

The American Journal of Human Genetics, Volume 95

Supplemental Data

Transcriptome Sequencing of a Large Human Family Identifies the Impact of Rare Noncoding Variants

Xin Li, Alexis Battle, Konrad J. Karczewski, Zach Zappala, David A. Knowles, Kevin S. Smith, Kim R. Kukurba, Eric Wu, Noah Simon, and Stephen B. Montgomery

FIGURE S1. FAMILY STRUCTURE.	4
FIGURE S2. FLOWCHART OF GENOTYPE CALLING AND RNA-SEQ QUALITY CONTROL STEPS.	5
FIGURE S3. HAPLOTYPE / IDENTITY-BY-DESCENT (IBD) INFERENCE.	6
FIGURE S4. HISTOGRAM OF HAPLOTYPE LENGTHS.....	7
FIGURE S5. DISTRIBUTION CIS-EQTL AND CIS-SQTL VARIANTS NEAR A GENE: LOCAL HAPLOTYPE BLOCKS HAVE THE LARGEST NUMBER OF EQTL (LEFT) OR SQTL (RIGHT) EFFECTS.	8
FIGURE S6. COMPARISON OF EQTL DISCOVERY BETWEEN RNA-SEQ AND MICROARRAY.	9
FIGURE S7. CONCORDANCE OF EQTL EFFECT SIZES (B) BETWEEN RNA-SEQ AND MICROARRAY.....	10
FIGURE S8. IDENTIFYING LARGE-EFFECT CIS-EQTL GENES IN FAMILY COMPARED TO POPULATION.....	11
FIGURE S9. EFFECT SIZE MEASURED BY THE ONE-REGRESSOR AND THE TWO-REGRESSOR MODEL.	13
FIGURE S10. LARGE-EFFECT FAMILY CIS-EQTL GENES.....	14
FIGURE S11. OVERLAP OF B AND FIT (R^2) EFFECT SIZE OUTLIERS.....	15
FIGURE S12. LARGE RELATIVE B VS. ABSOLUTE B.....	16
FIGURE S13. EFFECT OF DIFFERENT QUANTIFICATION PIPELINES: COMPARISONS OF EFFECT SIZE B BETWEEN TOPHAT + CUFFLINKS AND GEM + FLUX PIPELINES.	17
FIGURE S14. ENRICHMENT OF RARE VARIANTS AT LARGE EFFECT SIZE B.	18
FIGURE S15. INFLUENCE OF DISCOVERY SAMPLE SIZES IN TAGGING CAUSAL SNPs.	19
FIGURE S16. INFLUENCE OF DIFFERENT CRITERIA IN SELECTING BEST SNP: SMALLEST P-VALUE OR LARGEST EFFECT SIZE.....	20
FIGURE S17. ADJUSTMENT OF EFFECT SIZE EMPIRICAL P-VALUES: COMPARISON OF EFFECT SIZE CONFIDENCE INTERVALS (NOISE LEVELS) BETWEEN THE FAMILY AND THE POPULATION.....	21
FIGURE S18. ADJUSTMENT OF EFFECT SIZE EMPIRICAL P-VALUES: DISTRIBUTION OF P-VALUES OF FAMILY VERSUS POPULATION EFFECT SIZES.	23
FIGURE S19. ADJUSTMENT OF EFFECT SIZE EMPIRICAL P-VALUES: COMPARISONS OF EMPIRICAL P-VALUE AND WELCH'S T-TEST.	25
FIGURE S20. CORRELATION OF EQTL EFFECT SIZE B AND ASE EFFECT SIZE (ALLELIC IMBALANCE).....	26
FIGURE S21. ENRICHMENT OF ASE EFFECTS AT LARGE-EFFECT GENES.	27
FIGURE S22. RARE REGULATORY VARIANTS CONTRIBUTING TO LARGE-EFFECT EQTLs: ENRICHMENT OF RARE VARIANTS NEAR THE TSS OF LARGE-EFFECT (B) CIS-EQTL GENES, COMPARING ANNOTATED AND ALL RARE VARIANTS.	28
FIGURE S23. RARE REGULATORY VARIANTS CONTRIBUTING TO LARGE-EFFECT EQTLs.	29
FIGURE S24. ENRICHMENT OF RARE REGULATORY AT LARGE EFFECT GENES.	30
FIGURE S25. IDENTIFICATION OF LARGE ASE EFFECT.....	31
FIGURE S26. MENDELIAN SEGREGATION OF ALTERNATIVE SPLICING PATTERNS.....	33
FIGURE S27. EXAMPLES OF ALTERNATIVE SPLICING PATTERNS DETERMINED BY HAPLOTYPE GROUPS.....	34
FIGURE S28. ASE HERITABILITY ANALYSIS.....	35
FIGURE S29. ALLELIC RATIO CORRELATION WITH SIBLINGS, USING NA12879 AS REFERENCE.	36
FIGURE S30. ALLELIC RATIO CORRELATION BETWEEN DIFFERENT TYPES OF SIBLINGS.	37

TABLE S1. NUMBER OF VARIANTS SEGREGATING IN THE FAMILY.	38
TABLE S2. GENOTYPES CONFIRMED WITH COMPLETE GENOMICS LONG FRAGMENT READ ¹ (LFR).	39
TABLE S3. PHASING CONFIRMED WITH MOLECULAR HAPLOTYPE BY LFR.	40
TABLE S4. GENOTYPES CONFIRMED WITH ILLUMINA PLATINUM GENOMES.	41
TABLE S5. LINKAGE ANALYSIS OF <i>cis</i> -EQTL: SUMMARY OF EQTL AND SQTL GENES IDENTIFIED IN THE FAMILY..	42
TABLE S6. EFFECT OF DIFFERENT DISCOVERY PANEL SIZES: NUMBER OF LARGE EFFECT B GENES GIVEN DIFFERENT DISCOVERY PANEL SIZES.	44
TABLE S7. PREDICTION OF EQTLs AT RARE VARIANTS GIVEN ANNOTATION: PROPORTION OF GENES BEING AN EQTL GIVEN A REGULATORY VARIANT NEAR TSS.	46
TABLE S8. EXAMPLES OF SQTL GENES.	47
TABLE S9. FAMILY-SPECIFIC <i>cis</i> -EQTL MODIFYING COMPLEX TRAIT GENES.	49
TABLE S10. EXAMPLES OF RARE REGULATORY VARIANTS INFLUENCING GWA GENES.	50

CEPH/Utah Pedigree 1463

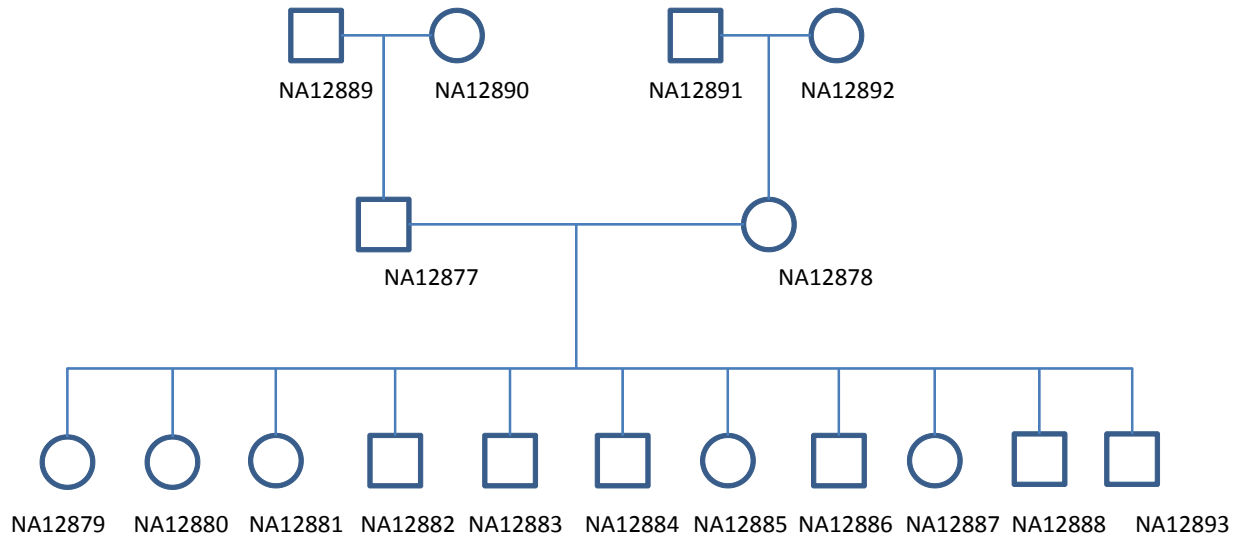


Figure S1. Family structure.

Four grandparents, two parents and eleven children. All family members are RNA-sequenced. Whole genome DNA-sequencing data of all family members were generated by Complete Genomics. Whole genome sequencing was also performed again by both Illumina Platinum Genomes and Complete Genomics Long Fragment Read¹ technology. All three sets of genome sequencing data are compared to confirm genotyping correctness (Table S2, Table S3, Table S4).

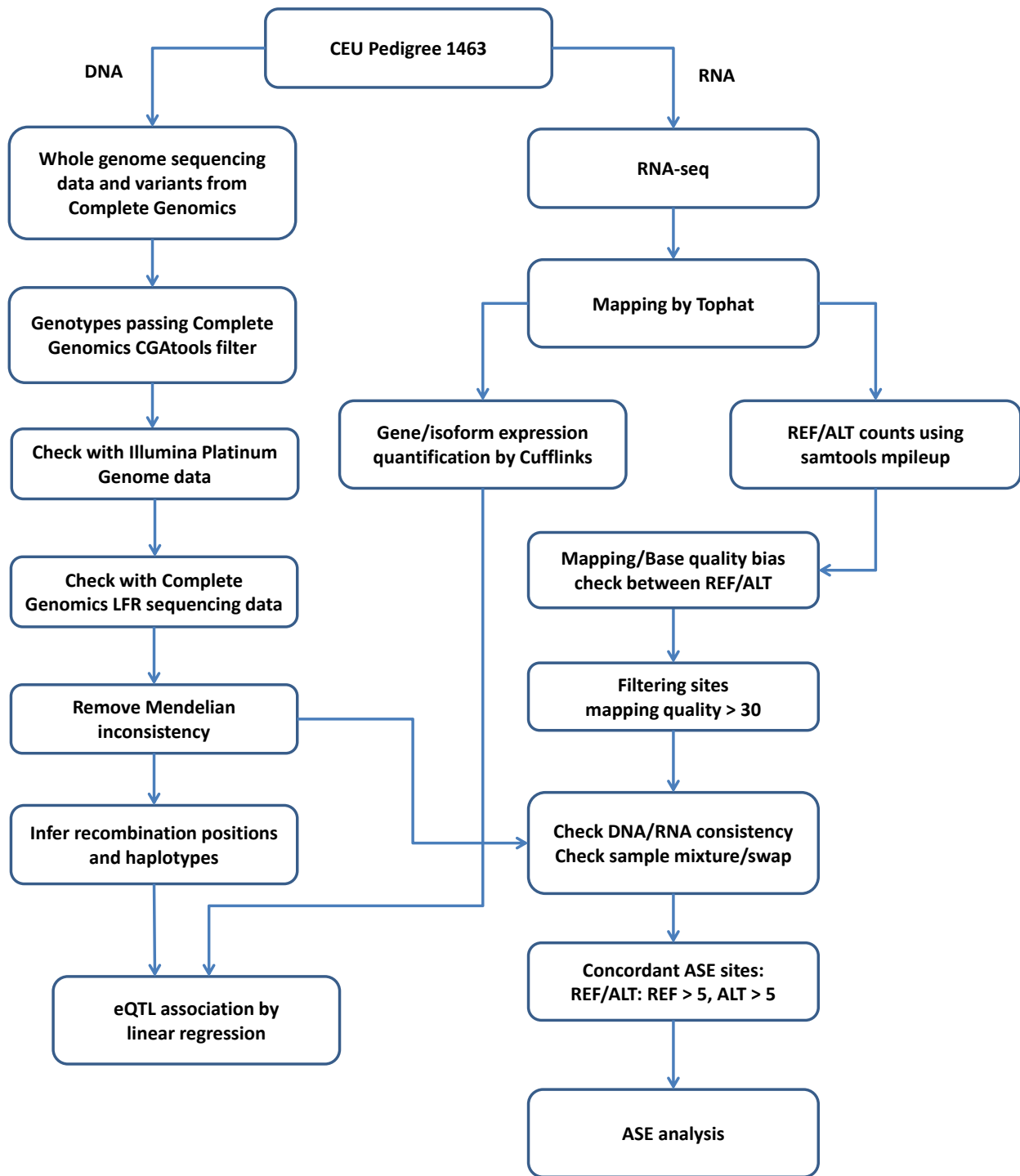


Figure S2. Flowchart of genotype calling and RNA-Seq quality control steps.

Genotyping data were confirmed across three sequencing platforms (Table S2, Table S3, Table S4) to guarantee correctness especially at rare variants. We further filtered variants by stringent Mendelian consistency throughout the whole family. RNA/DNA concordance was checked at heterozygous sites to avoid sample mixture/swap.

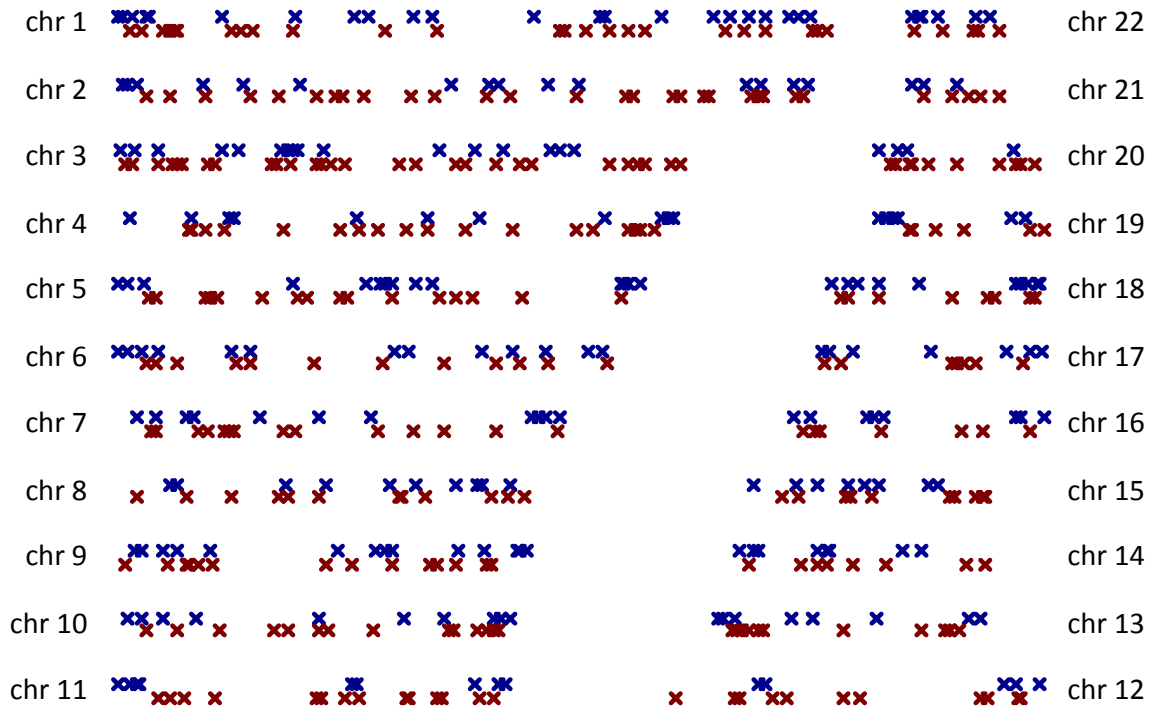


Figure S3. Haplotype / identity-by-descent (IBD) inference.

Distribution of recombination breakpoints. Red: maternal recombinations, Blue: paternal recombinations. We inferred 813 recombination positions in CEU family 1463. We partition chromosomes into haplotypes according to these recombination positions.

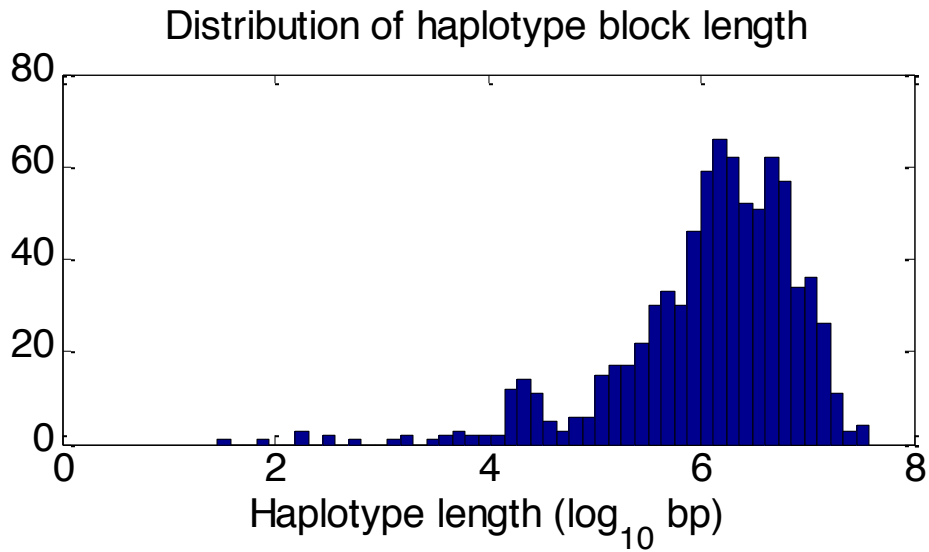


Figure S4. Histogram of haplotype lengths.

Haplotype blocks are defined by recombination positions as shown in Figure S3. Majority of haplotype blocks are long enough to include the most intensive cis-regulatory regions of a gene (100kb near TSS). The median haplotype length is 1.65Mb, and 90% of haplotype blocks range from 0.02Mb to 12Mb.

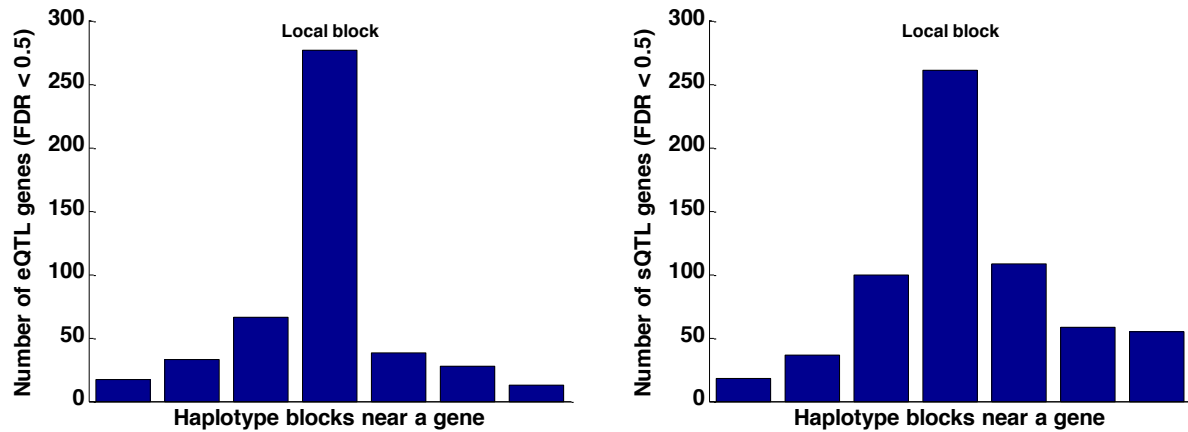


Figure S5. Distribution cis-eQTL and cis-sQTL variants near a gene: Local haplotype blocks have the largest number of eQTL (left) or sQTL (right) effects.

Compared with three nearby blocks, local haplotypes shows substantially larger number of eQTL / sQTL effects, compared to up and downstream haplotype blocks. We tested the local haplotype block containing each gene and three nearby haplotype blocks for eQTL linkage. As we expected, local haplotypes that contain the tested gene show the largest number of eQTL associations compared with nearby blocks. Local haplotypes also show largest number of sQTL associations compared with nearby blocks. The result suggests that most *cis*-acting expression or splicing QTL variants are located in the local haplotype blocks.

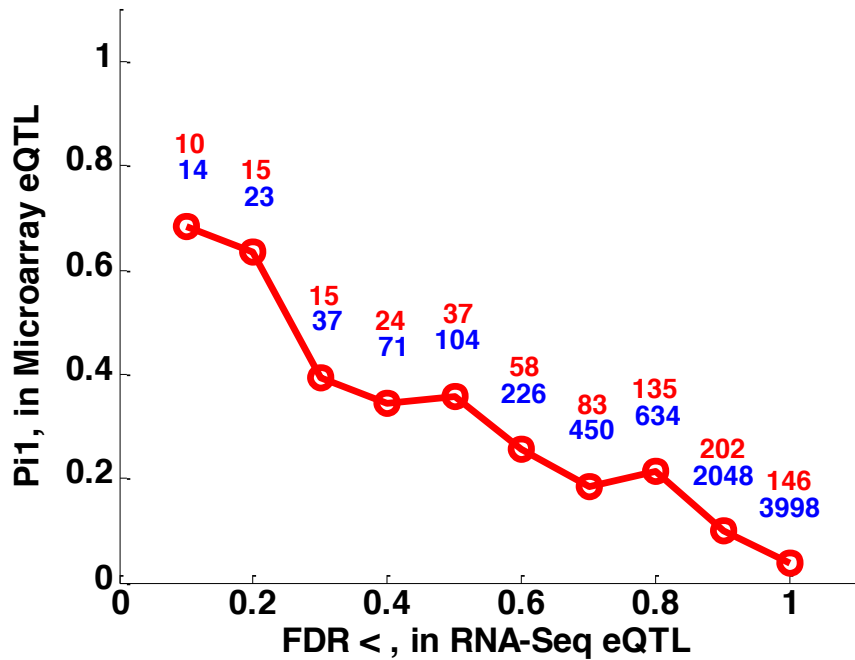


Figure S6. Comparison of eQTL discovery between RNA-Seq and microarray.

We tested eQTLs within the same family quantified in published microarray studies² (only seven of the siblings are available from this microarray data). We measured the number of eQTLs from RNA-Seq data that can also be detected using microarray. Blue numbers are total number of eQTL genes detected by RNA-Seq passing that FDR cutoff, red numbers are number of genes also showing eQTL effects by microarray as indicated by π_1 ³. We can observe that given more stringent FDRs that the two approaches give more concordant discoveries. Furthermore, both eQTL discoveries (Figure S6) and effect sizes (Figure S7) show concordant patterns between the two studies.

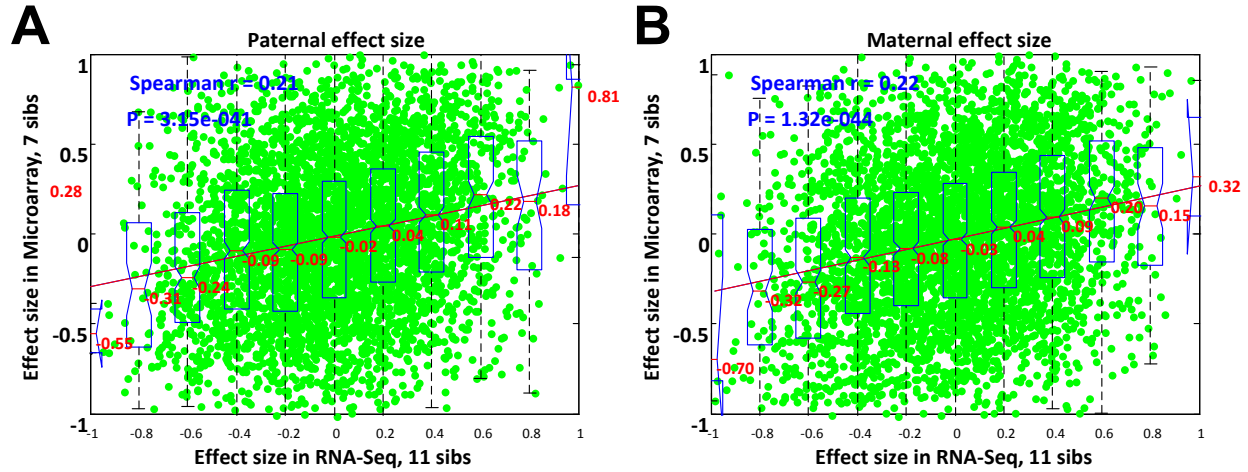


Figure S7. Concordance of eQTL effect sizes (β) between RNA-Seq and microarray.

(A) Paternal effect sizes. (B) Maternal effect sizes. We report effect sizes of eQTL as measured from RNA-Seq data or microarray data. Sign of effect size indicates whether the paternal haplotype of a parent (father or mother) increases or decreases expression in children. Red numbers are medians of each box. Effect sizes measured between microarray and RNA-Seq quantification are modestly concordant.

$$T_i = \mu + \beta_j p + \beta_k m$$

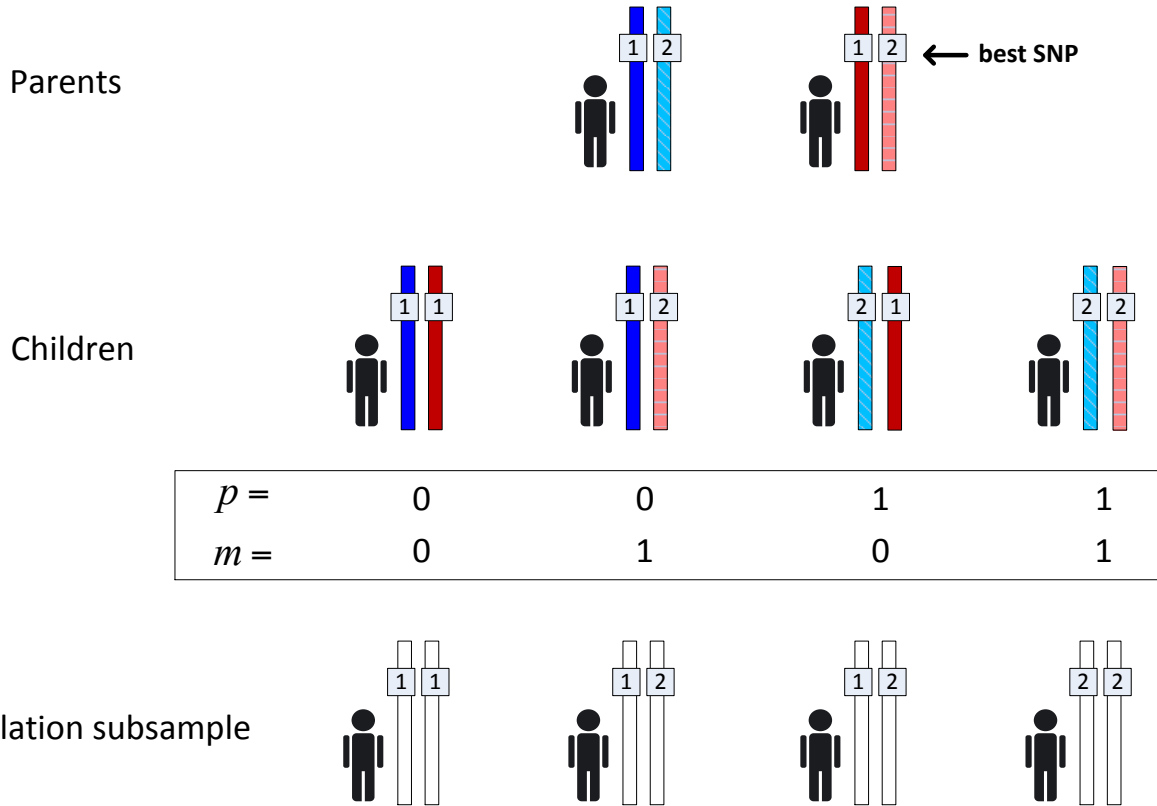


Figure S8. Identifying large-effect cis-eQTL genes in family compared to population

Effect sizes measured in the family are compared to those measured in genotype-matched population subsamples. We have two β s for both the family and the population data to avoid effect size inflation due to more regressors in the family than the population. We use the same regression to measure effect sizes for both the family and the population data: $T_i \sim \mu + \beta_j p + \beta_k m$, p, m are two regressors indicating paternal and maternal haplotypes in the family. We can use the same regression formula because genotypes are matched exactly between family and population subsample at the best associated SNP, so the two regressors p, m match segregating patterns of the best SNP in both the family and the population subsample. If we assume only the best SNP is functioning in both the family and the population samples, β_j, β_k are expected to be the same between the family and the population subsample. For population heterozygotes (with identical, unphased SNP genotypes) the maternal and paternal alleles are assigned arbitrarily from the two possible options, as needed, to match family genotypes. However this extra information does not influence the measure of effect sizes on either side.

Subsequently, large-effect outlier genes are identified by comparing effect size (β or R^2) of genes in the family to those in the population. Effect sizes (β) can be directly compared using analytical tests (Welch's t). However, to explore the behavior of effect sizes under different

sample sizes, we applied a subsampling approach among the population individuals to re-generate the expected effect size distribution of the best associated SNP among 11 individuals.

In specific, the best eQTL SNP is discovered in a separate 180-individual discovery panel to avoid bias of multiple selections (a phenomenon otherwise known as regression to the mean or winner's curse). Effect sizes of genes in the population are then assessed by subsampling the same number ($N=11$) of individuals from the 193-individual replication panel (of 373 Geuvadis European samples) to account for potential biases due to different sample sizes. β or R^2 are regression slope and coefficient of determination (fit) measured by linear regression. The method is illustrated in Figure S8. We then generated the population effect size distribution by subsampling down to 11 matched individuals multiple (100) times from the population data. For each gene, we generated an empirical p -value of observing a larger effect in the family, by counting how often the effect size in the family is larger than those from the population subsamples. We estimated the total number of genes exceeding population effect sizes using the π_1 statistic, based on empirical p -values of all genes. The estimated π_1 for large-effect cis-eQTLs is 0.0611 (for eQTLs selected by R^2). For β effect size estimates, we noticed there is difference in noise levels between the population and the family estimates which result in a p -value distribution significantly skewed towards 1; we discuss an adjustment for such differences in Figure S17, Figure S18 and Figure S19.

We do not use β and only use R^2 when comparing splicing QTL effect sizes, since the estimation noise of β is too large for transcript ratios in Geuvadis data. The estimated π_1 for large-effect cis-sQTL is 0.0749.

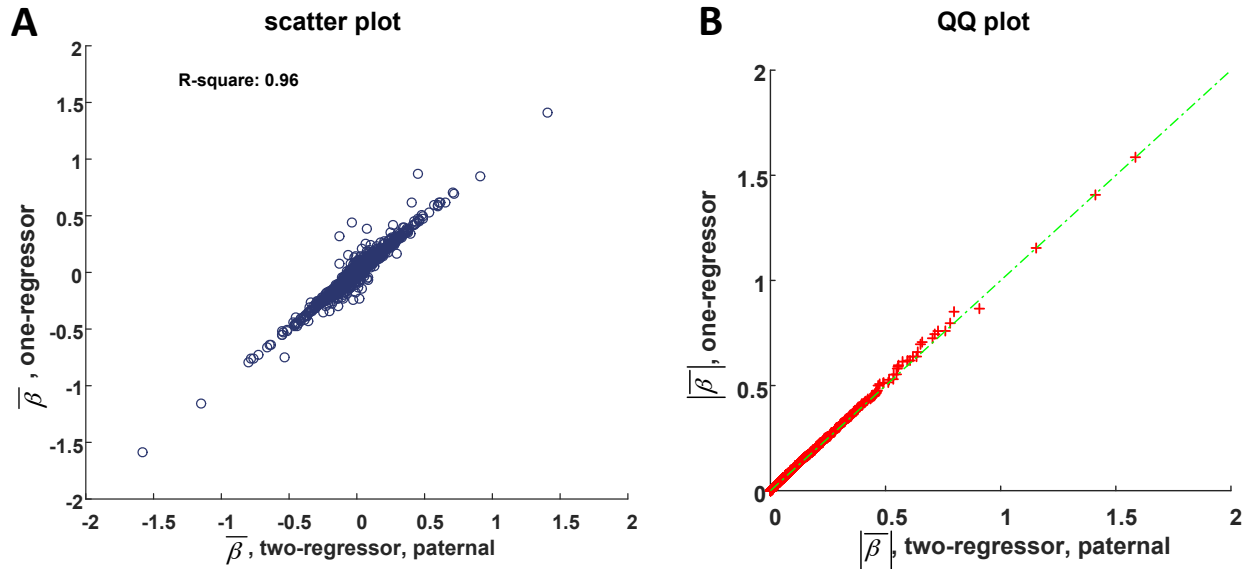


Figure S9. Effect size measured by the one-regressor and the two-regressor model.

For the family, in order to measure the effect of the whole haplotype, we need to use a two-regressor model. For the population subsample, however we chose to use the same two-regressor model to avoid possible effect size (β or R^2) inflation caused by the use of more regressors. This makes a fairest comparison between the family and the population subsample as they are now measured on exactly the same model with the same number of regressors, and effect sizes differences are truly due to biological factors specific to the family instead of different regression methods.

In order to match the regressors of the family, we actually implicitly phased the SNP of the population subsample according to the family (Figure S8), this information is arbitrary for the population subsample however this arbitrary splitting of one regressor into two regressors does not actually influence the measure of effect sizes. The two regressors, p and m , which indicate transmissions from either parents are statistically independent of each other or, in terms of linear relationships, orthogonal: $p \perp m$, $E(p \cdot m) = 0$, therefore each will capture their own effect without interfering with one another. Figure S9 shows the comparison of the actual effect sizes measured by the one-regressor and two-regressor models, which verifies that the two-regressor model unbiasedly captures the same effect sizes ($\bar{\beta}$) as the one-regressor model.

Panel A shows effect sizes measured using the mean from 100 subsamples out of the population. Using the one-regressor model, we simply regress on a single SNP (considering 00, 01, 10, 11 to be 0, 1, 1, 2). Using the two-regressor model, we regress on both the paternal haplotype and the maternal haplotype. Here, we show the effect size β measured on the paternal haplotype, oriented according to the SNP's phase on the father (to match the sign of the SNP β). Panel B shows the QQ plot of effect sizes comparing effect sizes (absolute values) measured by the one-regressor and the two-regressor models. However, the dispersion ($var(\hat{\beta})$) of a two-regressor model is expected to be substantially larger and is a reason why we emphasize the usage of the same model between the family and the population.

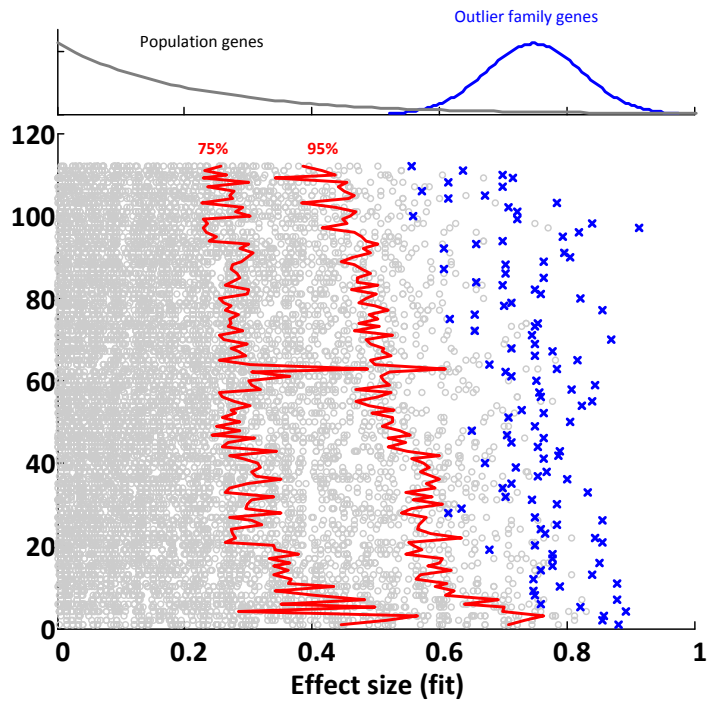


Figure S10. Large-effect family *cis*-eQTL genes.

Effect sizes are compared to population by fit (R^2) of linear regression. Shown are family eQTL genes (blue) with effect sizes greater than the 0.99 quantile (empirical p -value < 0.01) of population effect sizes (grey). The magnitude of the outlier proportion has been extended on the top to illustrate the range of effect sizes for measured large-effect *cis*-eQTLs.

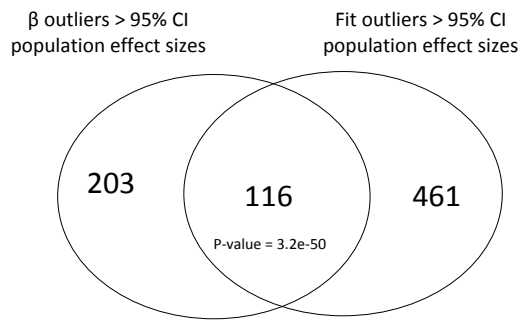
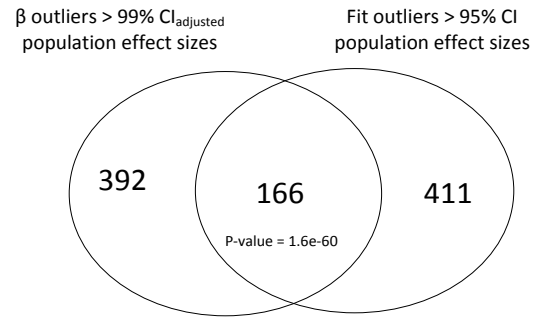
A**B**

Figure S11. Overlap of β and fit (R^2) effect size outliers.

(A) There are 319 β outlier genes and 577 fit outlier genes with effect sizes greater than 95% quantile (empirical p -value < 0.05) of the population. The overlap is 116 genes. The overlap is statistically significant by Fisher's exact test, indicating shared effects they are capturing. (B) After adjustment of confidence intervals of β (described in Figure S17-Figure S19), there are 558 β outlier genes with effect sizes greater than 99% CI_{adjusted} of the population. The overlap with 577 (> 95% CI) fit outlier genes is 166 genes.

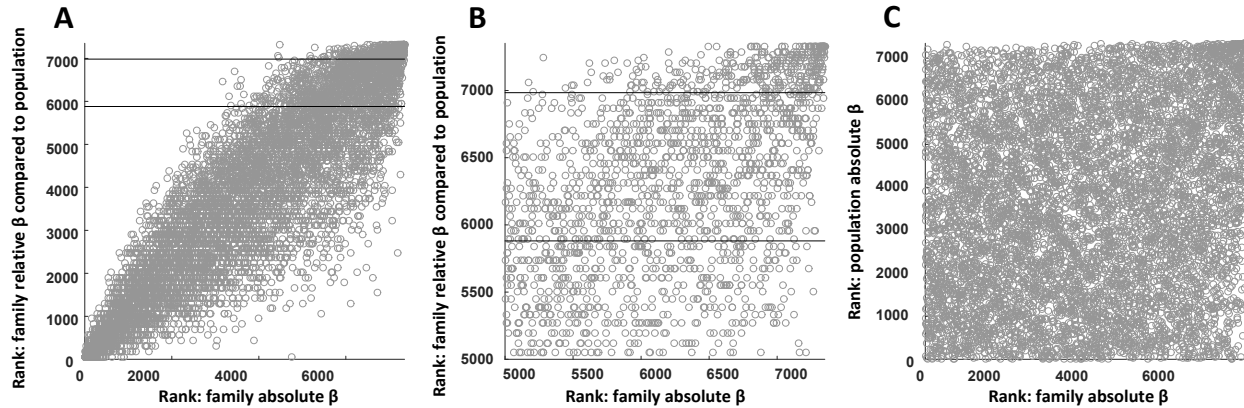


Figure S12. Large relative β vs. absolute β .

We tested the properties of large-effect genes compared to the population to see whether the comparison to the population adds additional information to the ranking of genes. Figure S12 shows that large absolute β are not necessarily highly ranked in the relative scale and vice versa. This does indicate that we are gaining novel information by ranking genes according to their relative effect sizes (empirical p -values) instead of just ranking them by their absolute β . (A) Ranks of absolute β , compared to ranks of relative β . Relative β is the empirical p -values comparing the family with the population effect sizes. Absolute β is just the original effect size β yielded by the linear regression. (B) Zoom-in of upper right rectangle of (A). The two lines indicate the top 5% and 20% of genes. The figure shows that by comparing to the population, the ranks of genes are not the same as simply ranking the genes by their original β . Of the top 5% of genes by each metric, the overlap is 52%. (C) Ranks of family absolute β compared to ranks of population absolute β . Family effect sizes and population effect sizes are not correlated in general.

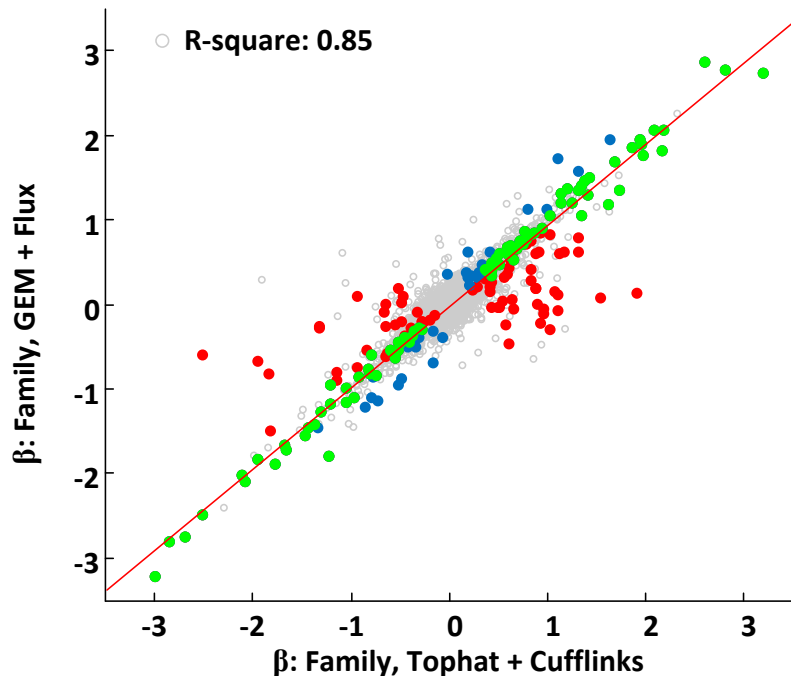


Figure S13. Effect of different quantification pipelines: comparisons of effect size β between Tophat + Cufflinks and GEM + Flux pipelines.

Effect sizes are highly correlated between two quantification methods. Here we plot only paternal side β , maternal side β patterns are very similar. Discovery of large effect size genes (> 95 CI of population, paternal side only) are: 165 genes by Tophat + Cufflinks (red), 125 genes by GEM + Flux (blue) and 90 genes of their intersection (green). Geuvadis expression values were based on a different quantification pipeline than used in the family data. To exclude the possibility that large-effect eQTL genes are due to technical differences between the family and population data, we compared discovery of large effect genes and enrichment of rare variants in the family using the same pipeline (GEM + Flux) as Geuvadis. For the family data, the effect size estimates are highly correlated between two pipelines. We observe a similar discovery set of large effect genes and also similar patterns of rare variant enrichment as Tophat + Cufflinks pipeline (Figure S13, Figure S14). The confidence intervals used for β effect sizes in Figure S13-Figure S16 are raw (if not otherwise specified), without further adjustment (described in Figure S17-Figure S19).

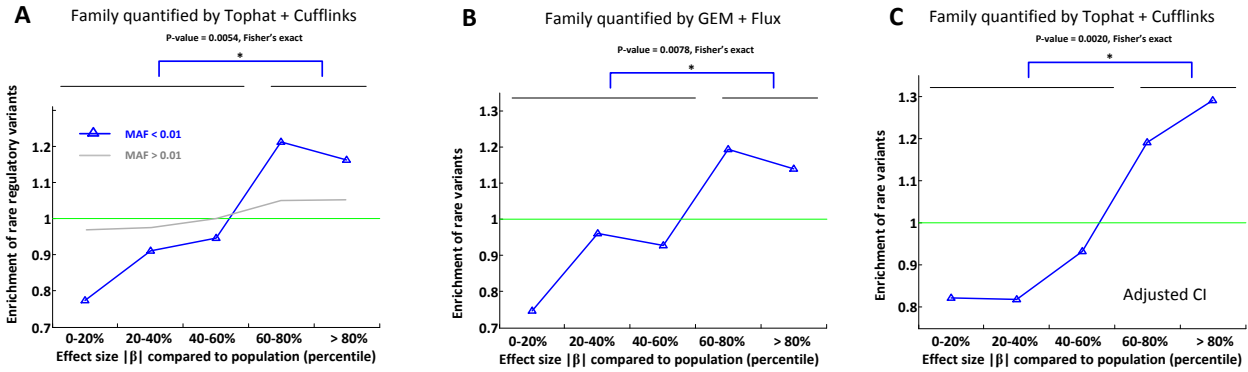


Figure S14. Enrichment of rare variants at large effect size β .

The enrichment pattern is similar when using GEM + Flux (A) pipeline compared to Tophat + cufflinks pipeline (B) and (C). X-axes in (A) and (B) are raw CIs, (C) is adjusted CI. We ranked effect sizes of genes based on 1 – their empirical p -values: how often their effect sizes in the family are larger than effect sizes among the population subsamples. The distribution of effect sizes in the population was generated by repeatedly subsampling 11 individuals from the population. Rare variants are defined as those with MAF < 0.01, within Encode TF binding + DNase peaks and PhyloP score > 1. Enrichment is defined as proportions of genes with such an annotated rare variant in each effect size bin divided by proportions of genes with such an annotated rare variant across all effect size bins.

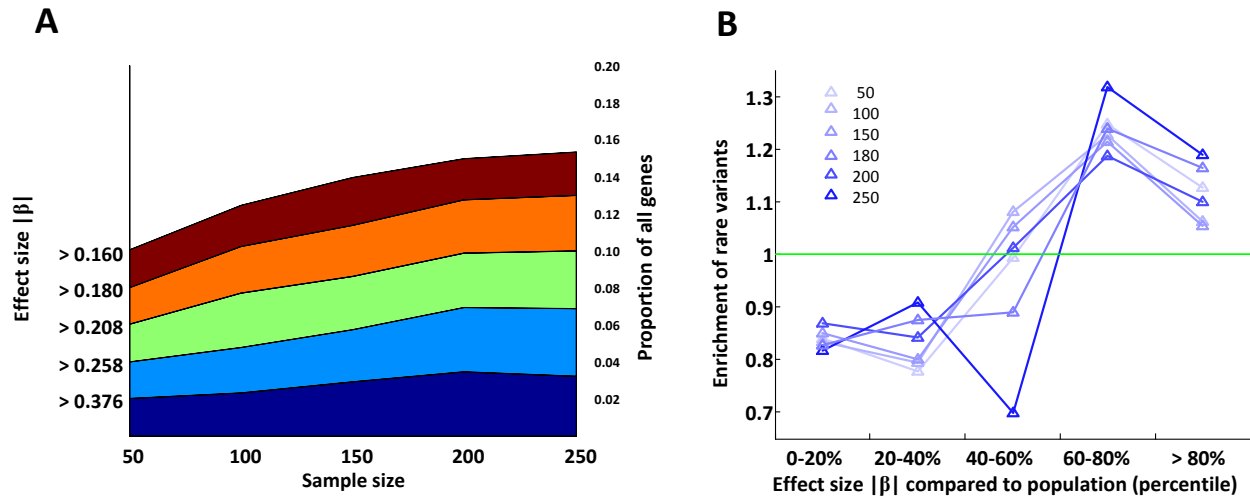


Figure S15. Influence of discovery sample sizes in tagging causal SNPs.

As we are using whole genomes for both the family and the population data, we have the ability to test all SNPs including the causal SNP in our association models. However, a general rule of statistics indicates that the power to capture the true largest effect or causal SNP in the population also depends on sample sizes. A small discovery panel may result in poor choice of SNP and deflation of effect sizes in the population, which can potentially over-estimate the number of large effect genes in the family. We analyze how sample sizes of the population discovery panel influence our identification of large effect genes in the family. We observe continuously increasing number of large effect eQTLs discovered in the population given larger sample sizes (Figure S15A), which indicates that large sample sizes do increase chance of tagging a true causal SNP. We consider this a very important effect suggesting the necessity of large sample sizes to accurately measure effect sizes. However, given our particular application, as largest effect genes are likely to saturate first, increasing sample sizes does not have a significant influence on our discovery of family large effect genes (Table S6). The enrichment of rare variants at large-effect genes is also comparable given different discovery panel sizes (Figure S15B).

(A) Number of large effect genes discovered in the population given larger sample sizes. Best SNPs are discovered in the discovery panel of varied size. We re-measured the effect sizes of those SNP in the replication panel. There are increased numbers of large effect SNPs discovered given a larger discovery panel size. Note that y-axis is piled inversely, with largest effect sizes stacked at the bottom. (B) Enrichment of rare variants at large effect family genes given different population discovery panel sizes. We ranked effect sizes of genes based on 1 - their empirical p -values: how often effect sizes in the family are larger than those of the population subsamples. Rare variants are defined as those with $MAF < 0.01$, within Encode TF binding + DNase peaks and PhyloP score > 1 . Enrichment is defined as proportions of genes with such an annotated rare variant in each effect size bin divided by proportions of genes with such an annotated rare variant across all effect size bins.

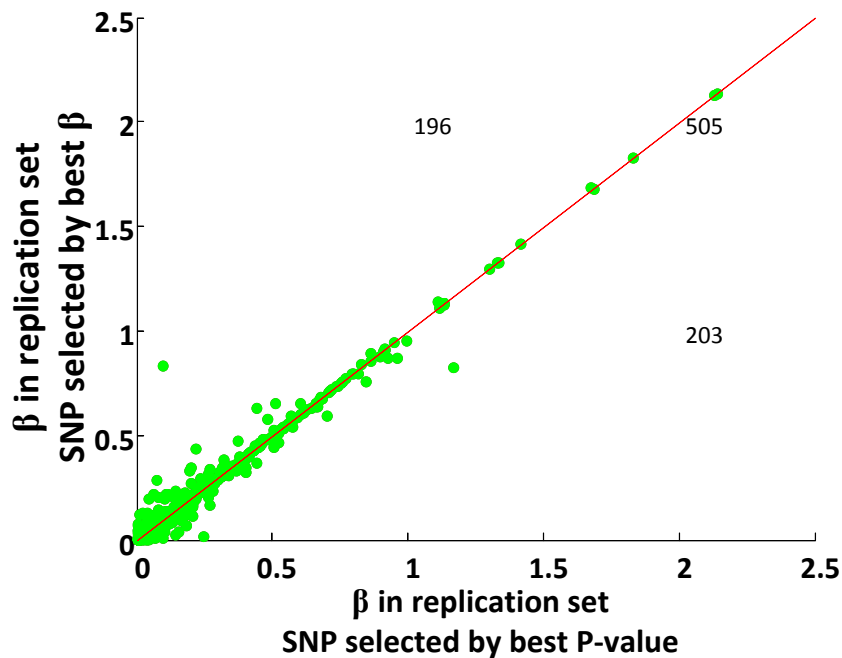


Figure S16. Influence of different criteria in selecting best SNP: smallest p -value or largest effect size.

X axis: re-measured β in replication panel if SNP is selected by best P value in the discovery panel. Y axis: re-measured β if SNP is selected by best β in the discovery panel. Shown here are 904 genes with a best SNP $< 1e-5$ in discovery panel. For each gene, we choose a best β among all SNPs within $10 * p$ -value of the best SNP. 203, 196 and 505 are genes with $X > Y$, $X < Y$ and $X = Y$.

For each gene, we select the best SNP based on p -value in the discovery panel. However, this SNP is not necessarily truly the largest effect SNP as p -value is an indicator of best fit (R^2) instead of largest β , such that we may possibly miss a secondary effect SNP with larger β . To test the possibility that we miss larger effect SNPs and under-estimate the effect sizes in the population, we analyzed the differences of choosing the best SNP by p -value or β . For each gene, we first find a best SNP with smallest p -value, we then pick another SNP of largest β (could be the same one) among all SNPs with a p -value no more than an order of magnitude less significant than that of the best p -value, we then re-measure their effect sizes in the replication set. We observe that approximately half of the time, the best p -value and best β SNP is the same SNP. Further, even when they are not the same SNP, the measured effect size in the replication set is very similar. This suggests that most effect size differences near the best SNP is due to random noise, the existence of a secondary effect SNP with even larger effect size is not significant.

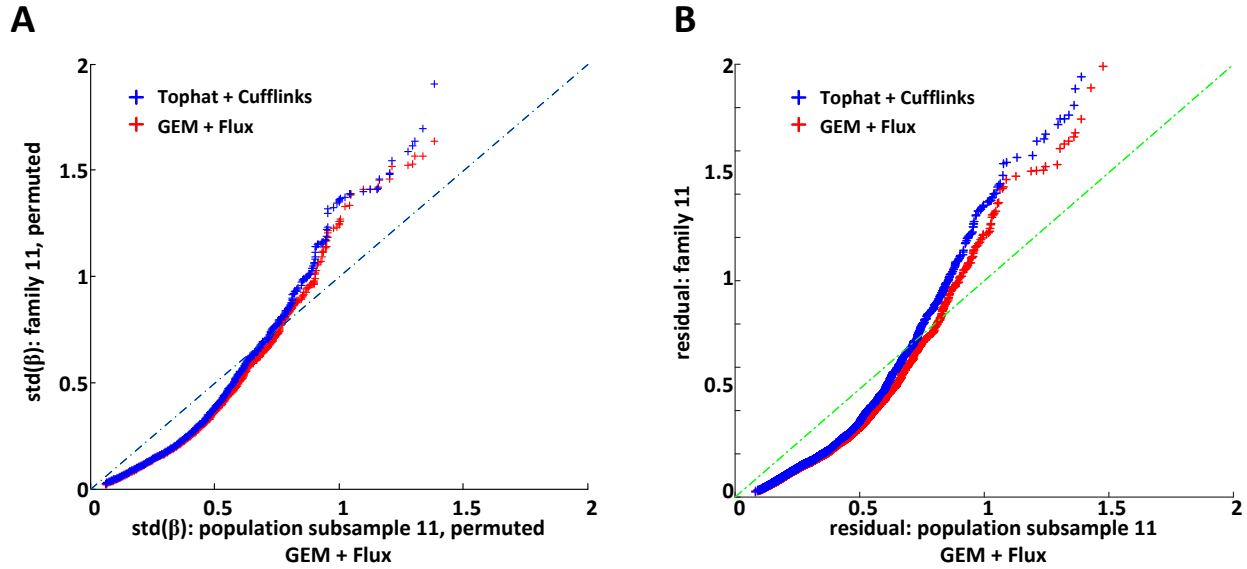


Figure S17. Adjustment of effect size empirical p -values: comparison of effect size confidence intervals (noise levels) between the family and the population.

(A) QQ plot of standard deviation of measured effect size $\hat{\beta}_i$, comparing the family and the population. Data are permuted such that $E(\hat{\beta}) = 0$, $std(\hat{\beta})$ measures the noise levels of measures of β . Population (Geuvadis) samples are quantified using GEM + Flux pipeline. Family data are quantified using both GEM + Flux and Tophat + Cufflinks pipelines respectively. Standard deviation of $\hat{\beta}$ is significantly larger in the population no matter which pipeline is used in family. (B) QQ plot of analytical estimation of standard error of β : \hat{S}_β , comparing the family and the population. \hat{S}_β is computed by $\frac{1}{\sqrt{\sum(x_i - \bar{x})}} \sqrt{\frac{\sum \varepsilon_i^2}{11-3}}$, where ε_i is the residuals in linear regression $T_i \sim \mu + \beta_j p + \beta_k m + \varepsilon_i$, x is either p or m .

The degree of noise in estimated β is different between the population and the family even if we match the sample size and quantification methods. The noise in estimated effect sizes is significantly smaller in the family than in the population. This difference reflects both the fact that family members are more homogeneous (sharing more covariates such as genetic, environment, lifestyles and etc., thus having tighter fit to the regression slopes) and also the possible existence of other technical factors, which we cannot tell apart.

The discovery of such differences is actually biologically informative, however our subsampling scheme is not intended to reflect and calibrate the noise of effect size estimates of the family members. Therefore, regardless of the source of these differences, they have undesirably shifted the empirical p -values (see comparisons to analytic p -values Figure S18A and C, Figure S19A).

Here, we explored two methods to adjust different noise levels between effect size estimations, which yield empirical p -values closer to analytical p -values and less conservative estimates of FDR. However as there is not a robust way to precisely calculate this FDR, we

leave the over-conservative empirical p -values unadjusted for all main analyses. It is important to note that “FDR” here measures overall excess (FDR < 1) of large effect sizes between the family and the population, it does not mean that the ordering of effect sizes (empirical p -values) are all due to random chance regardless of the outcome of this FDR. Though we did not estimate an accurate FDR here, the relative ranks of genes according to their effect sizes compared to the population (empirical p -values) are not affected, which are still valid and biologically meaningful.

To correct for additional noise in population subsamples, we measured the standard deviation of β of randomly permuted data ($\bar{\beta} = 0$) in both the family and population. We estimated that the standard deviation of β estimates of the family is 0.55 times the size of that of the population: $std(\hat{\beta}_{\text{family}}) = 0.55 * std(\hat{\beta}_{\text{population}})$ (Figure S17A). Such difference will make the confidence intervals which are measured from subsampled population to be larger than the actual noise of β estimates in the family. To adjust for such differences, we narrowed the empirical distribution of a gene by moving each subsampled effect size in the population towards their mean: $\beta_{\text{adjusted}} = \bar{\beta} + 0.55 * (\beta - \bar{\beta})$. This adjustment shrinks the distribution of population effect sizes (and consequently reduces the empirical p -values) but retains the estimates of $\bar{\beta}$ of that gene from the population. After adjustment, the distribution of empirical p -values testing whether the family effect size is bigger than the subsampled population is more uniform (Figure S18). As the adjustment of confidence intervals mainly influences calculation of the FDR, the enrichment pattern of rare variants was very similar under the adjusted confidence intervals (comparing Figure S14, B and C).

Alternatively, we also directly estimated analytic standard errors (confidence intervals) of $\hat{\beta}$ without using permutation. Assuming residuals are normally distributed, from linear model theory the standard error of $\hat{\beta}$ is estimated (MLE) by $\frac{1}{\sqrt{\sum(x_i - \bar{x})^2}} \sqrt{\frac{\sum \varepsilon_i^2}{11-3}}$, where ε_i is the residuals in linear regression $T_i \sim \mu + \beta_j p + \beta_k m + \varepsilon_i$, x is either p or m . We compared the difference of standard errors of $\hat{\beta}$ between the family and the population subsamples (Figure S17B). The estimated global difference of the standard error of $\hat{\beta}$ follows $\hat{S}_{\beta_{\text{family}}} = 0.55 * \hat{S}_{\beta_{\text{population}}}$. The scaling factor is very similar to that inferred by permutation.

Here, both standard error tuning methods make the assumption that the noise in the family for each gene is approximately a constant scaling factor less than the noise in the population. The tuning factor is obtained by matching ranks (U statistic of Wilcoxon rank-sum test) of two distributions of standard errors of $\hat{\beta}$, until the test is not significant.

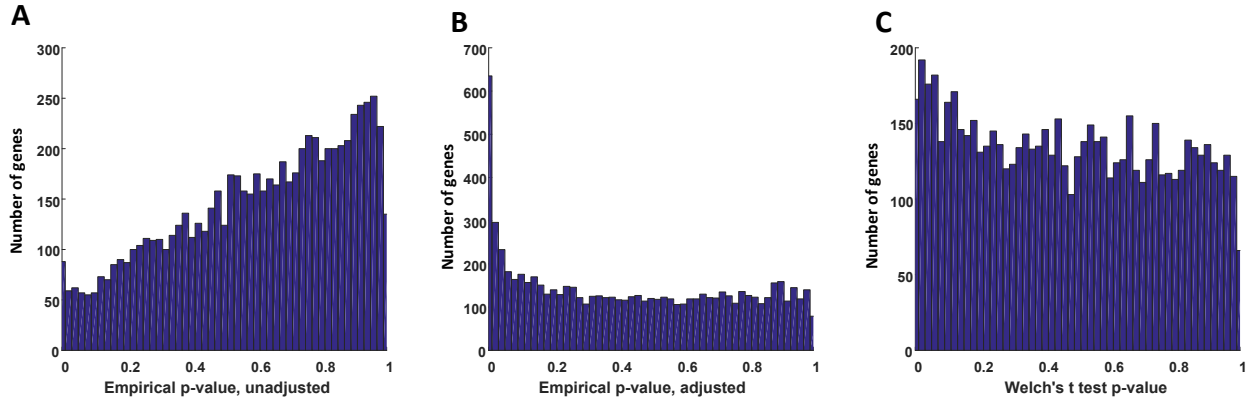


Figure S18. Adjustment of effect size empirical p -values: distribution of p -values of family versus population effect sizes.

Empirical p -values generated by subsampling: (A) before noise adjustment and (B) after noise adjustment. For each gene, we compute how often $(1 - \text{empirical } p\text{-values})$ the family effect size are larger than the effect sizes of the population subsamples. Distribution of effect sizes in the population is generated by subsampling 100 times from the population. After adjustment of their different noise levels, empirical p -values are more evenly distributed. (C) Welch's t -test p -values. Welch's t -test is performed by directly using a t -test between two regression slopes (two β 's), with standard errors estimated analytically.

We computed the straightforward analytical p -values (Welch's t -test) without using any subsampling (single-SNP regression over whole replication panel), which provide a bottom-line theoretical control of the empirical p -values. The adjusted empirical p -values lie much closer to theoretical p -values than the raw empirical p -values. Here, we can simply use a pure analytic test to compare regression slopes β_{family} and $\beta_{\text{population}}$ without either subsampling or

permutation by applying Welch's t test:
$$\frac{\beta_{\text{family}} - \beta_{\text{population}}}{\sqrt{(\hat{\sigma}_{\beta_{\text{family}}})^2 + (\hat{\sigma}_{\beta_{\text{population}}})^2}}$$
. Under normality assumption of

regression residuals, this test statistic follows t distribution, the standard errors are analytic

estimations from regression residuals:
$$\hat{\sigma}_{\beta_{\text{family}}} = \frac{1}{\sqrt{\sum(x_i - \bar{x})}} \sqrt{\frac{\sum \varepsilon_i^2}{11-3}}, \quad \hat{\sigma}_{\beta_{\text{population}}} = \frac{1}{\sqrt{\sum(x_i - \bar{x})}} \sqrt{\frac{\sum \varepsilon_i^2}{373-180-2}},$$

degree of freedom is
$$\frac{((\hat{\sigma}_{\beta_{\text{family}}})^2 + (\hat{\sigma}_{\beta_{\text{population}}})^2)^2}{(\hat{\sigma}_{\beta_{\text{family}}})^4 / (11-3) + (\hat{\sigma}_{\beta_{\text{population}}})^4 / (373-180-2)},$$
 373-180 = 193 is the size of

replication panel where effect size of the best associated SNP is re-measured.

$\text{FDR}_{\text{adjusted}}$ of those large effect genes at a given empirical p -value $(1 - \text{CI}_{\text{adjusted}})$ cutoff is calculated as (total number of genes * p -value cutoff) / number of discoveries. At p -value < 0.01 ($\text{CI}_{\text{adjusted}} > 0.99$), there are 558 larger effect β in the family compared to the population, $\text{FDR}_{\text{adjusted}} = 7341 * 0.01 * 2 / 558 \sim 0.26$ (* 2 because we combined paternal and maternal discoveries). It is important to note that while our reported FDR is conservative, the adjusted FDR may be permissive if the differences in the variance reflect meaningful biological differences. By Welch's t -test, there are 320 larger effect β (p -value < 0.01) in the family compared to the population, $\text{FDR} = 7341 * 0.01 * 2 / 320 \sim 0.46$. We conclude that there is

definitely a significant excess of large effect cis-eQTLs in the family than in the population, however as there is not yet a very robust estimation of this proportion, we choose to state the conservative FDR in the main text.

It is also important to note that this “FDR” measures whether there is overall excess of large effect genes in the family. The ranking of empirical p -values which reflects the positioning of effect sizes in a population spectrum is still biologically meaningful regardless of this excess. The downstream analysis based on rankings of effect sizes does not rely on this estimation of FDR.

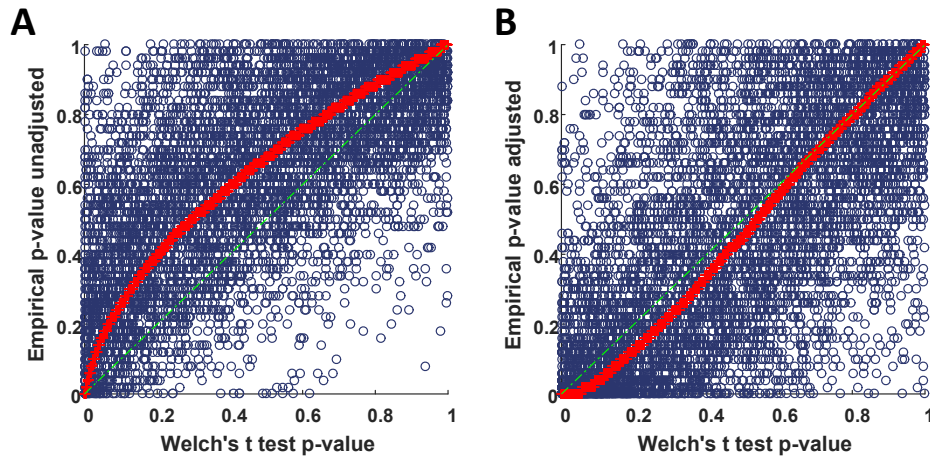


Figure S19. Adjustment of effect size empirical p -values: comparisons of empirical p -value and Welch's t -test.

Unadjusted empirical p -values (A) are significantly conservative than Welch's t -test, while adjusted p -values (B) are more optimistic than Welch's t -test. Here we only show empirical p -values and Welch's p -values measuring the difference of the paternal side β_j in the family and single-regressor $\beta_{\text{population}}$ of the population, the maternal side is similar.

As effect size (β) can be directly compared using analytical tests, to gain a theoretical control of the correctness of the subsampling scheme, we performed the conventional analytical test (Welch's t) to compare effect sizes. Here, population β is just a one-regressor (the best SNP) straightforward measurement of effect sizes over the whole replication panel without subsampling or implicit phasing. The analytical test can be used to gauge the overall soundness of empirical p -values. Comparing Welch's t -test with subsampling + permutation (empirical p -value) based test, these three p -values are mostly concordant with each other, however the empirical p -values are more conservative than Welch's t test before adjustment but more optimistic than Welch's t test after adjustment (Figure S19).

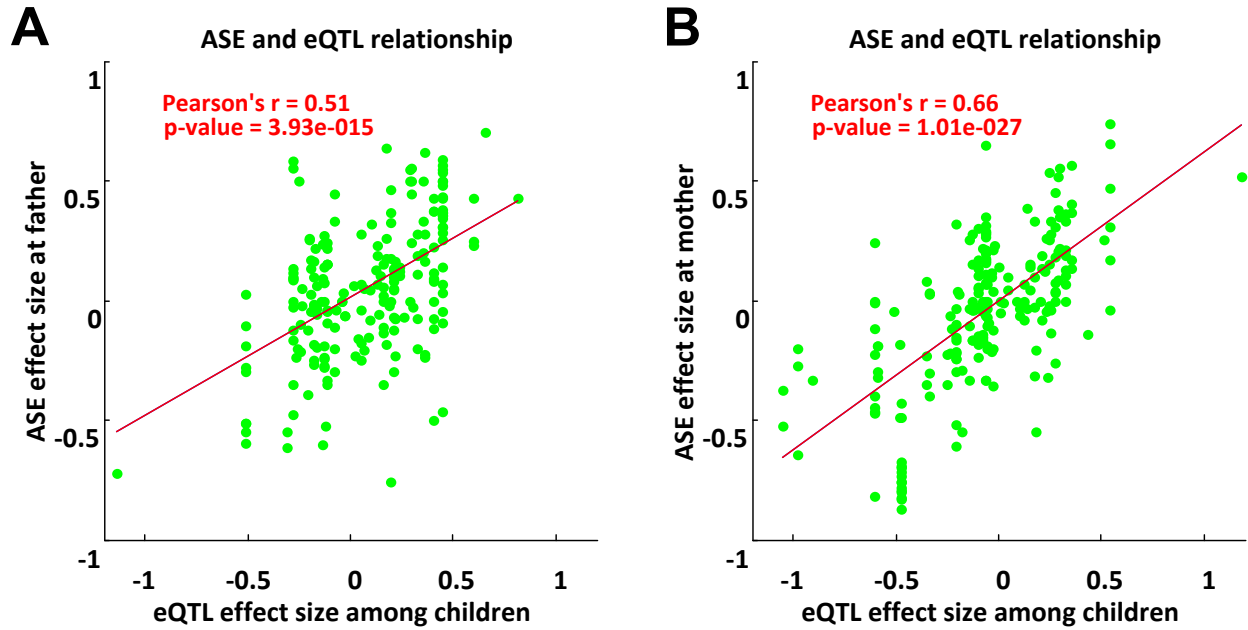


Figure S20. Correlation of eQTL effect size β and ASE effect size (allelic imbalance).

(A) Paternal β and allelic imbalance in father. (B) Maternal β and allelic imbalance in mother. ASE effect sizes in the parents and *cis*-eQTL effect sizes among the children should have a simple linear relationship. We computed correlation between ASE (quantified by allelic imbalance) in the parents and eQTL effect sizes (quantified by linear regression β) among the children. Indeed, we observed a linear relationship between ASE effect sizes in the parents and eQTL effect sizes among the children. In other words, the difference between two homologous alleles in a parent will exhibit as between-individual differences among the children, as expected by Mendelian segregation. For example, expression level difference between two homologous alleles (ASE) of the parent NA12878 (a_1, a_2): $a_1 - a_2$ is proportional to expression level differences between her offspring $(a_1, *) - (a_2, *)$ depending on which haplotype they inherit. We observed that when a haplotype is highly expressed in a parent as indicated by ASE, children inheriting that haplotype also have higher expression levels. ASE effect size at a heterozygous site is represented by $(\text{paternal reads} - \text{maternal reads}) / (\text{paternal reads} + \text{maternal reads})$, i.e., $2 * (\text{paternal allelic imbalance} - 0.5)$. *cis*-eQTL effect size is defined as the difference of gene expression levels between children inheriting different haplotypes (which is simply β of linear regression: $T_i \sim \mu + \beta_j p + \beta_k m$). We can observe that β is linearly determined by allelic imbalances.

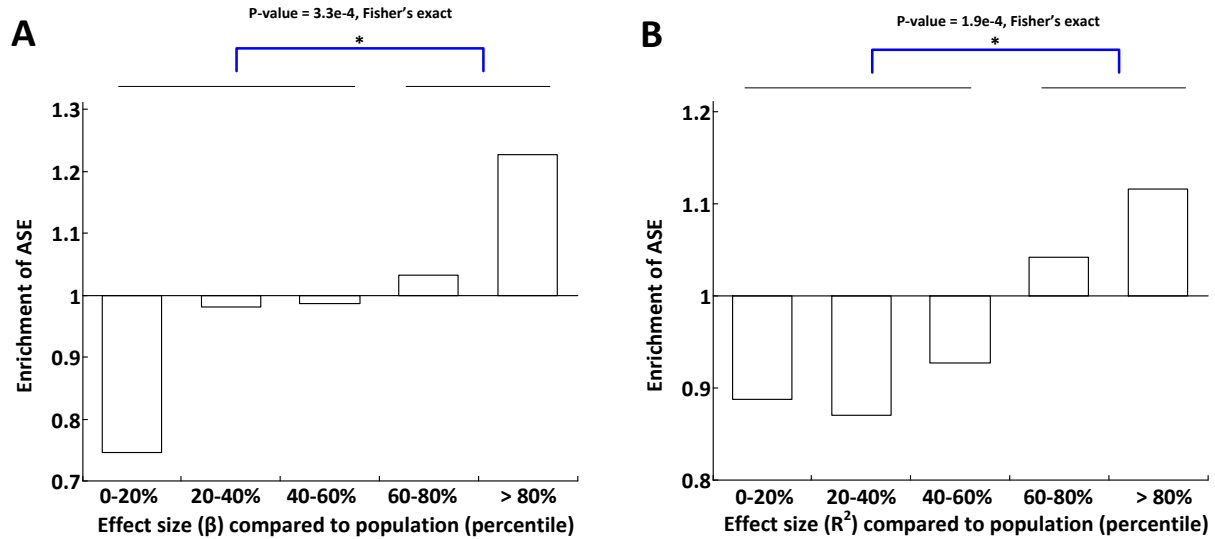


Figure S21. Enrichment of ASE effects at large-effect genes.

To further confirm that identified large-effect genes are potentially due to rare and heterozygous variants instead of technical artefacts, we assessed enrichment of ASE effects for large-effect eQTLs. ASE effects are evaluated in both parents of the family. In theory, large-effect eQTLs among siblings should also exhibit as ASE effects among at least one of the parents. For both β and R^2 , we observed increasing incidence of ASE effects at larger effect eQTLs. (A) Enrichment of ASE at large-effect (measured by β) genes. ASE effects (measured upon two parents) are defined by those passing binomial test p -value < 0.01 and allelic imbalance > 0.05 . (B) Enrichment of ASE at large-effect (measured by R^2) genes. We ranked effect sizes of genes in the family based on the how often (x -axis, $1 - \text{empirical } p$ -values) their measured effect in the family was greater than in the population subsamples. Enrichment is defined as the proportion of genes exhibiting ASE in each effect size bin divided by the proportion of genes exhibiting ASE across all effect size bins. We only consider genes testable for ASE, i.e., with heterozygous sites in RNA covered regions.

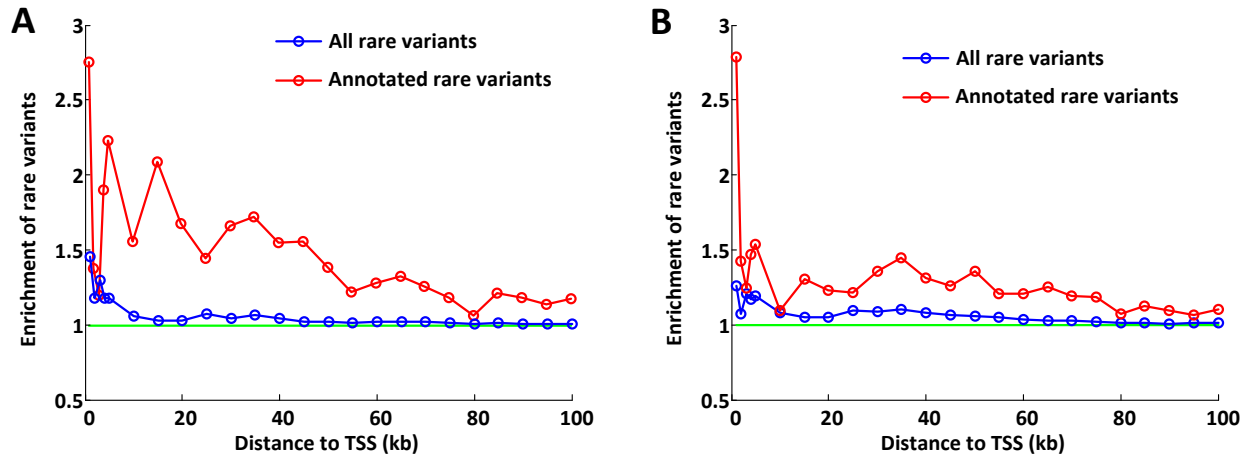


Figure S22. Rare regulatory variants contributing to large-effect eQTLs: enrichment of rare variants near the TSS of large-effect (β) cis-eQTL genes, comparing annotated and all rare variants.

We examined enrichment of rare variants near the transcription start site (TSS) of eQTL genes. Allele frequencies are based on the Phase 1 release of the 1000 Genomes Project European populations. We narrowed variants by multiple functional annotations such as conservation score (PhyloP) and regulatory features from annotations in RegulomeDB⁴ indicating Encode⁵ TF binding and DNaseI hypersensitivity peaks.

We observed an increasing enrichment of rare variants at larger effect size genes. Likewise, given a rare variant in an annotated regulatory region, we also see a significantly increased proportion of large effect genes. The enrichment is stronger in the immediately proximity of the TSS but also spreads across the 100kb regions. The enrichment is also much stronger among annotated regulatory sites than all other sites. Enrichment is defined as proportions of genes with such an annotated rare variant in each effect size bin divided by proportions of genes with such an annotated rare variant across all effect size bins.

(A) 319 genes $CI > 0.95$. (B) 558 genes $CI_{adjusted} > 0.99$. The plot shows enrichment for all rare variants (MAF < 0.01, 100kb near TSS) and annotated rare variants (MAF < 0.01, 100kb near TSS, within Encode TF binding and DNaseI hypersensitivity peaks and with PhyloP score > 1). We observed increased enrichment of rare variants near the TSS of larger family effect size genes. Enrichment is stronger for annotated rare variants.

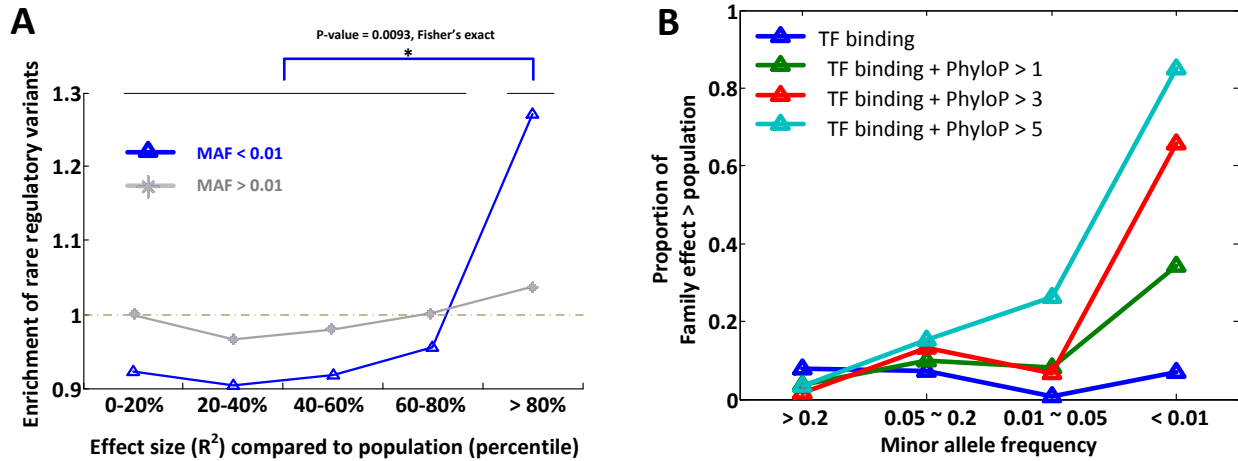


Figure S23. Rare regulatory variants contributing to large-effect eQTLs.

(A) Enrichment of rare variants at large-effect genes. Effect sizes are measured by fit (R^2) and binned by comparing to population effect sizes. We ranked genes according to their effect sizes in the family as percentile (x-axis, 1 – empirical p -values) in the population. Rare variants are defined as those of MAF < 0.01, Encode TF + DNase peak, PhyloP > 1 and within 100kb near TSS. (B) Utility of rare variants in predicting a larger effect in family than population. Enrichment is defined as proportions of genes with such an annotated rare variant in each effect size bin divided by proportions of genes with such an annotated rare variant across all effect size bins. (R^2). Rare variants are restricted to those in Encode TF + DNase peaks and different PhyloP score cutoffs. We estimate the proportion of family effects larger than population effects using π_1 statistics. π_1 is estimated from empirical p -values of whether a family effect size is larger than population by counting the number of times a family effect size is greater than subsampled population.

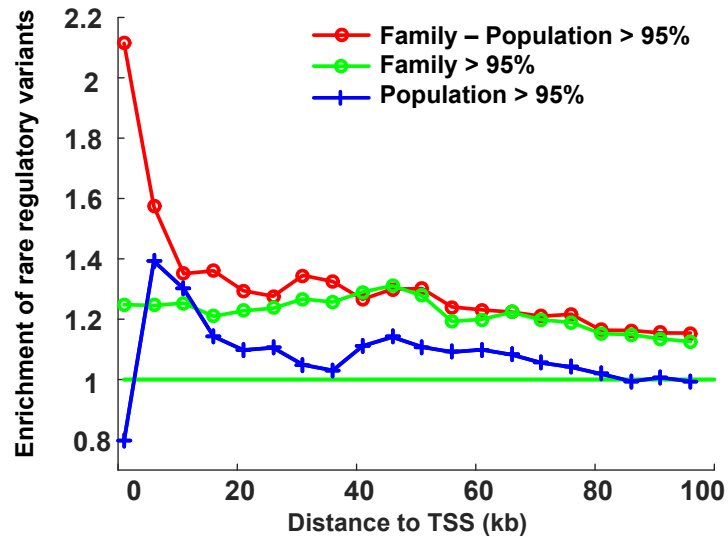


Figure S24. Enrichment of rare regulatory at large effect genes.

One might ask whether enrichment of rare variants is a general property of large-effect eQTL genes regardless of whether they are family specific or not (due to potentially larger number of regulatory elements near those genes). To explore such possibilities, we evaluated enrichment of rare regulatory variants at three categories of genes: genes whose effect sizes are larger in the family than in the population, genes whose effect sizes are large in the family regardless of whether they are larger than the population and genes whose effect sizes are large in the population. Here we consider three types of genes: genes with larger effect in the family than the population (red), genes with large-effect in the family regardless of effect sizes in the population (green) and genes with large effects in the population (blue). We consider the top 5% of the genes in each category. We only observe enrichment in the former two categories. This indicates that enrichment of rare variants is only at those family-specific large effect genes, it is not due to general enrichment of regulatory elements near large-effect genes.

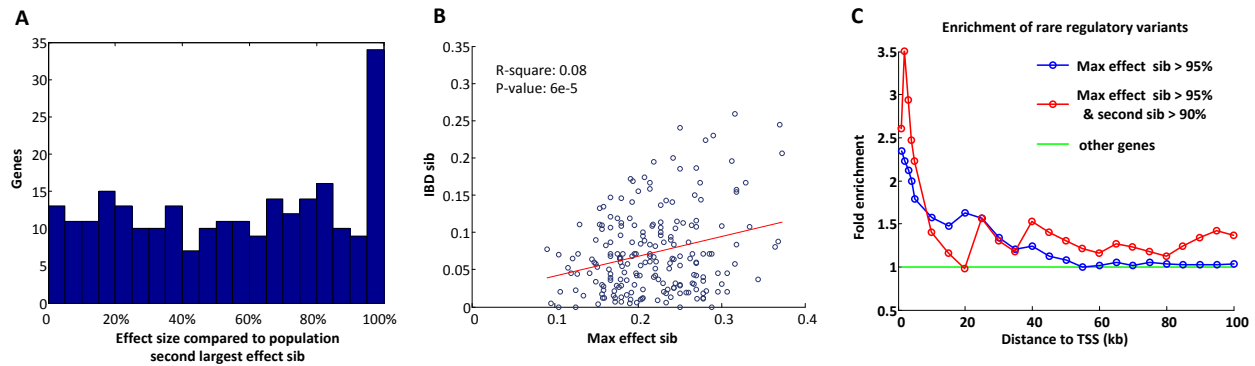


Figure S25. Identification of large ASE effect.

Genes are chosen where largest ASE effect sibling is > 95% quantile of population ASE effects. We checked ASE effects at a second largest effect sibling and IBD siblings to further confirm those large effects. (A) Here we show distribution of effect sizes (percentile as compared to population, 1 – empirical p -values) of second largest ASE effect sibling at outlier genes. (B) ASE effect size (allelic imbalance) between largest effect sibling and its double-IBD (identical-by-descent) sibling. (C) Enrichment of rare regulatory variants near TSS at those genes where both the first and second largest effect sibs are significantly larger than population. At each TSS cutoff, enrichment is defined as the proportion of large ASE effect genes with an annotated rare variant divided by proportions of genes with an annotated rare variant for all genes which are testable for ASE.

As ASE effects are evaluated at individual heterozygous sites, we wanted to exclude the possibility that large-effect ASE is due to technical artifacts such as mapping biases or sequencing errors. To achieve this, we looked at ASE for the second largest effect sibling. Our rationale being that the second sibling would be less likely to be a coincidental artifact than the first. We observe that the ASE effect at the second largest sibling is also significantly enriched for larger effect sizes (Figure S25A). Furthermore, we also looked at an identical-by-descent sibling of the maximum effect sibling. We observe that large effects are repeated at the IBD sibling (Figure S25B). When looking at genes with both first and second largest ASE effects greater than population, we observed strong enrichment of rare variants at those genes (Figure S25C).

ASE discovery. The following are additional notes on discovering large ASE effects and its FDR. We applied a similar method to identify large-effect ASE in the family. We use allelic ratio as a measure for ASE. Large-effect ASE genes are detected by comparing maximum allelic imbalance among 11 siblings and ASE in subsampled population data. We subsampled 11 individuals from the population (373 European individuals from Geuvadis data) and take the maximum allelic imbalance. We calculated empirical p -values of family effect sizes according to the effect size distribution of the population subsamples. To account for the differences in read depths between the family and population data, we further down-sampled the population data by a ratio of 1.97.

We discover 223 large effect genes at $CI > 0.95$ (with empirical p -value < 0.05), which yields an $FDR = 1777 * 0.05 / 223 = 40\%$. We do not expect to see more large ASE genes in the family than in a population subsample. Unlike eQTLs, there is equivalent statistical power in the family and a population subsample to detect ASE effects either arising from rare or common variants. The excessive number of large effect genes mainly reflects read depth differences (lower read depth leads to larger allelic imbalance) between two datasets we have not yet corrected out. We are trying to correct out this factor by using a uniform down-sampling factor of 1.97 which reflects the global read depth difference between two datasets. However as there are substantial variability of read depth between individuals and sites, this global correction cannot remove all the technical differences.

It is very important to mention that by theory FDR for ASE should be inherently 1. However this “FDR” is a measurement of the excess of large ASE in the family compared to the population (which should be zero), it does not mean that large effect sizes are out of random chance. Empirical p -values here are not just random noises; they have biological meanings individually, reflecting the positioning of ASE effect sizes of one individual among the natural spectrum of all individuals. Therefore, the ranking of those genes by their empirical p -values are still biologically meaningful regardless of whether there is overall excess of larger effect sizes. As our purpose is not to estimate whether there is excess or not but to obtain an ordering of genes by their relative ASE effect sizes, it is therefore critical to emphasize the meaning of this FDR here.

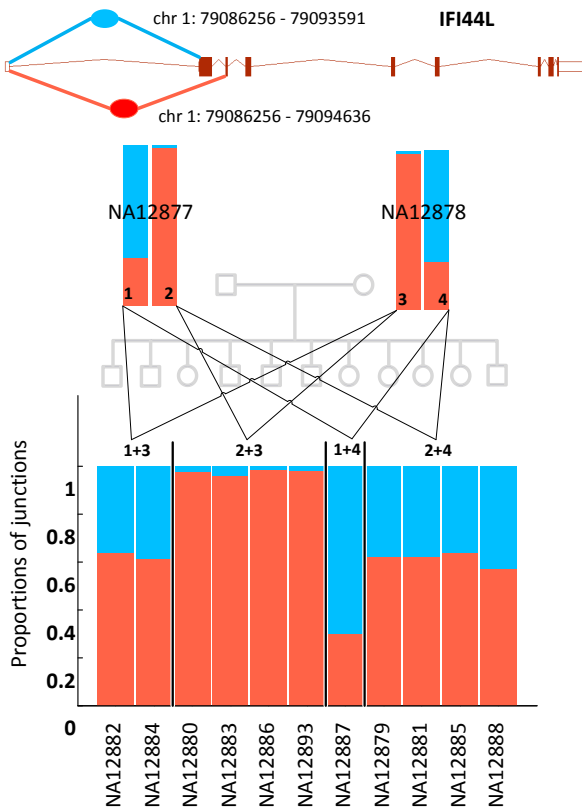


Figure S26. Mendelian segregation of alternative splicing patterns.

Alternative splicing patterns determined by haplotype groups. 1,2,3,4 are paternal grandfather, grandmother and maternal grandfather and grandmother haplotypes. We observed that transcript ratios can exhibit as Mendelian segregation in the family. We use JunctionsTK (junction toolkit, a tool developed by our group) to quantify such segregation patterns using splicing reads. JunctionsTK uses reads spanning splice junctions from junctions.bed files produced by TopHat. It calculates proportions of junction reads from one donor exon to different acceptors (or different donors to a same acceptor). Compared to transcript abundance reported by Cufflinks, splice junction reads, as they are specific to each alternative transcript, more directly inform alternative splicing differences between individuals. The figure shows segregation of splicing junction usage of these genes by different ancestral haplotypes. We show differentially expressed transcripts among four groups of siblings (depending on which two grandparental alleles are inherited), each group is divided by vertical bars. Here, the y-axis shows proportional usage of each junction site, from the same donor exon to different acceptor exons (or different donors to a same acceptor). We can observe usage of splicing junctions is highly consistent within same ancestral haplotypes, while distinct between different ancestral haplotypes.

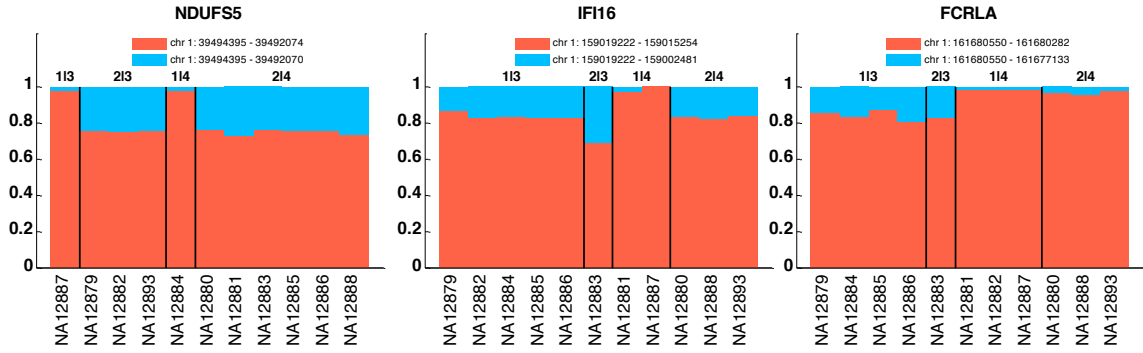


Figure S27. Examples of alternative splicing patterns determined by haplotype groups.

1, 2, 3 and 4 are paternal grandfather, grandmother and maternal grandfather and grandmother haplotypes, respectively. Explanations of segregation patterns are provided in Figure S26. Additional information about those genes is provided in Table S8.

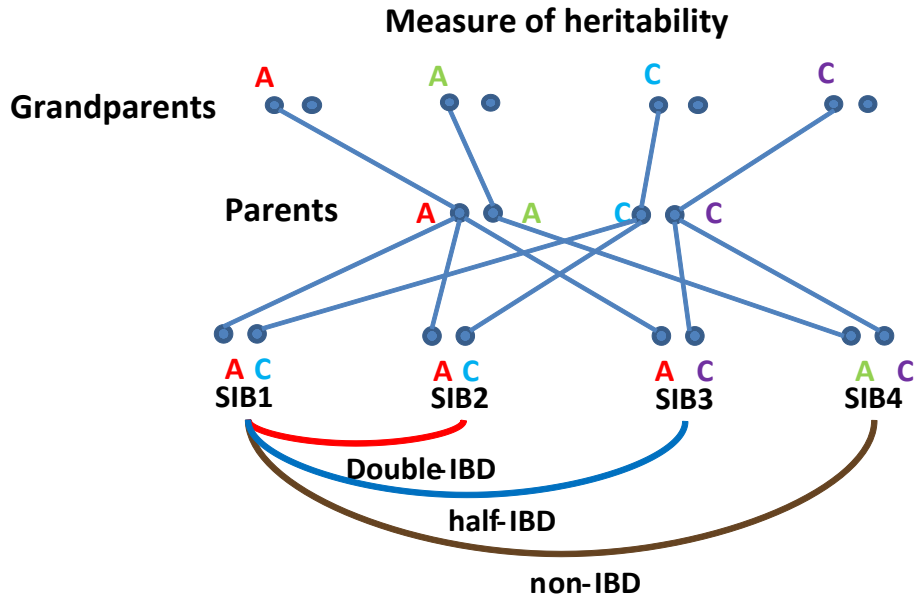


Figure S28. ASE heritability analysis.

Method of measuring heritability in family. IBD means Identical-by-descent, descending from the same ancestral haplotype. Sib1 and Sib2 are double-IBD siblings as they share both haplotypes. Sib1 and Sib3 are half-IBD as they share only one haplotype. Sib1 and Sib4 are non-IBD as they share neither of their haplotypes. We use each child as a reference and calculated the correlation of allelic ratios with their double-IBD, half-IBD and non-IBD siblings. We repeat this for each of 11 children.

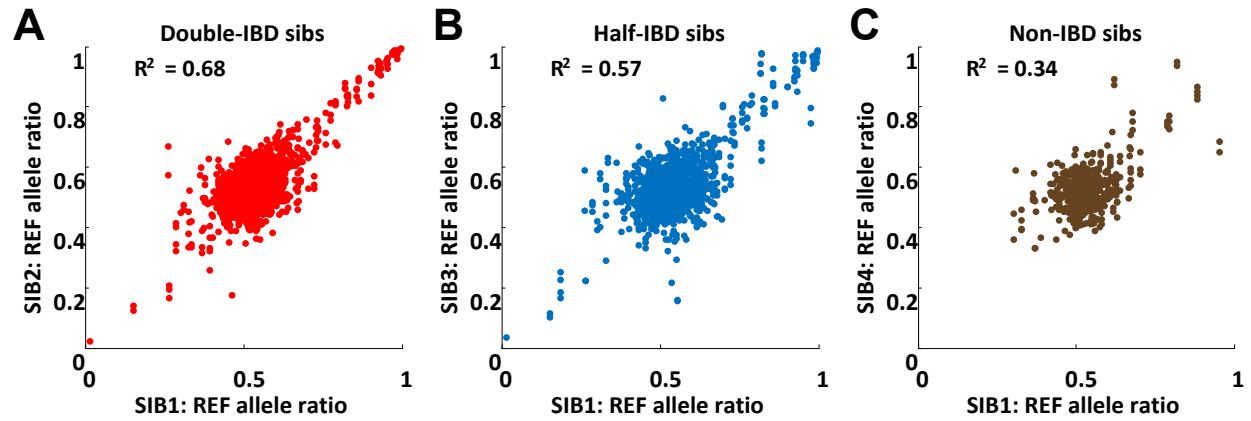


Figure S29. Allelic ratio correlation with siblings, using NA12879 as reference.

The figure shows measured correlation for each experiment. (A) Double-IBD, (B) half-IBD and (C) non-IBD are defined as sharing both, only one or neither haplotypes. To reduce random sampling noise, the result is based on sites of depths greater than 100.

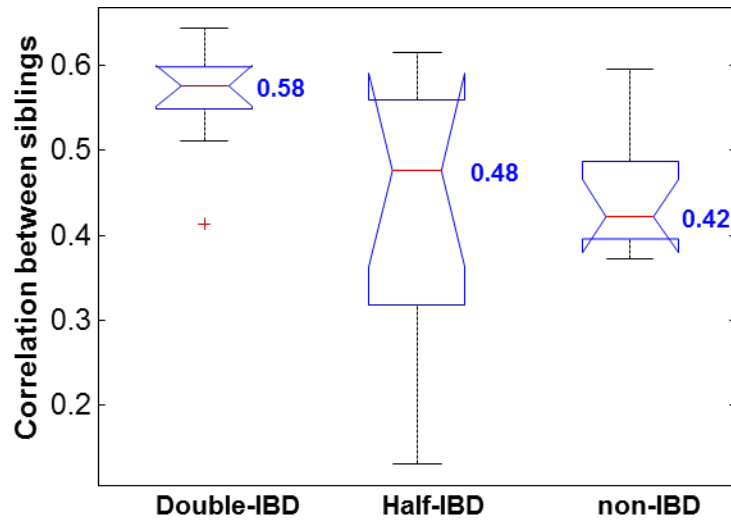


Figure S30. Allelic ratio correlation between different types of siblings.

From left to right, measured correlation between double-IBD, half-IBD and non-IBD siblings using each sibling as reference. Median correlation coefficients for double-IBD, half-IBD and non-IBD siblings are 0.58, 0.48, and 0.42 respectively.

Supplemental Tables

	Total segregating SNPs	MAF in 1000 Genomes European			
		≤ 0.05		≤ 0.01	
CEU family	2,936,403	345,232	11.76%	91,882	3.13%

Table S1. Number of variants segregating in the family.

Segregating variants among children are those variants that are heterozygous in at least one of the parents and both alleles are transmitted to the children. We define rare variants as variants with minor allele frequency below 0.01 (or otherwise specified in the paper) in the Phase 1 release of the 1000 Genomes Project European populations. As calling of rare variants is especially susceptible to genotyping errors, we use stringent Mendelian constraints to reduce these errors. We require all called variants to be completely consistent along the IBD inheritance across the whole family.

	NA12877	NA12885	NA12886	NA12891	NA12892
Sites genotyped in both LFR and original sequencing	1835141	1711513	1754935	1716804	1763045
Concordant sites	1834123	1710681	1753403	1716076	1762294
% concordant	0.999445	0.999514	0.999127	0.999576	0.999574

Table S2. Genotypes confirmed with Complete Genomics Long Fragment Read¹ (LFR).

Comparison of genotypes between original Complete Genomics sequencing and LFR technology. We have five individuals in the family that were sequenced again using LFR technology which can generate molecular haplotypes. LFR can place 92% of heterozygous SNP into long ~500kb contigs. It has very high genotyping accuracy, with an error rate of 1 in 10Mb. Genotypes used in this study were found to be more than 99.9% concordant between original sequencing data and LFR data. Further, we confirmed called genotypes with both new sequencing data from both Complete Genomics Long Fragment Read¹ technology and from Illumina Platinum Genomes (Table S4).

		NA12877	NA12885	NA12886
Total phased het sites	PEDIBD	1913154	1903914	1922039
	LFR	1876780	1747734	1795687
Overlap		1834122	1710681	1753402
Phase concordant % concordant		1831234 0.998425	1709182 0.999124	1752028 0.999216

Table S3. Phasing confirmed with molecular haplotype by LFR.

Comparison of inferred haplotypes by Ped-IBD and molecular haplotypes generated by LFR technology. Phasing results are confirmed to be 99.9% concordant with molecular haplotypes generated by LFR technology.

	Variant sites genotyped by Complete Genomics	Variant sites genotyped by Illumina	Genotyped by both	Concordance
NA12889	3232336	2862355	2703495	0.9983
NA12890	3237426	2798769	2672122	0.9981
NA12891	3216575	2816832	2678531	0.9982
NA12892	3231991	2869083	2711768	0.9983
NA12877	3228441	2828846	2681250	0.9976
NA12878	3220524	2802242	2667261	0.9979
NA12879	3220425	2796242	2662748	0.9977
NA12880	3231151	2806653	2670442	0.9980
NA12881	3236471	2802998	2670859	0.9979
NA12882	3270994	2825547	2684441	0.9881
NA12883	3234742	2796940	2667074	0.9982
NA12884	3226121	2785306	2653732	0.9966
NA12885	3233256	2781640	2657075	0.9982
NA12886	3202270	2847260	2668380	0.9961
NA12887	3209866	2825088	2672224	0.9972
NA12888	3234370	2816855	2675737	0.9977
NA12889	3245741	2825098	2686108	0.9976

Table S4. Genotypes confirmed with Illumina Platinum Genomes.

Comparison of genotypes between Illumina Platinum Genomes and Complete Genomics sequencing data. The same CEU family samples were also sequenced by Illumina as part of Illumina Platinum Genomes. All 17 members were sequenced to 50x depth on a HiSeq 2000 system. We compared genotypes called by Complete Genomics and Illumina (passing genotype filter and quality score = 99). On average, Complete Genomics data cover more sites than Illumina and include the majority (> 95%) of Illumina sites. Genotype concordance between the two platforms at overlapping sites ranged from 0.9966 to 0.9983 across individuals.

	Tested genes	Number of eQTL genes	
eQTL genes	8,974 genes	$\pi_1=0.078$	~698
		FDR < 50%	274
sQTL genes	7,954 genes	$\pi_1=0.079$	~624
		FDR < 50%	261

Table S5. Linkage analysis of *cis*-eQTL: summary of eQTL and sQTL genes identified in the family.

Total number of eQTL genes π_1^3 and numbers of genes below FDR 0.5. Numbers of haplotype blocks holding these eQTL or sQTL genes are also listed.

Gene expression quantification. We used the Tophat/Cufflinks to quantify expression levels of whole genes and each transcript from RNA-Seq data (Figure S2). We performed eQTL discovery using linear regression of gene expression levels with local haplotype blocks. We identified *cis* expression QTLs (*cis*-eQTLs) by restricting association to the haplotype block that contains the tested gene. We only considered protein-coding genes, and to minimize possible technical artifacts in quantification we also exclude all pseudogenes, all immunoglobulin and HLA genes where there is an increased potential for mapping biases and sequencing errors. We required an average FPKM greater than 2 and at least 3 individuals with FPKM greater than 1 for a gene to be tested. Setting this threshold, we tested 8,974 genes for eQTLs. For *cis*-splicing QTLs (sQTLs), we additionally require each gene to have two or more quantified alternative transcripts (N=7,954 genes).

eQTL and sQTL discoveries. To detect eQTLs in the family, we used a two-variable (paternal and maternal haplotypes) linear regression to test for gene expression \sim haplotype association. For each haplotype block, the two parental haplotypes of each child are encoded using two variables, p and m . The maternal haplotype m_i of a child i , for example, is either 0 or 1, depending on which of the two possible maternal alleles is present. Then, an expression trait is regressed as the summation of effects of two parental haplotypes, $T_i \sim \mu + \beta_j p_i + \beta_k m_i$, where T_i is the trait of individual i , the effects of two parental alleles k and j are expressed by β_j and β_k and μ is the intercept. For sQTL, we selected the most significant p -value among all transcripts for each gene. P -values are further adjusted using permutation as described below.

Empirical p -values were generated using permutation by swapping phenotypes across individuals. We performed 10000 permutations at each gene and computed p -values by counting how many times permuted p -values fell below the nominal p -value.

To quantify effects for common variants, we used linear regression to test common variants among 373 unrelated European individuals from Geuvadis study⁶. To ensure

discoveries in the population were relevant to the family, we only test variants which were also polymorphic in the family.

Sample size	50	100	150	200	250
50	312	257	243	247	252
100	257	321	243	245	254
150	243	243	308	241	245
200	247	245	241	315	249
250	252	254	245	249	366

Table S6. Effect of different discovery panel sizes: number of large effect β genes given different discovery panel sizes.

On the diagonal are numbers of genes with family effect sizes > 95% CI population effect sizes. Off diagonal cells show their intersections.

MAF	Distance to TSS	TF binding + DNase peak	PhyloP score	Motif	# of genes	% eQTL in family	% eQTL in population
< 0.01	100 kb	-	-	-	7912	0.0858	0.1662
< 0.01	100 kb	Yes	-	-	3123	0.0990	0.1775
< 0.01	100 kb	Yes	> 1	-	367	0.4577	0.1759
< 0.01	100 kb	Yes	> 1	Yes	135	0.5627	0.2049
< 0.01	5 kb	-	-	-	1525	0.1099	0.1807
< 0.01	5 kb	Yes	-	-	386	0.1815	0.1743
< 0.01	5 kb	Yes	> 1	-	41	0.5151	0.1622
< 0.01	5 kb	Yes	> 1	Yes	17	0.8999	0.2000
< 0.01	100 kb	-	-	-	7912	0.0869	0.1284
< 0.01	100 kb	Yes	-	-	3123	0.0968	0.1371
< 0.01	100 kb	Yes	> 3	-	88	0.8303	0.1688
< 0.01	100 kb	Yes	> 3	Yes	30	0.9528	0.2400
> 0.01	100 kb	-	-	-	8312	0.0785	0.1647
> 0.01	100 kb	Yes	-	-	8186	0.0757	0.1652
> 0.01	100 kb	Yes	> 1	-	6110	0.0820	0.1676
> 0.01	100 kb	Yes	> 1	Yes	2456	0.0965	0.1688
> 0.01	5 kb	-	-	-	7525	0.0801	0.1684
> 0.01	5 kb	Yes	-	-	6092	0.0941	0.1755
> 0.01	5 kb	Yes	> 1	-	1359	0.1676	0.1833

> 0.01	5 kb	Yes	> 1	Yes	413	0.1616	0.1777
> 0.01	100 kb	-	-	-	8312	0.0781	0.1647
> 0.01	100 kb	Yes	-	-	8186	0.0762	0.1652
> 0.01	100 kb	Yes	> 3	-	1542	0.1012	0.1717
> 0.01	100 kb	Yes	> 3	Yes	457	0.1504	0.1762

Table S7. Prediction of eQTLs at rare variants given annotation: proportion of genes being an eQTL given a regulatory variant near TSS.

We assessed the utility of different variant annotations in predicting eQTLs. We incrementally add annotations for minor allele frequency, distance to TSS, transcription factor binding, DNase sites, conservation score and transcription factor motif.

We selected those annotations which are previously found to be informative in predicting eQTLs: distance to TSS, transcription factor binding, DNase sites, conservation and transcription factor motifs. Encode transcription factor binding and DNase hypersensitivity peaks were obtained from RegulomeDB database⁴. Conservation scores using PhyloP⁷ (phyloP100way) software were downloaded from the UCSC genome browser (genome.ucsc.edu). Motif disrupting sites were downloaded from HaploReg database (v2)⁸.

Gene name	transcript ratio p -value	splice junction p -value	function
NDUFS5	0.0030	2.51E-09	neurological disorders
IFI44L	0.0020	1.33E-08	response to viral infection
IFI16	0.0010	1.91E-06	response to viral infection
FCRLA	0.0001	9.30E-05	B-cell development

Table S8. Examples of sQTL genes.

These genes are identified as sQTL genes by both transcript ratios and splice junction read ratios. Transcript ratio p -values are based on quantification of transcript abundances by Cufflinks, splice junction p -values are based on quantification of splice junction reads by JunctionTK. Two methods are in general concordant with each other. Segregating patterns of those genes are illustrated in Figure S26 and Figure S27.

GENE	GWA SNP ID*	chr	bps	Distance of nearest rare regulatory variant to TSS (bps)	Trait
B4GALT1	rs10813960-T	9	33180362	-	Urate levels
BAK1	rs210134-A	6	33540209	-	Platelet counts
PHTF2	rs12234571-C	7	77549906	-	Obesity-related traits
BAK1	rs9469457-A	6	33489882	-	Obesity-related traits
TCFL5	rs17854409-G	20	61491494	-	Obesity-related traits
TRAF3IP2	rs3851228-T	6	111848191	123850	Inflammatory bowel disease
PHTF2	rs848452-?	7	77596812	-	Dental caries
EPB41L5	rs13401620-A	2	120513133	-	Breast size
BAK1	rs210142-C	6	33546837	-	Chronic lymphocytic leukemia
INSIG1	rs10263087-C	7	154970469	49137	Formal thought disorder in schizophrenia
BAK1	rs210134-G	6	33540209	-	Mean platelet volume
BAK1	rs210134-G	6	33540209	-	Platelet counts
ENTPD6	rs1044573-A	20	25206654	128236	Allergic rhinitis
ABHD12	rs7267979-G	20	25298087	-	Liver enzyme levels (alkaline phosphatase)
PLCL2	rs9821630-G	3	16970938	-	Multiple sclerosis
ZKSCAN5	rs11761528-T	7	99118801	-	Dehydroepiandrosterone sulphate levels
PLCL2	rs1372072-A	3	16955259	-	Primary biliary cirrhosis
TRAF3IP2	rs33980500-T	6	111913262	123850	Psoriasis
TRAF3IP2	rs33980500-T	6	111913262	123850	Psoriatic arthritis
TRAF3IP2	rs240993-A	6	111673714	123850	Psoriasis
PRMT7	rs7197653-C	16	68383047	-	Magnesium levels
DHCR7	rs12785878-?	11	71167449	158269	Vitamin D insufficiency

NADSYN1	rs12785878-?	11	71167449	158510	Vitamin D insufficiency
CST3	rs911119-?	20	23612737	6080	Chronic kidney disease
PHTF2	rs6465825-C	7	77416439	-	Chronic kidney disease
ALDH7A1	rs13182402-G	5	125918148	-	Osteoporosis
BAK1	rs210138-G	6	33542538	-	Testicular germ cell tumor
IL16	rs7172689-?	15	81533695	-	Inattentive symptoms

*Identified as an eQTL in⁶, but not polymorphic in family

Table S9. Family-specific cis-eQTL modifying complex trait genes.

We assessed the number of GWA loci which were influenced by large-effect family eQTLs. Here, we identified family-specific eQTL for genes in the NHGRI GWA catalog¹⁷. We tested all those GWA SNPs in Geuvadis data and select those which are eQTLs ($\pi_1^3 \sim 0.3$, 315 genes at an FDR < 0.05). We then sub-selected those where the associated GWA SNP was not polymorphic in the family, so the GWA SNP is not causing the eQTL in the family, and another regulatory variant for the same gene is likely to be present. This highlights the potential for rare regulatory variants manifesting as family-specific eQTLs to be modifying important complex disease associated genes. The table lists examples of large-effect (CI > 0.80) family eQTL that influence GWA genes. The GWA SNP for the trait is determined to be an eQTL SNP (FDR 0.05) but not polymorphic in the family. Rare regulatory variants are defined as those within Encode TF binding + DNase peaks, MAF < 0.01 and PhyloP > 0, within 200kb near TSS.

GENE	GWA SNP ID	chr	bps	Distance of nearest rare regulatory variant to TSS (bps)	Trait
IRF5	rs12531711-G	7	128617466	50878	Primary biliary cirrhosis
IRF5	rs10488631-C	7	128594183	50878	Primary biliary cirrhosis
IRF5	rs10488631-C	7	128594183	50878	Rheumatoid arthritis
IRF5	rs729302-A	7	128568960	50878	Systemic lupus erythematosus
IRF5	rs12531711-G	7	128617466	50878	Systemic lupus erythematosus
IRF5	rs10488631-C	7	128594183	50878	Systemic lupus erythematosus
IRF5	rs4728142-A	7	128573967	50878	Systemic lupus erythematosus
IRF5	rs12537284-A	7	128717906	50878	Systemic lupus erythematosus
IRF5	rs10488631-C	7	128594183	50878	Systemic lupus erythematosus
IRF5	rs10488631-C	7	128594183	50878	Systemic sclerosis
IRF5	rs4728142-A	7	128573967	50878	Ulcerative colitis
IRF5	rs4728142-A	7	128573967	50878	Ulcerative colitis
NAPRT1	rs2290416-?	8	144657600	46423	Attention deficit hyperactivity disorder
NT5E	rs494562-G	6	86117129	40975	Metabolic traits
TCF7	rs756699-A	5	133446575	63462	Multiple sclerosis

Table S10. Examples of rare regulatory variants influencing GWA genes.

The GWA SNP for the trait is determined to be an eQTL SNP (FDR 0.05). Rare regulatory variants are defined as those within Encode TF binding + DNase peaks, MAF < 0.01 and PhyloP > 3, within 200kb near TSS.

References

1. Peters, B.A., Kermani, B.G., Sparks, A.B., Alferov, O., Hong, P., Alexeev, A., Jiang, Y., Dahl, F., Tang, Y.T., Haas, J., et al. (2012). Accurate whole-genome sequencing and haplotyping from 10 to 20 human cells. *Nature* 487, 190-195.
2. Cheung, V.G., Nayak, R.R., Wang, I.X., Elwyn, S., Cousins, S.M., Morley, M., and Spielman, R.S. (2010). Polymorphic cis- and trans-regulation of human gene expression. *PLoS biology* 8, 14.
3. Storey, J.D. (2002). A direct approach to false discovery rates. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 64, 479-498.
4. Boyle, A.P., Hong, E.L., Hariharan, M., Cheng, Y., Schaub, M.A., Kasowski, M., Karczewski, K.J., Park, J., Hitz, B.C., Weng, S., et al. (2012). Annotation of functional variation in personal genomes using RegulomeDB. *Genome research* 22, 1790-1797.
5. Consortium, E.P., Dunham, I., Kundaje, A., Aldred, S.F., Collins, P.J., Davis, C.A., Doyle, F., Epstein, C.B., Frietze, S., Harrow, J., et al. (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature* 489, 57-74.
6. Lappalainen, T., Sammeth, M., Friedlander, M.R., t Hoen, P.A., Monlong, J., Rivas, M.A., Gonzalez-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G., et al. (2013). Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* 501, 506-511.
7. Cooper, G.M., Stone, E.A., Asimenos, G., Program, N.C.S., Green, E.D., Batzoglou, S., and Sidow, A. (2005). Distribution and intensity of constraint in mammalian genomic sequence. *Genome research* 15, 901-913.
8. Ward, L.D., and Kellis, M. (2012). HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic acids research* 40, D930-934.