all CpG-gene pairs at Bonferroni-adjusted $P < 0.05$, resulting in a non-redundant 12,652 eQTM set across tissues. This set, defined as ‘single-tissue eQTM set’ was complemented with significant cases derived from a cross-tissue approach to constitute the complete set of significant eQTMs reported, defined as ‘complete eQTM set’. This approach²⁹ enables joint modeling of cross-tissue effects, and it is implemented in the R package *mashr*. For the cross-tissue analysis, we considered the set of 12,652 eQTMs, i.e. CpG-gene pairs, significant in at least one tissue derived from the single-tissue eQTM-mapping approach. We applied Fisher's z transformation to Spearman correlation coefficients with the function *fisherz* of the R package *psych*, and calculated corresponding standard errors. For all tissues, we selected Fisher-transformed Spearman coefficients and corresponding standard errors for each of the 12,652 eQTMs, as well as for 100,000 randomly selected CpG-gene pairs that were tested across all tissues. Those estimates were used to fit the *mashr* model. The local false sign rate (LFSR) generated by *mashr* was used to identify significant (LFSR < 0.05) eQTMs. The complete set of significant eQTMs, referred to as ‘complete eQTM set’, was defined by the union of significant cases derived from the single-tissue (FDR < 0.05) and multi-tissue (LFSR < 0.05) eQTM-mapping approaches (FDR < 0.05 or LFSR < 0.05). This resulted in an expansion of the 15,839 eQTM-tissue significant eQTM set - corresponding to

12,652 eQTM - to 47,783 eQTM-tissue significant cases (Supplementary Figure 1a, Supplementary Table 2). Across the article, we refer to CpGs with at least one significant eQTM as eCpGs. Considering eQTM-based downstream analyses, the ‘single-tissue eQTM set’ or the ‘complete eQTM set’ were used depending on the particular analysis, as noted.

The number of significant eQTMs was strongly correlated with per-tissue sample size (Spearman’s $\rho = 0.86$), indicating power limitations to detect eQTMs in low-sampled tissue sets, and we observed that eQTMs tend to be either tissue-specific or shared across most tissue types (Supplementary Figure 1b). However, the limited sample size of the sets utilized for eQTM mapping ($N < 40$ in 4 tissues, Supplementary Figure 1a) and the limited number of tissues ($N = 8$) impose limitations in accurately estimating how much eQTMs are shared across tissues. This analysis would benefit from a more powered and exhaustive eQTM catalog.

Replication of eQTMs in external cohorts

We assessed eQTM replication for all tissues in the FUSION Skeletal Muscle Study cohort, where $N = 265$ individuals were utilized for eQTM mapping. The FUSION cohort characteristics, and the eQTM-mapping procedure, which is similar to the one employed herein, are described in ¹⁹; we employed available summary statistics. For a particular tissue, we considered for replication analyses CpG-gene pairs from the single-tissue eQTM set that passed $P < 0.05$ in the tissue tested. Across tested tissues, we observed a high replication rate (cross-tissue average $\pi_1 = 0.75$), being the highest for muscle ($\pi_1 = 0.84$), possibly due to muscle-specific eQTMs contribution (Supplementary Table 2).

Characterization of eQTM predictors

To annotate mCpGs for gene regulatory elements, we extended the span of their genomic location by ± 100 bps, and checked for overlap (≥ 1 bp) with regulatory regions. For each regulatory element class, a Fisher’s exact test was conducted to determine enrichment of CpGs included in the single-tissue eQTM set in gene regulatory elements, and significance was defined at Bonferroni-adjusted $P < 0.01$. To assess the relative contribution of molecular signatures to eQTMs linked to different gene regulatory elements, we stratified eQTM tests by CpG-overlap with gene regulatory element classes. We considered one eQTM test per CpG per tissue, corresponding to the CpG-gene test with the smallest p-value. For each element class, in each tissue, a logistic regression model of eQTM likelihood was built to

predict whether an eQTM test was significant given the CpG-Gene TSS distance (in Kb), direction of the eQTM effect ('1' for negative correlation between methylation and expression, '0' otherwise), the tissue-averaged methylation (in M-value units) and gene expression abundance (in $\log_2(\text{TPM}+1)$ units), as predictors. For all predictors in all logit models, multicollinearity was tested with the R package *vif*. None of the predictors displayed problematic variance inflation; *vif* score for any predictor was below 1.5. Cross-tissues meta-effect was evaluated by modeling single-tissue effect estimates (log of odds ratio) with a random-effects model (*rma* function, *metafor* R package). Summary statistics relative to the characterization of eQTM are provided in Supplementary Table 2. Additionally, to evaluate the contribution of e/mQTLs to eQTM, we aligned our e/mQTL colocalization results (see Characterization of mQTL-eQTL shared signal) with eQTM identified at $\text{FDR} < 0.05$, and observed that 37% of eQTM correspond to e/mQTL colocalizations, highlighting a substantial contribution of genetics to CpG-gene expression correlation.

Comparison of empirical associations of DNA methylation with gene expression to array annotations

To investigate how accurately the CpG-gene assignments provided by Illumina reflect the eQTM results observed in GTEx data, we contrasted our eQTM findings with the EPIC array CpG annotation file provided by Illumina. We observed that while 76% (4,489/5,898) of the eCpGs we identify have an assigned gene provided in the annotation file, for only 45% (2,641/5,898) of these eCpGs does the annotated gene match a CpG-gene association detected through eQTM mapping. Moreover, only 22% (2,828/12,652) of the eQTM we identify in at least one tissue match any annotated CpG-gene pair. Overall, our eQTM results can enhance our ability to assign CpGs to gene(s) with which they are biologically linked, therefore facilitating the interpretation of methylation-derived analyses, as methylation patterns are often interpreted in light of their predicted impact on gene regulation.

Mapping of mQTLs and eQTLs

To define mQTLs we analyzed all samples with available methylation and genotype data (Supplementary Table 1), comprising a total of 856 samples, from 42 - muscle skeletal - to 190 - lung - per tissue, derived from 367 subjects⁴ and interrogated a total of 754,054 CpGs across all tissues. For each variant-CpG pair, we fit a linear regression model separately in each tissue, and tested for significance of genotype on methylation estimates while adjusting for additional known and unknown factors:

$$Y = \beta_0 + \beta_G \text{ Genotype} + \beta_{(1,...,m)} C + \beta_{(1,...,n)} \text{ PEER} + \varepsilon$$

where,

Y is the inverse-normal-transformed DNAm levels

β_0 is the intercept

β are the corresponding effect sizes. β_G is the effect size of genotype on DNAm.

C represents a subset of covariates that were used in *cis*-eQTL mapping⁴. These covariates include 5 genotype principal components, 2 covariates derived from the generation of genotype data by whole genome sequencing - described in⁴ - and biological sex status.

PEER represents PEER factors⁵¹ derived from DNAm. The number of PEER factors was selected to maximize mQTL discovery, across two sample size bins: tissues with < 50 samples and tissues with ≥ 100 samples. The optimization was performed similarly to⁴, and resulted in the selection of 5 and 20 PEER factors, respectively, for the two sample size bins. In the optimization step, PEERs were calculated from inverse-normalized DNAm β values from CpGs in chromosome 1 (~70K CpGs) and significant mQTLs were defined at nominal $P < 1e-05$. To correct for multiple testing of variants per CpG, we permuted DNAm estimates 1,000 times, adjusting p-values with a beta distribution approximation⁵⁴. Genome-wide CpG multiple testing correction was performed on top-significant CpG-variant beta-adjusted p-values using Storey *qvalue*⁵³. The set of significant mQTL CpGs (mCpGs) identified at $FDR < 0.05$ was defined as ‘single-tissue mQTL set’, and complemented by significant cases derived from cross-tissue QTL mapping (see ‘Definition of QTL sets’).

To identify independent mQTLs, we started from the set of mCpGs discovered in the first pass of association analysis (complete mQTL set: $FDR < 0.05$ or $LFSR < 0.05$). Then, the maximum beta-adjusted p-value (correcting for multiple testing across the variants) over these CpGs was taken as the CpG-level threshold. The next stage proceeded iteratively for each CpG and threshold. A *cis*-scan of the window was performed in each iteration, using 1,000 permutations and correcting for all previously discovered variants. If the beta-adjusted p-value for the most significant CpG-variant, i.e. best association, was not significant at the CpG-level threshold, the forward stage was complete and the procedure moved on to the backward step. If this p-value was significant, the best association was added to the list of

discovered mQTLs as an independent signal and the forward step proceeded to the next iteration. Once the forward stage was complete for a given CpG, a list of associated variants was produced which we refer to as forward signals. The backward stage consisted of testing each forward signal separately, controlling for all other discovered signals. To do this, for each forward signal we ran a *cis* scan over all variants in the window using *FastQTL*, fitting all other discovered signals as covariates. If no variant was significant at the CpG-level threshold the signal being tested was dropped, otherwise the best association from the scan was chosen as the variant that represented the signal best in the full model.

We define eQTLs as *cis*-gene variants with a significant genotype effect on gene expression, utilizing a single-tissue approach analogous to the mQTL-mapping one. We included the same covariates and variant set ($\pm 1\text{Mb}$ from gene transcription start site, $\text{MAF} < 0.01$) employed for eQTL mapping in ⁴. A total of 3,438 samples was considered, from 73 - kidney cortex - to 706 - muscle skeletal - samples per tissue, from a total of 829 subjects. Analogously to mQTLs, we identified multiple independent eQTLs, and the complete set of significant eQTLs was obtained by complementing the single-tissue mQTL set with significant cases derived from the cross-tissue approach (see ‘Definition of QTL sets’).

Definition of mQTL and eQTL sets

To overcome QTL-mapping limited power due to per-tissue available sample sizes, and to determine QTL tissue-specific patterns, we used an approach to perform a cross-tissue QTL analysis by leveraging QTL signal across tissues ²⁹, implemented in the R package *mashr*. Considering the set of 286,153 mCpGs significant in at least one tissue derived from the single-tissue mQTL-mapping approach, i.e. the ‘single-tissue mQTL set’, for every top mQTL per CpG per tissue, mQTL effect sizes, corresponding standard errors and 301,801 randomly selected variant-CpG pairs that were tested across all tissues were used to fit the *mashr* model. The *mashr* version employed herein (0.2.6) sets missing effect size values to 0 and corresponding standard error to 1,000,000. The local false sign rate (LFSR) generated by *mashr* was used to define significant ($\text{LFSR} < 0.05$) mQTLs.

The ‘complete mQTL set’ set was defined by the union of significant cases derived from the single-tissue ($\text{FDR} < 0.05$) and cross-tissue ($\text{LFSR} < 0.05$) mQTL-mapping approaches ($\text{FDR} < 0.05$ or $\text{LFSR} < 0.05$). This resulted in an expansion of the 607,987 mCpG-tissue significant mQTL set - corresponding to 286,152 mCpGs - to 1,385,225 mCpG-tissue significant cases. An equivalent approach was employed to perform cross-tissue eQTL

meta-analysis. Considering mQTL- or eQTL-based (e/mQTL) analyses, ‘single-tissue e/mQTL set’ or the ‘complete e/mQTL set’ are used depending on the particular analysis, as noted.

Replication of mQTLs in external cohorts

We assessed mQTL replication in the BEST blood cohort (N = 337, A = HumanMethylation 450K) ²⁶ and the FUSION Skeletal Muscle Study cohort (N = 282, A = HumanMethylation EPIC) ¹⁹; ‘N’ and ‘A’ define the number of individuals utilized for mQTL-mapping and the Illumina array used to profile methylation, respectively. In all cases, we tested for replication the mCpG lead variants (best variant per mCpG) from the single-tissue mQTL set. The mQTL-mapping procedure and the BEST cohort characteristics are described in ²⁶; we employed available summary statistics. In brief, DNAm β values were logit-transformed and adjusted for potential batch effects. A linear model was fit with genotype, age, sex and 10 methylation PCs, considering variants in a ± 500 Kb window from the CpG locus. A relatively similar procedure was employed to map mQTLs in the FUSION Skeletal Muscle cohort, described in ¹⁹. Replication was assessed by means of π_1 , which measures the estimated true positive rate ⁵³. We observed high true positive rate values in the blood ($\pi_1 = 0.91$) and muscle ($\pi_1 = 0.93$) cohorts, based on 7,245 and 3,548 tested variant-gene pairs, respectively.

Characterization of tissue specificity patterns of mQTLs and eQTLs

The number of mCpGs detected per tissue was strongly correlated with per-tissue sample size (Spearman’s $\rho = 0.92$). The overall tissue specificity of mQTLs follows a skewed U-shaped curve, i.e. for a particular CpG, genetic regulation of DNAm tends to be either highly tissue-specific or highly shared across tissue types (Supplementary Figure 2a). The fraction of mCpGs identified that were detected as mCpGs exclusively in a single tissue (Fig. 2b) is assumed to be a lower bound, as we observe that the abundance of tissue-specific mCpGs is strongly correlated (Spearman’s $\rho = 0.80$) with sample size, indicating power limitations to detect tissue-specific QTLs in low-sampled tissue sets. This assumption is compatible with the larger tissue-specific eGene fractions (Supplementary Figure 2b) observed for eQTLs, mapped in larger sample sets. Differential tissue-sharing distribution of comparing eQTLs to mQTLs was tested by means of a Wilcoxon rank-sum test, and the null hypothesis was rejected ($P = 2.9\text{e-}03$). That implies that, compared to eQTLs, mQTLs appear to be more tissue-shared. This pattern could be due to more stable cross-tissue DNAm QTL effects

compared to expression QTLs, but also to substantially lower mQTL sample sizes (compared to eQTL sample sizes).

Functional genomic characterization of mQTLs and eQTLs

We observed eQTLs to be more strongly enriched in open chromatin sites than mQTLs (Supplementary Figure 3a). Additionally, mQTLs appear to be depleted in transcribed genes and genic enhancers but enriched in distal, active enhancers (Supplementary Figure 3b).

mQTL-eQTL colocalization

We investigated the associations between mQTLs and eQTLs (single-tissue QTL set: $FDR < 0.05$) by means of QTL effect size colocalization with *coloc*⁵⁶ using default priors. For both QTL types, we considered unconditional QTL mappings, i.e. agnostic to multiple independent QTLs, due to computational limitations of performing colocalization on the complete combinatorial space considering multiple independent QTLs for both QTL types. A mQTL locus was defined as overlapping with an eQTL locus, and subsequently tested for colocalization, if the mQTL-eQTL region a) had at least 50 variants in common and b) included potentially causal mQTL and eQTL variants. That is, it included at least one fine-mapped and/or conditional QTL mapping lead variant, for both mQTL and eQTL signals. Fine-mapped QTL variants were estimated with *dap-g*⁶², and QTL credible sets were defined at 90% confidence.

To classify mQTL loci (mCpGs) into the mutually exclusive colocalization categories depicted in Fig. 2e, we first annotated mQTL loci that did not overlap any eQTL locus in any tissue. For the remaining eQTL-overlapping mQTL loci set, which was tested for eQTL colocalization, a mQTL locus was annotated as involved in a eQTL-mQTL colocalization if it colocalized ($PP4 > 0.5$) with at least one eQTL in at least one tissue. For the remaining set, a mQTL locus was annotated as independent to eQTL signal if it exhibited a $PP3 > 0.5$ for at least one eQTL colocalization test in at least one tissue. The remaining set was considered to exhibit inconclusive colocalization signal ($PP0 + PP1 + PP2 > 0.5$).

Across tissues, 93% (266,239/286,152) of mQTL loci do overlap with an eQTL. Using a moderately permissive threshold for the posterior probability of sharing the same causal variant ($PP4 > 0.5$), only 21% of mQTL loci are suggestively colocalized ($PP4 > 0.5$) with at least one eQTL, whereas for 38% of cases there is evidence of independent variants driving

mQTL and eQTL signals ($PP3 > 0.5$). For the remaining 34% of the cases analyzed, we lack adequate power to conduct meaningful colocalization analyses ($PP0 + PP1 + PP2 > 0.5$). Our results indicate that a considerable fraction of mQTLs do not show clear associations with local gene expression in the same tissue type, but we acknowledge that limited power for e/mQTL detection, allelic heterogeneity, and colocalization assumptions may limit our ability to accurately estimate this fraction.

mQTL-eQTL concordance in direction of effects across regulatory regions

We assessed whether the discordance/concordance rate varied as a function of mCpG location in gene regulatory regions, considering promoters and proximal enhancers jointly, distal enhancers and insulators. Gene regulatory element annotations were derived from ENCODE5 cCREs catalog (see eQTM section above). To annotate the mCpGs, we extended the span of their genomic location by ± 100 bps, and checked for overlap (≥ 1 bp) with regulatory regions. A mCpG was annotated with promoter/proximal enhancer status if, in addition to overlapping with an ENCODE-predicted promoter or proximal enhancer, it overlapped the 2kb region upstream from the TSS of the corresponding colocalized eGene. To compare the discordance/concordance rate across regulatory regions, we performed a multi-sample test for equality of proportions without continuity correction (*prop.test* function, *stats* R package).

Characterization of mQTL-eQTL regulatory pleiotropy

Regulatory pleiotropy categories are defined in Methods and illustrated in Supplementary Figure 4a. The most common scenario, comprising 54% of the eQTL-colocalized mCpGs (Supplementary Figure 4b), corresponds to eQTL-mQTL colocalizations involving multiple mCpGs and a single eGene (Tier 3 in Supplementary Figure 4a). Overall, we observe a higher mCpGs per eGene than eGenes per mCpG ratio (Supplementary Figure 4c). The mCpGs per eGene ratio is correlated with sample size (Spearman's $\rho = -0.75$), suggesting that the observed value is a lower bound estimate due to mQTL-mapping power limitations. Across tissues, the largest pleiotropic sets tend to involve mCpGs and eGenes located in the Major Histocompatibility Complex (MHC).

Colocalization of GWAS with QTL signal

The approach to identify colocalization of GWAS with QTL signal is described in Methods. In total, 6,720 GWAS-GWAS-hit tuples were considered for downstream analyses, from 1 - mothers's age at death, epilepsy, self-reported schizophrenia, intracranial volume, insomnia - to 733 - standing height - GWAS hits depending on the GWAS trait.

To evaluate the conservativeness of selected priors, we compared *coloc* mQTL-GWAS colocalization results to those generated with default priors ($p_1 = 1e-04$, $p_2 = 1e-04$, $p_{12} = 1e-05$) for the GWAS with largest amount of signal, i.e., UKB standing height GWAS. We observed that results are strongly correlated (Spearman's $\rho = 0.93$), but colocalization probabilities derived from *fastenloc*-derived priors tend to be more conservative (Supplementary Figure 5a). That is, at $PP_4 > 0.5$, considering *fastenloc*-derived priors, we identify 53% less colocalized cases than with the default-priors approach. That is expected, given the higher ratio between *fastenloc*-derived p_2 (GWAS association) and p_{12} (mQTL and GWAS association) priors compared to corresponding default-priors one (Supplementary Table 5).

Evaluation of mQTL-GWAS colocalization approach

Considering results with suggestive colocalization probability (the intersection set of *coloc* $PP_4 > 0.1$ and *fastenloc* $RCP > 0.1$), we observe a strong correlation (Spearman's $\rho = 0.79$) between results from both methods (Supplementary Figure 5b). We identified as significantly colocalized those GWAS-GWAS hit-mCpG/eGene-tissue tuples with both corresponding *coloc* $PP_4 > 0.3$ and *fastenloc* $RCP > 0.3$, i.e. the intersection of cases with moderate colocalization signal derived from both methods ($RCP > 0.3$ and $PP_4 > 0.3$). Across the article, *coloc* PP_4 is provided as the reference colocalization probability. We identify 55% of GWAS hits (1,505/2,734) colocalizing with at least one mQTL but with no eQTLs at $RCP > 0.3$ and $PP_4 > 0.3$; this estimate can range from 44 to 66% depending on the combination of PP_4 and RCP thresholds selected. Colocalization cases involving at least one mQTL but no eQTLs are defined as 'mQTL-specific' colocalizations, as opposed to 'eQTL-specific' colocalizations, which comprise cases involving at least one eQTL but no mQTLs. Colocalization cases involving at least one mQTL and one eQTL are defined as 'e/mQTL-shared' colocalizations.

Scope of mQTL-GWAS colocalizations

Instances of mQTL-GWAS colocalizations were observed among 81% (67/83) of tested GWASs and involved 41% (2,734/6,720) of GWAS hits, and 3,381 and 940 colocalized (trait-linked) mCpGs and eGenes, respectively. For 4.5% (102/2,254) of GWAS hits involved in mQTL-GWAS colocalizations, the colocalizing signal corresponded to a secondary mQTL. For nine GWAS traits, colocalizations were only detected for mQTLs, including osteoporosis and certain balding and metabolic phenotypes, among other traits.

Characterization of DNAm signatures of trait-linked mCpGs

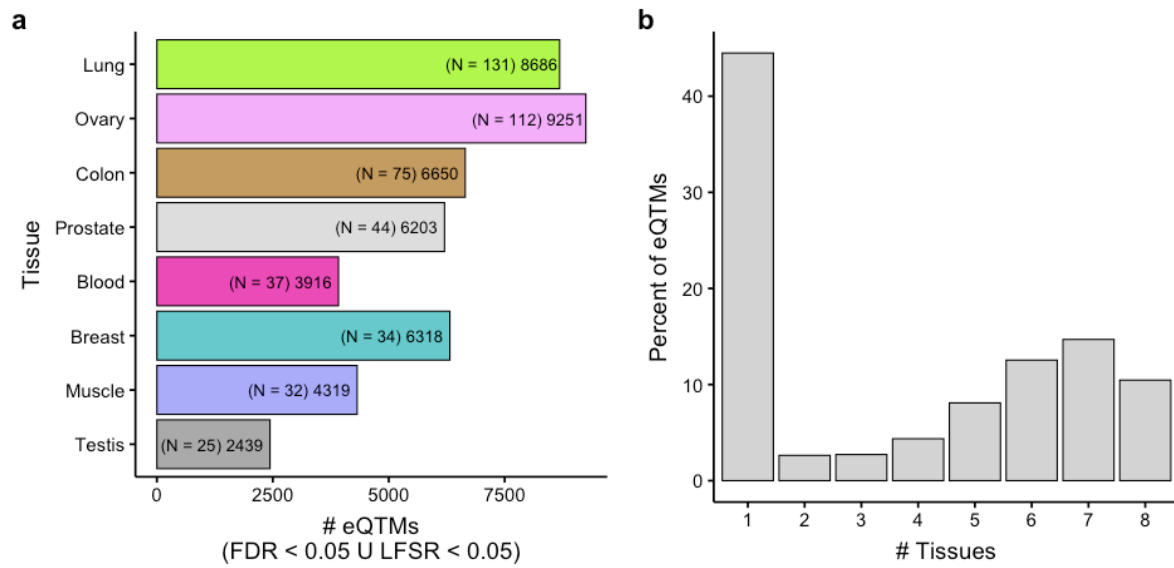
To evaluate the overlap of trait-linked mCpGs with open chromatin regions, we extended the genomic location span of mCpGs tested for GWAS colocalization by ± 100 bps, and checked for overlap (≥ 1 bp) with the aggregated set of DNase-seq derived ENCODE5 open chromatin regions utilized to characterize mQTL signatures. Trait-linked mCpGs were classified as eQTL-shared or mQTL-specific (see Fig. 4). Enrichment significance of eQTL-shared or mQTL-specific trait-linked mCpGs in open chromatin regions was estimated at Fisher's exact test $P < 0.05$. To evaluate the methylation signatures of mQTL-GWAS colocalizations, for each tissue, we performed a Wilcoxon rank-sum test comparing DNAm levels of mCpGs tested for colocalization to those significantly colocalized ($RCP > 0.3$ and $PP4 > 0.3$). For the majority (8/9) of tissues, DNAm levels of colocalized mCpGs were significantly (Wilcoxon $P < 0.05$) lower than tested ones. We applied an analogous approach to eQTL-GWAS colocalizations, and observed an inverse pattern: for the majority (6/9) of tissues, expression levels of colocalized eGenes were significantly (Wilcoxon $P < 0.05$) higher than tested ones. Bootstrapped ($N = 5,000$ replicates) values for DNAm and gene expression means - averaging mCpGs and eGenes within each tissue - for all QTL-GWAS colocalization groups are displayed in Supplementary Figure 6; confidence intervals were computed using bootstrapping with replacement.

Integration of trait-linked genetically-regulated methylated loci with functional maps

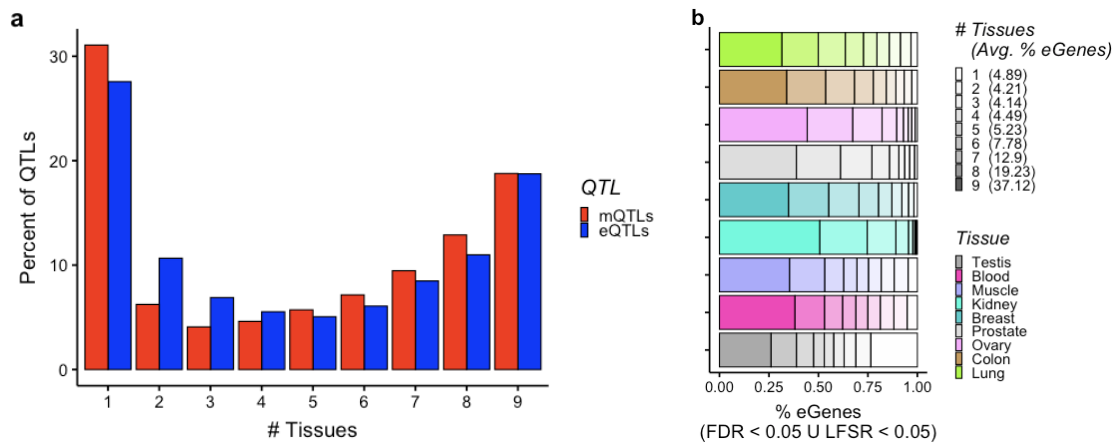
To identify genes involved in trait-linked mCpGs that co-located with gene regulatory elements, we integrated mQTL-derived colocalization results with curated promoter- and enhancer-gene target predictions^{38,39} and eQTM associations generated herein (Methods). We identified 68% (1,307/1,911) of mQTL-specific trait-linked mCpGs as co-located with enhancers and/or as eCpGs, across 61 GWASs and 1,129 GWAS hits and reported findings in

Supplementary Table 6. For 35% (400/1,129) of these loci, multiple mCpGs consistently support the same gene candidate(s). Among highly supported (by ≥ 3 mCpGs) cases, we identify poorly or not characterized gene-trait associations. For instance, the topmost supported instance corresponds to the RUNX1 locus associated with asthma. For the asthma GWASs analyzed, we observe 12 distinct mCpGs linked to RUNX1 regulatory regions. Given that other members of the RUNX transcription factor family are reported to play a role in asthma ⁶⁴, RUNX1 is a strong candidate to be involved in the etiology of the trait. Another well-supported case corresponds to the TMEM72 locus, associated with red blood cell counts, for which we identify 6 mCpGs linked to TMEM72 regulatory regions. The TMEM72 transmembrane protein is strongly and differentially expressed in ductal cells of the kidney ⁶⁵, which plays a major role in red blood cell homeostasis ⁶⁶.

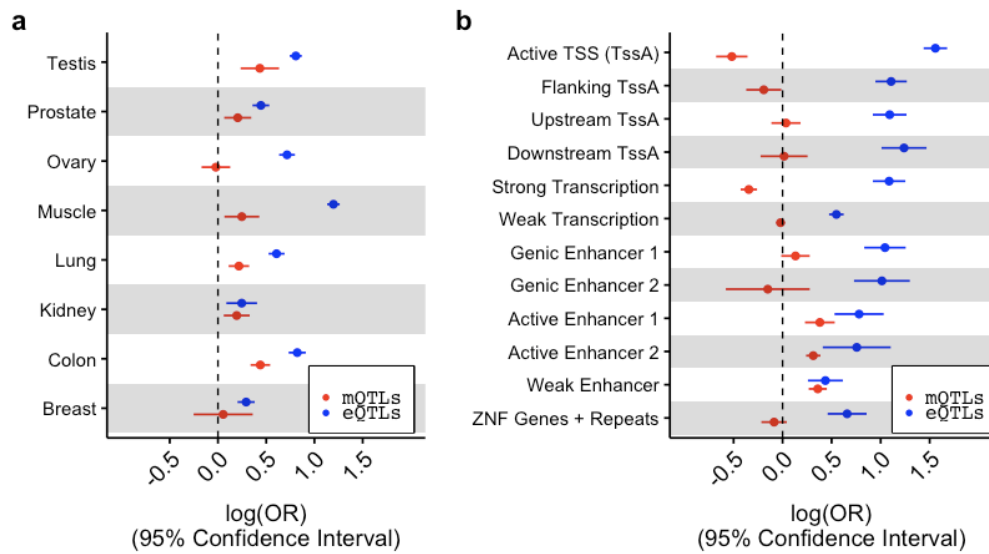
Supplementary Figures



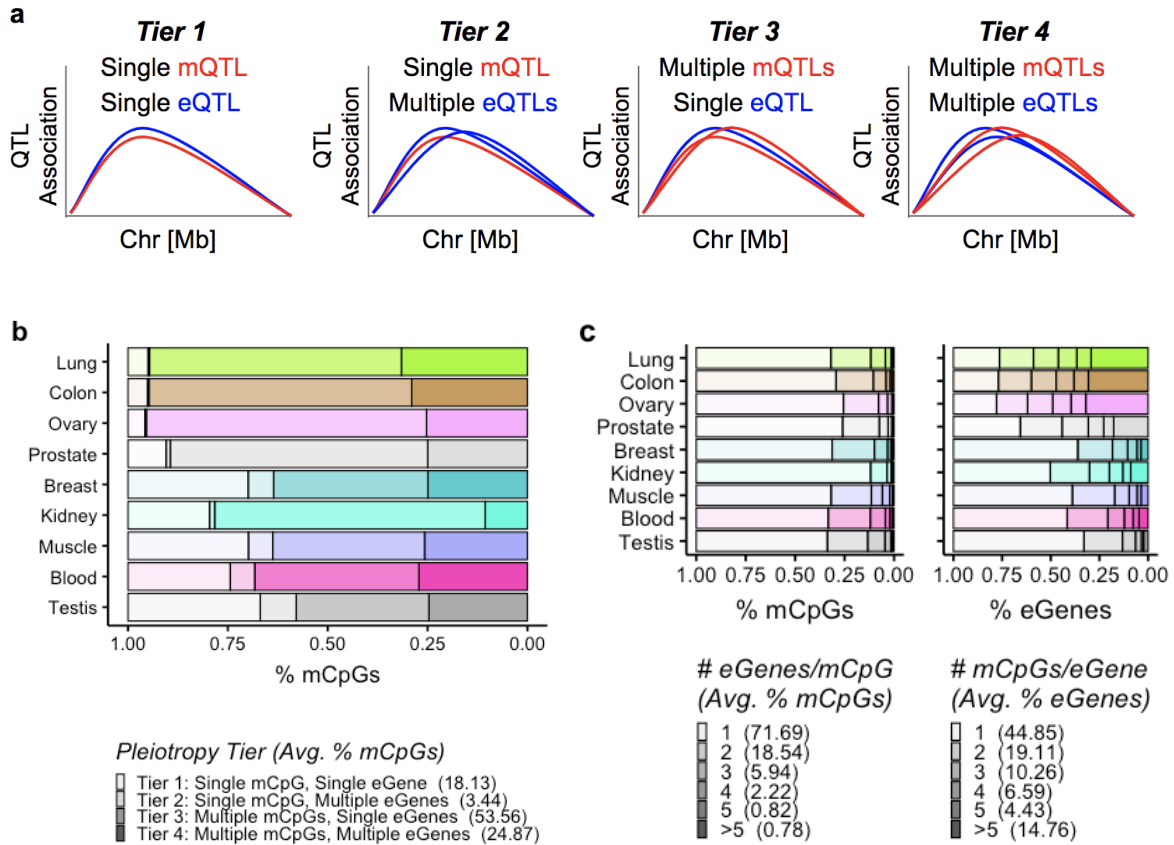
Supplementary Figure 1. eQTM discovery and tissue specificity patterns. (a) Number of eQTMs per tissue, shown with per-tissue eQTM-mapping sample sizes in parentheses. **(b)** Tissue sharing profile of eQTMs.



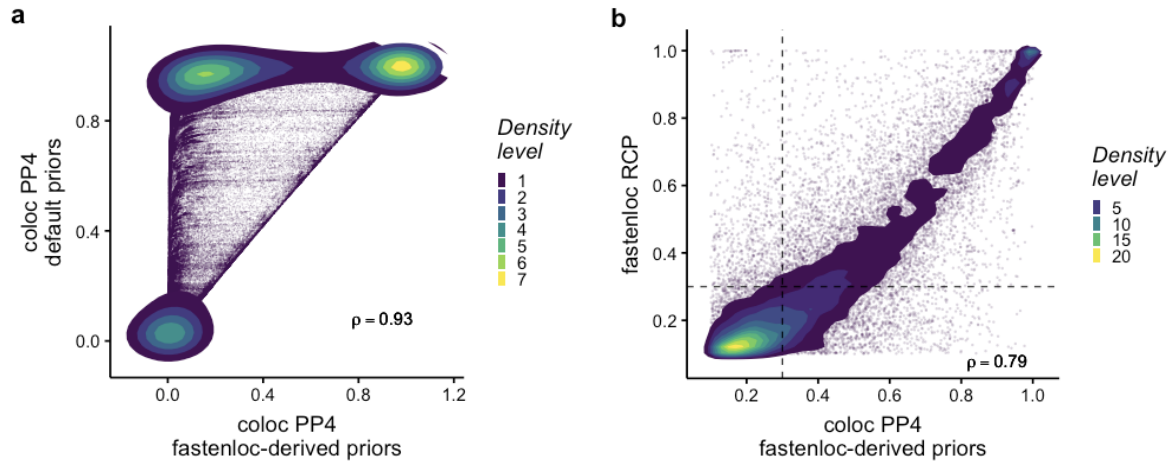
Supplementary Figure 2. Tissue specificity of QTLs. (a) Tissue sharing profile of mQTLs and eQTLs. **(b)** Cross-tissue sharing of eGenes. Cross-tissue average percent of eGenes per tissue-sharing category is shown in parentheses. Of note, testis is an outlier for tissue specificity, as 23.5% of eGenes were not detected in any other tissue.



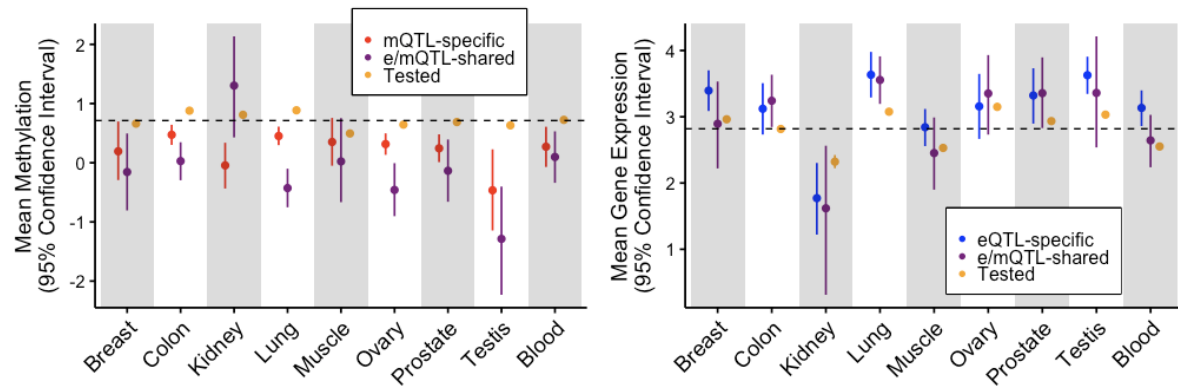
Supplementary Figure 3. Enrichment of QTLs in chromatin states. (a) QTL enrichment (x-axis) in tissue-matching open chromatin regions derived from ENCODE DNase-seq profiles per tissue (y-axis). Whole blood is excluded due to lack of a tissue-matching DNase-seq profile. Enrichment differences between tissues may be due in part to per-tissue DNase-seq data quality. **(b)** QTL enrichment (x-axis) in active chromatin states. OR: Odds Ratio.



Supplementary Figure 4. Characterization of mQTL pleiotropy. (a) Scheme of possible scenarios of eQTL-mQTL colocalization regarding QTL variants' pleiotropic effect on multiple mCpGs and eGenes. (b) Quantification of mQTL-eQTL pleiotropy per tier per tissue, in percent of mCpGs belonging to each tier. Tier details illustrated in (a). (c) Distribution of the number of eGenes per mCpG (left panel) and mCpGs per eGene (right panel) involved in mQTL-eQTL colocalization events, stratified by tissue.



Supplementary Figure 5. Evaluation of mQTL-GWAS colocalization approach. (a) Density plot of mQTL-GWAS colocalization scores based on *coloc* run with default (y axis) and *fastenloc*-derived priors (x axis) on UKB standing height GWAS; Spearman's rho is shown. Each dot corresponds to a colocalization test for a particular GWAS hit, independent mQTL and tissue combination. (b) Density plot of mQTL-GWAS colocalization scores based on *coloc* (x axis) and *fastenloc* (y axis) approaches on all GWASs; Spearman's rho is shown. Each dot corresponds to a colocalization test for a particular GWAS, GWAS hit, independent mQTL and tissue combination. Dots within the top-right quadrant correspond to significant (RCP > 0.3 and PP4 > 0.3) colocalizations.



Supplementary Figure 6. Signatures of trait-linked QTLs. DNAm - in M-values - of mCpGs (left panel) and gene expression - in $\log_2(\text{TPM}+1)$ - of eGenes (right panel) tested for colocalization, stratified by tissue and colocalization group (see Fig. 4). Average DNAm and gene expression across tissues is indicated by dashed line.