

# The impact of rare variation on gene expression across tissues

Xin Li<sup>1\*</sup>, Yungil Kim<sup>2\*</sup>, Emily K. Tsang<sup>1,3\*</sup>, Joe R. Davis<sup>1,4\*</sup>, Farhan N. Damani<sup>2</sup>, Colby Chiang<sup>5</sup>, Gaelen T. Hess<sup>4</sup>, Zachary Zappala<sup>1,4</sup>, Benjamin J. Strober<sup>6</sup>, Alexandra J. Scott<sup>5</sup>, Amy Li<sup>4</sup>, Andrea Ganna<sup>7,8,9</sup>, Michael C. Bassik<sup>4</sup>, Jason D. Merker<sup>1</sup>, GTEx Consortium†, Ira M. Hall<sup>5,10,11</sup>, Alexis Battle<sup>2§</sup> & Stephen B. Montgomery<sup>1,4§</sup>

**Rare genetic variants are abundant in humans and are expected to contribute to individual disease risk<sup>1–4</sup>. While genetic association studies have successfully identified common genetic variants associated with susceptibility, these studies are not practical for identifying rare variants<sup>1,5</sup>. Efforts to distinguish pathogenic variants from benign rare variants have leveraged the genetic code to identify deleterious protein-coding alleles<sup>1,6,7</sup>, but no analogous code exists for non-coding variants. Therefore, ascertaining which rare variants have phenotypic effects remains a major challenge. Rare non-coding variants have been associated with extreme gene expression in studies using single tissues<sup>8–11</sup>, but their effects across tissues are unknown. Here we identify gene expression outliers, or individuals showing extreme expression levels for a particular gene, across 44 human tissues by using combined analyses of whole genomes and multi-tissue RNA-sequencing data from the Genotype-Tissue Expression (GTEx) project v6p release<sup>12</sup>. We find that 58% of underexpression and 28% of overexpression outliers have nearby conserved rare variants compared to 8% of non-outliers. Additionally, we developed RIVER (RNA-informed variant effect on regulation), a Bayesian statistical model that incorporates expression data to predict a regulatory effect for rare variants with higher accuracy than models using genomic annotations alone. Overall, we demonstrate that rare variants contribute to large gene expression changes across tissues and provide an integrative method for interpretation of rare variants in individual genomes.**

Our analysis focused on individuals with extremely high or extremely low expression of a particular gene compared with the population, using the GTEx v6p release data, which include RNA-sequencing data for 449 individuals and 44 tissues. We refer to these individuals as gene expression outliers. The GTEx data enable the identification of both single-tissue and multi-tissue expression outliers (Fig. 1a), with the latter defined by consistent extreme expression across many tissues (see Methods). To account for broad environmental and technical confounders, we removed hidden factors estimated by PEER (probabilistic estimation of expression residuals)<sup>13</sup> from each tissue before outlier discovery (Extended Data Figs 1, 2 and Supplementary Tables 1, 2).

We identified a single-tissue expression outlier for  $\geq 99\%$  of expressed genes in each tissue and a multi-tissue outlier for 4,919 out of 18,380 genes that were tested (27%). Each individual was a single-tissue outlier for a median of 83 genes per tissue and a multi-tissue outlier for a median of 10 genes. Single-tissue outliers that were found in one tissue replicated in other tissues at rates of up to 33%, with higher rates among related tissues (Fig. 1b and Extended Data Fig. 3). The

replication rate for multi-tissue outliers was much higher and increased with the number of tissues used for discovery (Fig. 1c).

We investigated the influence of rare genetic variation on extreme expression levels, focusing on the individuals of European ancestry with whole-genome sequencing data (1,144 multi-tissue outliers). Multi-tissue outliers were strongly enriched for nearby rare variants. The enrichment was most pronounced for structural variants, as previously described<sup>14</sup>, and greater for short insertions and deletions (indels) than for single-nucleotide variants (SNVs) (Fig. 2a and Extended Data Fig. 4). Because most rare variants occur as heterozygotes, expression outliers driven by rare variants in *cis* should exhibit allele-specific expression (ASE). Both single-tissue and multi-tissue outliers were significantly enriched for ASE compared to non-outliers (see Methods; two-sided Wilcoxon rank-sum tests, each nominal  $P < 2.2 \times 10^{-16}$ ; Fig. 2c). For underexpression outliers with exonic rare variants, the rare allele was generally underexpressed with respect to the common allele and conversely so for overexpression outliers, consistent with the rare variant causing the effect (two-sided Wilcoxon rank-sum tests, each nominal  $P < 4.0 \times 10^{-8}$ ; Extended Data Fig. 5a). The enrichment for rare variants and ASE was stronger for multi-tissue outliers than for single-tissue outliers (Fig. 2b, c and Extended Data Fig. 6a), especially at higher Z-score thresholds.

To characterize the properties of rare variants that correlated with large changes in gene expression, we assessed the enrichment of different classes of variants in outliers compared to non-outliers (Supplementary Table 3a). Outliers were enriched, in order of significance, for structural variants, variants near splice sites, introducing frameshifts, at start or stop codons, near the transcription start site and in conserved regions (Fig. 3a). Variants in coding regions contributed disproportionately to outlier expression; enrichments weakened for all variants types (SNVs, indels and structural variants) when excluding exonic regions (Extended Data Fig. 6b). Additionally, 90% of stop-gain and frameshift variants were predicted to trigger nonsense-mediated decay in outliers (see Methods), suggesting a biological mechanism for these cases.

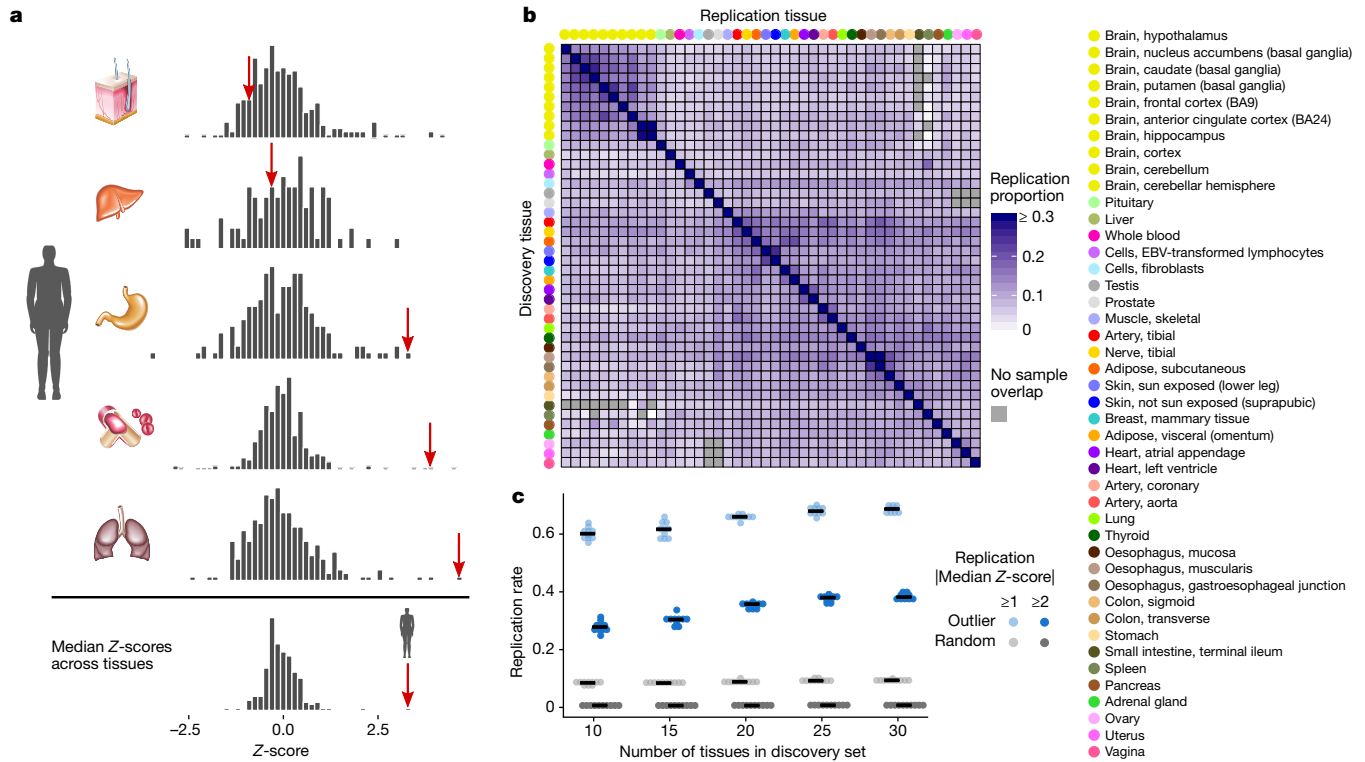
We also tested the relationship between outlier gene expression and functional annotations. Multi-tissue outliers were strongly enriched for variants in promoter or CpG-rich regions and had variants with higher conservation<sup>15–18</sup> and CADD (combined annotation-dependent depletion)<sup>19</sup> scores than non-outliers. We observed weaker enrichment in enhancers and transcription-factor-binding sites (Fig. 3b and Extended Data Fig. 7). Combining all classes of variation, other than non-conserved, non-coding, rare variants (excluded as less likely candidates for causal effects),

<sup>1</sup>Department of Pathology, Stanford University, Stanford, California 94305, USA. <sup>2</sup>Department of Computer Science, Johns Hopkins University, Baltimore 21218, Maryland, USA. <sup>3</sup>Biomedical Informatics Program, Stanford University, Stanford, California 94305, USA. <sup>4</sup>Department of Genetics, Stanford University, Stanford, California 94305, USA. <sup>5</sup>McDonnell Genome Institute, Washington University School of Medicine, St Louis, Missouri 63108, USA. <sup>6</sup>Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland 21218, USA. <sup>7</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. <sup>8</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA. <sup>9</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA. <sup>10</sup>Department of Medicine, Washington University School of Medicine, St Louis, Missouri 63110, USA. <sup>11</sup>Department of Genetics, Washington University School of Medicine, St Louis, Missouri 63110, USA.

†Lists of participants and their affiliations appear in the online version of the paper.

\*These authors contributed equally to this work.

§These authors jointly supervised this work.



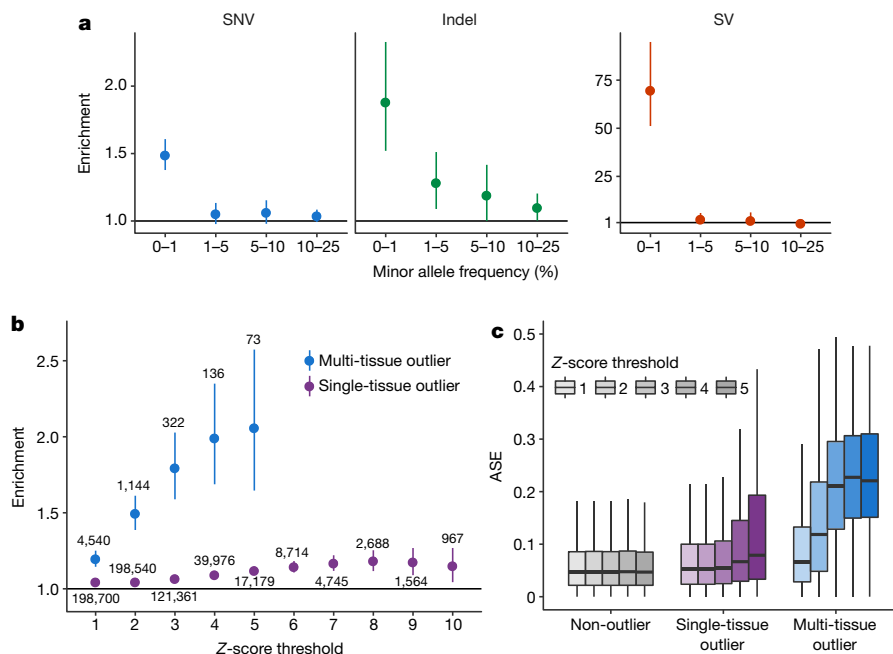
**Figure 1 | Gene expression outliers and sharing between tissues.**

**a**, A multi-tissue outlier. The individual has extreme expression values for the gene *AKRIC4* in multiple tissues (red arrows) and the most extreme median expression value across tissues. **b**, Outlier expression sharing between tissues, as measured by the proportion of single-tissue

outliers that have a  $|Z\text{-score}| \geq 2$  with the same effect direction for the corresponding genes in each replication tissue. Tissues are hierarchically clustered by gene expression. **c**, Estimated replication rate of multi-tissue outliers in a constant held-out set of tissues for different sets of discovery tissues.

we observed that 58% of underexpression and 28% of overexpression outliers had rare variants near the relevant gene, compared to 8% for non-outliers (Fig. 3c). Overexpression outliers were more

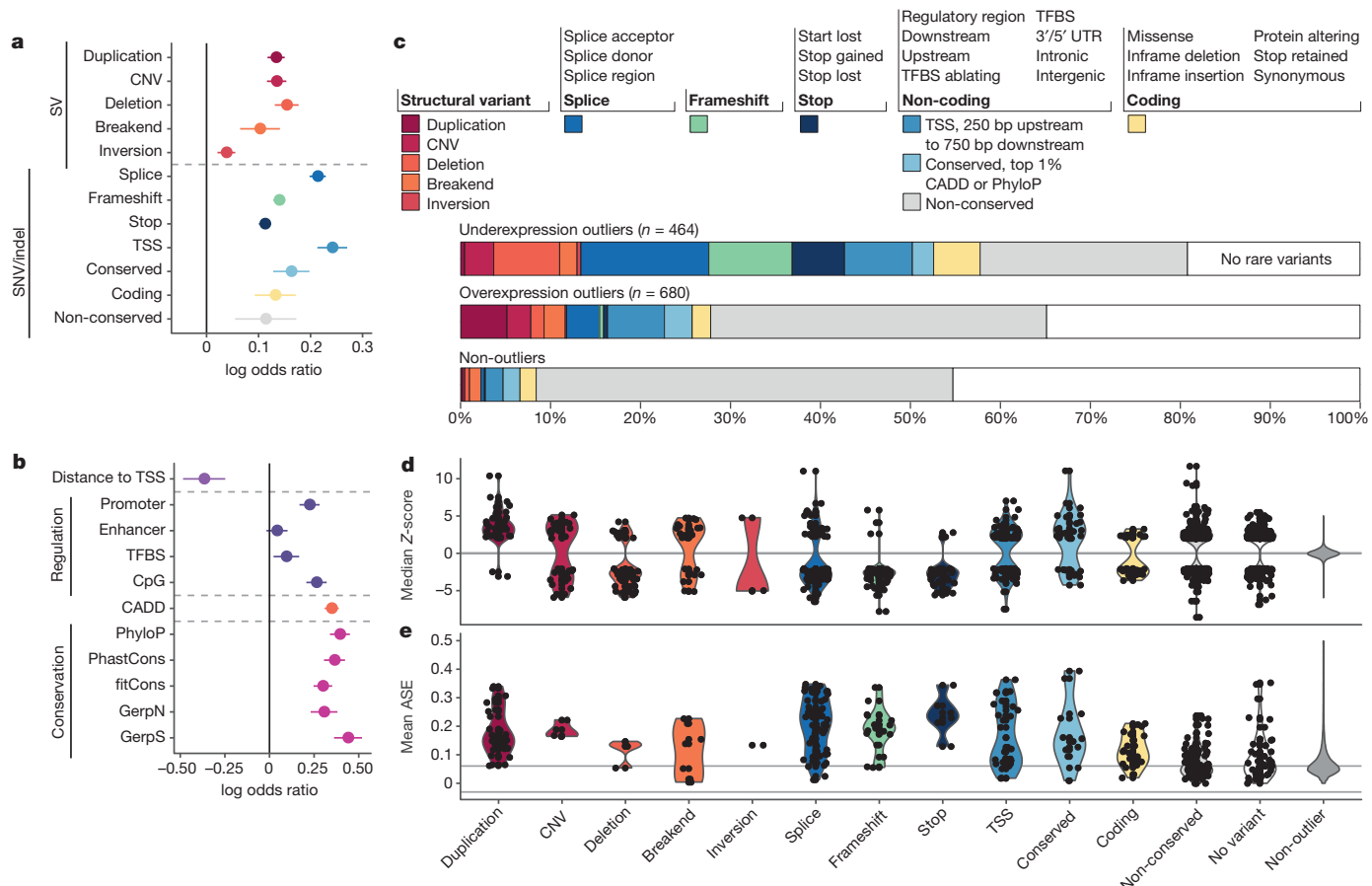
common overall, potentially because detection of underexpression outliers for very low expression genes is inherently limited (Extended Data Fig. 5b). Overexpression outliers were also less enriched



**Figure 2 | Enrichment of rare variants and ASE in outliers.**

**a**, Enrichment of SNVs, indels and structural variants (SVs) within 10 kb of the transcription start site (TSS) among outliers. For each frequency stratum, we calculated enrichment as the relative risk of having a nearby rare variant given the outlier status (see Methods). Lines indicate 95% Wald confidence intervals of the relative risk estimates. **b**, Rare SNV

enrichments at increasing Z-score thresholds. Text labels indicate the number of outliers at each threshold. **c**, ASE at increasing Z-score thresholds. ASE is measured as the magnitude of the difference between the reference-allele ratio and the null expectation of 0.5. The non-outlier category is defined in the Methods.



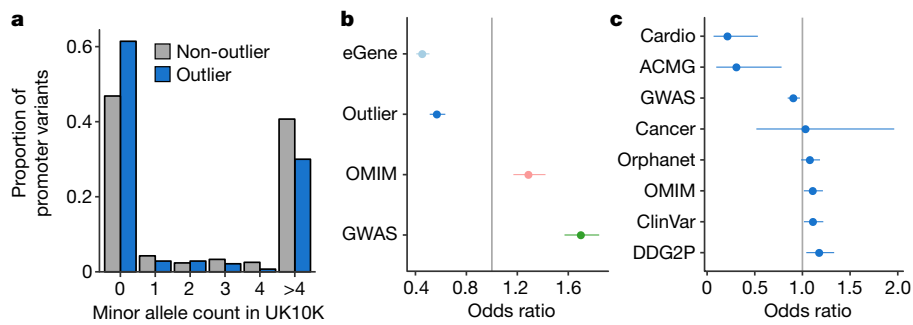
**Figure 3 | Stratification of multi-tissue outliers by rare variant classes.** We considered rare variants in the gene body and within 10 kb of the gene (200 kb for structural variants and enhancers). **a**, Enrichment of disjoint variant classes among outliers calculated as the log odds ratio with 95% Wald confidence intervals. **b**, Enrichment of functional annotations for

rare SNVs. **c**, Proportion of genes with an outlier potentially explained by each rare variant class. **d**, Distribution of median Z-scores for each variant class. **e**, For each variant class, distribution of ASE (see Methods) averaged across tissues. Grey lines mark the median values among non-outliers.

for functionally annotated rare variants (Extended Data Fig. 5c). Some variant classes had strong directionality concordant with their expected impact: duplications caused overexpression, whereas deletions, start- and stop-codon variants and frameshifts coincided with underexpression (Fig. 3d). We also observed strong ASE for outliers carrying all classes of variants, except non-conserved variants (Fig. 3e).

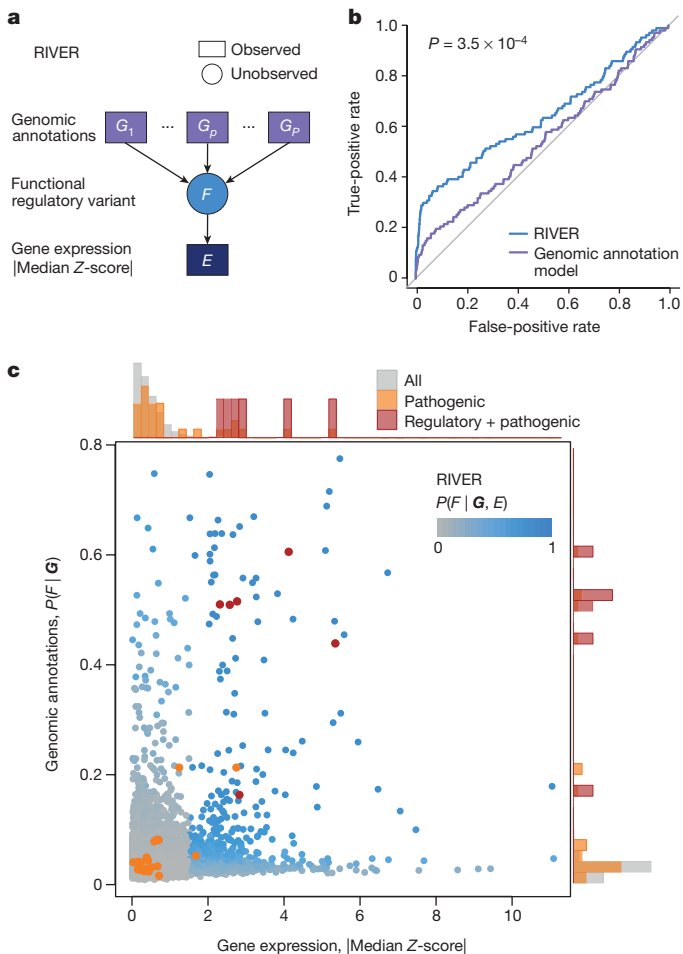
We hypothesized that functional, large-effect rare variants have been under recent selective pressure. As expected, we found that rare promoter variants of outliers were significantly less frequent in the UK10K cohort of 3,781 individuals<sup>3</sup> than rare promoter variants of

non-outliers for the same genes (two-sided Wilcoxon rank-sum test,  $P = 0.0060$ ; Fig. 4a). Additionally, genes intolerant to loss-of-function and missense mutations were depleted of both multi-tissue outliers and multi-tissue expression quantitative trait loci (eQTLs; Fisher's exact test, all  $P < 2 \times 10^{-15}$ ; Fig. 4b and Extended Data Fig. 8a). We observed a similar depletion in two curated disease gene lists—genes involved in heritable cardiovascular disease and genes in the guidelines of the American College of Medical Genetics and Genomics for incidental findings<sup>20</sup>—but not in broader gene lists (Fig. 4c and Extended Data Fig. 8b, c). Genes with a multi-tissue outlier were more likely to have a



**Figure 4 | Evolutionary constraint of genes with multi-tissue outliers.** **a**, Distributions of UK10K minor allele frequencies for promoter SNVs in outlier and non-outlier individuals at genes with multi-tissue outliers. **b**, Odds ratio of being intolerant to loss-of-function variants for genes with multi-tissue outliers, genes with shared eQTLs (eGenes), genes reported in the genome-wide association study (GWAS) catalogue and

Online Mendelian Inheritance in Man (OMIM) genes. **c**, Odds ratio of a gene having a multi-tissue outlier for each of eight sets of genes involved in complex traits or diseases (gene lists are described in the Methods; DDG2P: Developmental Disorders Genotype-to-Phenotype). In **b** and **c**, lines represent 95% confidence intervals (Fisher's exact test).



**Figure 5 | Performance of RIVER for prioritizing functional regulatory variants.** **a**, RIVER probabilistic graphical model (see Methods). **b**, Predictive power of RIVER compared to an L2-regularized logistic regression model using only genomic annotations. Accuracy was assessed using held-out individuals, who shared the same rare SNVs as observed individuals (AUCs compared with DeLong's approach<sup>29</sup>). **c**, Distribution of RIVER scores (shades of blue) as a function of expression and genomic annotation scores. The distributions of variant categories across expression and genomic annotation scores are shown as histograms aligned opposite the corresponding axes.

multi-tissue eQTL (two-sided Wilcoxon rank-sum test,  $P < 2.2 \times 10^{-16}$ ; Extended Data Fig. 8d, e), suggesting that rare and common regulatory variation influence similar genes. However, we found evidence that genes with outliers were more constrained than genes with multi-tissue eQTLs, because genes with outliers had less missense and loss-of-function variation (Tukey's range test, missense Z-score  $P = 0.0070$ , probability of loss-of-function intolerance score  $P = 0.032$ ; Fig. 4b and Extended Data Fig. 8a). This suggests that outlier expression analysis can yield unique insights into constraints on gene regulation.

Next, we sought to prioritize rare variants in each individual genome by their predicted impact on gene expression. We developed RIVER (RNA-informed variant effect on regulation), a Bayesian statistical model that jointly analyses genome and transcriptome data from the same individual to estimate the probability that a variant has regulatory impact (<https://bioconductor.org/packages/release/bioc/html/RIVER.html>, see Methods). RIVER uses a generative model that assumes that genomic annotations (Supplementary Table 3b) determine the prior probability that a variant is a functional regulatory variant, in terms of influence on gene expression, which in turn affects whether nearby genes are likely to display outlier levels of expression (Fig. 5a). RIVER does not require a labelled set of functional/non-functional variants;

rather it derives its power from identifying expression patterns that coincide with predictive genomic annotations.

We trained RIVER on the GTEx v6p cohort, and evaluated the model on held-out pairs of individuals who shared the same rare variants. We then computed the RIVER score (the posterior probability of having a functional regulatory variant) for one individual, using both expression and genomic data, and assessed the accuracy with respect to the expression levels of the second individual that had been held out (see Methods). Incorporating expression data significantly improved prediction compared with a model that uses genomic annotations alone (area under the curve (AUC) of 0.64 and 0.54, respectively,  $P = 3.5 \times 10^{-4}$ ; Fig. 5b and Extended Data Fig. 9a, b), and RIVER learned, unsupervised, to prioritize variants supported by both genomic annotations and extreme expression levels across tissues (Fig. 5c and Extended Data Fig. 9c). ASE was also enriched among the top RIVER hits compared with the genomic annotation model (Extended Data Fig. 9d). Finally, even after accounting for the most informative genomic annotations or summary scores, personal expression data were highly informative of rare variant effects (average log odds ratio, 2.76; Extended Data Fig. 9e, f).

RIVER can be used to predict regulatory effects on gene expression of disease-associated variants and aid in prioritization of rare variants in disease studies. To investigate this potential, we evaluated 27 pathogenic variants from ClinVar<sup>21</sup> present in 21 GTEx donors (Fig. 5c and Extended Data Fig. 10a). Overall, pathogenic variants had RIVER scores that were higher than background variants (two-sided Wilcoxon rank-sum test,  $P = 3.3 \times 10^{-9}$ ; Extended Data Fig. 10b–d), and the six that were probably regulatory variants (those not annotated as missense or as an indel within a coding region) scored in the 99.9th percentile. Several cases, which we evaluated in detail, illustrated that rare disease-causing variants can have a regulatory impact evident from RNA-sequencing data, even from healthy individuals that have those variants (in whom the variants are often heterozygous; Extended Data Fig. 10e, f). Note that RIVER trained on healthy cohorts, such as GTEx, can then be directly applied to new cohorts that include disease samples.

To experimentally validate a subset of the variants that were identified through outlier analysis, we used CRISPR–Cas9-mediated genome editing<sup>22,23</sup>. In K562 cells, we tested six SNVs and matched controls in transcribed regions of genes with an outlier (see Methods and Extended Data Fig. 11a, b), and compared the allelic ratios between mRNA and genomic DNA (gDNA), which was used as an internal control. All variants that were tested were SNVs in underexpression outliers and were therefore expected to decrease expression. Two variants were excluded owing to low cDNA and gDNA total reads counts. The four remaining SNVs in outliers all showed lower proportions of the alternate (installed) allele in the cDNA compared to the gDNA, confirming that these variants decreased expression (Extended Data Fig. 11c).

In summary, by combining data across multiple tissues, we curated a set of gene expression outliers that replicated at higher rates and showed stronger enrichment of rare variants than those from any single tissue. We found that rare structural variants, frameshift indels, coding variants and variants near the transcription start site were most likely to have large effects on expression. However, our ability to characterize the genetic basis of multi-tissue outliers remains incomplete. Outliers without an underlying rare variant in our analysis may be due to variants in more distal regions or in annotations we did not consider, or may be attributable to residual technical or environmental effects.

Although variant interpretation remains challenging, RIVER demonstrates the value of incorporating personal gene expression data to examine the consequences of rare variants that may be uncertain based on the sequence alone. Our results suggest that a general approach can be applied to studies that supplement genome sequencing with other molecular phenotypes, such as methylation<sup>24–26</sup> and histone modification<sup>27,28</sup>. We anticipate that such integrative approaches will be

essential for effective interpretation of genome-wide genetic variation on a personalized level.

**Online Content** Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

**Received 8 September 2016; accepted 13 September 2017.**

1. Tennesen, J. A. *et al.* Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *Science* **337**, 64–69 (2012).
2. Nelson, M. R. *et al.* An abundance of rare functional variants in 202 drug target genes sequenced in 14,002 people. *Science* **337**, 100–104 (2012).
3. The UK10K Consortium. The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82–90 (2015).
4. Keinan, A. & Clark, A. G. Recent explosive human population growth has resulted in an excess of rare genetic variants. *Science* **336**, 740–743 (2012).
5. Uricchio, L. H., Zaitlen, N. A., Ye, C. J., Witte, J. S. & Hernandez, R. D. Selection and explosive growth alter genetic architecture and hamper the detection of causal rare variants. *Genome Res.* **26**, 863–873 (2016).
6. MacArthur, D. G. *et al.* A systematic survey of loss-of-function variants in human protein-coding genes. *Science* **335**, 823–828 (2012).
7. Narasimhan, V. M. *et al.* Health and population effects of rare gene knockouts in adult humans with related parents. *Science* **352**, 474–477 (2016).
8. Montgomery, S. B., Lappalainen, T., Gutierrez-Arcelus, M. & Dermitzakis, E. T. Rare and common regulatory variation in population-scale sequenced human genomes. *PLoS Genet.* **7**, e1002144 (2011).
9. Zhao, J. *et al.* A burden of rare variants associated with extremes of gene expression in human peripheral blood. *Am. J. Hum. Genet.* **98**, 299–309 (2016).
10. Zeng, Y. *et al.* Aberrant gene expression in humans. *PLoS Genet.* **11**, e1004942 (2015).
11. Li, X. *et al.* Transcriptome sequencing of a large human family identifies the impact of rare noncoding variants. *Am. J. Hum. Genet.* **95**, 245–256 (2014).
12. The GTEx Consortium. Genetic effects on gene expression across human tissues. <https://doi.org/10.1038/nature24277> (2017).
13. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nat. Protoc.* **7**, 500–507 (2012).
14. Chiang, C. *et al.* The impact of structural variation on human gene expression. *Nat. Genet.* **49**, 692–699 (2017).
15. Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* **20**, 110–121 (2010).
16. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).
17. Cooper, G. M. *et al.* Distribution and intensity of constraint in mammalian genomic sequence. *Genome Res.* **15**, 901–913 (2005).
18. Arbiza, L. *et al.* Genome-wide inference of natural selection on human transcription factor binding sites. *Nat. Genet.* **45**, 723–729 (2013).
19. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
20. Green, R. C. *et al.* ACMG recommendations for reporting of incidental findings in clinical exome and genome sequencing. *Genet. Med.* **15**, 565–574 (2013).
21. Landrum, M. J. *et al.* ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic Acids Res.* **44**, D862–D868 (2016).
22. Hendel, A. *et al.* Chemically modified guide RNAs enhance CRISPR–Cas genome editing in human primary cells. *Nat. Biotechnol.* **33**, 985–989 (2015).
23. Hess, G. T. *et al.* Directed evolution using dCas9-targeted somatic hypermutation in mammalian cells. *Nat. Methods* **13**, 1036–1042 (2016).
24. Grundberg, E. *et al.* Global analysis of DNA methylation variation in adipose tissue from twins reveals links to disease-associated variants in distal regulatory elements. *Am. J. Hum. Genet.* **93**, 876–890 (2013).
25. Gamazon, E. R. *et al.* Enrichment of *cis*-regulatory gene expression SNPs and methylation quantitative trait loci among bipolar disorder susceptibility variants. *Mol. Psychiatry* **18**, 340–346 (2013).
26. Bell, J. T. *et al.* DNA methylation patterns associate with genetic and gene expression variation in HapMap cell lines. *Genome Biol.* **12**, R10 (2011).

27. Waszak, S. M. *et al.* Population variation and genetic control of modular chromatin architecture in humans. *Cell* **162**, 1039–1050 (2015).
28. Grubert, F. *et al.* Genetic control of chromatin states in humans involves local and distal chromosomal interactions. *Cell* **162**, 1051–1065 (2015).
29. DeLong, E. R., DeLong, D. M. & Clarke-Pearson, D. L. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics* **44**, 837–845 (1988).

**Supplementary Information** is available in the online version of the paper.

**Acknowledgements** We thank members of the MacArthur laboratory and the Laboratory, Data Analysis, and Coordinating Center (LDACC) for performing the quality control of the whole genome sequencing data, D. Conrad for help with the structural variant calls, D. A. Knowles for code review, J. T. Leek and C. D. Brown for feedback on the manuscript, and the artists of the graphics that we modified in Fig. 1 (<https://pixabay.com/en/man-silhouette-stand-straight-308387/> and <http://www.allvectors.com/human-organs/>). The Genotype-Tissue Expression (GTEx) project was supported by the Common Fund of the Office of the Director of the National Institutes of Health (NIH). Additional funds were provided by the National Cancer Institute; National Human Genome Research Institute (NHGRI); National Heart, Lung, and Blood Institute; National Institute on Drug Abuse; National Institute of Mental Health; and National Institute of Neurological Disorders and Stroke. Donors were enrolled at the Biospecimen Source Sites funded by Leidos Biomedical, Inc. (Leidos) subcontracts to the National Disease Research Interchange (10XS170) and Roswell Park Cancer Institute (10XS171). The LDACC was funded through a contract (HHSN268201000029C) to The Broad Institute. Biorepository operations were funded through a Leidos subcontract to the Van Andel Institute (10ST1035). Additional data repository and project management were provided by Leidos (HHSN261200800001E). The Brain Bank was supported by a supplement to University of Miami grant DA006227. We are grateful for support from a Hewlett-Packard Stanford Graduate Fellowship (E.K.T.), a doctoral scholarship from the Natural Science and Engineering Council of Canada (E.K.T.), a Lucille P. Markey Biomedical Research Stanford Graduate Fellowship (J.R.D.), the Stanford Genome Training Program (SGTP; NHGRI T32HG000044) (J.R.D., Z.Z.), the National Science Foundation GRFP (DGE-114747) (Z.Z.), the Joseph C. Pistrutto Research Fellowship (F.N.D.), NIH training grant T32 GM007057 (B.J.S.), a Mr and Mrs Spencer T. Olin Fellowship for Women in Graduate Study (A.J.S.), the Searle Scholars Program (A.B.), NIH grants 1R01MH109905-01 (A.B.), R01MH101814 (NIH Common Fund; GTEx Program) (A.B. and S.B.M.), R01HG008150 (NHGRI; Non-Coding Variants Program) (A.B., S.B.M.), and NHGRI grants U01HG007436 and U01HG009080 (S.B.M.).

**Author Contributions** X.L., Y.K., E.K.T., J.R.D., A.B. and S.B.M. designed the study, performed analyses and wrote the manuscript. Y.K., F.N.D. and A.B. developed RIVER. G.T.H., A.L. and M.C.B. designed and executed the validation using CRISPR–Cas9. C.C., A.J.S. and I.M.H. provided the set of structural variants. J.D.M. provided the lists of curated cancer and cardiovascular disease genes. Z.Z., B.J.S. and A.G. contributed to analysis and feedback.

**Author Information** Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations. Correspondence and requests for materials should be addressed to A.B. ([ajbatt@cs.jhu.edu](mailto:ajbatt@cs.jhu.edu)) or S.B.M. ([smontgom@stanford.edu](mailto:smontgom@stanford.edu)).

**Reviewer Information** *Nature* thanks E. Birney, A. Clark and Y. Gilad for their contribution to the peer review of this work.



This work is licensed under a Creative Commons Attribution 4.0 International (CC BY 4.0) licence. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons licence, users will need to obtain permission from the licence holder to reproduce the material. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

**GTEx Consortium****Laboratory, Data Analysis & Coordinating Center (LDACC)—Analysis**

**Working Group** François Aguet<sup>1</sup>, Kristin G. Ardlie<sup>1</sup>, Beryl B. Cummings<sup>1,2</sup>, Ellen T. Gelfand<sup>1</sup>, Gad Getz<sup>1,3</sup>, Kane Hadley<sup>1</sup>, Robert E. Handsaker<sup>1,4</sup>, Katherine H. Huang<sup>1</sup>, Seva Kashin<sup>4</sup>, Konrad J. Karczewski<sup>1,2</sup>, Monkol Lek<sup>1,2</sup>, Xiao Li<sup>1</sup>, Daniel G. MacArthur<sup>1,2</sup>, Jared L. Nedzel<sup>1</sup>, Duyen T. Nguyen<sup>1</sup>, Michael S. Noble<sup>1</sup>, Ayyellet V. Segrè<sup>1</sup>, Casandra A. Trowbridge<sup>1</sup>, Taru Tukiainen<sup>1,2</sup>

**Statistical Methods groups—Analysis Working Group** Nathan S. Abell<sup>5,6</sup>, Brunilda Balliu<sup>6</sup>, Ruth Barshir<sup>7</sup>, Omer Basha<sup>7</sup>, Alexis Battle<sup>8</sup>, Gireesh K. Bogu<sup>9,10</sup>, Andrew Brown<sup>11,12,13</sup>, Christopher D. Brown<sup>14</sup>, Stephane E. Castel<sup>15,16</sup>, Lin S. Chen<sup>17</sup>, Colby Chiang<sup>18</sup>, Donald F. Conrad<sup>19,20</sup>, Nancy J. Cox<sup>21</sup>, Farhan N. Daman<sup>8</sup>, Joe R. Davis<sup>5,6</sup>, Olivier Delaneau<sup>11,12,13</sup>, Emmanouil T. Dermitzakis<sup>11,12,13</sup>, Barbara E. Engelhardt<sup>22</sup>, Eleazar Eskin<sup>23,24</sup>, Pedro G. Ferreira<sup>25,26</sup>, Laure Frésard<sup>5,6</sup>, Eric R. Gamazon<sup>21,27,28</sup>, Diego Garrido-Martin<sup>9,10</sup>, Ariel D.H. Gewirtz<sup>29</sup>, Genna Gliner<sup>30</sup>, Michael J. Gludemans<sup>5,6,31</sup>, Roderic Guigo<sup>9,10,32</sup>, Ira M. Hall<sup>18,19,33</sup>, Buhm Han<sup>34</sup>, Yuan He<sup>35</sup>, Farhad Hormozdizadeh<sup>23</sup>, Cedric Howald<sup>11,12,13</sup>, Hae Kyung Im<sup>36</sup>, Brian Jo<sup>29</sup>, Eun Yong Kang<sup>23</sup>, Yungil Kim<sup>8</sup>, Sarah Kim-Hellmuth<sup>15,16</sup>, Tuuli Lappalainen<sup>15,16</sup>, Gen Li<sup>37</sup>, Xin Li<sup>6</sup>, Boxiang Liu<sup>5,6,38</sup>, Serghei Mangui<sup>23</sup>, Mark I. McCarthy<sup>39,40,41</sup>, Ian C. McDowell<sup>42</sup>, Pejman Mohammadi<sup>15,16</sup>, Jean Monlong<sup>9,10,43</sup>, Stephen B. Montgomery<sup>5,6</sup>, Manuel Muñoz-Aguirre<sup>9,10,44</sup>, Anne W. Ndungu<sup>39</sup>, Dan L. Nicolae<sup>36,45,46</sup>, Andrew B. Nobel<sup>47,48</sup>, Meritxell Oliva<sup>36,49</sup>, Halit Ongen<sup>11,12,13</sup>, John J. Palowitch<sup>47</sup>, Nikolaos Panousis<sup>11,12,13</sup>, Panagiotis Papanikolaou<sup>9,10</sup>, YoSon Park<sup>14</sup>, Princy Parsana<sup>8</sup>, Anthony J. Payne<sup>39</sup>, Christine B. Peterson<sup>50</sup>, Jie Quan<sup>51</sup>, Ferran Reverter<sup>9,10,52</sup>, Chiara Sabatti<sup>53,54</sup>, Ashis Saha<sup>8</sup>, Michael Sammeth<sup>55</sup>, Alexandra J. Scott<sup>18</sup>, Andrew A. Shabalin<sup>56</sup>, Reza Sodaei<sup>9,10</sup>, Matthew Stephens<sup>45,46</sup>, Barbara E. Stranger<sup>36,49,57</sup>, Benjamin J. Strober<sup>35</sup>, Jae Hoon Sul<sup>58</sup>, Emily K. Tsang<sup>5,31</sup>, Sarah Urbut<sup>46</sup>, Martijn van de Bunt<sup>39,40</sup>, Gao Wang<sup>46</sup>, Xiaquan Wen<sup>59</sup>, Fred A. Wright<sup>60</sup>, Hualin S. Xi<sup>51</sup>, Esti Yeger-Lotem<sup>7,61</sup>, Zachary Zappala<sup>5,6</sup>, Judith B. Zaugg<sup>62</sup>, Yi-Hui Zhou<sup>60</sup>

**Enhancing GTEx (eGTEx) groups** Joshua M. Akey<sup>29,63</sup>, Daniel Bates<sup>64</sup>, Joanne Chan<sup>5</sup>, Lin S. Chen<sup>17</sup>, Melina Claussnitzer<sup>1,65,66</sup>, Kathryn Demanelis<sup>17</sup>, Morgan Diegel<sup>64</sup>, Jennifer A. Doherty<sup>67</sup>, Andrew P. Feinberg<sup>35,68,69,70</sup>, Marian S. Fernando<sup>36,49</sup>, Jessica Halow<sup>64</sup>, Kasper D. Hansen<sup>68,71,72</sup>, Eric Haugen<sup>64</sup>, Peter F. Hickey<sup>72</sup>, Lei Hou<sup>1,73</sup>, Farzana Jasmine<sup>17</sup>, Ruiqi Jian<sup>5</sup>, Lihua Jiang<sup>5</sup>, Audra Johnson<sup>64</sup>, Rajinder Kaul<sup>64</sup>, Manolis Kellis<sup>1,73</sup>, Muhammad G. Kibriya<sup>17</sup>, Kristen Lee<sup>64</sup>, Jin Billy Li<sup>5</sup>, Qin Li<sup>5</sup>, Xiao Li<sup>5</sup>, Jessica Lin<sup>5,74</sup>, Shin Lin<sup>5,75</sup>, Sandra Linde<sup>5,6</sup>, Caroline Linke<sup>36,49</sup>, Yaping Liu<sup>1,73</sup>, Matthew T. Maurano<sup>76</sup>, Benoit Molinier<sup>1</sup>, Stephen B. Montgomery<sup>5,6</sup>, Jemma Nelson<sup>64</sup>, Fidencio J. Neri<sup>64</sup>, Meritxell Oliva<sup>36,49</sup>, Yongjin Park<sup>1,73</sup>, Brandon L. Pierce<sup>17</sup>, Nicola J. Rinaldi<sup>1,73</sup>, Lindsay F. Rizzardi<sup>68</sup>, Richard Sandstrom<sup>64</sup>, Andrew Skol<sup>36,49,57</sup>, Kevin S. Smith<sup>5,6</sup>, Michael P. Snyder<sup>5</sup>, John Stamatoyannopoulos<sup>64,74,77</sup>, Barbara E. Stranger<sup>36,49,57</sup>, Hua Tang<sup>5</sup>, Emily K. Tsang<sup>5,31</sup>, Li Wang<sup>1</sup>, Meng Wang<sup>5</sup>, Nicholas Van Wittenbergh<sup>1</sup>, Fan Wu<sup>36,49</sup>, Rui Zhang<sup>5</sup>

**NIH Common Fund** Concepcion R. Nierras<sup>78</sup>

**NIH/NCI** Philip A. Branton<sup>79</sup>, Latarsha J. Carithers<sup>79,80</sup>, Ping Guan<sup>79</sup>, Helen M. Moore<sup>79</sup>, Abhi Rao<sup>79</sup>, Jimmie B. Vaught<sup>79</sup>

**NIH/NHGRI** Sarah E. Gould<sup>81</sup>, Nicole C. Lockart<sup>81</sup>, Casey Martin<sup>81</sup>, Jeffery P. Struwing<sup>81</sup>, Simona Volpi<sup>81</sup>

**NIH/NIMH** Anjene M. Addington<sup>82</sup>, Susan E. Koester<sup>82</sup>

**NIH/NIDA** A. Roger Little<sup>83</sup>

**Biospecimen Collection Source Site—NDRI** Lori E. Brigham<sup>84</sup>, Richard Hasz<sup>85</sup>, Marcus Hunter<sup>86</sup>, Christopher Johns<sup>87</sup>, Mark Johnson<sup>88</sup>, Gene Kopen<sup>89</sup>, William F. Leinweber<sup>89</sup>, John T. Lonsdale<sup>89</sup>, Alisa McDonald<sup>89</sup>, Bernadette Mestichelli<sup>89</sup>, Kevin Myer<sup>86</sup>, Brian Roe<sup>86</sup>, Michael Salvatore<sup>89</sup>, Saboor Shad<sup>89</sup>, Jeffrey A. Thomas<sup>89</sup>, Gary Walters<sup>88</sup>, Michael Washington<sup>88</sup>, Joseph Wheeler<sup>87</sup>

**Biospecimen Collection Source Site—RPCI** Jason Bridge<sup>90</sup>, Barbara A. Foster<sup>91</sup>, Bryan M. Gillard<sup>91</sup>, Ellen Karasik<sup>91</sup>, Rachna Kumar<sup>91</sup>, Mark Miklos<sup>90</sup>, Michael T. Moser<sup>91</sup>

**Biospecimen Core Resource—VARI** Scott D. Jewell<sup>92</sup>, Robert G. Montroy<sup>92</sup>, Daniel C. Rohrer<sup>92</sup>, Dana R. Valley<sup>92</sup>

**Brain Bank Repository—University of Miami Brain Endowment Bank** David A. Davis<sup>93</sup>, Deborah C. Mash<sup>93</sup>

**Leidos Biomedical—Project Management** Anita H. Undale<sup>94</sup>, Anna M. Smith<sup>95</sup>, David E. Tabo<sup>95</sup>, Nancy V. Roche<sup>95</sup>, Jeffrey A. McLean<sup>95</sup>, Negin Vatanian<sup>95</sup>, Karna L. Robinson<sup>95</sup>, Leslie Sobin<sup>95</sup>, Mary E. Barcus<sup>96</sup>, Kimberly M. Valentino<sup>95</sup>, Liqun Qi<sup>95</sup>, Steven Hunter<sup>95</sup>, Pushpa Hariharan<sup>95</sup>, Shilpi Singh<sup>95</sup>, Ki Sung Um<sup>95</sup>, Takunda Matose<sup>95</sup>, Maria M. Tomaszewski<sup>95</sup>

**ELSI Study** Laura K. Barker<sup>97</sup>, Maghboeba Mosavel<sup>98</sup>, Laura A. Siminoff<sup>97</sup>, Heather M. Traino<sup>97</sup>

**Genome Browser Data Integration & Visualization—EBI** Paul Flicek<sup>99</sup>, Thomas Juettemann<sup>99</sup>, Magali Ruffier<sup>99</sup>, Dan Sheppard<sup>99</sup>, Kieron Taylor<sup>99</sup>, Stephen J. Trevanion<sup>99</sup>, Daniel R. Zerbino<sup>99</sup>

**Genome Browser Data Integration & Visualization—UCSC Genomics Institute, University of California Santa Cruz** Brian Craft<sup>100</sup>, Mary Goldman<sup>100</sup>, Maximilian Haussler<sup>100</sup>, W. James Kent<sup>100</sup>, Christopher M. Lee<sup>100</sup>, Benedict Paten<sup>100</sup>, Kate R. Rosenbloom<sup>100</sup>, John Vivian<sup>100</sup>, Jingchun Zhu<sup>100</sup>

<sup>1</sup>The Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts 02142, USA. <sup>2</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. <sup>3</sup>Massachusetts General Hospital Cancer Center and Department of Pathology, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. <sup>4</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts 02114, USA. <sup>5</sup>Department of Genetics, Stanford University, Stanford, California 94305, USA. <sup>6</sup>Department of Pathology, Stanford University, Stanford, California 94305, USA. <sup>7</sup>Department of Clinical Biochemistry and Pharmacology, Faculty of Health Sciences, Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel. <sup>8</sup>Department of Computer Science, Johns Hopkins University, Baltimore, Maryland 21218, USA. <sup>9</sup>Centre for Genomic Regulation (CRG), The Barcelona Institute for Science and Technology, Dr Aiguader 88, 08003 Barcelona, Spain. <sup>10</sup>Universitat Pompeu Fabra (UPF), 08002 Barcelona, Spain. <sup>11</sup>Department of Genetic Medicine and Development, University of Geneva Medical School, 1211 Geneva, Switzerland. <sup>12</sup>Institute for Genetics and Genomics in Geneva (iG3), University of Geneva, 1211 Geneva, Switzerland. <sup>13</sup>Swiss Institute of Bioinformatics, 1211 Geneva, Switzerland. <sup>14</sup>Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania 19104, USA. <sup>15</sup>New York Genome Center, New York, New York 10013, USA. <sup>16</sup>Department of Systems Biology, Columbia University Medical Center, New York, New York 10032, USA. <sup>17</sup>Department of Public Health Sciences, The University of Chicago, Chicago, Illinois 60637, USA. <sup>18</sup>McDonnell Genome Institute, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>19</sup>Department of Genetics, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>20</sup>Department of Pathology & Immunology, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>21</sup>Division of Genetic Medicine, Department of Medicine, Vanderbilt University Medical Center, Nashville, Tennessee 37232, USA. <sup>22</sup>Department of Computer Science, Center for Statistics and Machine Learning, Princeton University, Princeton, New Jersey 08540, USA. <sup>23</sup>Department of Computer Science, University of California, Los Angeles, California 90095, USA. <sup>24</sup>Department of Human Genetics, University of California, Los Angeles, California 90095, USA. <sup>25</sup>Instituto de Investigação e Inovação em Saúde (i3S), Universidade do Porto, 4200-135 Porto, Portugal. <sup>26</sup>Institute of Molecular Pathology and Immunology (IPATIMUP), University of Porto, 4200-625 Porto, Portugal. <sup>27</sup>Department of Clinical Epidemiology, Biostatistics and Bioinformatics, Academic Medical Center, University of Amsterdam, 1105 AZ Amsterdam, The Netherlands. <sup>28</sup>Department of Psychiatry, Academic Medical Center, University of Amsterdam, 1105 AZ Amsterdam, The Netherlands. <sup>29</sup>Lewis Sigler Institute, Princeton University, Princeton, New Jersey 08540, USA. <sup>30</sup>Department of Operations Research and Financial Engineering, Princeton University, Princeton, New Jersey 08540, USA. <sup>31</sup>Biomedical Informatics Program, Stanford University, Stanford, California 94305, USA. <sup>32</sup>Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), 08003 Barcelona, Spain. <sup>33</sup>Department of Medicine, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>34</sup>Department of Convergence Medicine, University of Ulsan College of Medicine, Asan Medical Center, Seoul 138-736, South Korea. <sup>35</sup>Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland 21218, USA. <sup>36</sup>Section of Genetic Medicine, Department of Medicine, The University of Chicago, Chicago, Illinois 60637, USA. <sup>37</sup>Department of Biostatistics, Mailman School of Public Health, Columbia University, New York, New York 10032, USA. <sup>38</sup>Department of Biology, Stanford University, Stanford, California 94305, USA. <sup>39</sup>Wellcome Trust Centre for Human Genetics, Nuffield Department of Medicine, University of Oxford, Oxford, OX3 7BN, UK. <sup>40</sup>Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford, Churchill Hospital, Oxford, OX3 7LE, UK. <sup>41</sup>Oxford NIHR Biomedical Research Centre, Churchill Hospital, Oxford, OX3 7LJ, UK. <sup>42</sup>Computational Biology & Bioinformatics Graduate Program, Duke University, Durham, North Carolina 27708, USA. <sup>43</sup>Human Genetics Department, McGill University, Montreal, Quebec H3A 0G1, Canada. <sup>44</sup>Departament d'Estadística i Investigació Operativa, Universitat Politècnica de Catalunya, 08034 Barcelona, Spain. <sup>45</sup>Department of Statistics, The University of Chicago, Chicago, Illinois 60637, USA. <sup>46</sup>Department of Human Genetics, The University of Chicago, Chicago, Illinois 60637, USA. <sup>47</sup>Department of Statistics and Operations Research, University of North Carolina, Chapel Hill, North Carolina 27599, USA. <sup>48</sup>Department of Biostatistics, University of North Carolina, Chapel Hill, North Carolina 27599, USA. <sup>49</sup>Institute for Genomics and Systems Biology, The University of Chicago, Chicago, Illinois 60637, USA. <sup>50</sup>Department of Biostatistics, The University of Texas MD Anderson Cancer Center, Houston, Texas 77030, USA. <sup>51</sup>Computational Sciences, Pfizer Inc, Cambridge, Massachusetts 02139, USA. <sup>52</sup>Universitat de Barcelona, 08028 Barcelona, Spain. <sup>53</sup>Department of Biomedical Data Science, Stanford University, Stanford, California 94305, USA. <sup>54</sup>Department of Statistics, Stanford University, Stanford, California 94305, USA. <sup>55</sup>Institute of Biophysics Carlos Chagas Filho (IBCCF), Federal University of Rio de Janeiro (UFRJ), 21941902 Rio de Janeiro, Brazil. <sup>56</sup>Department of Psychiatry, University of Utah, Salt Lake City, Utah 84108, USA. <sup>57</sup>Center for Data Intensive Science, The University of Chicago, Chicago, Illinois 60637, USA. <sup>58</sup>Department of Psychiatry and Biobehavioral Sciences, University of California, Los Angeles, California 90095, USA. <sup>59</sup>Department of Biostatistics, University of Michigan, Ann Arbor, Michigan 48109, USA. <sup>60</sup>Bioinformatics Research Center and Departments of Statistics and Biological Sciences, North Carolina State University, Raleigh, North Carolina 27695, USA. <sup>61</sup>National Institute for Biotechnology in the Negev, Beer-Sheva 84105 Israel. <sup>62</sup>European Molecular Biology Laboratory, 69117 Heidelberg, Germany. <sup>63</sup>Department of Ecology and Evolutionary Biology, Princeton University, Princeton, New Jersey 08540, USA. <sup>64</sup>Altius Institute for Biomedical Sciences, Seattle, Washington 98121, USA. <sup>65</sup>Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, Massachusetts 02215, USA. <sup>66</sup>University of Hohenheim, 70599 Stuttgart, Germany. <sup>67</sup>Huntsman Cancer Institute, Department of Population Health Sciences, University of Utah, Salt Lake City, Utah 84112, USA. <sup>68</sup>Center for Epigenetics, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA. <sup>69</sup>Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA. <sup>70</sup>Department of Mental Health, Johns Hopkins University School of Public Health, Baltimore, Maryland 21205, USA. <sup>71</sup>McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins School of Medicine,

Baltimore, Maryland 21205, USA. <sup>72</sup>Department of Biostatistics, Johns Hopkins University, Baltimore, Maryland 21205, USA. <sup>73</sup>Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA. <sup>74</sup>Department of Medicine, University of Washington, Seattle, Washington 98195, USA. <sup>75</sup>Division of Cardiology, University of Washington, Seattle, Washington 98195, USA. <sup>76</sup>Institute for Systems Genetics, New York University Langone Medical Center, New York, New York 10016, USA. <sup>77</sup>Department of Genome Sciences, University of Washington, Seattle, Washington 98195, USA. <sup>78</sup>Office of Strategic Coordination, Division of Program Coordination, Planning and Strategic Initiatives, Office of the Director, NIH, Rockville, Maryland 20852, USA. <sup>79</sup>Biorepositories and Biospecimen Research Branch, Division of Cancer Treatment and Diagnosis, National Cancer Institute, Bethesda, Maryland 20892, USA. <sup>80</sup>National Institute of Dental and Craniofacial Research, Bethesda, Maryland 20892, USA. <sup>81</sup>Division of Genomic Medicine, National Human Genome Research Institute, Rockville, Maryland 20852, USA. <sup>82</sup>Division of Neuroscience and Basic Behavioral Science, National Institute of Mental Health, NIH, Bethesda, Maryland 20892, USA. <sup>83</sup>Division of Neuroscience and Behavior, National Institute on Drug Abuse, NIH, Bethesda, Maryland 20892, USA. <sup>84</sup>Washington Regional Transplant Community, Falls Church, Virginia 22003, USA. <sup>85</sup>Gift of Life Donor Program, Philadelphia, Pennsylvania 19103, USA. <sup>86</sup>LifeGift, Houston, Texas 77055, USA. <sup>87</sup>Center for Organ Recovery and Education, Pittsburgh, Pennsylvania 15238, USA. <sup>88</sup>LifeNet Health, Virginia Beach, Virginia 23453, USA. <sup>89</sup>National Disease Research Interchange, Philadelphia, Pennsylvania 19103, USA. <sup>90</sup>Unyts, Buffalo, New York 14203, USA. <sup>91</sup>Pharmacology and Therapeutics, Roswell Park Cancer Institute, Buffalo, New York 14263, USA. <sup>92</sup>Van Andel Research Institute, Grand Rapids, Michigan 49503, USA. <sup>93</sup>Brain Endowment Bank, Miller School of Medicine, University of Miami, Miami, Florida 33136, USA. <sup>94</sup>National Institute of Allergy and Infectious Diseases, NIH, Rockville, Maryland 20852, USA. <sup>95</sup>Biospecimen Research Group, Clinical Research Directorate, Leidos Biomedical Research, Inc., Rockville, Maryland 20852, USA. <sup>96</sup>Leidos Biomedical Research, Inc., Frederick, Maryland 21701, USA. <sup>97</sup>Temple University, Philadelphia, Pennsylvania 19122, USA. <sup>98</sup>Department of Health Behavior and Policy, School of Medicine, Virginia Commonwealth University, Richmond, Virginia 23298, USA. <sup>99</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton CB10 1SD, UK. <sup>100</sup>UCSC Genomics Institute, University of California Santa Cruz, Santa Cruz, California 95064, USA.

## METHODS

**Study population.** All human subjects were deceased donors. Informed consent was obtained for all donors via next-of-kin consent to permit the collection and banking of de-identified tissue samples for scientific research. The research protocol was reviewed by Chesapeake Research Review Inc., Roswell Park Cancer Institute's Office of Research Subject Protection, and the institutional review board of the University of Pennsylvania. We used the RNA-seq, allele-specific expression, and whole-genome sequencing (WGS) data from the v6p release of the GTEx project. The generation of these data are described in the supplementary information of ref. 12.

**Correction for technical confounders.** We restricted our expression analyses to the 449 individuals and 44 tissues for which sex and the top three genotype principal components, which capture major population stratification, were available. For each tissue, we  $\log_2$ -transformed all expression values ( $\log_2(\text{RPKM} + 2)$ ), where RPKM is the number of reads per kilobase of transcript per million mapped reads. We then standardized the expression of each gene to prevent shrinkage of outlier expression values caused by quantile normalization. To remove unmeasured batch effects and other confounders, for each tissue separately, we estimated hidden factors using PEER<sup>13</sup> on the transformed expression values. In each tissue, we defined expressed genes and corrected for the same number of PEER factors as in the GTEx eQTL analyses (see supplementary information of ref. 12). We regressed out the PEER factors, the top three genotype principal components and sex (where appropriate) from the transformed expression data for each tissue using the following linear model:

$$Y_g = \mu_g \mathbf{1} + \sum_{n=1}^N \alpha_{g,n} \mathbf{P}_n + \sum_{k=1}^3 \beta_{g,k} \mathbf{G}_k + \gamma_g \mathbf{S} + \varepsilon_g$$

where  $Y_g$  is the transformed expression of a given gene  $g$ ,  $\mu_g$  is the mean expression level for the gene,  $\mathbf{P}_n$  is the  $n$ th PEER factor,  $\mathbf{G}_1$ ,  $\mathbf{G}_2$ ,  $\mathbf{G}_3$  are the top three genotype principal components, and  $\mathbf{S}$  is the sex covariate. We assumed the residual vector  $\varepsilon_g$  follows the multivariate normal distribution  $\varepsilon_g \sim N(0, \sigma^2 \mathbf{I})$ . Finally, we standardized the expression residuals  $\varepsilon_g$  for each gene, which yielded  $Z$ -scores.

To better understand the effect of PEER correction on the removal of technical and biological confounders, we compared the PEER factors in each tissue separately to pre-collected sample and subject covariates. We considered the subset of covariates with  $>50$  observations in at least 31 tissues, where we first selected covariates with more than one unique entry in each tissue. For categorical covariates, we only considered categories with more than 20 observations. For each PEER factor and each covariate, we fit a linear model with the PEER factor as the response and the covariate as the predictor. From this model, we computed the proportion of that PEER factor's variance explained by the covariate as the adjusted  $R^2$ :

$$\text{Adjusted } R^2 = R^2 - \left[ (1 - R^2) \cdot \frac{p}{n - p - 1} \right]$$

where  $p$  and  $n$  are the number of parameters and samples, respectively, and

$$R^2 = \frac{SS_T - SS_R}{SS_T}$$

$SS_T$  and  $SS_R$  refer to the total and residual sums of squares, respectively.

To quantify the degree to which each covariate was captured by the combination of all PEER factors, genotype principal components and sex (where appropriate) for each tissue, we considered the expression component regressed out from the uncorrected data:

$$W_g = Y_g - \varepsilon_g$$

For each covariate, we then fit a linear model with  $W_g$  as the response and the covariate as the predictor. We assessed the proportion of the variance of  $W_g$  explained by each covariate by computing the adjusted  $R^2$  for the covariate across all genes. We used the formula above, but summed across all genes to compute  $SS_T$  and  $SS_R$ .

To assess the impact of PEER correction on rare variant enrichment, we also tried removing either the top five PEER factors for each tissue or no PEER factors. We then performed multi-tissue outlier calling and tested the enrichment of rare and common variants in the two partially corrected datasets (see 'Enrichment of rare and common variants near outlier genes').

**Single-tissue and multi-tissue outlier discovery.** Single-tissue and multi-tissue outlier calling was restricted to autosomal lincRNA and protein-coding genes. For each tissue, an individual was called a single-tissue outlier for a particular gene if that individual had the largest absolute  $Z$ -score and the absolute value was

at least 2. For each gene, the individual with the most extreme median  $Z$ -score taken across tissues was identified as a multi-tissue outlier for that gene provided the absolute median  $Z$ -score was at least 2. Therefore, each gene had at most one single-tissue outlier per tissue and one multi-tissue outlier. Under this definition an individual could be an outlier for multiple genes. In addition, we only tested for multi-tissue outliers among individuals with expression measurements for the gene in at least five tissues. To reduce cases where non-genetic factors may cause widespread extreme expression, we removed eight individuals that were multi-tissue outliers for 50 or more genes from all downstream analyses, including before single-tissue outlier discovery. Removing these individuals with extreme expression across many genes improved our rare variant enrichments, but the precise threshold mattered less (Extended Data Fig. 2g). We chose the threshold of 50 to strike a balance between removing extreme individuals while not excluding a large proportion of our cohort.

**Replication of expression outliers.** We calculated the proportion of single-tissue outliers discovered in one tissue that had  $|Z\text{-score}| \geq 2$  with the same direction of effect for the same gene in the replication tissue. Since certain groups of tissues were sampled in a specific subset of individuals, we evaluated the extent to which replication was influenced by the size and the overlap of the discovery and replication sets. We repeated the replication analysis with the discovery and replication in exactly 70 overlapping individuals for each pair of tissues with enough samples and compared the replication patterns to those obtained by using all individuals. To estimate the extent to which individual overlap biased replication estimates, for each pair of tissues with sufficient samples, we defined three disjoint groups of individuals: 70 individuals with data for both tissues, 69 distinct individuals with data in the first tissue, and 69 distinct individuals with data in the second tissue. We discovered outliers in the first tissue using the shared set of individuals then tested for replication using the same individuals in the second tissue. Then, for each gene, we added the identified outlier to the distinct set of individuals and tested the replication again in the second tissue. We repeated the process running the discovery in the second tissue and the replication in the first one. We compared the replication rates when using the same or different individuals for the discovery and replication.

We assessed the confidence of our multi-tissue outliers using cross-validation. We separated the tissue expression data randomly into two groups: a discovery set of 34 tissues and a replication set of 10 tissues. For  $t = 10, 15, 20, 25$ , and 30, we randomly sampled  $t$  tissues from the discovery set and performed outlier calling as described above. Owing to incomplete tissue sampling, the number of tissues supporting each outlier is at least five but less than  $t$ . We computed the replication rate as the proportion of outliers in the discovery set with  $|\text{median } Z\text{-score}| \geq 1$  or 2 in the replication set. We set no restriction on the number of tissues required for testing in the replication set. To calculate the expected replication rate, we randomly selected individuals in the discovery set with at least five tissues that expressed the gene and computed the replication rate. We repeated this process 10 times for each discovery set size.

**Quality control of genotypes and rare variant definition.** We restricted our rare variant analyses to individuals of European descent, as they constituted the largest homogenous population within our dataset. We considered only autosomal variants that passed all filters in the VCF (those marked as PASS in the Filter column). Minor allele frequencies (MAFs) within the GTEx data were calculated from the 123 individuals of European ancestry with WGS data (average coverage  $30\times$ ). The MAF was the minimum of the reference and the alternate allele frequency where the allele frequencies of all alternate alleles were summed together. Rare variants were defined as having  $\text{MAF} \leq 0.01$  in GTEx, and for SNVs and indels we also required  $\text{MAF} \leq 0.01$  in the European population of the 1000 Genomes Project Phase 3 data<sup>30</sup>. To ensure that population structure among the individuals of European descent was unlikely to confound our results, we verified that the allele frequency distribution of rare variants included in our analysis (within 10 kb of a protein-coding or lincRNA gene, see below) was similar for the five European populations in the 1000 Genomes Project (Extended Data Fig. 4d).

**Enrichment of rare and common variants near outlier genes.** We assessed the enrichment of rare SNVs, indels and structural variants near outlier genes. Proximity was defined as within 10 kb of the transcription start site for most analyses. For Fig. 3 and Extended Data Figs 5, 7, 8, we included all variants within 10 kb of the gene, including the gene body, to also capture coding variants. In Fig. 3 and Extended Data Figs 5, 8, we extended the window to 200 kb for enhancers and structural variants. For each gene with an outlier, we chose the remaining set of individuals tested for outliers at the same gene as non-outlier controls. We only considered genes that had both an outlier and at least one control. We stratified variants of each class into four minor allele frequency bins (0–1%, 1–5%, 5–10%, 10–25%) to compare the relative enrichments of rare and common variants. We also assessed the enrichment of SNVs at different  $Z$ -score cutoffs. Enrichment was defined as the ratio of the proportion of outliers with a variant whose frequency lies within the range to the corresponding proportion for non-outliers. This



enrichment analysis is equivalent to the relative risk of having a nearby rare variant given outlier status. We used the asymptotic distribution of the log relative risk to obtain 95% Wald confidence intervals. Within our set of European individuals, we observed some individuals with minor admixture that had relatively more rare variants than the rest (Extended Data Fig. 1b). We confirmed that inclusion of these admixed individuals did not substantially affect our results (Extended Data Fig. 1c). We also calculated rare variant enrichments when restricting to variants outside protein-coding and lincRNA exons in the Gencode v.19 annotation (extending internal exons by 5 bp to capture canonical splice regions).

To measure the informativeness of variant annotations, we used logistic regression to model outlier status as a function of the feature of interest; this yielded log odds ratios with 95% Wald confidence intervals. Note that for the feature enrichment analysis in Fig. 3b and Extended Data Fig. 7, we required that outliers and their gene-matched non-outlier controls have at least one rare variant near the gene. We standardized all features, including binary features, to facilitate comparison between features of different scale. We also calculated the proportion of overexpression outliers, underexpression outliers and non-outliers with a rare variant near the gene (within 10 kb for SNVs and indels and 200 kb for structural variants). To each outlier instance, we assigned at most one of the 12 rare variant classes that we considered (Supplementary Table 3a). If an outlier had rare variants from multiple classes near the relevant genes, we selected the class that was most significantly enriched among outliers.

**Annotation of variants.** We obtained structural variant annotations from ref. 14 and computed features for rare SNVs and indels using three primary data sources: Roadmap Epigenomics<sup>31</sup>, CADD v.1.2 (ref. 19) and VEP v.80 (ref. 32). Promoter and enhancer annotation tracks were obtained from the Roadmap Epigenomics Project ([http://www.broadinstitute.org/~meuleman/reg2map/HoneyBadger2\\_release/](http://www.broadinstitute.org/~meuleman/reg2map/HoneyBadger2_release/)). We mapped 28 unique tissues in the GTEx project to 19 tissue groups in the Roadmap Project. Using these annotations, for each individual, we assessed whether each SNV or indel overlapped a promoter or enhancer region in at least one of the 19 Roadmap tissue groups. Features, including conservation<sup>15–18</sup>, transcription factor binding and deleteriousness, were extracted from the full annotation tracks of the CADD v.1.2 release (downloaded 15 May 2015; <http://cadd.gs.washington.edu/download>). Finally, we obtained protein-coding and transcription-related annotations from VEP and LOFTEE. This information was provided in the GTEx v6p VCF file (described in ref. 12). Stop-gain and frameshift variants annotated as high-confidence loss-of-function variants by LOFTEE were assumed to trigger nonsense-mediated decay. We generated gene-level features described in Supplementary Table 3.

**Allele-specific expression (ASE).** We only considered sites with at least 30 total reads and at least five reads supporting each of the reference and alternate alleles. To minimize the effect of mapping bias, we filtered out sites that showed mapping bias in simulations<sup>33</sup>, that were in low mappability regions (<ftp://hgdownload.cse.ucsc.edu/gbdb/hg19/bbi/wgEncodeCrgMapabilityAlign50mer.bw>) or that were rare variants or within 1 kb of a rare variant in the given individual (the variants were extracted from the GTEx exome-sequencing data described in ref. 12). The first two filters were provided in the GTEx ASE data release. The third filter was applied to eliminate potential mapping artefacts that mimic genetic effects from rare variants. We measured ASE at each testable site as the absolute deviation of the reference-allele ratio from 0.5. For each gene, all testable sites in all tissues were included. We compared ASE in single-tissue and multi-tissue outliers at different Z-score thresholds to non-outliers using two-sided Wilcoxon rank-sum tests. To obtain a matched background, we only included a gene in the comparison when ASE data existed for both the outlier individual and at least one non-outlier. In the case of single-tissue outliers, we also required the tissue to match between the outlier and the non-outlier. All individuals that were neither multi-tissue outliers for the given gene nor single-tissue outliers for the gene in the corresponding tissue were included as non-outliers.

In cases where outliers had rare coding variants in the gene, if the rare variants were causing the extreme expression in *cis*, we expected to see ASE at the rare variant matching the direction of the effect. For underexpression outliers, we expected the (rare) minor allele to be underexpressed compared to the major allele. For overexpression outliers, we expected the minor allele to be overexpressed. To test this, we used the same filters as above, but looked exclusively at rare variants (instead of excluding them). We measured ASE as the minor-allele ratio: the number of reads supporting the minor allele over the total number of reads.

We also used ASE to evaluate the performance of both the genomic annotation model and RIVER (see below) by testing the association between allelic imbalance and model predictions using Fisher's exact test. Here, we defined allelic imbalance as the top 10% of the median absolute deviation, across tissues, of the reference-allele ratio from 0.5.

**Allele frequency measurements in UK10K.** UK10K<sup>3</sup> VCF files of whole-genome cohorts were downloaded from <https://www.ebi.ac.uk>. We merged the

Avon Longitudinal Study of Parents and Children (ALSPAC) EGAS00001000090 and the Department of Twin Research and Genetic Epidemiology (TWINsUK) EGAS00001000108 datasets for a total of 3,781 individuals. We counted the occurrence of all rare GTEx SNVs in Roadmap Epigenomics-annotated promoter regions among the UK10K samples. GTEx variants absent from the UK10K cohorts were assigned a count of 0.

**Definition of multi-tissue eGenes.** We defined multi-tissue eGenes using two approaches. For the tissue-by-tissue approach, we obtained lists of significant eGenes ( $q$  value  $\leq 0.05$ ) for each of the 44 tissues from the GTEx v6p release. The second approach used *cis*-eQTLs with shared effects across tissues estimated by the RE2 model of the Meta-Tissue software<sup>34</sup>, as described in ref. 12. We chose, for each gene, the variant with the lowest nominal  $P$  value from the RE2 model. We then determined the number of tissues in which this variant-gene pair showed a *cis*-eQTL effect ( $m$  value  $\geq 0.9$  (ref. 34)). For each of the 18,380 genes tested for multi-tissue outliers, we calculated the number of tissues in which the gene appeared as a significant eGene (tissue-by-tissue approach) or had a shared eQTL effect (Meta-Tissue approach). To show that the enrichment of outlier genes as multi-tissue eGenes was not confounded by gene expression level, using the Meta-Tissue results, we stratified genes tested for multi-tissue outliers into RPKM deciles and repeated the comparison between genes with and without a multi-tissue outlier. When comparing the enrichment for eGenes among constrained and disease gene lists, we classified the top  $n$  Meta-Tissue eGenes (ranked by nominal  $P$  value from the RE2 model) as multi-tissue eGenes and considered the remaining genes as background. We selected  $n$  to match the number of multi-tissue outliers in the comparison.

**Evolutionary constraint of genes with multi-tissue outliers.** We obtained gene-level estimates of evolutionary constraint from the Exome Aggregation Consortium<sup>35</sup> (<http://exac.broadinstitute.org/>, ExAC release v.0.3). We intersected the 17,351 autosomal lincRNA and protein-coding genes with constraint data from ExAC with the 18,380 genes tested for multi-tissue outliers from GTEx, yielding 14,379 genes for further analysis (3,897 and 10,482 genes with and without a multi-tissue outlier, respectively). We examined three functional constraint scores from the ExAC database: synonymous Z-score, missense Z-score and probability of loss-of-function intolerance (pLI). Synonymous- and missense-intolerant genes were defined as those with corresponding Z-scores above the 90th percentile. We defined loss-of-function intolerant genes as those with a pLI score above 0.9, following the guidelines provided by ExAC. We calculated odds ratios and 95% confidence intervals for the enrichment of genes with multi-tissue outliers in these lists using a Fisher's exact test. We repeated this analysis for three other gene sets: 19,182 multi-tissue eGenes from GTEx v6p defined using Meta-Tissue, 9,480 reported GWAS genes from the NHGRI-EBI catalogue<sup>36</sup> (<http://www.ebi.ac.uk/gwas>, accessed 30 November 2015) and 3,576 OMIM genes (<http://omim.org/>, accessed 26 May 2016).

We tested for a difference in the mean constraint for genes with multi-tissue outliers and genes with multi-tissue eQTLs using ANOVA. For each constraint score in ExAC, we treated the score for each gene as the response and the status of the gene as having a multi-tissue outlier and/or a multi-tissue eQTL as a categorical predictor with four classes. After fitting the model, we performed a Tukey's range test to determine whether there was a significant difference in the mean constraint between genes with a multi-tissue outlier but no multi-tissue eQTL and genes with a multi-tissue eQTL but no multi-tissue outlier.

**Overlap of genes with multi-tissue outliers and disease genes.** We examined the enrichment of genes with multi-tissue outliers in eight disease gene lists: the GWAS catalogue and OMIM (described above), as well as ClinVar (6,279 genes; <http://www.ncbi.nlm.nih.gov/clinvar/>), OrphaNet (3,451 genes; <http://www.orpha.net/>), ACMG<sup>20</sup> (58 genes; <http://www.ncbi.nlm.nih.gov/clinvar/docs/acmg/>), Developmental Disorders Genotype-to-Phenotype<sup>37</sup> (DDG2P; 1,693 genes; <http://www.ebi.ac.uk/gene2phenotype/>), and two curated gene lists of 86 cardiovascular disease genes and 55 cancer genes (described below). We computed odds ratios and 95% confidence intervals using a Fisher's exact test to compare each disease gene list to the genes with multi-tissue outliers and repeated the comparison for genes with multi-tissue eQTLs.

Heritable cancer predisposition and heritable cardiovascular disease gene lists were curated by local experts in clinical and laboratory-based genetics in the two respective areas (Stanford Medicine Clinical Genomics Service, Stanford Cancer Center's Cancer Genetics Clinic and Stanford Center for Inherited Cardiovascular Disease). Genes were included if both the clinical and laboratory-based teams agreed there was sufficient published evidence to support using variants in these genes in clinical decision making.

For each of the eight disease gene lists above and for genes with multi-tissue outliers or multi-tissue eQTLs, we computed the number of variants (SNVs and indels within 10 kb and structural variants within 200 kb of the gene, including the gene body) at each gene in the 123 individuals of European ancestry with WGS

data. For each gene list and for each MAF bin (0–1%, 1–5%, 5–10%, 10–25%), we compared the mean number of variants near genes in the list to the mean number near all other annotated autosomal protein-coding and lincRNA genes using a two-sided *t*-test.

**The RIVER integrative model for predicting regulatory effects of rare variants.** RIVER (RNA-informed variant effect on regulation) is a hierarchical Bayesian model that predicts the regulatory effects of rare variants by integrating gene expression with genomic annotations. The RIVER model consists of three layers: a set of nodes  $\mathbf{G} = G_1, \dots, G_P$  in the topmost layer representing  $P$  observed genomic annotations over all rare variants near a particular gene; a latent binary variable  $F$  in the middle layer representing the unobserved functional regulatory status of the rare variants; and one binary node  $E$  in the final layer representing expression outlier status of the nearby gene. We model each conditional probability distribution as follows:

$$F|\mathbf{G} \sim \text{Bernoulli}(\psi), \quad \psi = \text{logit}^{-1}(\beta'\mathbf{G})$$

$$E|F \sim \text{categorical}(\theta)$$

$$\beta_i \sim \mathcal{N}\left(0, \frac{1}{\lambda}\right)$$

$$\theta \sim \text{Beta}(C, C)$$

with parameters  $\beta$  and  $\theta$  and hyper-parameters  $\lambda$  and  $C$ .

Because  $F$  is unobserved, the RIVER log-likelihood objective over instances  $n = 1, \dots, N$   $\sum_{n=1}^N \log \sum_{F_n=0}^1 P(E_n, \mathbf{G}_n, F_n | \beta, \theta)$  is non-convex. We therefore optimize model parameters using Expectation–Maximization<sup>38</sup> (EM) as follows:

In the E-step, we compute the posterior probabilities ( $\omega_n^{(i)}$ ) of the latent variables  $F_n$  given current parameters and observed data. For example, at the  $i$ th iteration, the posterior probability of  $F_n = 1$  for the  $n$ th instance is

$$\begin{aligned} \omega_n^{(i)} &= P(F_n = 1 | \mathbf{G}_n, \beta^{(i)}, E_n, \theta^{(i)}) \\ &= \frac{P(F_n = 1 | \mathbf{G}_n, \beta^{(i)}) P(E_n | F_n = 1, \theta^{(i)})}{\sum_{F_n=0}^1 P(F_n | \mathbf{G}_n, \beta^{(i)}) \cdot P(E_n | F_n, \theta^{(i)})} \\ \omega_{0n}^{(i)} &= 1 - \omega_n^{(i)} \end{aligned}$$

In the M-step, at the  $i$ th iteration, given the current estimates  $\omega^{(i)}$ , the parameters ( $\beta^{(i+1)}$ ) are estimated as

$$\text{argmax}_{\beta^{(i+1)}} \sum_{n=1}^N \sum_{F_n=0}^1 \log(P(F_n | \mathbf{G}_n, \beta^{(i+1)})) \cdot \omega_{F_n, n}^{(i)} - \frac{\lambda}{2} \|\beta^{(i+1)}\|_2$$

where  $\lambda$  is an L2 penalty hyper-parameter derived from the Gaussian prior on  $\beta$ .

The parameter  $\theta$  gets updated as:

$$\theta_{st}^{(i+1)} = \sum_{n=1}^N I(E_n = t) \cdot \omega_{s, n}^{(i)} + C$$

where  $I$  is an indicator operator,  $t$  is the binary value of expression  $E_n$ ,  $s$  is the possible binary values of  $F_n$ , and  $C$  is a pseudo count derived from the Beta prior on  $\theta$ . The E and M steps are applied iteratively until convergence.

**RIVER application to the GTEx cohort.** As input, RIVER requires a set of genomic features  $\mathbf{G}$  and a set of corresponding expression outlier observations  $\mathbf{E}$ , each over instances of individual and gene pairs. Using the variant annotations described above, we generated site-level genomic features for the 116 European individuals with GTEx WGS data that had fewer than 50 multi-tissue outliers. We then collapsed these features for all rare SNVs within 10 kb of each transcription start site to generate the gene-level features that are described in Supplementary Table 3b. This produced a matrix of genomic features  $\mathbf{G}$  of size (116 individuals  $\times$  1,736 genes)  $\times$  (112 genomic features), where we standardized features before use. For the values of  $\mathbf{E}$ , we defined any individual with  $|\text{median } Z\text{-score}| \geq 1.5$  as an outlier if expression was observed in at least five tissues; the remaining individuals were labelled as non-outliers for the gene. We used this more lenient threshold in order to obtain a sufficiently large set of outliers for robust training and testing. In total, we extracted 48,575 instances where an individual had at least one rare variant within 10 kb of the transcription start site of a gene.

To train and evaluate RIVER on the GTEx cohort, we used the 3,766 instances of individual and gene pairs where two individuals had the same rare SNVs near a particular gene. We held out those instances and trained RIVER parameters with the remaining instances. RIVER requires two hyper-parameters  $\lambda$  and  $C$ . To select  $\lambda$ , we first applied an L2-regularized multivariate logistic regression with features  $\mathbf{G}$  and response variable  $\mathbf{E}$ , selecting  $\lambda$  with the minimum squared error via tenfold cross-validation (we selected  $\lambda = 0.01$ ). We selected  $C = 50$ , informed simply by the total number of training instances available, as validation data were not available for extensive cross-validation. Initial parameters for EM were set to  $\theta = (P(E = 0 | F = 0), P(E = 1 | F = 0), P(E = 0 | F = 1), P(E = 1 | F = 1)) = (0.99, 0.01, 0.3, 0.7)$  and  $\beta$  from the multivariate logistic regression above, although different initializations did not significantly change the final parameters (Extended Data Fig. 9b).

The 3,766 held-out pairs of instances were used to create a labelled evaluation set. For one of the two individuals from each pair, we estimated the posterior probability of a functional rare variant  $P(F | \mathbf{G}, E, \beta, \theta)$ . The outlier status of the second individual, whose data were not observed either during training or prediction, was then treated as a 'label' of the true status of functional effect  $F$ . Using this labelled set, we compared the RIVER score to the posterior  $P(F | \mathbf{G}, \beta)$  estimated from the plain L2-regularized multivariate logistic regression model with genomic annotations alone. We produced receiver operating characteristic curves and computed areas under the curve (AUCs) for both models, testing for significant differences using DeLong's method<sup>29</sup>. This analysis relied on outlier status reflecting the consequences of rare variants. Indeed, pairs of individuals who shared rare variants tended to have highly similar outlier status even after regressing out effects of common variants (Kendall's  $\tau$  rank correlation,  $P < 2.2 \times 10^{-16}$ ). We repeated this evaluation, varying the median  $Z$ -score threshold used to define outliers, and we also compared RIVER to individual features that were strongly enriched among outliers as well as PolyPhen<sup>39</sup> and SIFT<sup>40</sup>.

**Supervised model integrating expression and genomic annotation.** To assess the information gained by incorporating gene expression data in the prediction of functional rare variants, we applied a simplified supervised approach to a limited dataset. We used the instances where two individuals had the same rare SNVs to create a labelled training set where the outlier status of the second individual was used as the response variable. We then trained a logistic regression model with only two features: (1) the outlier status of the first individual and (2) a single genomic feature value, such as CADD or deleterious annotation of genetic variants using neural networks (DANN). We estimated parameters from the entire set of rare-variant-matched pairs using logistic regression to determine the log odds ratio and corresponding  $P$  value of expression status as a predictor. While this approach was not amenable to training a full predictive model over all genomic annotations jointly given the limited number of instances, it provided a consistent estimate of the log odds ratio of outlier status. We tested five genomic predictors: CADD<sup>19</sup>, DANN<sup>41</sup>, transcription-factor-binding site annotations, PhyloP scores<sup>15</sup> and one aggregated feature: the posterior probability from a multivariate logistic regression model learned with all genomic annotations.

**RIVER assessment of pathogenic ClinVar variants.** We downloaded variants from the ClinVar database<sup>21</sup> (accessed 04 May 2015) and searched for these disease variants within the set of rare variants segregating in the GTEx cohort. Any disease variant reported as pathogenic, likely pathogenic or a risk factor for disease was considered pathogenic. We further categorized the pathogenic variants as likely regulatory if they were annotated as splice-site variants, synonymous or nonsense, whereas missense variants were considered unlikely to have a regulatory effect. To explore RIVER scores for those pathogenic variants, all instances were used for training RIVER. We then computed a posterior probability  $P(F | \mathbf{G}, E, \beta, \theta)$  for each instance coinciding with a pathogenic ClinVar variant.

**Stability of estimated parameters with different parameter initializations.** We tried several different initialization parameters for  $\beta$  and  $\theta$  to explore how this affected the estimated parameters. We initialized a noisy  $\beta$  by adding  $K\%$  Gaussian noise compared to the mean of  $\beta$  with fixed  $\theta$  (for  $K = 10, 20, 50, 100, 200, 400, 800$ ). For  $\theta$ , we fixed  $P(E = 1 | F = 0)$  and  $P(E = 0 | F = 0)$  as 0.01 and 0.99, respectively, and initialized  $(P(E = 1 | F = 1), P(E = 0 | F = 1))$  as (0.1, 0.9), (0.4, 0.6) and (0.45, 0.55) instead of (0.3, 0.7) with  $\beta$  fixed. For each parameter initialization, we computed Spearman rank correlations between parameters from RIVER using the original initialization and the alternative initializations. We also investigated how many instances within top 10% of posterior probabilities from RIVER under the original settings were replicated in the top 10% of posterior probabilities under the alternative initializations (replication accuracy in Extended Data Fig. 9b).

**Validation of large-effect rare variants using CRISPR–Cas9 genome editing.** To select rare, coding SNVs for validation by CRISPR–Cas9 editing, we first restricted to the (gene, individual, variant) tuples identified in multi-tissue outliers without a rare structural variant or a rare indel within 200 kb or 10 kb of the gene, respectively.

We considered the 116 rare SNVs with a coding consequence for the corresponding gene as annotated by VEP<sup>32</sup>; coding annotations included stop gained, stop lost, splice acceptor variant, splice donor variant, start lost, missense variant, splice region variant, stop retained variant, synonymous variant, coding sequence variant and 5'/3' UTR variant. Using RNA-seq data from ENCODE, we further restricted our variant list to the 59 SNVs occurring in genes with an average FPKM (fragments per kilobase per million reads) of at least 10 in K562 cells (ENCODE experiment accession numbers ENCSR000AEL and ENCSR000AEN)<sup>42</sup>. Finally, we filtered for rare, coding SNVs in (gene, individual) pairs with  $|\text{median } Z\text{-score}| > 4$  and a RIVER score above the 99.5th percentile. These filters yielded a final set of 13 rare SNVs from which we chose the six exonic SNVs for testing.

As controls, we selected SNVs present within the same cDNA amplicon region as the corresponding outlier SNV (see details on targeted sequencing below). We first searched for coding SNVs present within these regions in the GTEx cohort that did not occur in the outlier individual. If no SNV could be found satisfying these criteria, we expanded our search for SNVs using the ExAC database (ExAC release v.0.3)<sup>35</sup>. If multiple possible control variants existed for an outlier SNV, we ranked the controls by CADD score<sup>19</sup> and prioritized synonymous variants.

Sequences of single-guide RNAs (sgRNAs) used in the study are listed in Extended Data Fig. 11b. For each variant, a sgRNA and two donor oligonucleotides (with the reference and alternative alleles) were designed such that the PAM was located as close to the variant as possible. The donors were 99 bp long centred on the variant being installed. The variants were installed into K562 cells as previously described<sup>22,23</sup>. The K562 cells were those generated previously<sup>23</sup> and were regularly tested for mycoplasma infection. sgRNAs were expressed in the pGH020 (Addgene plasmid 85405) expression vector. For each donor oligonucleotide, K562 cells constitutively expressing a Cas9–BFP fusion protein were electroporated with 3  $\mu\text{g}$  of sgRNA plasmid DNA and 1  $\mu\text{l}$  of 100  $\mu\text{M}$  donor oligonucleotide using the T-016 program on a Lonza Nucleofector 2b. After electroporation, cells were allowed to recover for five days. Cells electroporated with the reference and alternative allele donor oligonucleotides were mixed in a 1:1 ratio and grown together for three more days to control for differences in culturing conditions. We included cells electroporated with the reference allele to ensure that any changes in expression we observed were not due to the editing process itself. Because the editing efficiency is not 100% and varies between loci, we expected fewer than half the cells to carry the alternative allele and for this proportion to vary by locus. One to two million cells were collected for RNA and genomic DNA extraction.

Genomic DNA (gDNA) was extracted using the QiaAmp DNA mini kit (Qiagen). Total RNA was extracted using QiaShredder and RNeasy Mini kit (Qiagen). Subsequently, 6  $\mu\text{g}$  of RNA was converted into cDNA using AMV reverse transcriptase (Promega). cDNA was purified and concentrated with the PCR Purification Kit (Qiagen). PCR primers were designed to generate 300–400-bp amplicons including the variant in either the gDNA or cDNA locus. For both gDNA and cDNA samples, 400 ng of DNA was amplified in triplicate (technical replicates) using Phusion High-Fidelity polymerase (Fisher) and the amplicon was purified on a 1% TAE agarose gel. The amplicons were then prepared for sequencing using the Nextera XT kit (Illumina) and sequenced together on a NextSeq 500.

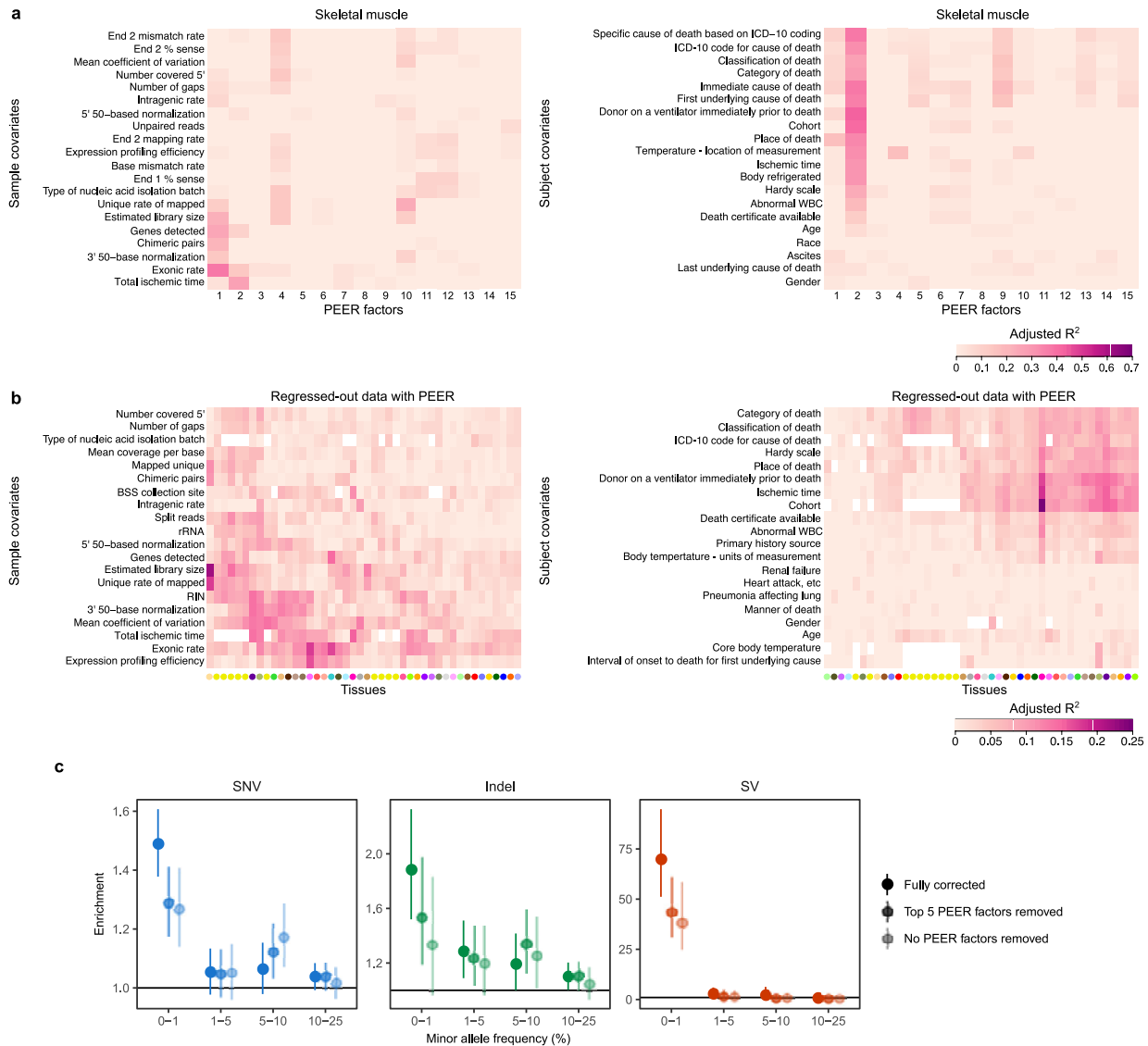
Reads were trimmed with cutadapt<sup>43</sup> (v.1.13) and aligned using bwa<sup>44</sup> (v.0.7.12-r1039) allowing no mismatches (bwa aln -n 0), which excluded any reads with indels created during editing. We used custom reference sequences, one each for the reference and alternate alleles of the targeted cDNA and gDNA amplicon regions. Allele counts at the target locus were computed for each sample using samtools pileup as implemented in the R package Rsamtools<sup>45</sup> (v.1.22.0). Only reads with a minimum mapping quality of 20 were considered. Two of the tested loci amplified poorly in preparation for sequencing, and they had extremely low

mapping rates and total read counts over the target locus (median read count across replicates <400 compared to 281,000 and 397,000 for gDNA and cDNA, respectively, for the remaining loci). As such, we removed these two loci from further analysis. Finally, to assess the effect of each variant on expression, we tested for a significant difference between the cDNA and gDNA alternate allele proportions with a two-sided *t*-test. We corrected for multiple testing using the Bonferroni procedure.

**Code availability.** RIVER is available at <https://bioconductor.org/packages/release/bioc/html/RIVER.html>. Additionally, the code for running analyses and producing the figures throughout this manuscript is available separately (<https://github.com/joed3/GTExV6PRareVariation>).

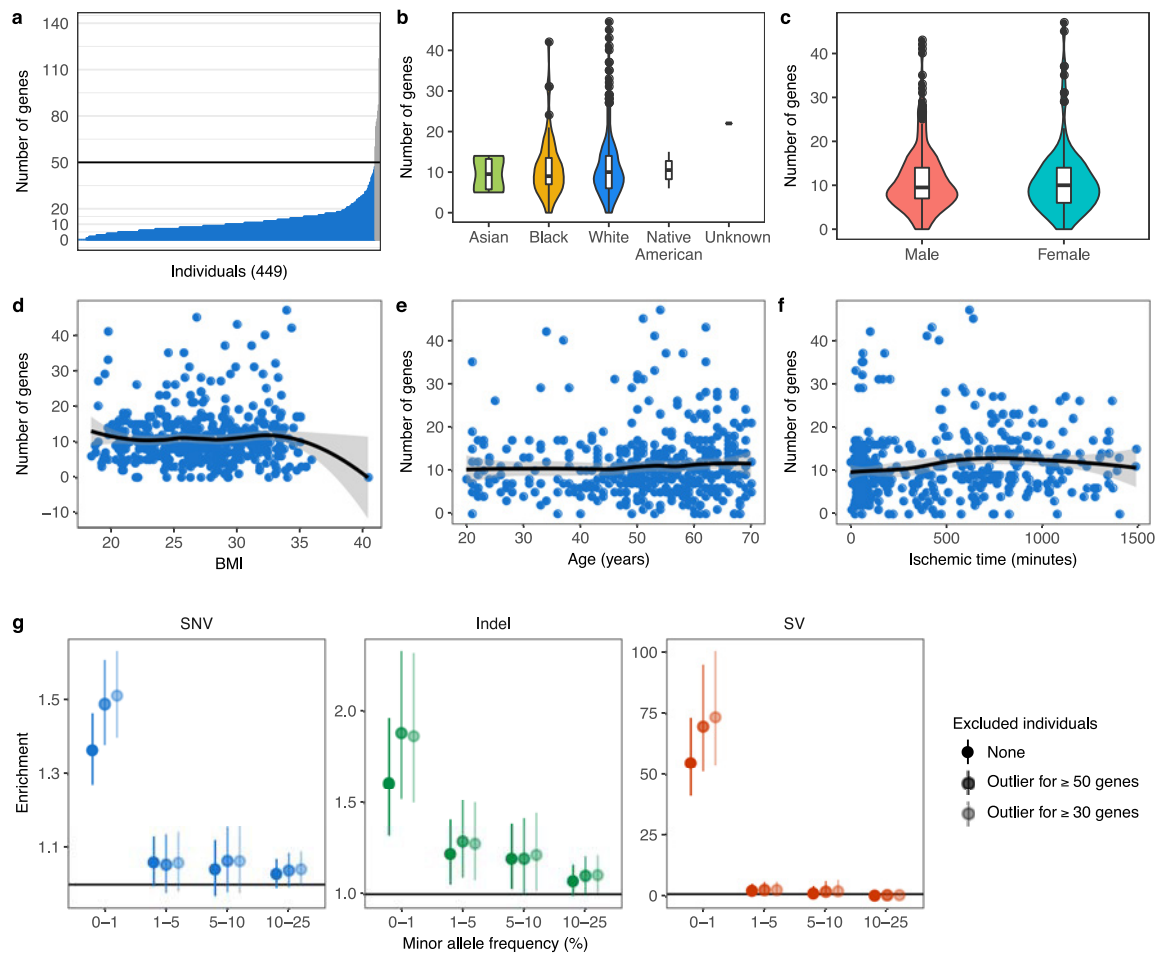
**Data availability.** The GTEx v6p release genotype and allele-specific expression data are available from dbGaP (study accession phs000424.v6.p1; [http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000424.v6.p1](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000424.v6.p1)). Expression data from the v6p release and eQTL results are available from the GTEx portal (<http://gtexportal.org>).

30. 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
31. The Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
32. McLaren, W. *et al.* The Ensembl variant effect predictor. *Genome Biol.* **17**, 122 (2016).
33. Panousis, N. I., Gutierrez-Arcelus, M., Dermizakis, E. T. & Lappalainen, T. Allelic mapping bias in RNA-sequencing is not a major confounder in eQTL studies. *Genome Biol.* **15**, 467 (2014).
34. Sul, J. H., Han, B., Ye, C., Choi, T. & Eskin, E. Effectively identifying eQTLs from multiple tissues by combining mixed model and meta-analytic approaches. *PLoS Genet.* **9**, e1003491 (2013).
35. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
36. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* **42**, D1001–D1006 (2014).
37. Deciphering Developmental Disorders Study. Large-scale discovery of novel genetic causes of developmental disorders. *Nature* **519**, 223–228 (2015).
38. Dempster, A. P., Laird, N. M. & Rubin, D. B. Maximum likelihood from incomplete data via the EM algorithm. *J. R. Stat. Soc.* **39**, 1–38 (1977).
39. Adzhubei, I. A. *et al.* A method and server for predicting damaging missense mutations. *Nat. Methods* **7**, 248–249 (2010).
40. Ng, P. C. & Henikoff, S. Predicting deleterious amino acid substitutions. *Genome Res.* **11**, 863–874 (2001).
41. Quang, D., Chen, Y. & Xie, X. DANN: a deep learning approach for annotating the pathogenicity of genetic variants. *Bioinformatics* **31**, 761–763 (2015).
42. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
43. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal* **17**, 10–12 (2011).
44. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
45. Morgan, M., Pagès, H., Obenchain, V. & Hayden, N. Rsamtools: Binary alignment (BAM), FASTA, variant call (BCF), and tabix file import. R package v.1.28.0 <http://bioconductor.org/packages/release/bioc/html/Rsamtools.html> (2017).
46. Dror, Y. & Freedman, M. H. Shwachman–Diamond syndrome. *Br. J. Haematol.* **118**, 701–713 (2002).
47. Austin, K. M. *et al.* Mitotic spindle destabilization and genomic instability in Shwachman–Diamond syndrome. *J. Clin. Invest.* **118**, 1511–1518 (2008).
48. Schmidt, A. *et al.* Severely altered guanidino compound levels, disturbed body weight homeostasis and impaired fertility in a mouse model of guanidinoacetate *N*-methyltransferase (GAMT) deficiency. *Hum. Mol. Genet.* **13**, 905–921 (2004).



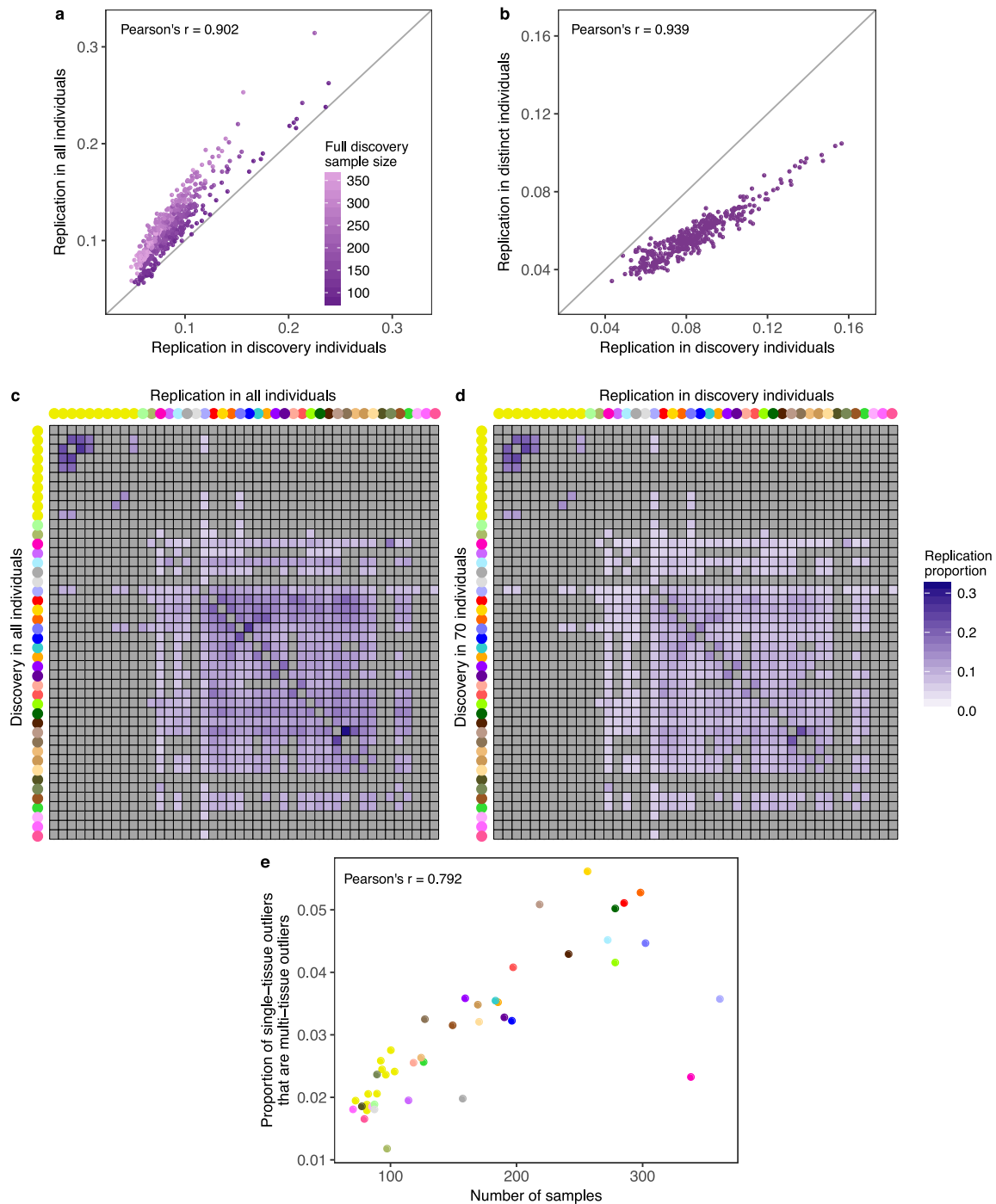
**Extended Data Figure 1 | PEER correction.** **a**, Adjusted  $R^2$  between top 15 PEER factors and top 20 sample (left) and subject (right) covariates in an example tissue, skeletal muscle. Covariates were ranked by the average adjusted  $R^2$  across all PEER factors and hierarchically clustered. The corresponding data for all tissues are provided in Supplementary Tables 1, 2. **b**, Adjusted  $R^2$  between the total expression component removed by PEER in each tissue and the top 20 sample (left) and subject (right) covariates. The covariates were ranked by the average adjusted  $R^2$  across

all tissues, and both axes were hierarchically clustered. White denotes missing values, and tissues are coloured as in Fig. 1. PEER factors captured slightly different covariates across tissues, with a noticeable difference between the brain and other tissues. **c**, Rare variant enrichments as in Fig. 2a for different levels of PEER correction. The fully corrected data show substantially stronger rare variant enrichments than the two partially corrected datasets.



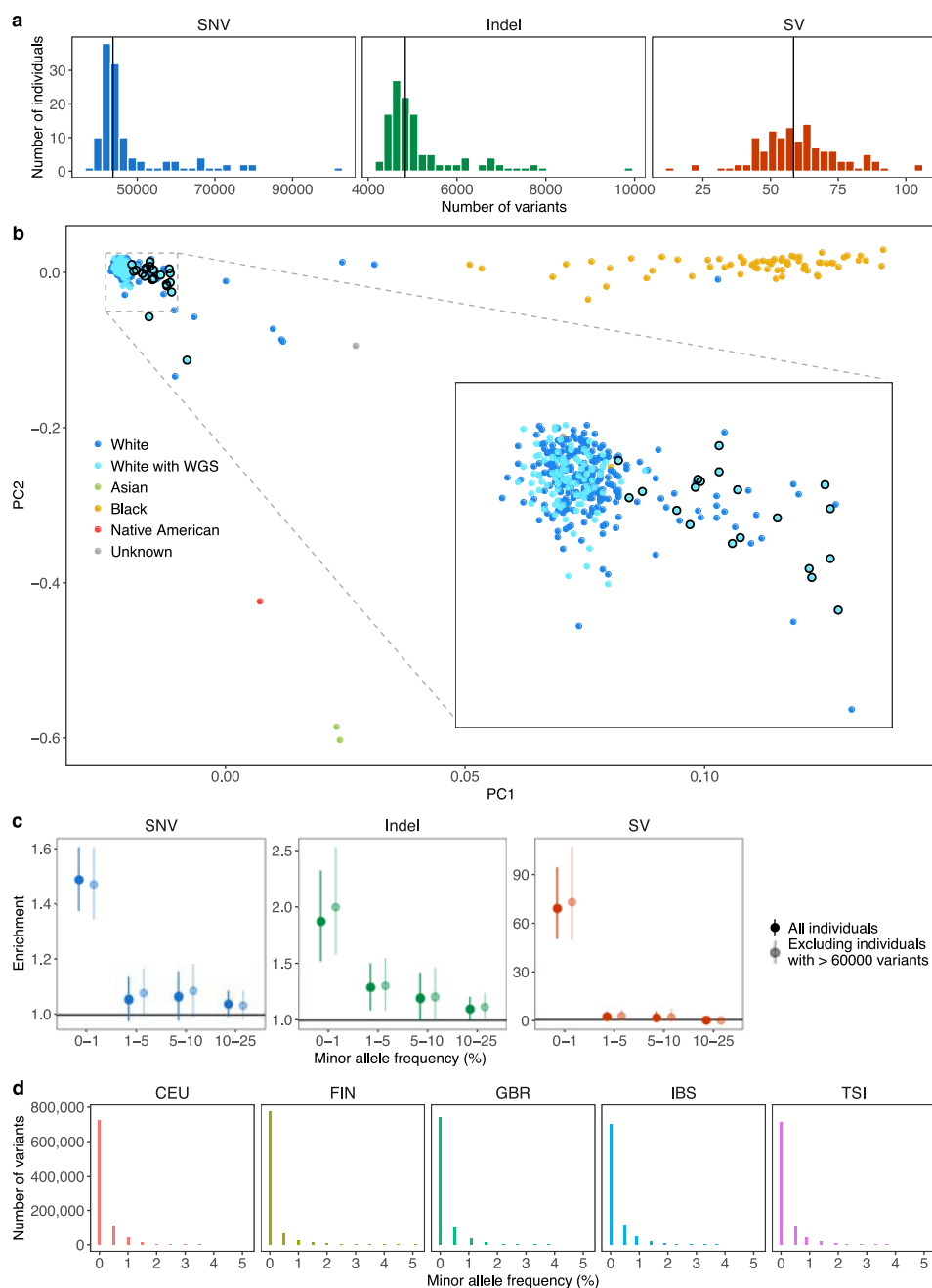
**Extended Data Figure 2 | Distribution of the number of genes with a multi-tissue outlier.** **a**, Distribution of the number of genes for which each individual was a multi-tissue outlier. Each individual was an outlier for a median of 10 genes. Individuals with 50 or more outliers are coloured in grey and were excluded from downstream analyses. **b–f**, Distribution of the number of genes for which individuals, stratified by common covariates, were multi-tissue outliers. For race and sex, we compared the distributions using an unsigned Wilcoxon rank-sum test, whereas we used Spearman's  $\rho$  to test for association with the remaining covariates.

Only age (Spearman's  $\rho = 0.10$ ,  $P = 0.033$ ) and ischaemic time (Spearman's  $\rho = 0.18$ ,  $P = 0.00022$ ) were nominally associated with the number of outlier genes per individual. The association with age fails to achieve significance after correcting for multiple testing using the Bonferroni method. Note that in **b** we only tested for a significant difference in the distribution of the number of outlier genes between white and black individuals, because there were too few individuals in the other groups. **g**, Enrichments as shown in Fig. 2a either including all individuals, or excluding individuals that are outliers for 50 (matches Fig. 2a) or 30 genes.



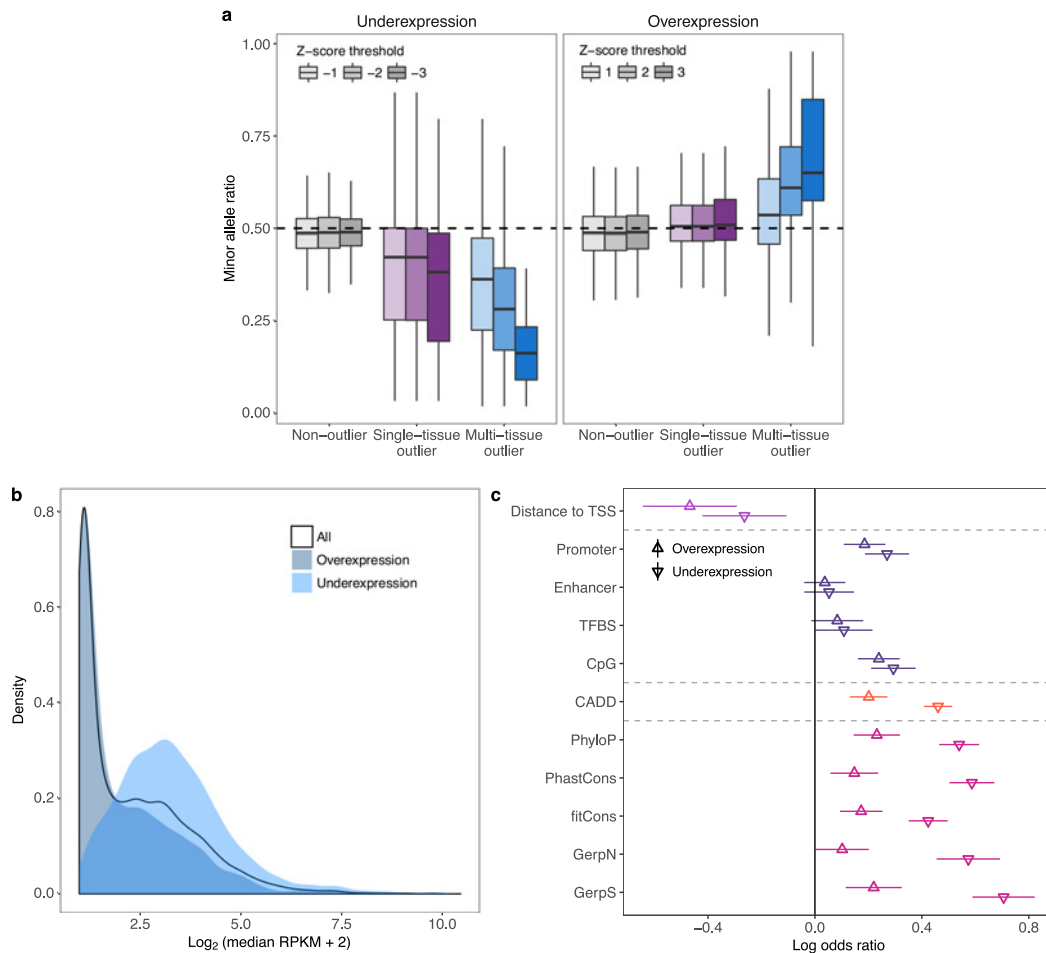
**Extended Data Figure 3 | Single-tissue outlier replication.** **a**, Correlation between the replication proportions (see Methods) obtained from all samples and from a subset of 70 overlapping individuals per tissue pair (Pearson's correlation,  $P < 2.2 \times 10^{-16}$ ). When restricting to 70 individuals, the replication rates decreased more for discovery tissues with larger sample sizes in the full dataset, indicating that replication rates were underestimated for tissues with small sample sizes. **b**, Correlation between replication in the 70 individuals used for discovery and replication assessed in a set of 70 individuals that included the outlier individual and 69 individuals excluded from the discovery set (Pearson's correlation,  $P < 2.2 \times 10^{-16}$ ). Replication was higher when computed in the discovery individuals rather than in a distinct set of individuals. **c**, Single-tissue

outlier replication using all individuals, as in Fig. 1b, but data are only shown for pairs with at least 70 overlapping individuals. Tissue pairs with insufficient overlap are in grey. **d**, For each pair of tissues with sufficient samples, outlier discovery and replication using 70 individuals sampled in both tissues. The replication values decreased compared with replication performed in all individuals (**c**), particularly for tissues with large sample sizes in the complete dataset. However, the pattern of replication, with more similar tissues having higher replication rates, is maintained. **e**, For each tissue, the proportion of (individual, gene) outlier pairs where the individual was also a multi-tissue outlier for the gene. This proportion was positively correlated with the tissue sample size ( $P = 1.4 \times 10^{-10}$ ). Points are coloured by tissue as in Fig. 1.



**Extended Data Figure 4 | Number of rare variants per individual and population structure.** **a**, The distribution of the number of rare variants of each type for individuals of European descent (reported as white). Certain individuals had many more rare variants than the population median (vertical black line). **b**, Principal component analysis of all individuals. Individuals are plotted according to their first two genotype principal components (PCs) and coloured by their reported ancestry. White individuals with WGS data, included in **a**, are coloured in a lighter shade of blue and those with 60,000 or more rare variants are circled in black.

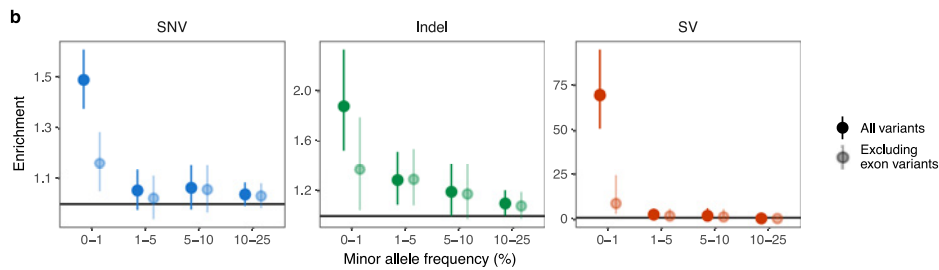
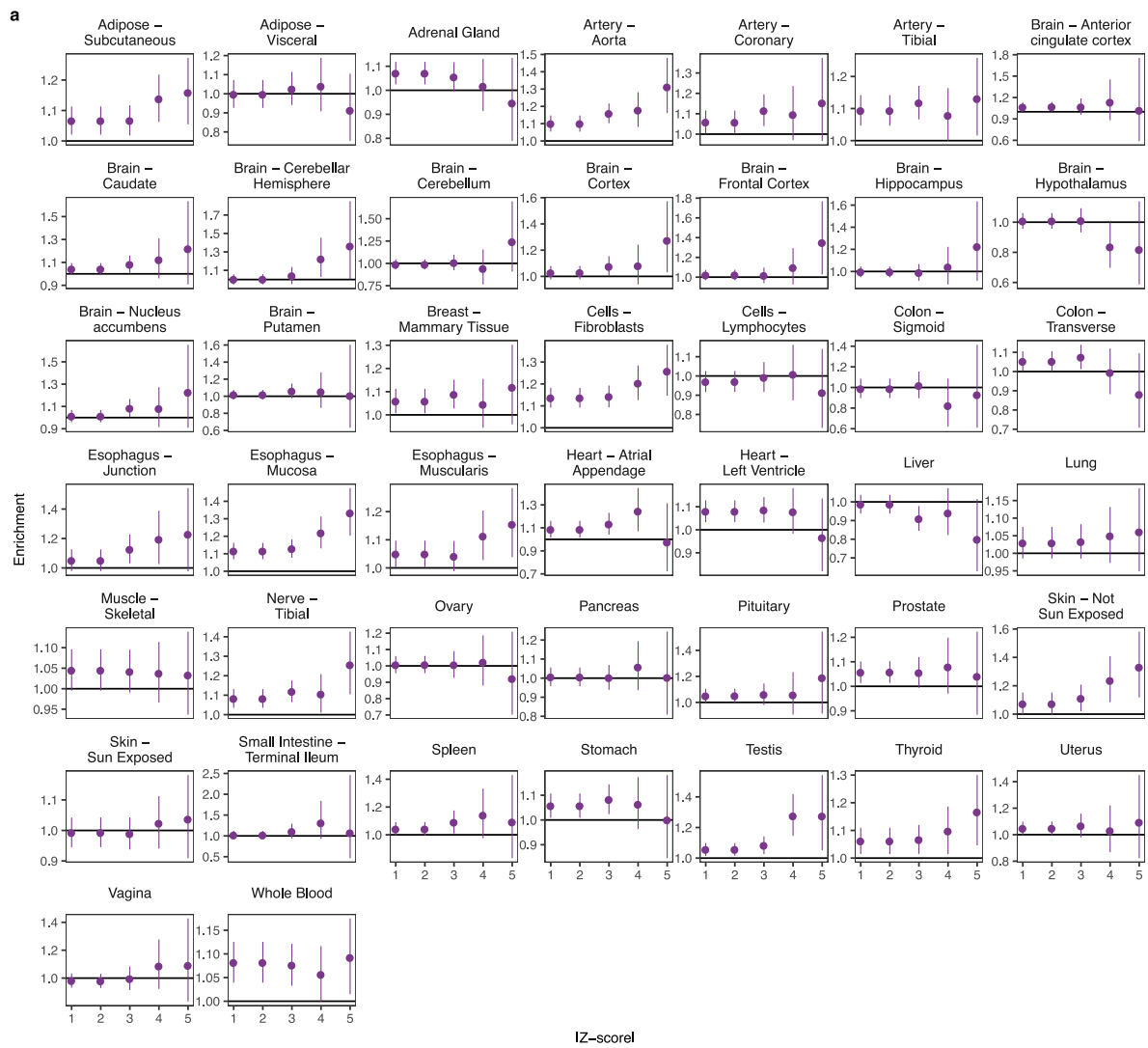
The individuals with an excess of rare variants probably had African or Asian admixture. **c**, Enrichments as in Fig. 2a and excluding individuals with >60,000 rare variants (circled in **b**), which did not substantially affect the enrichment patterns. **d**, European population allele frequency distributions in the 1000 Genomes Project of rare SNVs and indels used in our analysis. The rare variants included in our analysis were constrained to have  $MAF \leq 0.01$  in the 1000 Genomes European super population, but they were also relatively rare in each of the individual European populations.



**Extended Data Figure 5 | Comparison of overexpression and underexpression outliers.** **a**, ASE at rare exonic variants. ASE is shown as the ratio of the number of reads supporting the minor allele to the total number of reads at the site. If the rare variant is driving the extreme expression, we expect this ratio to be below 0.5 for underexpression outliers and above 0.5 for overexpression outliers. Rare coding variants were enriched for ASE in the direction of the extreme expression effect (two-sided Wilcoxon rank-sum tests, each nominal  $P < 4.0 \times 10^{-8}$ ). **b**, Expression level distribution of all genes and genes with overexpression

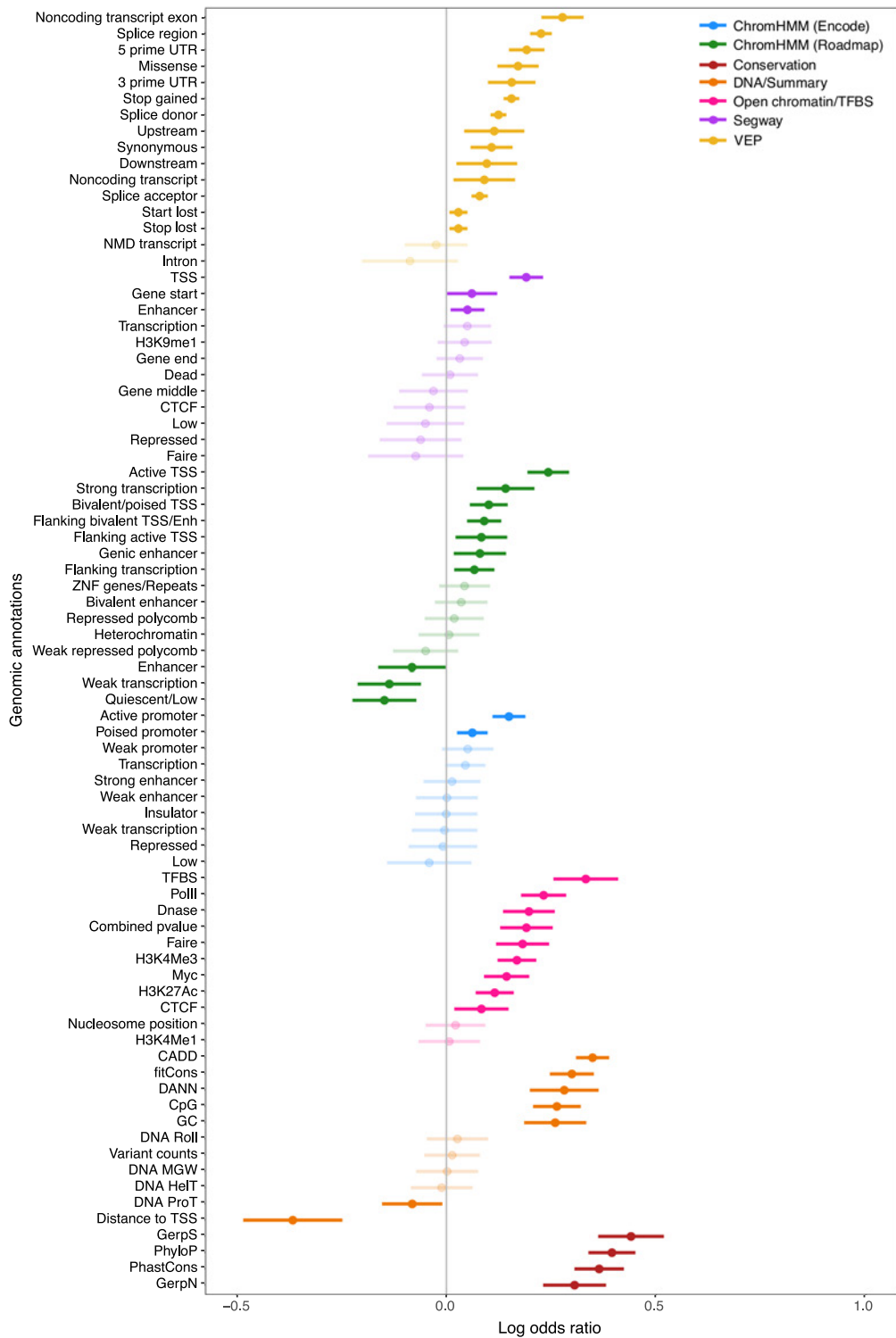
or underexpression outliers. Expression is shown as the  $\log_2$  of the median (RPKM + 2), where the median was first taken across individuals in each tissue then across expressed tissues for each gene. For genes with low expression, even an RPKM of 0 may not yield a Z-score  $\leq -2$ . Indeed, underexpression outliers were depleted among low expressed genes whereas the opposite was true of overexpression outliers (two-sided Wilcoxon rank-sum test comparing to all genes,  $P < 2.2 \times 10^{-16}$  for both overexpression and underexpression). **c**, Feature enrichments (as in Fig. 3b) shown separately for over and underexpression outliers.





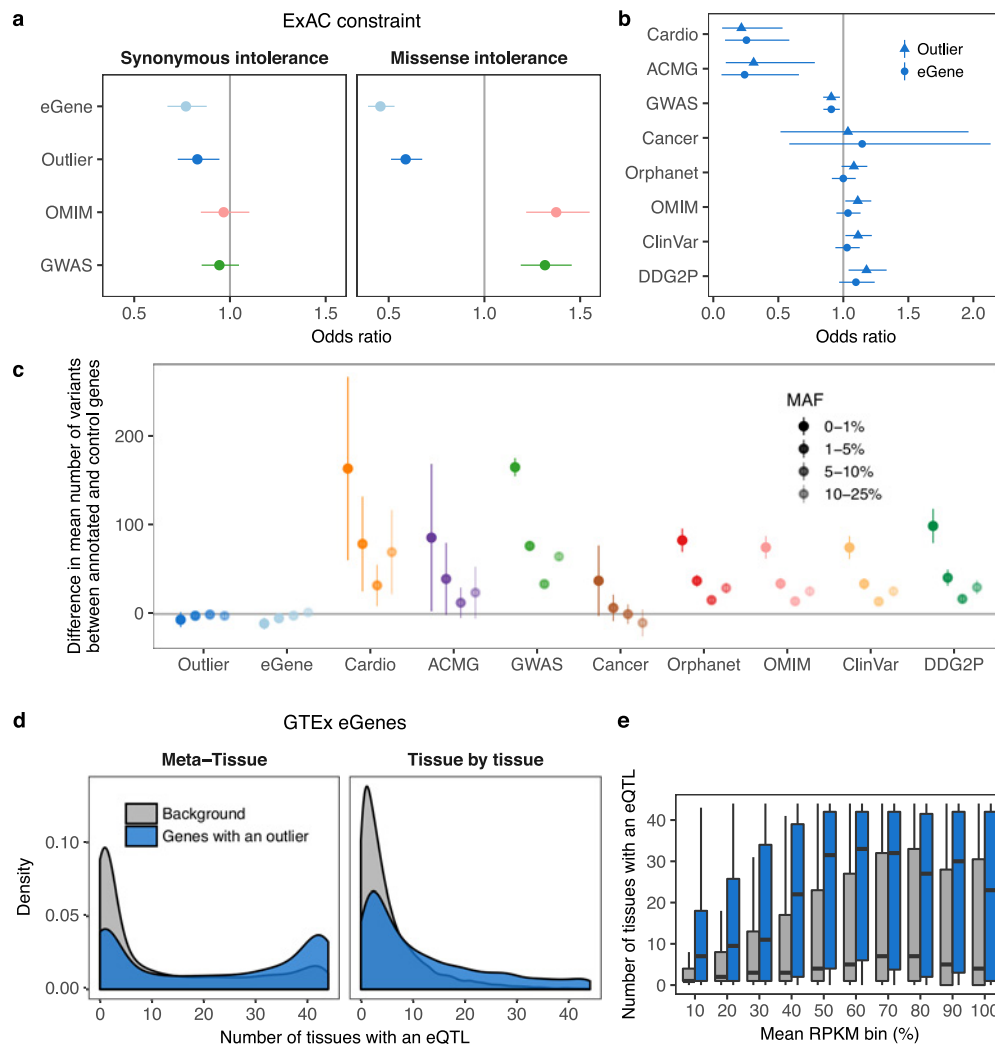
**Extended Data Figure 6 | Extended rare variant enrichments. a**, For each tissue, rare SNV enrichment in single-tissue outliers compared with non-outliers at the same genes for increasing Z-score thresholds. Enrichments calculated as in Fig. 2. The rare variant enrichments varied between tissues though the overall pattern mirrored that of multi-tissue outliers when combining all the tissues (Fig. 2b). The high variance in the enrichments

underscores the noise in single-tissue outlier discovery. **b**, As in Fig. 2a, enrichment for SNVs, indels and structural variants in outliers compared with the same genes in non-outliers, either including all rare variants or only those outside protein-coding or lincRNA exons in Gencode v.19. The enrichment of rare variants was weaker, but still significant, for all variant types when excluding exonic regions.



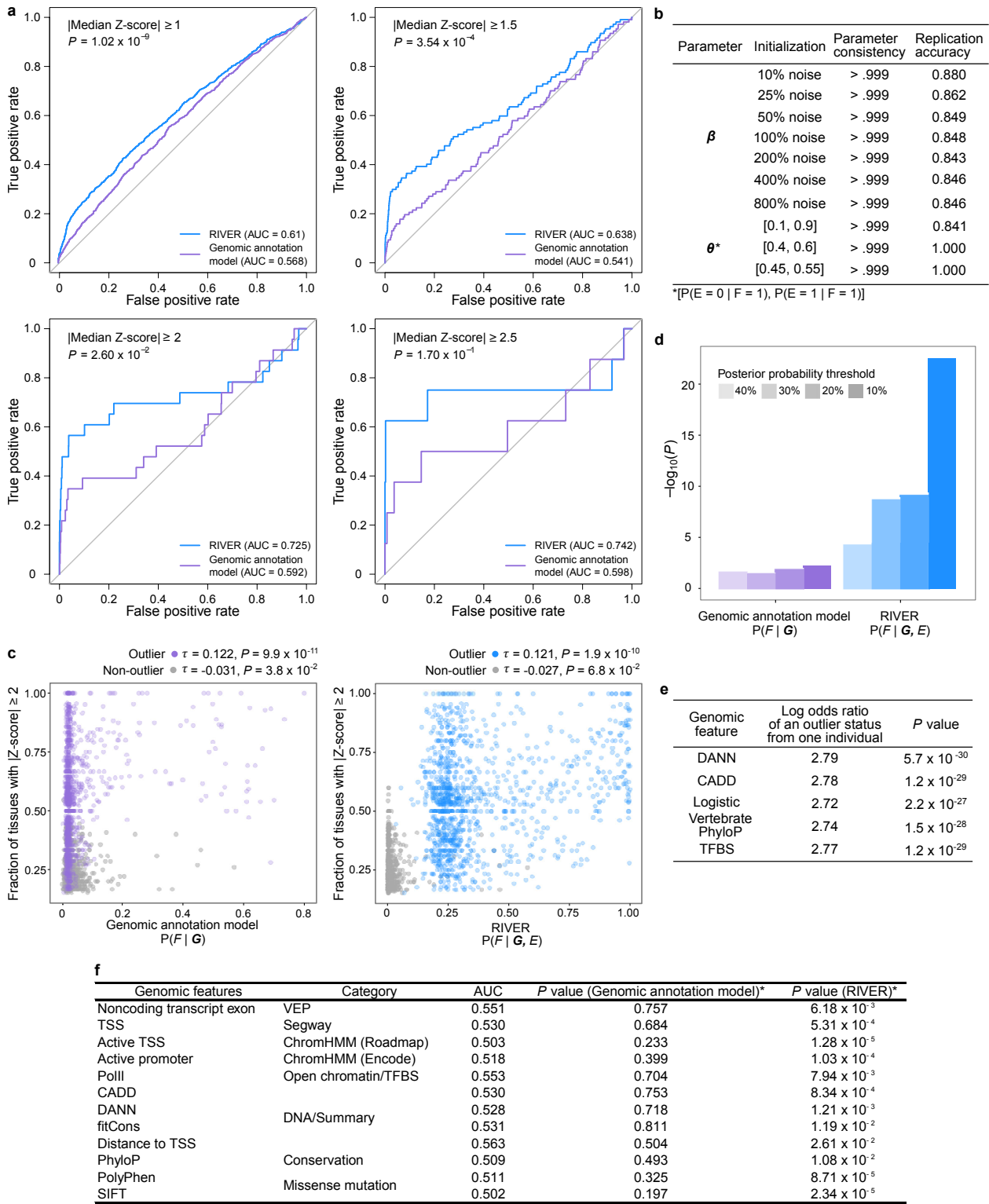
**Extended Data Figure 7 | Enrichment of an extended list of functional genomic annotations.** log odds ratios and 95% Wald confidence intervals from logistic regression models of outlier status as a function of each genomic feature. Features were calculated among rare SNVs within 10 kb of the gene. When more than one feature corresponded to the same

genomic annotation (for example, the number or the presence of rare variants in a splice region; Supplementary Table 3b), the feature with the highest enrichment is shown. Lighter shading indicates a non-significant log odds ratio (nominal  $P > 0.05$ ).



**Extended Data Figure 8 | Evolutionary constraint and regulatory control of multi-tissue outlier genes.** **a**, Odds ratio of being intolerant to synonymous and missense variants for genes with multi-tissue eQTLs (eGenes), genes with multi-tissue outliers, OMIM and GWAS genes (see Methods). As expected, GWAS and OMIM genes showed no enrichment or depletion for synonymous variation intolerant genes. Genes with multi-tissue outliers and eGenes showed slight depletion for these genes. Genes with multi-tissue outliers and eGenes were strongly depleted for genes intolerant to missense variation compared with OMIM and GWAS genes. **b**, Comparison of the depletion of disease genes among genes with a multi-tissue outlier and eGenes. Similar to Fig. 4c, bars represent 95% confidence intervals from Fisher's exact test. **c**, For each of ten gene lists, the difference in the mean number of variants near genes in the list compared with the mean for all other annotated genes. Results are stratified by minor allele frequency, and bars indicate the 95% confidence interval for the

difference from a two-sided *t*-test. Disease genes had more variants than control genes in general, and the difference was particularly striking for rare variants. This suggests that the depletion of outliers and eQTLs for certain groups of disease genes is not due to less rare variation near these genes. Instead, we hypothesize that the variation around these genes in our healthy cohort is less likely to have large regulatory effects. **d**, Distribution of the number of tissues with an eQTL for genes with and without outliers. Genes with multi-tissue outliers had eQTLs in more tissues than genes without. This suggests that they are more susceptible to shared regulatory control. This result held for both multi-tissue eQTL definitions (see Methods; Meta-Tissue: 23 versus 3 tissues, Wilcoxon rank-sum test  $P < 2.2 \times 10^{-16}$ ; tissue-by-tissue: 7 versus 3 tissues,  $P < 2.2 \times 10^{-16}$ ). **e**, This eGene enrichment was robust across different mean expression levels across tissues (two-sided Wilcoxon rank-sum tests, Bonferroni-adjusted  $P < 1 \times 10^{-11}$ ).



Extended Data Figure 9 | See next page for caption.

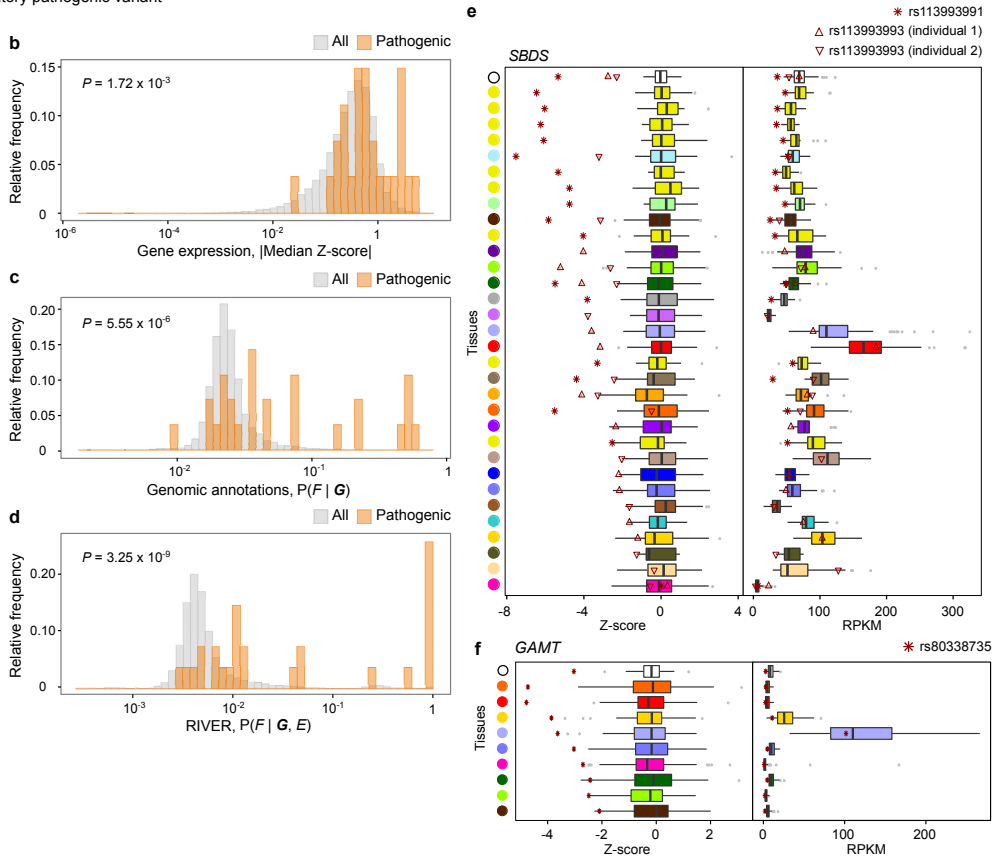
**Extended Data Figure 9 | RIVER performance.** **a**, Comparison between the predictive power of RIVER and that of the genomic annotation model, as in Fig. 5a, across different Z-score thresholds for outlier calling. Increasing the Z-score threshold improved AUC values, but reduced the number of outlier examples, which led to noisy receiver operating characteristic curves. **b**, Stability analysis of estimated parameters with different parameter initializations (see Methods). **c**, Correlations, using Kendall's  $\tau$ , between the fraction of tissues with  $|Z\text{-score}| \geq 2$  and the test probabilities from the genomic annotation model (left) and RIVER (right). We calculated test posterior probabilities using tenfold cross-validation and only considered individual and gene pairs with a fraction of tissues with  $|Z\text{-score}| \geq 2$  that was significantly different from 0.05 (one-sided binomial exact test, Benjamini–Hochberg adjusted  $P < 0.05$ ). **d**,  $P$  values from a one-sided Fisher's exact test measuring the association

between allelic imbalance (see Methods) and the posterior probability of a functional rare variant according to the genomic annotation model and RIVER. The posterior probabilities from RIVER were more strongly associated with allelic imbalance across all four thresholds tested.

**e**, Assessment of the advantage of incorporating gene expression with genomic annotations for predicting outlier status using simplified supervised models (see Methods). All models showed consistent improvement of the log odds ratio of outlier status when incorporating expression. **f**, Performance of models with 12 individual genomic features compared with the genomic annotation model and RIVER. Some models with single genomic features provided slightly better AUCs compared with the genomic annotation model, but they were not statistically different. On the other hand, RIVER predicted the effects of rare variants significantly better than each of the models that included a single feature.

Gene	Variant ID	P(F   G)	P(F   G, E)	Median Z-score	Disease	Variant type
<i>SBDS</i>	rs113993991*	0.447	0.985	-5.337	Shwachman-Diamond syndrome	nonsense
<i>TPP1</i>	rs119455955*	0.619	0.995	-4.11	Ceroid lipofuscinosis neuronal 2, Neuronal ceroid lipofuscinosis, Inborn genetic diseases	nonsense
<i>GAMT</i>	rs80338735*	0.162	0.929	-2.813	Deficiency of guanidinoacetate methyltransferase	synonymous
<i>SBDS</i>	rs113993993*	0.526	0.989	-2.753	Shwachman-Diamond syndrome, susceptibility to aplastic anemia	splice donor
<i>OGG1</i>	rs104893751	0.213	0.963	-2.733	Clear cell carcinoma of kidney	missense
<i>BBS2</i>	rs121908176*	0.519	0.992	-2.56	Bardet-Biedl syndrome 2	nonsense
<i>SBDS</i>	rs113993993*	0.52	0.988	-2.301	Shwachman-Diamond syndrome, susceptibility to aplastic anemia	splice donor
<i>NAGA</i>	rs121434529	0.047	0.063	-1.663	Schindler disease, type 1	missense
<i>OGG1</i>	rs104893751	0.213	0.239	-1.231	Clear cell carcinoma of kidney	missense
<i>SLC25A11</i>	rs140547520	0.009	0.004	-0.7	Amyotrophic lateral sclerosis 18	missense
<i>DSTYK</i>	rs200780796	0.077	0.049	-0.694	Susceptibility to congenital anomalies of the kidney and urinary tract 1	missense
<i>CLPTM1</i>	rs120074114	0.027	0.006	-0.66	Apolipoprotein c-ii variant	missense
<i>MUTYH</i>	rs34612342	0.078	0.038	0.65	Endometrial carcinoma, MYH-associated polyposis, Carcinoma of colon, Hereditary cancer-predisposing syndrome	missense
<i>IVD</i>	rs28940889	0.074	0.045	0.573	Isovaleryl-CoA dehydrogenase deficiency	missense
<i>GPR97</i>	rs121908464	0.025	0.009	0.508	Bilateral frontoparietal polymicrogyria	missense
<i>ZNF200</i>	rs61732874	0.017	0.003	-0.431	Familial Mediterranean fever	missense, 3' UTR
<i>APOC4</i>	rs120074114	0.038	0.012	0.411	Apolipoprotein c-ii variant	missense
<i>SLC7A9</i>	rs79389353	0.044	0.014	-0.375	Cystinuria	missense
<i>RPL29</i>	rs121912698	0.023	0.008	-0.371	Aminoacylase 1 deficiency	missense
<i>RPS19</i>	rs147508369	0.018	0.013	0.304	Diamond-Blackfan anemia 1	missense
<i>ABHD14B</i>	rs121912698	0.035	0.011	0.224	Aminoacylase 1 deficiency	missense
<i>ZNF200</i>	rs104895091	0.022	0.005	0.218	Autosomal dominant familial Mediterranean fever	inframe, 3' UTR
<i>ABHD14B</i>	rs121912701	0.02	0.004	0.206	Aminoacylase 1 deficiency	missense
<i>ZNF200</i>	rs28940579	0.025	0.006	0.175	Familial Mediterranean fever	missense, 3' UTR
<i>RPL29</i>	rs121912698	0.036	0.012	0.153	Aminoacylase 1 deficiency	missense
<i>RPL29</i>	rs121912701	0.021	0.005	0.142	Aminoacylase 1 deficiency	missense
<i>ABHD14B</i>	rs121912698	0.035	0.011	0.025	Aminoacylase 1 deficiency	missense

\* Regulatory pathogenic variant



Extended Data Figure 10 | See next page for caption.

**Extended Data Figure 10 | Evaluation of known pathogenic variants using RIVER.** **a**, The 27 GTEx rare SNVs reported as disease variants in ClinVar. **b–d**, Relative frequency of the  $|\text{median } Z\text{-score}|$  (**b**), posterior probabilities from the genomic annotation model (**c**) and posterior probabilities from RIVER (**d**) for all (individual, gene) pairs (grey) and 27 pairs with pathogenic variants from ClinVar (orange).  $P$  values were computed using two-sided Wilcoxon rank-sum tests. We note that rare indels and structural variants were not found nearby the genes in the individuals carrying these pathogenic variants. **e, f**, The  $Z$ -score and RPKM distributions for *SBDS* (**e**) and *GAMT* (**f**) were compared with the values from four individuals carrying regulatory pathogenic variation (red asterisks and triangles). The median  $Z$ -score and RPKM values across tissues are shown at the top of each plot (black circle). Tissues are coloured as in Fig. 1 and sorted in decreasing order of the difference

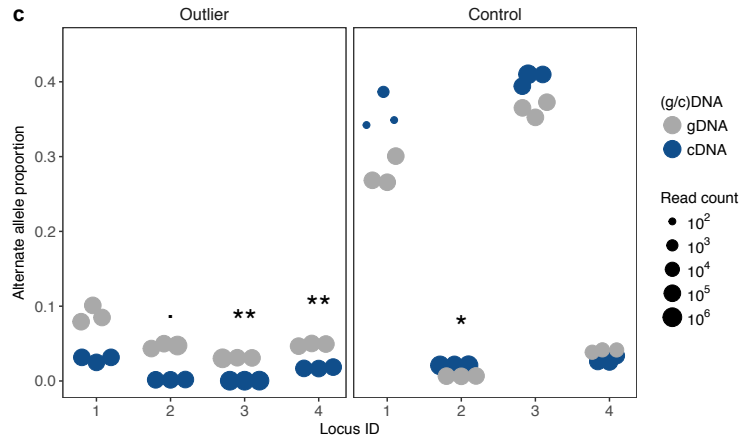
between the average  $Z$ -score of individual(s) with a regulatory pathogenic variant and the median  $Z$ -score for the tissue. Three individuals carrying a total of two unique rare variants are shown for *SBDS*. Both variants are associated with the recessive Shwachman–Diamond syndrome, which causes systemic symptoms that include pancreatic, neurological and haematologic abnormalities<sup>46</sup> and can disrupt fibroblast function<sup>47</sup>. The individuals, being heterozygous for these variants, lacked the disease phenotype. Nonetheless, we saw extreme underexpression of *SBDS* across almost all tissues in these individuals, including brain tissues, fibroblasts and pancreas. One individual had a rare variant for *GAMT* associated with cerebral creatine deficiency syndrome 2, shown to cause neurological deficiencies and also lead to low body fat<sup>48</sup>. The individual had the most extreme underexpression in (subcutaneous) adipose tissue.

**a**

Locus ID	Chr:Position	Ref/Alt	GTEx MAF	Gene	Median Z-score	RIVER score	CADD score	Coding consequence
Outlier 1	7:66459273	T/A	0.004	<i>SBDS</i>	-5.337	0.985	2.821	Stop gained
Control 1	7:66459256	T/C	0.130	<i>SBDS</i>	[-2.753, 0.773]	[0.003, 0.989]	2.191	Synonymous
Outlier 2	12:4766944	C/T	0.004	<i>NDUFA9</i>	-5.569	0.982	0.609	Stop gained
Control 2	12:4766925	G/T	0	<i>NDUFA9</i>	N/A	N/A	-0.198	Synonymous
Outlier 3	7:102944937	G/A	0.004	<i>PMPCB</i>	-5.936	0.969	5.789	Missense; Splice region; 3' UTR
Control 3	7:102948074	A/G	0	<i>PMPCB</i>	N/A	N/A	1.995	Synonymous; 3' UTR
Outlier 4	19:13885293	T/A	0.004	<i>C19orf53</i>	-4.229	0.956	2.184	Start lost
Control 4	19:13885309	C/T	0.256	<i>C19orf53</i>	[-2.496, 0.919]	[0.004, 0.400]	2.172	Synonymous

**b**

Locus ID	sgRNA
Outlier 1	GTGTTGTAAATGTTTCTAA
Control 1	ACTGATGAGATCTTCCTTTT
Outlier 2	TGCTGTGTGTACTACTCGT
Control 2	CTTCTGCTATTATAGGAAT
Outlier 3	ATAGTGCTGCTGCTGCTGG
Control 3	GACTTAGCAAAGTTTCATTT
Outlier 4	TTCCGCTGCGTGCCGGACCA
Control 4	GCAGGGGCAGCGCAAGTTTC



**Extended Data Figure 11 | Validation of large-effect rare variants using CRISPR–Cas9 genome editing.** **a**, SNVs in outliers and controls assayed for expression effects using CRISPR–Cas9 genome editing. For common SNVs in controls (MAF >1% in the GTEx cohort), the range of median Z-scores and RIVER scores are given for all individuals with the minor allele. Missing values indicate that the variant was absent from our cohort. **b**, sgRNAs for four SNVs found in outliers and four control SNVs in the same genes. **c**, Alternate (installed) gDNA and cDNA allele

proportions for four rare, coding SNVs in outliers (left) and four matched control SNVs (right). Each gDNA and cDNA sample was sequenced in triplicate (technical replicates). Asterisks denote the Bonferroni-adjusted significance level from a two-sided *t*-test of the difference between the gDNA and cDNA alternate allele proportions:  $\cdot P < 0.05$ ,  $*P < 0.01$ ,  $***P < 0.001$ . Although one control SNV showed a significant difference in the alternate allele proportion between cDNA and gDNA, it displayed an increase rather than a decrease in expression.



## Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form is intended for publication with all accepted life science papers and provides structure for consistency and transparency in reporting. Every life science submission will use this form; some list items might not apply to an individual manuscript, but all fields must be completed for clarity.

For further information on the points included in this form, see [Reporting Life Sciences Research](#). For further information on Nature Research policies, including our [data availability policy](#), see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

### ► Experimental design

#### 1. Sample size

Describe how sample size was determined.

See Figure 1 of a previous GTEx manuscript (PMID: 23715323) for a detailed description of statistical power and sample size quantification.

#### 2. Data exclusions

Describe any data exclusions.

We used the GTEx samples included in the eQTL analyses. The data inclusion/exclusion criteria are described in that manuscript (submitted concurrently). We further excluded any individuals with more than 50 multi-tissue outliers, as described in the methods.

#### 3. Replication

Describe whether the experimental findings were reliably reproduced.

We validated the effects on expression of four rare variants via CRISPR/Cas9 genome editing. We attempted to validate two additional variants but excluded them from our analysis because they failed to amplify successfully, so we could not accurately measure their effect.

#### 4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

Order of sample processing for library preparation and sequencing was randomized to avoid batch effects.

#### 5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

No blinding took place.

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

#### 6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or in the Methods section if additional space is needed).

n/a Confirmed

- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.)
- A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- A statement indicating how many times each experiment was replicated
- The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section)
- A description of any assumptions or corrections, such as an adjustment for multiple comparisons
- The test results (e.g.  $P$  values) given as exact values whenever possible and with confidence intervals noted
- A clear description of statistics including central tendency (e.g. median, mean) and variation (e.g. standard deviation, interquartile range)
- Clearly defined error bars

See the web collection on [statistics for biologists](#) for further resources and guidance.

## ► Software

Policy information about [availability of computer code](#)

### 7. Software

Describe the software used to analyze the data in this study.

RIVER is available at <https://bioconductor.org/packages/release/bioc/html/RIVER.html>. Additionally, the code for running analyses and producing the figures throughout the manuscript is available separately (<https://github.com/joed3/GTEXV6PRareVariation>).

For manuscripts utilizing custom algorithms or software that are central to the paper but not yet described in the published literature, software must be made available to editors and reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). *Nature Methods* [guidance for providing algorithms and software for publication](#) provides further information on this topic.

## ► Materials and reagents

Policy information about [availability of materials](#)

### 8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

Residual biospecimens are available to all researchers according to the Genotype-Tissue Expression (GTEx) project biospecimens access policy. The policy and related forms can be found on the GTEx Portal ([gtexportal.org](http://gtexportal.org)) under the Biobank Tab.

### 9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

N/A

### 10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

For the experimental validation, we used a K562 cell line previously described (PMID:27798611).

b. Describe the method of cell line authentication used.

None.

c. Report whether the cell lines were tested for mycoplasma contamination.

Yes, the K562 cells were tested for mycoplasma.

d. If any of the cell lines used are listed in the database of commonly misidentified cell lines maintained by [ICLAC](#), provide a scientific rationale for their use.

N/A

## ► Animals and human research participants

Policy information about [studies involving animals](#); when reporting animal research, follow the [ARRIVE guidelines](#)

### 11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

N/A

Policy information about [studies involving human research participants](#)

### 12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

The human research participants and relevant covariates are described in this manuscript and the GTEx eQTL manuscript submitted concurrently.