

A genome-wide mutational constraint map quantified from variation in 76,156 human genomes

Siwei Chen^{1,2,†}, Laurent C. Francioli^{1,2,†}, Julia K. Goodrich^{1,2}, Ryan L. Collins^{1,3,4}, Masahiro Kanai^{1,2}, Qingbo Wang^{1,5}, Jessica Alföldi^{1,2}, Nicholas A. Watts^{1,2}, Christopher Vittal^{1,2}, Laura D. Gauthier⁶, Timothy Poterba^{1,2,7}, Michael W. Wilson^{1,2}, Yekaterina Tarasova¹, William Phu^{1,2}, Mary T. Yohannes¹, Zan Koenig¹, Yossi Farjoun⁶, Eric Banks⁶, Stacey Donnelly⁷, Stacey Gabriel^{1,7}, Namrata Gupta^{1,7}, Steven Ferreira⁷, Charlotte Tolonen⁶, Sam Novod⁶, Louis Bergelson⁶, David Roazen⁶, Valentin Ruano-Rubio⁶, Miguel Covarrubias⁶, Christopher Llanwarne⁶, Nikelle Petrillo⁶, Gordon Wade⁶, Thibault Jeandet⁶, Ruchi Munshi⁶, Kathleen Tibbetts⁶, gnomAD Project Consortium*, Anne O'Donnell-Luria^{1,2,8}, Matthew Solomonson^{1,2}, Cotton Seed^{2,9}, Alicia R. Martin^{1,2}, Michael E. Talkowski^{1,3}, Heidi L. Rehm^{1,3}, Mark J. Daly^{1,2,10}, Grace Tiao^{1,2}, Benjamin M. Neale^{1,2,9,‡}, Daniel G. MacArthur^{1,2,11,12,‡}, Konrad J. Karczewski^{1,2,9}

¹Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

²Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

³Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA

⁴Division of Medical Sciences, Harvard Medical School, Boston, MA, USA

⁵Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita, Japan

⁶Data Science Platform, Broad Institute of MIT and Harvard, Cambridge, MA, USA

⁷Broad Genomics, Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA

⁸Division of Genetics and Genomics, Boston Children's Hospital, Boston, MA

⁹Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA

¹⁰Institute for Molecular Medicine Finland, (FIMM) Helsinki, Finland

¹¹Centre for Population Genomics, Garvan Institute of Medical Research and UNSW Sydney, Sydney, Australia

¹²Centre for Population Genomics, Murdoch Children's Research Institute, Melbourne, Australia

*Lists of authors and their affiliations appear at the end of the paper

‡These authors contributed equally: Benjamin M. Neale, Daniel G. MacArthur.

†These authors contributed equally: Siwei Chen, Laurent C. Francioli.

Correspondence should be addressed to K.J.K (konradk@broadinstitute.org) and S.C (siwei@broadinstitute.org)

Abstract

The depletion of disruptive variation caused by purifying natural selection (constraint) has been widely used to investigate protein-coding genes underlying human disorders, but attempts to assess constraint for non-protein-coding regions have proven more difficult. Here we aggregate, process, and release a dataset of 76,156 human genomes from the Genome Aggregation Database (gnomAD), the largest public open-access human genome reference dataset, and use this dataset to build a mutational constraint map for the whole genome. We present a refined mutational model that incorporates local sequence context and regional genomic features to detect depletions of variation across the genome. As expected, protein-coding sequences overall are under stronger constraint than non-coding regions. Within the non-coding genome, constrained regions are enriched for known regulatory elements and variants implicated in complex human diseases and traits, facilitating the triangulation of biological annotation, disease association, and natural selection to non-coding DNA analysis. More constrained regulatory elements tend to regulate more constrained protein-coding genes, while non-coding constraint captures additional

functional information underrecognized by gene constraint metrics. We demonstrate that this genome-wide constraint map provides an effective approach for characterizing the non-coding genome and improving the identification and interpretation of functional human genetic variation.

Introduction

The expansion in the scale of human whole-genome or exome reference data has allowed characterization of the patterns of variation in human genes. With these data it is possible to directly assess the strength of negative selection on loss-of-function (LoF) and missense variation by modeling “constraint,” the depletion of variation in a gene compared to an expectation conditioned on that gene’s mutability. Using coding variant data from sequencing of more than 125K humans¹, we previously developed a constraint metric that classifies each protein-coding gene along a spectrum of LoF intolerance¹, providing a valuable resource for studying the functional significance of human genes²⁻⁵. Although of outsized biological importance, protein-coding regions comprise less than 2% of the human genome, and the vast non-coding genome has been much less characterized, even though the importance of non-coding variation in human complex diseases has been long recognized⁶⁻¹⁰.

Several challenges arise when extending the gene constraint model to the non-coding space. First, the sample size of human whole-genome reference data has been relatively small compared to the exome, limiting the power of detecting depletions of variation at a fine scale. Second, our detailed understanding of coding region exon structure and effect of specific variants on amino acid translation enables a precision not available in non-coding analysis. Third, there is a strong expectation from Mendelian genetics and existing constraint analyses that the coding regions, while a small fraction of the genome, harbor a massively disproportionate amount of rare and common disease mutations under selection. Fourth, the mutation rate in non-coding regions is highly heterogeneous, and can be affected not only by local sequence context as commonly modeled in gene constraint metrics but also by a variety of genomic features at larger scales^{11,12}.

Current methods attempting to evaluate non-coding constraint can be broadly divided into three categories: 1) context-dependent mutational models that assess the deviation of observed variation from an expectation based on the sequence composition of *k*-mers (e.g., Orion¹³, CDTs¹⁴, DR¹⁵); 2) machine-learning classifiers that are trained to differentiate between disease-associated variants and benign variants (e.g., CADD¹⁶, GWAVA¹⁷, JARVIS¹⁸); and 3) phylogenetic conservation scores that use comparative genomics data to infer evolutionary constraint (e.g., phastCons¹⁹, phyloP²⁰). While all these methods aid in our understanding of the non-coding genome, each suffer from limitations/biases, respectively as 1) overlooking the influence of regional genomic features beyond the scale of flanking nucleotides on mutation rate; 2) a strong dependence on the availability of well-characterized functional mutations as training data; and 3) compromised power to detect regions that have only recently been under selection in the human lineage and may have a functional impact on human-specific traits or diseases.

Here we present a genome-wide map of human constraint, generated from a high-quality set of variant calls from 76,156 whole-genome sequences (gnomAD v3.1.2 <https://gnomad.broadinstitute.org>). We describe an improved model of human mutation rates that jointly analyzes local sequence context and regional genomic features and quantifies the depletion of variation in tiled windows across the entire genome. By building a more comprehensive picture of genic constraint rather than solely focusing on coding variation, we facilitate the functional interpretation of non-coding regions and improve the characterization of gene function in the context of the regulatory network. Our study aims to depict a

genome-wide view of how natural selection shapes patterns of human genetic variation and generate a more comprehensive catalog of functional genomic elements with potential clinical significance.

Results

Aggregation and quality control of genome sequence data

We aggregated, reprocessed, and performed joint variant-calling on 153,030 whole genomes mapped to human genome reference build GRCh38, of which 76,156 samples were retained as high-quality sequences from unrelated individuals, without known severe pediatric disease, and with appropriate consent and data use permissions for the sharing of aggregate variant data. Among these samples, 36,811 (48.3%) are of non-European ancestry, including 20,744 individuals with African ancestries and 7,647 individuals with admixed Amerindigenous ancestries. After stringent quality control (see Supplementary Information), we discovered a set of 644,267,978 high-confidence short nuclear variants (single nucleotide/indel variants; gnomAD v3.1.2), of which 390,393,900 low-frequency (allele frequency [AF] $\leq 1\%$), high-quality single nucleotide variants were used for building the genome-wide constraint map. These correspond to approximately one variant every 4.9 bp (one low-frequency variant every 8 bp) of the genome, providing a high density of variation.

Quantifying mutational constraint across the genome

To construct a genome-wide mutational constraint map, we divided the genome into continuous non-overlapping 1kb windows, and quantified constraint for each window by comparing the expected and the observed variation in our gnomAD dataset. Here, we implemented a refined mutational model, which incorporates trinucleotide sequence context, base-level methylation, and regional genomic features to predict expected levels of variation under neutrality. In brief, we estimated the relative mutability for each single nucleotide substitution with one base of adjacent nucleotide context (e.g., ACG \rightarrow ATG), with adjustment for the effect of methylation on mutation rate at CpG sites, which become saturated for mutation at sample sizes above $\sim 10K$ genomes²¹ (Extended Fig. 1a,b; Methods). Meanwhile, we adjusted the effects of regional genomic features for each trinucleotide mutation rate based on the occurrence of *de novo* mutations ($N=413,304$ previously detected in family-based whole-genome sequencing studies^{22,23}; Extended Fig. 1c), and then applied it to establish the expected number of variants per 1kb across the entire genome (Methods).

We quantified the deviation from expectation for each 1kb window using a Z score²⁴ (Methods; Extended Fig. 1d,e), which was centered around zero for non-coding regions (median=0.08), and was significantly higher (more constrained) for windows containing any protein-coding sequences (median=1.47, Wilcoxon $P < 10^{-200}$; Fig. 1a). The constraint Z score is positively correlated with the percentage of coding bases in a window and presented a substantial shift towards higher constraint for exonic sequences from directly concatenating coding exons into 1kb windows (median=3.17; Extended Fig. 2a-c). About 3.12% and 0.05% of the non-coding windows exhibited constraint as strong as the 50th and 90th percentile of exonic regions (Extended Fig. 2d). Comparing our Z score against the adjusted proportion of singletons (APS) score, a measure of constraint developed for structural variation (SV)²⁵, we found a significant correlation (linear regression $\beta = 0.01$, $P = 4.3 \times 10^{-65}$, Fig. 1b; Methods), providing an internal validation of our approach.

Investigating genomic properties of non-coding regions under constraint

To further validate our constraint metric and investigate the functional relevance of non-coding regions under selection, we examined the correlation between our constraint Z score and several annotations of functional non-coding sequences (Fig. 2a). First, we found that candidate cis-regulatory elements (cCREs, derived from ENCODE²⁶ integrated DNase- and ChIP-seq data) are significantly enriched in the most constrained percentile of the genome ($Z \geq 4$, OR=2.77 compared to the genome-wide average, Fisher's exact $P < 10^{-200}$); cCREs with a promoter-like signature (cCRE-PLS) presented the strongest enrichment (OR=7.28), followed by elements with a proximal/distal enhancer-like signature (pELS OR=4.35, dELS

OR=2.14), and as a negative control, elements bound by CTCF but not associated with a regulatory signature showed no enrichment (CTCF-only OR=0.82). These patterns indicate that a large fraction of the constrained non-coding regions may serve a regulatory role, in line with previous findings^{13,14,18}. Similarly, significant enrichment was found for an independent set of active, *in vivo*-transcribed enhancers (identified by FANTOM CAGE analyses²⁷; OR=3.58) and super enhancers²⁸ (OR=3.41), which are groups of enhancers in close genomic proximity regulating genes important for cell type specification²⁹. By aggregating the regulatory annotations, we estimated that ~10.4% and ~6.3% of promoters and enhancers, respectively, are under selection as strong as coding exons on average (Extended Fig. 3a; Methods). A much higher proportion, 22.2%, was found for sequences encoding microRNAs (miRNAs), which are increasingly recognized as key mediators in various developmental and physiological processes³⁰. In contrast, only 3.7% of long non-coding RNAs (lncRNAs) exhibited such strong constraint, similar to that of non-coding regions overall (3.1%; Extended Fig. 2d and 3b).

We next examined the distribution of potentially functional non-coding variants on the constraint spectrum. There was significant enrichment for non-coding variants implicated by genome-wide association studies (GWAS) in the constrained end of the genome: 837/19,471 constrained windows [$Z \geq 4$] overlapped with GWAS Catalog³¹ annotations (OR=1.57 compared to the genome-wide average of 51,430/1,843,559, Fisher's exact $P=2.5 \times 10^{-32}$, **Fig. 2b**; Methods). The enrichment became stronger when restricted to the subset of variants that had been replicated by an independent study (OR=2.08, $P=4.1 \times 10^{-13}$). Moreover, further strong signals were found for likely causal GWAS variants fine-mapped for 148 complex diseases and traits in large-scale biobanks³² (OR=3.24, $P=3.0 \times 10^{-10}$; Methods). Across the 95% credible set (CS)-trait pairs, strong enrichment was predominantly seen in disease phenotypes, including coronary artery disease (CAD), inguinal hernia, fibroblastic disorders, and glaucoma (ORs 3.31-6.02, **Fig. 2c**; Methods). In the 95% CS of CAD, for instance, the highest constraint score was found for rs1897107 and rs1897109 (both within the same genomic window chr6:160725000-160726000, $Z=6.32$); high constraint ($Z \geq 4$) was also found for 26 variants from the same CS (totaling 28/52), which together spanned a ~153 kb sequence downstream of the gene *PLG* (**Fig. 2d**). *PLG* encodes the plasminogen protein that circulates in blood plasma and is converted to plasmin to dissolve the fibrin of blood clots. While dysregulation of the PLG-plasmin system has been frequently associated with CAD³³⁻³⁸, no specific variants in *PLG* have been implicated. Our results prioritized a set of non-coding variants in highly constrained regions of *PLG*, which adds quantitative evidence to the implication of *PLG* in CAD and may help direct or prioritize follow-up functional experiments.

Collectively, these results demonstrated a significant positive correlation between constraint and existing functional non-coding annotations, validating our non-coding constraint metric. Yet, we suggest that our metric provides additional information for the characterization of non-coding regions. For instance, prioritizing ENCODE cCREs by constraint Z score revealed increasingly stronger GWAS enrichment in the more constrained cCREs (Extended Fig. 4a), and constrained regions outside cCREs also captured significant signals, reflecting the value of non-coding constraint independent of regulatory annotations. Moreover, besides prioritizing existing GWAS results, constraint can be used as a prior for statistical fine-mapping. Using UK Biobank (UKBB) traits as examples, incorporating constraint Z score into the functionally informed fine-mapping model³⁹ predicted ~13K variant-trait pairs to have an increased posterior inclusion probability of causality ($\Delta\text{PIP} \geq 0.01$), in which 164 likely causal associations were newly identified at $\text{PIP} \geq 0.8$ (Extended Fig. 4b; Methods). While only functional tests can ultimately validate the underlying causality, our constraint map presents a valuable resource for expanding or refining the catalog of functional non-coding variants in the human genome.

Comparing constraint Z score to other genome-wide predictive scores

To benchmark the performance of our constraint map in prioritizing non-coding variants, we extended the analyses of GWAS variants to compare constraint Z score to other population genetics-based constraint metrics (Orion¹³, CDS¹⁴, gWRVIS¹⁸, and DR¹⁵). Specifically, we assessed the performance of different metrics in identifying putative functional non-coding variants – as aforementioned, a) GWAS Catalog³¹ variants (N=9,229 with an independent replication); b) GWAS fine-mapping³² variants (N=2,191), and additionally, c) a subset of high-confidence causal variants from b (N=140); and d) pathogenic Mendelian variants (N=288 from ClinVar⁴⁰) – against background variants in the population with a similar allele frequency (hereafter referred to as “positive” and “negative” variant set, respectively; Methods). Overall, constraint Z score achieved the highest performance across all comparisons, as measured by the area under curve (AUC) statistic (**Fig. 3a,b** and Extended Fig. 5). The performance was also more stable than others when varying the allele frequency threshold for the negative variant set (Extended Fig. 5). This may be due to other metrics being informed by the site frequency spectrum, which made the classification performance sensitive to differences in allele frequency between the positive and negative variants. We also showed that our performance was robust to the artificial break of genomic windows (non-overlapping 1kb) by reconstructing constraint Z scores in a sliding-window (1kb stepped by 100bp) approach as adopted by other metrics (Extended Fig. 6).

Extending the comparison to include phylogeny-based conservation scores (phyloP²⁰, phastCons¹⁹, and GERP⁴¹) revealed relatively low performance compared to the population genetics-based constraint metrics (**Fig. 3a,b**). The conservation scores were weakly correlated with constraint (Spearman's rank correlation coefficient 0.017-0.19, Extended Fig. 7), suggesting that intraspecies (human lineage-specific) constrained regions complement, rather than being a subset of, regions that are conserved across species. Each individual metric also contributed to the classification when modeled as independent predictive variables (**Fig. 3c,d**; Methods), reinforcing the complementary nature of different approaches. Variants that were uniquely captured by constraint Z score, for instance, tended to be in regions with high recombination rates (3.45-fold the rest of the positive variant set) and high DNA methylation (2.74-fold; Methods), both associated with an increased mutation rate that had been adjusted in our refined mutational model. To further illustrate this improvement, we rebuilt our constraint model from solely the local sequence context, i.e., without adjustment on mutation rate by regional genomic features, and confirmed that constraint Z score outperformed such metrics (Extended Fig. 6). Altogether, we demonstrate that constraint Z score is an effective metric for identifying functional variants in the non-coding genome; at the same time, we suggest that a combination of different metrics is likely to provide the most informative results.

Exploring non-coding dosage sensitivity in the constrained genome

Besides single nucleotide variants (SNVs) that have been extensively studied in GWAS, copy number variants (CNVs) causing dosage alterations (deletions/loss or duplications/gain) of DNA represent another significant risk factor for human disease⁴²⁻⁴⁷. Yet, unlike SNVs, CNVs can be large and determining the “minimal critical region”⁴⁸ with a pathogenic effect has been a major challenge. Although CNVs primarily affect non-coding sequences, the most commonly studied mechanism is still the dosage alteration of overlapping protein-coding genes⁴⁹. Using our genome-wide constraint map, we explored the possibility that constrained non-coding regions are also sensitive to a dosage effect, which may underlie the pathogenicity of corresponding CNVs.

We surveyed a collection of ~100K CNVs from a genome-wide CNV morbidity map of developmental delay and congenital birth defects^{50,51}. There was a substantial excess of CNVs that affected constrained non-coding regions ($Z \geq 4$) among individuals with developmental disorders (DD cases) in comparison to healthy

controls (42.6% versus 12.5%, OR=5.21, Fisher's exact $P < 10^{-200}$, **Fig. 4a**; Methods). Moreover, of the 19 loci that had been previously identified as pathogenic⁵⁰, all but one (94.7%) affected constrained non-coding regions; the high incidence was recapitulated in a curated set of ~4K putative pathogenic CNVs (85.5% in ClinVar⁴⁰, **Fig. 4a**). Importantly, the case-control enrichment remained significant, albeit attenuated, after adjusting for the size and gene content of each CNV and when being tested in the subset of CNVs that are exclusively non-coding (**Fig. 4b**; Methods). Non-coding constraint presented high association with DD CNVs conditioning on gene constraint (log[OR]=1.06, logistic regression $P < 10^{-100}$), lending support to the possibility that dosage alteration of constrained non-coding regions may be an alternative explanation for the mechanism of CNVs underlying DDs.

One known example of pathogenic non-coding dosage alteration is the duplication of *IHH* regulatory domain in synpolydactyly and craniosynostosis⁵²⁻⁵⁴. The four implicated duplications covered a ~102kb sequence upstream of *IHH*, with a ~10kb overlapping region ("critical region"⁴⁸; **Fig. 4c**). The region contained no genes but exhibited high levels of constraint (median $Z = 2.52$, Wilcoxon $P = 1.3 \times 10^{-3}$ compared to the rest of the genome). The most constrained window (chr2:219111000-219112000, $Z = 4.12$) overlapped with the major enhancer of *IHH*, the duplication of which has been shown to result in dosage-dependent *IHH* misexpression and consequently syndactyly and malformation of the skull⁵⁴. This result highlights a potential use of the constraint metric to prioritize non-coding regions within large CNVs. As a further illustration, we examined a set of non-coding CNVs that had the highest constraint score among the DD cases. The most constrained genomic window (chr11:133208000-133209000, $Z = 8.87$) was affected by 12 deletions spanning a ~400kb non-coding sequence (**Fig. 4d**). While of varying size, the deletions shared a common region of ~20kb (potential "critical region"), which encompassed the most constrained window and overall, showed a significantly higher constraint than the other affected regions (median $Z = 1.63$ vs 0.84, Wilcoxon $P = 1.6 \times 10^{-3}$; **Fig. 4d**). In addition, the ~400kb sequence also harbored two deletions from healthy controls, which interestingly, overlapped with the two lowest constraint scores within the region and were significantly less constrained than those from DD cases (median $Z = 1.07$ vs 0.62, Wilcoxon $P = 4.74 \times 10^{-4}$). These findings suggest that our constraint metric can be a useful indicator of critical regions affected by large CNVs, facilitating the identification of non-coding risk factors in CNV disease association studies.

Leveraging non-coding constraint to improve gene function characterization

Given the significant role of constrained non-coding regions in gene regulation, it is natural to expect that more constrained regulatory elements would regulate more constrained genes. To test this, we analyzed the constraint for enhancers that had been linked to specific genes⁵⁵ (Methods). More constrained non-coding regions were more frequently linked to regulating a gene (**Fig. 5a**), and as expected, enhancers linked to constrained genes (predicted by loss-of-function observed/expected upper bound fraction [LOEUF]¹, or curated disease genes from ⁵⁶⁻⁵⁸; Methods) were significantly more constrained than those linked to presumably less constrained genes (median $Z = 2.71$ versus 1.99, Wilcoxon $P = 1.3 \times 10^{-26}$, **Fig. 5b**; Methods), thus supporting a correlated constraint between genes and their regulatory elements.

On the other hand, a particularly interesting set of associations are the links between constrained enhancers and the "unconstrained" genes predicted by LOEUF, because these links may reflect functional significance of the "unconstrained" genes that had been previously unrecognized. The lack of predicted gene constraint can be explained by the design of LOEUF as a measure of intolerance to rare LoF variation, where small genes with few expected LoF variants are likely underpowered. Indeed, stratifying genes by the number of expected LoF variants showed a significantly higher enhancer constraint for genes that were underpowered (≤ 5 expected LoF variants)¹ compared to genes that were sufficiently powered while

scored as unconstrained (median $Z=2.64$ versus 2.27, Wilcoxon $P=9.8\times 10^{-4}$, **Fig. 5a**). This suggests that certain underpowered genes may be functionally important but were not recognized in gene constraint evaluation. For instance, *ASCL2*, a basic helix-loop-helix (bHLH) transcription factor, had only 0.57 expected LoFs (versus 0 observed) across >125K exomes¹; although being depleted for LoF variation, the absolute difference was too small to obtain a precise estimate of LoF intolerance. Yet, we found *ASCL2* had a highly constrained enhancer ($Z=5.58$), located ~16kb upstream of the gene, where >40% of the expected variants were depleted (188.6 expected versus 112 observed, chr11:2286000-2287000). The same genomic window also contained an eQTL chr11:2286192:G>T that was predicted to be significantly associated with *ASCL2* expression⁵⁹; elevated *ASCL2* expression has been implicated in the development and progression of several human cancers⁶⁰⁻⁶². This example highlights the value of non-coding constraint – as a complementary metric to gene constraint – for identifying functionally important genes.

A practical implementation of this finding is to integrate the constraint of regulatory elements into the modeling of gene constraint, which essentially borrows power from extending the functional unit of a gene to encompass its regulatory components. As a proof-of-principle, we tested whether adding the enhancer constraint Z score to LOEUF improves the prioritization of underpowered genes. The enhancer constraint was found a significant predictor of constrained genes (logistic regression $P=7.4\times 10^{-11}$ conditioning on LOEUF) and significantly improved the performance of LOEUF in identifying constrained genes that were underpowered (AUC = 0.80 versus 0.73, bootstrap $P=0.03$, **Fig. 5b**; Methods). Moreover, such approaches would allow incorporation of tissue/cell-type specific information into gene constraint modeling given the diverse range of epigenomic data. We explored this by testing whether the constraint of tissue-specific enhancers is predictive of tissue-specific gene expression (as a proxy for tissue-specific gene function). The enhancer constraint Z score, again conditioning on LOEUF, was a significant predictor of the expression level of target genes in matched tissue types (**Fig. 5c**; Methods). These results further support the application of our constraint metric for improving the characterization of gene function. While we acknowledge that the biological consequences of mutations in enhancers are not clearly understood and thus natural selection may differ in strength depending on mechanistic consequence, an extended model to incorporate non-coding variation information in a biologically-informed way holds promise to facilitate our understanding of the molecular mechanisms underlying selection.

Discussion

We have previously developed constraint metrics that leverage population-scale exome and genome sequencing data to evaluate genic intolerance to coding variation for each protein-coding gene^{1,21}. Here, we adopted the same principle with an extended mutational model to assess constraint across the entire genome, using our latest release of gnomAD (v3.1.2), a dataset of harmonized high-quality whole-genome sequences from 76,156 individuals of diverse ancestries. Improvements to constraint modeling include unified fitting of the mutation rate for all substitution and trinucleotide contexts and inclusion of regional genomic features to refine the expected variation in non-coding regions (Methods). We validated our metric using a series of external functional annotations, with a focus on the non-coding genome, and demonstrated the value of our metric for prioritizing non-coding elements and identifying functionally important genes. We have made the constraint scores publicly accessible via the gnomAD browser (<https://gnomad.broadinstitute.org>).

One key challenge in quantifying non-coding constraint is the estimation of mutation rate under neutrality, which can be affected by various genomic features at different scales. To this end, we extended our previous mutational model, which computed the relative mutability of each substitution in a

trinucleotide context, to include adjustment for the effects of regional genomic features. The adjustment was applied to each specific trinucleotide context and allowed a varying genomic scale for each specific feature (Methods). The added value of this adjustment was demonstrated by the improved performance of the constraint Z score in identifying functional variants (Extended Fig. 6). Our constraint metric also outperformed other genome-wide predictive scores, while each metric tended to provide complementary information. We note that all comparisons were restricted to non-coding regions for explicitly evaluating the metrics in prioritizing non-coding variants, and we further eliminated potential bias from nearby genes by recapitulating the results within regions >10kb away from any protein-coding exons (Supplementary Information). Overall, our constraint metric presented consistent, high performance in identifying functional non-coding variants in the human genome.

Despite the clear constraint signal identified for non-coding regions, many limitations exist. First, the lack of prior classification of the molecular consequences of non-coding variants, as analogous to “nonsynonymous” versus “synonymous” informed by the genetic code in coding regions, limits the resolution of non-coding constraint assessment (e.g., to measure constraint against “LoF” variation). While there are rich resources defining regulatory elements in the non-coding genome, no method is available for determining the impact of each possible variant on gene expression and the distribution of their effect sizes genome-wide. Further, the interpretation of non-coding constraint, especially in the context of gene regulation, can only be informative when considered in a particular context, such as a tissue/cell type, developmental stage, or environment. Such information is not inherently built into our constraint metric nor in the mutational dataset; thus *ad hoc* integration of external annotations (e.g., tissue-specific enhancers as analyzed in this study) is often necessary for justifying specific biological implications. Also, since the detection of depletion of variation is immune to negative selection after reproductive age, genomic regions involved in late-onset phenotypes are likely to go underrecognized.

Finally, while this is the largest dataset of human genomes examined to date for non-coding constraint, our method will substantially increase in power and resolution as sample sizes increase. Benchmarking on the depletion of variation seen in coding regions, we are currently well-powered to detect extreme non-coding constraint as strong as the 90th percentile of coding exons of similar size, and we estimate a sample size of ~340K genomes to detect constraint as to the 50th percentile (Extended Fig. 8a; Methods). Much larger sample sizes will be required for further increasing the resolution, for instance from 1kb to a 100bp scale, we would need ~5.3M samples (Extended Fig. 8b); under the current sample size, 1kb presented optimal performance when compared to a various window size tested from 100bp-3kb (Extended Fig. 8c). Meanwhile, we emphasize the importance of increasing ancestral diversity in population-scale datasets like gnomAD. A more diverse population would identify a larger number of rare variants, thereby increasing the power of detecting depletions of variation. We explicitly demonstrated this by reconstructing our constraint metric from the subset of European population and comparing it to that from an equal-sized subset containing all diverse populations – the latter was proven to achieve a higher predictive power (Extended Fig. 8d). Future efforts towards a larger, more diverse human reference dataset would empower finer studies of the influence of human demography on constraint metrics, facilitating a fuller understanding of the distribution and effect of human genetic variation.

Overall, our study demonstrates the value of the genome-wide constraint map in characterizing both non-coding regions and protein-coding genes, providing a significant step towards a comprehensive catalog of functional genomic elements for humans.

Methods

Aggregation, variant-calling, and quality control of gnomAD genome data

We aggregated whole genome sequence data from 153,030 individuals spanning projects from case-control consortia and population cohorts, in a similar fashion to previous efforts¹. We harmonized these data using the GATK Best Practices pipeline and joint-called all samples using Hail⁶³, and developed and utilized an updated pipeline of sample, variant, and genotype quality control to create a high-quality callset of 76,156 individuals, computing frequency information for several strata of this dataset based on attributes such as ancestry and sex for each of 644,267,978 short nuclear variants (see Supplementary Information).

Estimation of trinucleotide context-specific mutation rates

We estimated the probability of a given nucleotide mutating to one of the three other possible bases in a trinucleotide context ($XY_1Z \rightarrow XY_2Z$), by computing the proportion of all possible variants observed per context in the human genome. Since CpG transitions begin to saturate (proportion observed approaching 1) at a sample size of ~10K genomes, we downsampled the gnomAD dataset to 1,000 genomes for this calculation. The computed proportion observed values, which represent the relative mutability of each trinucleotide context, were further scaled so that the weighted genome-wide average is the human per-base, per-generation mutation rate (1.2×10^{-8}) to obtain the absolute mutation rates μ . To estimate the proportion of variants expected to be observed in the full gnomAD dataset of 76,156 genomes, we fitted the actual proportion observed in the dataset against μ , using an exponential regression that caps at 1 for refining the estimates of (near-)saturated variant types ($R^2=0.999$, Extended Fig. 1a,b; Supplementary Data 1).

A total of 390,393,900 high-quality, rare ($AF \leq 0.1\%$) variants observed in 76,156 gnomAD genomes, a dataset of 6,079,733,538 possible variants at 2,026,577,846 autosomal sites (30-32X coverage), were used in the calculation of trinucleotide context-specific mutation rates. The estimates are well-correlated with the mutation rates reported in previous independent studies and are highly stable across different AF thresholds in gnomAD (Supplementary Information).

Adjustment of the effect of DNA methylation on CpG mutation rates

Given the strong effect of DNA methylation on increasing the mutation rate at CpG sites, we stratified all CpG sites by their methylation levels and computed the proportion observed within each context and methylation level. As an improvement to our previous methylation annotation (by averaging different tissues¹), we analyzed methylation data from germ cells across 14 developmental stages, comprising eight from preimplantation embryos (sperm, oocyte, pronucleus, two-cell-, four-cell-, eight-cell-, morula-, and blastocyst-stage embryos)⁶⁴ and six from primordial germ cells (7Wk, 10Wk, 11Wk, 13Wk, 17Wk, and 19Wk)⁶⁵. For each stage, we computed methylation level at each CpG site as the proportion of whole-genome bisulfite sequencing reads corresponding to the methylated allele. To derive a composite score from the 14 stages, we regressed the observation of a CpG variant in gnomAD (0 or 1) on the methylation computed at the corresponding site (a vector of 14), and we used the coefficients from the regression model as weights to compute a composite methylation score for each CpG site. This metric was further discretized into 16 levels (by a minimum step of 0.05: [0,0.05], (0.05,0.1], (0.1,0.15], (0.15,0.2], (0.2,0.25], (0.25,0.3], (0.3,0.35], (0.35,0.4], (0.4,0.45], (0.45,0.5], (0.5,0.55], (0.55,0.6], (0.6,0.65], (0.65,0.7], (0.7,0.75], (0.75,0.8], (0.8,0.85], (0.85,0.9], (0.9,1.0]) to stratify CpG variants in the mutation rate analysis.

Adjustment of the effects of regional genomic features on mutation rates

To estimate the effects of regional genomic features on mutation rates under neutrality, we utilized *de novo* mutations (DNMs), as a proxy of spontaneous mutations, and fitted logistic regression models using the genomic features as predictive variables. A set of 413,304 unique DNMs were compiled from two large-scale family-based whole-genome sequencing studies^{22,23}, and an exclusive set of 4,104,879 genomic sites (~10× the DNMs) randomly drew from the genome was used as the “nonmutated” background. For each DNM or background site, we computed 13 genomic features (see Collection of genomic features) at four scales by taking the mean value of 1kb, 10kb, 100kb, and 1Mb windows centering at the site. This generated a feature matrix of 13×4=52 columns and 413,304+4,104,879=4,518,183 rows. The matrix was further divided based on the trinucleotide context of each DNM or background site (by row) to assess the effects of genomic features on context-specific mutation rates. In particular, for CpG contexts, features that were correlated with DNA methylation (GC content, CpG_island, short interspersed nuclear element, and nucleosome density), which had been used for adjusting CpG mutation rates, were excluded from the analysis.

For each trinucleotide context, we first performed univariable logistic regression to select features that are significantly associated with an increased/decreased probability of observing a DNM. Features with a significant association surpassing the Bonferroni correction for 13×4=52 tests were selected; if a feature was significant at multiple genomic scales, the smallest window size was selected for the highest resolution (Extended Fig. 1c). Next, we fitted multivariable logistic regression using the selected features to predict DNMs from the background. To control for multicollinearity, we transformed the input feature matrix using principal components analysis (PCA⁶⁶) to generate decorrelated predictive variables (i.e., the principal components or PCs). The regression coefficients were the primary output of interest, which represent the effects of genomic features on increasing (a positive coefficient) or decreasing (a negative coefficient) the mutation rate, and were used for adjusting the expected number of variants in a given region. The selected features, the PCs, and the coefficients are summarized in Extended Fig. 1c and are available as pickle files for implementation (see Code availability in Supplementary Information).

Prediction of expected number of variants per 1kb

Using the trinucleotide mutation rate estimates and the above adjustments, we computed the expected number of variants in a given 1kb genomic window as follow:

$$Exp(w) = \sum_i^{64} r(w)_i \sum_{j=1}^3 \sum_{m=1}^k n(w)_{i,j,m} \times p_{i,j,m}$$

where i denotes one of the 64 trinucleotide contexts; j denotes one of the three bases substituting the central nucleotide; m denotes one of the k DNA methylation levels, where $k=16$ for CpG sites (see Adjustment of the effect of DNA methylation on CpG mutation rates) and $k=1$ for non-CpG sites (i.e., no stratification). Essentially, the expected value of variants in a genomic window w is calculated by multiplying the number of possible variants (n) in w by the probability of a variant (p) and summing across all trinucleotide contexts (i), substitutions (j), and methylation levels (m); $p_{i,j,m}$ is the trinucleotide mutation rate estimated in this study (as described in Estimation of trinucleotide context-specific mutation rates).

Additionally, Exp is adjusted by a factor r , which represents the effect of regional genomic features of w on mutation rate. For each i , specific features have been pre-selected and their effects on mutation rate have been estimated using logistic regression models (see Adjustment of the effects of regional genomic features on mutation rates). Denote the feature values, computed centering w and decorrelated by PCA, and the regression coefficients by $\mathbf{x} = \{x_1, x_2, \dots, x_t\}$ and $\boldsymbol{\beta} = \{\beta_1, \beta_2, \dots, \beta_t\}$, respectively, where t is the number of selected features for i , the adjustment factor r is defined as the ratio of logit given $\mathbf{x}(w)$ to

that of the genome-wide average \bar{x} : $r = \beta \cdot x(w) / \beta \cdot \bar{x}$; since the adjustment is specific to each trinucleotide context, r is further subscribed by i .

Construction of constraint Z score

We created a signed Z score to quantify the depletion of variation (constraint) at a 1kb scale by comparing the observed variation to an expectation:

$$\chi^2 = (Obs - Exp)^2 / Exp$$

$$Z = \begin{cases} \sqrt{\chi^2} & \text{if } Obs < Exp \\ -\sqrt{\chi^2} & \text{if } Obs \geq Exp \end{cases}$$

The observed variant count (*Obs*) is the number of unique rare ($AF \leq 0.1\%$) variants in a 1kb window identified in the gnomAD dataset of 76,156 genomes, and the expected number of variants (*Exp*) is established as described above based on the sequence context and the regional genomic features of the 1kb window.

Constraint Z scores were created for 2,689,987 non-overlapping 1kb windows across the human genome, comprising 2,561,056 on autosomes and 128,931 on chromosome X. Due to the lack of DNM data on chromosome X, the genomic feature adjustment factor r was assessed using autosomal regions and extrapolated to chromosome X. We performed downstream analyses separately for autosomes and chromosome X and presented the former as primary, with the latter provided in Supplementary Information. For the analyses, we filtered the dataset to windows where 1) the sites contained at least 1,000 possible variants, 2) at least 80% of the observed variants passed all variant call filters (INFO/FILTER equals to "PASS"), and 3) the mean coverage in the gnomAD genomes was between 25-35X (or 20-25X for chromosome X). This resulted in 1,984,900 autosomal windows (77.5% of initial) for the primary analyses, of which 141,341 overlapped with coding regions and 1,843,559 were exclusively non-coding. The computed constraint Z scores are available in Supplementary Data 2. We also computed the scores in a sliding window approach (1kb stepped by 100bp) and provided them in Supplementary Data 3.

Collection of genomic features

The 13 regional genomic features used for adjusting trinucleotide mutation rate are 1) GC content⁶⁷, 2) low-complexity region⁶⁸, 3) short and 4) long interspersed nuclear element⁶⁷, distance from the 5) telomere and the 6) centromere⁶⁷, 7) male and 8) female recombination rate²², 9) DNA methylation, 10) CpG island⁶⁷, 11) nucleosome density⁶⁹, 12) maternal and 13) paternal DNM cluster⁷⁰. Data were downloaded from the referenced resources, lifted over to GRCh38 coordinates when needed using CrossMap⁷¹, and files in .bed or .BigWig format were processed using bedtools⁷² and bigWigAverageOverBed⁶⁹ to obtain feature values within specific genomic windows.

Correlation between constraint Z score and APS

As an internal validation, we compared our constraint Z score against the SV constraint score APS²⁵. For each SV from the original study²⁵, we assessed its constraint by assigning the highest Z score among all overlapping 1kb windows. The correlation between constraint Z and APS was evaluated across 116,184 high-quality autosomal SVs scored by both metrics, using a linear regression test. In Fig. 1b, the correlation was presented by the mean value of APS across ascending constraint Z score bins, with 95% confidence intervals computed from 100-fold bootstrapping.

Correlation between constraint Z score and putative functional non-coding annotations

We validated the constraint metric using a number of external functional annotations, including 926,535 ENCODE cCREs²⁶ (34,803 promoter-like [PLS], 141,830 proximal enhancer-like [pELS], 667,599 distal enhancer-like [dELS], and 56,766 CTCF-only elements), 63,285 FANTOM5²⁷ enhancers, 331,601 super enhancers (SEdb²⁸), 111,308 GWAS Catalog³¹ variants (with an association $P \leq 5.0 \times 10^{-8}$; 9,229 with an independent replication), 2,191 GWAS variants fine-mapped across population biobanks with a posterior inclusion probability of causality ≥ 0.9 ³², and 100,530 CNVs from a CNV morbidity map of developmental delay^{50,51}.

To assess the correlation between constraint Z score and the collected functional elements, we intersected each annotation with the scored 1kb windows binned by constraint Z score (<-4 , $[-4,-3]$, $[-3,-2]$, $[-2,-1]$, $[-1,0]$, $[0,1]$, $[1,2]$, $[2,3]$, $[3,4]$, ≥ 4), and counted the frequency of overlapping windows within each bin. The enrichment of a given annotation (except CNVs) at a constraint level was evaluated by comparing the corresponding frequency to the genome-wide average using a Fisher's exact test. In the analysis of CNVs, we assessed their enrichment in constrained regions by assigning each CNV the highest Z score among its overlapping windows and comparing the proportions of constrained CNVs ($Z \geq 4$) from cases of developmental delay and healthy controls. The enrichment was further examined using a logistic regression model to adjust for the size and gene content (gene constraint¹ and gene number) of each CNV. We note that we performed all above analyses restricting to exclusively non-coding windows to evaluate the use of our constraint metric in characterizing the non-coding genome.

Estimation of constraint for aggregated regulatory annotations

We estimated how constrained the sequences encoding regulatory elements overall compared to coding exons by aggregating the regulatory annotations at a 1kb scale. These included 7,246 promoter-, 154,003 enhancer-, 117 microRNA (miRNA)-, and 414,084 long non-coding RNA (lncRNA)-1kb elements, created from concatenating ENCODE cCREs-PLS, cCREs-dELS, GENCODE⁷³ miRNA, and FANTOM6⁷⁴ lncRNA annotations, respectively, into 1kb windows. Similarly, 27,875 exonic 1kb elements were created from aggregating all protein-coding exons. Constraint Z scores were computed for the created 1kb elements and the percentiles of each regulatory annotation were compared against the exonic region. Benchmarking on the 50th percentile (median) of exonic regions, we estimated the proportion of the regulatory elements that are under selection as strong as the coding exons.

Incorporation of constraint Z score into GWAS fine-mapping

To demonstrate the use of our constraint metric in statistical fine-mapping, we performed approximate functionally informed fine-mapping³⁹ incorporating constraint Z score and our previous fine-mapping results for 119 UK Biobank (UKBB) traits³². The constraint Z scores were normalized and used as functional prior probabilities to update the posterior inclusion probabilities (PIPs; denoted as PIP_z) based on the previous UKBB fine-mapping (using a uniform prior, PIP_{unif}) and SuSiE⁷⁵. To exclude signals that potentially correspond to coding variants, we restricted our analysis to 60,121 non-coding variants in 6,592 SuSiE 95% credible set (CS)-trait pairs that do not contain variants within 1 kb of exonic regions. A total of 13,069 variant-trait pairs were predicted to have an increased PIP ($\Delta PIP \geq 0.01$) of causality. The variants, associated traits, and PIP scores (PIP_{unif} and PIP_z) are provided in Supplementary Data 4.

Comparison of constraint Z score and other predictive scores

We compared our constraint metric with other seven genome-wide predictive scores – Orion¹³, CDTs¹⁴, gwrVIS¹⁸, DR¹⁵, phyloP²⁰, phastCons¹⁹, and GERP⁴¹. Each score was downloaded from the original study, lifted over to GRCh38 coordinates (for Orion) and multiplied by -1 (for CDTs, gwrVIS, and DR) when needed so that a higher value represents a higher constraint/conservation for all metrics. Pairwise

correlation between the scores was assessed by comparing the mean value of each score on 1kb windows, using a Spearman's rank correlation test.

We evaluated the predictive performance of each metric in distinguishing functional non-coding variants ("positive" variant set) from background variants ("negative" variant set). Four positive variant sets were compiled from public databases: 1) 9,229 variants from GWAS Catalog³¹ (with an independent replication), 2) 2,191 variants from a recent fine-mapping study³² (with a posterior inclusion probability of causality ≥ 0.9), 3) 140 high-confidence variants from 2), and 4) 288 variants from ClinVar⁴⁰ (annotated as "pathogenic"). All variants were filtered to non-coding regions; in particular, ClinVar variants were more strictly filtered to intergenic/intron variants given its strong predominance of variants close to protein-coding exons (>90% were splice site/region variants). A further stringent non-coding subset was generated by excluding variants within 10kb to any exons, which resulted in 1) 4,379, 2) 967, 3) 59, and 4) 7 variants. For each positive variant set, a negative variant set was created by randomly drawing variants from gnomAD (to $\sim 10\times$ the size of corresponding positive variant set), of which the most severe molecular consequence is intergenic or intron and the AF approximates the positive variant set; AF > 5% and allele count (AC) = 1 were applied respectively for matching positive variant set 1)-3) and 4), based on their AF distributions in gnomAD (Fig. 3b). The selected variants were scored by each of the eight metrics, using bedtools⁷² (for .bed files) and bigWigAverageOverBed⁶⁹ (for .BigWig files), and the performance of each metric in classifying positive and negative variants was assessed by the area under curve (AUC) statistic, as presented by the receiver operating characteristic (ROC) curve.

To investigate whether different metrics capture complementary information in the classification, we fitted logistic regression models using all eight metrics as independent variables. The relative contribution of each metric was evaluated by the dominance analysis^{76,77}, which estimates the dominance of one predictor over another by comparing their additional R^2 contributions across all subset models. We further explored whether specific features were particularly captured by (and may have contributed to the performance of) our metric. We merged all positive variant sets and focused on a set of variants (N=204) that were uniquely prioritized by our metric, defined as being captured in the 99th percentile of constraint Z score but not in that of any other scores. Specific features associated with these variants were evaluated by comparing values of the 13 genomic features of these variants to the rest of the positive variant set. The fold change was used to indicate the extent to which a feature is distinguished in variants captured by constraint Z score from others.

Correlation of constraint between non-coding regulatory elements and protein-coding genes

To examine whether constraint of non-coding regulatory elements informs the constraint of their target genes, we compared constraint Z scores of enhancers linked to constrained genes and unconstrained genes. The former included well-established gene sets of 189 ClinGen⁵⁶ haploinsufficient genes, 2,454 MGI⁵⁷ essential genes mapped to human orthologs, 1,771 OMIM⁵⁸ autosomal dominant genes, and 1,920 LOEUF¹ first-decile genes; and the latter included a curated list of 356 olfactory receptor genes and 189 LOEUF last-decile genes with at least 10 expected LoF variants (which are sufficiently powered to be classified into the most constrained decile¹). The LOEUF underpowered list included 1,117 genes with ≤ 5 expected LoF variants. Enhancers linked to each gene were obtained from the Roadmap Epigenomics Enhancer-Gene Linking database, which used correlated patterns of activity between histone modifications and gene expression to predict enhancer-gene links^{78,79}. For each gene, we aggregated and merged enhancers predicted from all 127 reference epigenomes and assigned the most constrained enhancer to each gene for the analysis of enhancer-gene constraint correlation (Supplementary Data 5).

In the analysis of correlation between tissue-specific enhancer constraint and tissue-specific gene expression, we processed the enhancer-gene links with the same principle as described above but within specific tissue types (as defined in the Roadmap Epigenomics metadata⁵⁵). For each gene and tissue type, we searched for tissue-specific gene expression in the Genotype-Tissue Expression (GTEx⁵⁹) database (RNASeQCv1.1.9) and computed a normalized median expression for each gene ($\log_2(\text{TPM}+1)$). Enhancer constraint and gene expression values were calculated for 11 matched tissue types, and the correlation within each tissue type was evaluated by regressing gene expression on enhancer constraint, including gene constraint (LOEUF score) as a covariate.

Incorporation of non-coding constraint of regulatory elements into gene constraint modeling

To demonstrate the practical value of non-coding constraint in improving gene constraint modeling, we compared two models – using 1) LOEUF and 2) LOEUF+enhancer constraint Z score (as described in Correlation of constraint between non-coding regulatory elements and protein-coding genes) – in predicting constrained genes, with a particular focus on genes that were underpowered in LOEUF. A set of 3,220 unique constrained genes were curated from ClinGen⁵⁶, MGI⁵⁷, and OMIM⁵⁸ (see Correlation of constraint between non-coding regulatory elements and protein-coding genes), and a set of 356 olfactory receptor genes was used as the unconstrained genes. We trained logistic regression models on 50% of the genes and tested the performance on 77 underpowered genes in the remaining 50%. The predictive performance of the two models were measured by AUC, and the significance of the difference in AUCs was assessed using a bootstrap test⁸⁰.

Power of constraint detection

We estimated the power of our metric in detecting non-coding constraint as the percentage of the non-coding genome to obtain a high constraint Z score ($Z \geq 4$) under a certain strength of negative selection, which was quantified by the level of depletion of variation (i.e., 1-observed/expected). For a given depletion of variation, the minimum number of expected variants to achieve a $Z \geq 4$ was determined, and the number of samples required to achieve the expected number of variants was estimated using a linear model of $\log(\text{number of expected variants}) \sim \log(\text{number of samples})$ from downsampling the gnomAD dataset. The power was estimated at two scales – 1kb (used in this study) and 100bp – and benchmarked by the depletion of variation observed in coding exons of similar size.

Figure legends

Fig. 1: Distribution of constraint Z scores across the genome. **a**, Histograms of constraint Z scores for 1,984,900 1kb windows across the human autosomes. Windows overlapping coding regions (N=141,341 with ≥ 1 bp coding sequence; red) overall exhibit a higher constraint Z (stronger negative selection) than windows that are exclusively non-coding (N=1,843,559; blue); dashed lines indicate the medians. **b**, The correlation between constraint Z score and the adjusted proportion of singletons (APS) score developed for structural variation (SV) constraint. A collection of 116,184 autosomal SVs were assessed using constraint Z score by assigning each SV the highest Z among all overlapping 1kb windows, which shows a significant positive correlation with the SV constraint metric APS. Error bars indicate 100-fold bootstrapped 95% confidence intervals of the mean values.

Fig. 2: Correlation between constraint Z score and functional non-coding annotations. **a,b**, Distributions of candidate regulatory elements (**a**) and GWAS variants (**b**) along the spectrum of non-coding constraint. Enrichment was evaluated by comparing the proportion of non-coding 1kb windows, binned by constraint Z, that overlap with a given functional annotation to the genome-wide average. Error bars indicate 95% confidence intervals of the odds ratios. **c**, Enrichment of fine-mapped variants in constrained non-coding regions ($Z \geq 4$). Credible set (CS)-trait pairs with a significant enrichment are shown, ordered by the lower bound of 95% confidence interval; only lower bounds are shown for presentation purposes. **d**, The distribution of variants fine-mapped for coronary artery disease (CAD) in constrained regions ($Z \geq 4$) of *PLG*. Each bar shows the constraint Z score of a 1kb window (gaps indicate windows removed by quality filters); windows containing fine-mapped variants are colored by purple, and the number of variants in each window is annotated on top of the bar correspondingly. Ten variants are located within *PLG* introns, four are mapped to the antisense gene of *PLG* (ENSG00000287558), and 14 reside in the downstream intergenic regions.

Fig. 3: Performance of constraint Z score and other predictive scores in prioritizing non-coding variants. **a,b**, Receiver operating characteristic (ROC) curves of constraint Z score and other seven metrics in classifying putative functional non-coding variants – 2,191 GWAS fine-mapping variants (**a**) and 288 ClinVar pathogenic variants (**b**) – against background variants in the population. The performance of each metric was measured and ranked by the area under curve (AUC) statistic. **c,d**, The relative contribution of different metrics in classifying GWAS variants (**b**) and ClinVar variants (**c**). The eight metrics were modeled as eight independent predictors for the classification, and the relative contribution of one predictor over another was evaluated by estimating their additional R^2 contributions across all subset models.

Fig. 4: Contribution of non-coding constraint in evaluating copy number variants (CNVs). **a**, Proportions of constrained CNVs ($Z \geq 4$) identified in individuals with developmental delay (DD cases) versus healthy controls. Constrained CNVs are more common in DD cases than controls (7,239/17,004=42.6% versus 10,403/83,526=12.5%) and are most frequent for CNVs previously implicated as pathogenic (18/19=94.7% by DD and 3,433/4,014=85.5% by ClinVar). **b**, Contribution of non-coding constraint to predicting CNVs in DD cases versus controls. Non-coding constraint remains a significant predictor for the case/control status of CNVs after adjusting for gene constraint (LOEUF score), gene number, and size of CNVs (purple), as well as being tested in the subset of non-coding CNVs (blue). Error bars indicate 95% confidence intervals of the log odds ratios. **c**, CNVs at the *IHH* locus associated with synpolydactyly and craniosynostosis. The four implicated duplications (grey horizontal bars) span a ~102kb sequence upstream of *IHH*. Each vertical bar shows the constraint Z score of a 1kb window within the locus, with the highest score overlapping the *IHH* gene (red) and the highest non-coding score overlapping the major *IHH* enhancers (purple); gaps indicate

windows removed by quality filters. **d**, Non-coding CNVs with the highest constraint Z score identified in DD cases. The highest-scored window is located within the potential “critical region”⁴⁸ (purple vertical bars) shared by 12 DD deletions (red horizontal bars; grey indicates two deletions observed in controls). The critical region overall, has a significantly higher constraint Z score than the other regions affected by DD or control deletions, as shown in the kernel density estimate (KDE) plot on the right.

Fig. 5: Correlation of constraint between non-coding regulatory elements and protein-coding genes. **a**, The proportion of non-coding 1kb windows overlapping with enhancers that were predicted to regulate specific genes, as a function of their constraint Z scores. More constrained non-coding regions are more frequently linked to a gene. Error bars indicate standard errors of the proportions. **b**, Comparison of the constraint Z scores of enhancers linked to constrained and unconstrained genes. Enhancers of established sets of constrained genes (four blue boxes) are more constrained than enhancers of presumably less constrained genes (two grey boxes). Enhancers of genes that are underpowered for gene constraint detection (“LOEUF underpowered”) present a higher constraint than those powered yet unconstrained genes (“LOEUF unconstrained”). **c**, Improvement of incorporating enhancer constraint into LOEUF in prioritizing underpowered genes. ROC curves and AUCs show the performance of two models using LOEUF (blue) and LOEUF+enhancer constraint Z score (purple) to classify constrained and unconstrained genes, tested on a set of 77 underpowered genes. **d**, Contribution of enhancer constraint to predicting gene expression in specific tissue types. The x-axis shows the linear regression coefficient of tissue-specific enhancer constraint predicting the expression level of target genes in matched tissue types, conditioning on gene constraint (LOEUF score). Error bars indicate 95% confidence intervals of the coefficient estimates.

Extended Data Fig. 1: Construction of mutational model and constraint Z score. **a,b**, Estimation of trinucleotide context-specific mutation rates. The proportion of possible variants observed for each substitution and context in 76,156 gnomAD genomes (y-axis) is exponentially correlated with the absolute mutation rate estimated from 1,000 downsampled genomes (x-axis). Fit lines were modeled separately for human autosomes (**a**) and chromosome X (**b**). **c**, Estimation of the effects of regional genomic features on mutation rates. The effects of 13 genomic features at four scales (window sizes 1kb-1Mb; x-axis) on the mutation rate of 32 trinucleotide contexts (y-axis) are shown, colored by the coefficient from regressing *de novo* mutations (DNMs) on each specific feature and window size. Red/Blue color indicates a positive/negative effect of increasing the feature value on mutation rates; grey crosses indicate significant features at the smallest possible window size after Bonferroni correction for $13 \times 4 = 52$ tests. Abbreviations: LCR=low-complexity region, SINE/LINE=short/long interspersed nuclear element, Dist=Distance, Recomb=Recombination, Methyl=Methylation. **d,e**, The distribution of constraint Z score as a function of expected and observed variation. Each point represents the Z score of a 1kb window on the genome (N=1,984,900 on autosomes (**d**) and N=57,729 on chromosome X (**e**)), which quantifies the deviation of observed variation from expectation. A positive Z score (red) indicates depletion of variation (observed<expected) and the higher the Z score the stronger the depletion; the red dashed line indicates the 99th percentile of Z scores across the autosomes (**d**) or chromosome X (**e**).

Extended Data Fig. 2: Comparison of constraint Z score between coding and non-coding regions. **a**, The proportion of highly constrained windows ($Z \geq 4$) as a function of the percentage of coding sequences in a window. The intervals (x-axis) are left exclusive and right inclusive. “Exonic only” refers to the 1kb windows created from directly concatenating coding exons into 1kb sequences. **b**, The exonic-only regions (N=27,875; purple) present a significantly higher constraint Z score than regions that are exclusively non-coding (N=1,843,559; blue). Dashed lines indicate the medians. **c**, The proportion of highly constrained

windows ($Z \geq 4$) as a function of the proportion of exonic windows being added to the dataset of non-coding windows. **d**, Constraint Z score percentiles of non-coding versus exonic windows. About 0.05% (100-99.95%) and 3.12% (100-96.88%) of the non-coding windows exhibit similar constraint to the 90th and 50th of exonic regions, respectively.

Extended Data Fig. 3: Estimation of constraint for aggregated regulatory annotations. **a,b**, Constraint Z scores of aggregated promoter (dark purple), enhancer (light purple), microRNA (miRNA; dark blue), and long non-coding RNA (lncRNA; light blue) annotations are compared against those of exonic (**a**) and non-coding (**b**) regions at a 1kb scale. The constraint Z score percentiles of each annotation (y-axis) are benchmarked by the score deciles of exonic or non-coding regions (10-100 percentiles; x-axis); the grey dashed vertical line indicates the median (50th percentile).

Extended Data Fig. 4: Applications of constraint Z score for characterizing non-coding regions in addition to existing functional annotations. **a**, Use of constraint Z score for prioritizing non-coding regions with or without a regulatory annotation. Constrained non-coding regions are enriched for GWAS variants, independent of the candidate cis-regulatory element (cCRE) annotation from ENCODE. **b**, Use of constraint Z score in statistical fine-mapping. The increase in posterior inclusion probability (PIP) when incorporating constraint Z score as a functional prior into previous fine-mapping results (that used a uniform prior; denoted as PIP_z and PIP_{unif} , respectively) is shown for 164 new likely causal associations with a $PIP_z \geq 0.8$ as a function of PIP_z .

Extended Data Fig. 5: Comparison of constraint Z score and other predictive scores in prioritizing non-coding variants. **a**, Receiver operating characteristic (ROC) curves of constraint Z score and other seven metrics in classifying putative functional non-coding variants (“positive” variant set) – left to right: 9,229 GWAS Catalog variants, 2,191 GWAS fine-mapping variants, a subset of 140 high-confidence fine-mapped variants, and 288 ClinVar pathogenic variants – against “negative” variant set randomly drew from the population with a similar allele frequency (AF). $AF > 5\%$ and allele count (AC)=1 were applied respectively for matching the three GWAS variant sets and the ClinVar variant set, based on their AF distributions in gnomAD (shown in **b**). **b**, AUCs of the classification with a varying AF threshold for the negative variant set. As most GWAS variants are common and most ClinVar variants are very rare (not seen in the population), $AF > 5\%$ and $AC=1$ were applied respectively in the primary analyses shown in **a**.

Extended Data Fig. 6: Comparison of constraint Z scores built from different mutational models and genomic windows. Our constraint Z score (presented in this study) outperforms the scores rebuilt from mutational models that only consider local sequence context – trinucleotide (trimer-only) or heptanucleotide (heptamer-only) – without adjustment on mutation rate by regional genomic features, and the performance is robust to the artificial break of genomic windows when computed at a 1kb sliding by 100bp scale.

Extended Data Fig. 7: Pairwise correlations between different constraint/conservation metrics. The Spearman's rank correlation between each pair of the eight metrics was computed based on the mean value of each score on 1kb windows across the genome.

Extended Data Fig. 8: Power of constraint detection. **a,b**, The sample size required for well-powered non-coding constraint detection. The percentage of non-coding regions powered to detect constraint ($Z \geq 4$) at a 1kb (**a**) and 100bp (**b**) scale under varying levels of selection (depletion of variation) is shown as a function of log-scaled sample size. Lighter color indicates milder deletion of variation (weaker selection), which requires a larger sample size to detect constraint; the grey dashed vertical line indicates the current

sample size of 76,156 genomes. Dotted curves (left to right) benchmark the 95th, 90th, and 50th percentile of depletion of variation observed in coding exons of similar size. The number of samples required to obtain an 80% detection power is labeled at corresponding benchmarks. **c**, AUCs of constraint Z scores computed on different window sizes in identifying putative functional non-coding variants. 1kb (used in this study) presents the optimal window size with high performance while maintaining reasonable resolution. **d**, AUCs of constraint Z scores computed from different subsets of gnomAD in identifying putative functional non-coding variants. While with an equal sample size, the downsampled dataset with diverse ancestries presents higher performance than the Non-Finnish European (NFE)-only dataset.

References

- 1 Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434-443, doi:10.1038/s41586-020-2308-7 (2020).
- 2 Short, P. J. *et al.* De novo mutations in regulatory elements in neurodevelopmental disorders. *Nature* **555**, 611-616, doi:10.1038/nature25983 (2018).
- 3 Satterstrom, F. K. *et al.* Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism. *Cell* **180**, 568-584 e523, doi:10.1016/j.cell.2019.12.036 (2020).
- 4 Singh, T. *et al.* The contribution of rare variants to risk of schizophrenia in individuals with and without intellectual disability. *Nat Genet* **49**, 1167-1173, doi:10.1038/ng.3903 (2017).
- 5 Ganna, A. *et al.* Quantifying the Impact of Rare and Ultra-rare Coding Variation across the Phenotypic Spectrum. *Am J Hum Genet* **102**, 1204-1211, doi:10.1016/j.ajhg.2018.05.002 (2018).
- 6 Hindorff, L. A. *et al.* Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc Natl Acad Sci U S A* **106**, 9362-9367, doi:10.1073/pnas.0903103106 (2009).
- 7 Lanyi, J. K. Photochromism of halorhodopsin. cis/trans isomerization of the retinal around the 13-14 double bond. *J Biol Chem* **261**, 14025-14030 (1986).
- 8 Mathelier, A., Shi, W. & Wasserman, W. W. Identification of altered cis-regulatory elements in human disease. *Trends Genet* **31**, 67-76, doi:10.1016/j.tig.2014.12.003 (2015).
- 9 Spielmann, M. & Mundlos, S. Looking beyond the genes: the role of non-coding variants in human disease. *Hum Mol Genet* **25**, R157-R165, doi:10.1093/hmg/ddw205 (2016).
- 10 Zhang, F. & Lupski, J. R. Non-coding genetic variants in human disease. *Hum Mol Genet* **24**, R102-110, doi:10.1093/hmg/ddv259 (2015).
- 11 Seplyarskiy, V. B. & Sunyaev, S. The origin of human mutation in light of genomic data. *Nat Rev Genet* **22**, 672-686, doi:10.1038/s41576-021-00376-2 (2021).
- 12 Seplyarskiy, V. B. *et al.* Population sequencing data reveal a compendium of mutational processes in the human germ line. *Science* **373**, 1030-1035, doi:10.1126/science.aba7408 (2021).
- 13 Gussow, A. B. *et al.* Orion: Detecting regions of the human non-coding genome that are intolerant to variation using population genetics. *PLoS One* **12**, e0181604, doi:10.1371/journal.pone.0181604 (2017).
- 14 di Iulio, J. *et al.* The human noncoding genome defined by genetic diversity. *Nat Genet* **50**, 333-337, doi:10.1038/s41588-018-0062-7 (2018).
- 15 Halldorsson, B. V. *et al.* The sequences of 150,119 genomes in the UK Biobank. *Nature* **607**, 732-740, doi:10.1038/s41586-022-04965-x (2022).
- 16 Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat Genet* **46**, 310-315, doi:10.1038/ng.2892 (2014).
- 17 Yousefian-Jazi, A., Jung, J., Choi, J. K. & Choi, J. Functional annotation of noncoding causal variants in autoimmune diseases. *Genomics* **112**, 1208-1213, doi:10.1016/j.ygeno.2019.07.006 (2020).
- 18 Vitsios, D., Dhindsa, R. S., Middleton, L., Gussow, A. B. & Petrovski, S. Prioritizing non-coding regions based on human genomic constraint and sequence context with deep learning. *Nat Commun* **12**, 1504, doi:10.1038/s41467-021-21790-4 (2021).
- 19 Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res* **15**, 1034-1050, doi:10.1101/gr.3715005 (2005).
- 20 Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res* **20**, 110-121, doi:10.1101/gr.097857.109 (2010).
- 21 Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285-291, doi:10.1038/nature19057 (2016).
- 22 Halldorsson, B. V. *et al.* Characterizing mutagenic effects of recombination through a sequence-level genetic map. *Science* **363**, doi:10.1126/science.aau1043 (2019).
- 23 An, J. Y. *et al.* Genome-wide de novo risk score implicates promoter variation in autism spectrum disorder. *Science* **362**, doi:10.1126/science.aat6576 (2018).

- 24 Samocha, K. E. *et al.* A framework for the interpretation of de novo mutation in human disease. *Nat Genet* **46**, 944-950, doi:10.1038/ng.3050 (2014).
- 25 Collins, R. L. *et al.* A structural variation reference for medical and population genetics. *Nature* **581**, 444-451, doi:10.1038/s41586-020-2287-8 (2020).
- 26 Consortium, E. P. *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699-710, doi:10.1038/s41586-020-2493-4 (2020).
- 27 Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455-461, doi:10.1038/nature12787 (2014).
- 28 Jiang, Y. *et al.* SEdb: a comprehensive human super-enhancer database. *Nucleic Acids Res* **47**, D235-D243, doi:10.1093/nar/gky1025 (2019).
- 29 Pott, S. & Lieb, J. D. What are super-enhancers? *Nat Genet* **47**, 8-12, doi:10.1038/ng.3167 (2015).
- 30 Bartel, D. P. Metazoan MicroRNAs. *Cell* **173**, 20-51, doi:10.1016/j.cell.2018.03.006 (2018).
- 31 Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res* **42**, D1001-1006, doi:10.1093/nar/gkt1229 (2014).
- 32 Kanai, M. *et al.* Insights from complex trait fine-mapping across diverse populations. *medRxiv*, 2021.2009.2003.21262975, doi:10.1101/2021.09.03.21262975 (2021).
- 33 Jung, R. G. *et al.* Association between plasminogen activator inhibitor-1 and cardiovascular events: a systematic review and meta-analysis. *Thromb J* **16**, 12, doi:10.1186/s12959-018-0166-4 (2018).
- 34 Song, C., Burgess, S., Eicher, J. D., O'Donnell, C. J. & Johnson, A. D. Causal Effect of Plasminogen Activator Inhibitor Type 1 on Coronary Heart Disease. *J Am Heart Assoc* **6**, doi:10.1161/JAHA.116.004918 (2017).
- 35 Schaefer, A. S. *et al.* Genetic evidence for PLASMINOGEN as a shared genetic risk factor of coronary artery disease and periodontitis. *Circ Cardiovasc Genet* **8**, 159-167, doi:10.1161/CIRCGENETICS.114.000554 (2015).
- 36 Li, Y. Y. Plasminogen activator inhibitor-1 4G/5G gene polymorphism and coronary artery disease in the Chinese Han population: a meta-analysis. *PLoS One* **7**, e33511, doi:10.1371/journal.pone.0033511 (2012).
- 37 Drinane, M. C., Sherman, J. A., Hall, A. E., Simons, M. & Mulligan-Kehoe, M. J. Plasminogen and plasmin activity in patients with coronary artery disease. *J Thromb Haemost* **4**, 1288-1295, doi:10.1111/j.1538-7836.2006.01979.x (2006).
- 38 Lowe, G. D. *et al.* Tissue plasminogen activator antigen and coronary heart disease. Prospective study and meta-analysis. *Eur Heart J* **25**, 252-259, doi:10.1016/j.ehj.2003.11.004 (2004).
- 39 Wang, Q. S. *et al.* Leveraging supervised learning for functionally informed fine-mapping of cis-eQTLs identifies an additional 20,913 putative causal eQTLs. *Nat Commun* **12**, 3394, doi:10.1038/s41467-021-23134-8 (2021).
- 40 Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res* **46**, D1062-D1067, doi:10.1093/nar/gkx1153 (2018).
- 41 Davydov, E. V. *et al.* Identifying a high fraction of the human genome to be under selective constraint using GERP++. *PLoS Comput Biol* **6**, e1001025, doi:10.1371/journal.pcbi.1001025 (2010).
- 42 Greenway, S. C. *et al.* De novo copy number variants identify new genes and loci in isolated sporadic tetralogy of Fallot. *Nat Genet* **41**, 931-935, doi:10.1038/ng.415 (2009).
- 43 Mefford, H. C. *et al.* Recurrent reciprocal genomic rearrangements of 17q12 are associated with renal disease, diabetes, and epilepsy. *Am J Hum Genet* **81**, 1057-1069, doi:10.1086/522591 (2007).
- 44 Sebat, J. *et al.* Strong association of de novo copy number mutations with autism. *Science* **316**, 445-449, doi:10.1126/science.1138659 (2007).
- 45 Stefansson, H. *et al.* Large recurrent microdeletions associated with schizophrenia. *Nature* **455**, 232-236, doi:10.1038/nature07229 (2008).
- 46 Walsh, T. *et al.* Rare structural variants disrupt multiple genes in neurodevelopmental pathways in schizophrenia. *Science* **320**, 539-543, doi:10.1126/science.1155174 (2008).
- 47 Wright, C. F. *et al.* Genetic diagnosis of developmental disorders in the DDD study: a scalable analysis of genome-wide research data. *Lancet* **385**, 1305-1314, doi:10.1016/S0140-6736(14)61705-0 (2015).
- 48 Spielmann, M., Lupianez, D. G. & Mundlos, S. Structural variation in the 3D genome. *Nat Rev Genet* **19**, 453-467, doi:10.1038/s41576-018-0007-0 (2018).
- 49 Spielmann, M. & Mundlos, S. Structural variations, the regulatory landscape of the genome and their alteration in human disease. *Bioessays* **35**, 533-543, doi:10.1002/bies.201200178 (2013).

- 50 Coe, B. P. *et al.* Refining analyses of copy number variation identifies specific genes associated with developmental delay. *Nat Genet* **46**, 1063-1071, doi:10.1038/ng.3092 (2014).
- 51 Cooper, G. M. *et al.* A copy number variation morbidity map of developmental delay. *Nat Genet* **43**, 838-846, doi:10.1038/ng.909 (2011).
- 52 Klopocki, E. *et al.* Copy-number variations involving the IHH locus are associated with syndactyly and craniosynostosis. *Am J Hum Genet* **88**, 70-75, doi:10.1016/j.ajhg.2010.11.006 (2011).
- 53 Barroso, E. *et al.* Identification of the fourth duplication of upstream IHH regulatory elements, in a family with craniosynostosis Philadelphia type, helps to define the phenotypic characterization of these regulatory elements. *Am J Med Genet A* **167A**, 902-906, doi:10.1002/ajmg.a.36811 (2015).
- 54 Will, A. J. *et al.* Composition and dosage of a multipartite enhancer cluster control developmental expression of *Ihh* (Indian hedgehog). *Nat Genet* **49**, 1539-1545, doi:10.1038/ng.3939 (2017).
- 55 Roadmap Epigenomics, C. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317-330, doi:10.1038/nature14248 (2015).
- 56 Rehm, H. L. *et al.* ClinGen--the Clinical Genome Resource. *N Engl J Med* **372**, 2235-2242, doi:10.1056/NEJMSr1406261 (2015).
- 57 Blake, J. A. *et al.* The Mouse Genome Database (MGD): premier model organism resource for mammalian genomics and genetics. *Nucleic Acids Res* **39**, D842-848, doi:10.1093/nar/gkq1008 (2011).
- 58 McKusick, V. A. Mendelian Inheritance in Man and its online version, OMIM. *Am J Hum Genet* **80**, 588-604, doi:10.1086/514346 (2007).
- 59 Consortium, G. T. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580-585, doi:10.1038/ng.2653 (2013).
- 60 Xu, H. *et al.* Elevated ASCL2 expression in breast cancer is associated with the poor prognosis of patients. *Am J Cancer Res* **7**, 955-961 (2017).
- 61 Jubb, A. M. *et al.* Achaete-scute like 2 (*ascl2*) is a target of Wnt signalling and is upregulated in intestinal neoplasia. *Oncogene* **25**, 3445-3457, doi:10.1038/sj.onc.1209382 (2006).
- 62 Tian, Y. *et al.* MicroRNA-200 (*miR-200*) cluster regulation by achaete scute-like 2 (*Ascl2*): impact on the epithelial-mesenchymal transition in colon cancer cells. *J Biol Chem* **289**, 36101-36115, doi:10.1074/jbc.M114.598383 (2014).
- 63 Hail v. 0.2.62-84fa81b9ea3d. <https://github.com/hail-is/hail/commit/84fa81b9ea3d>.
- 64 Zhu, P. *et al.* Single-cell DNA methylome sequencing of human preimplantation embryos. *Nat Genet* **50**, 12-19, doi:10.1038/s41588-017-0007-6 (2018).
- 65 Tang, W. W. *et al.* A Unique Gene Regulatory Network Resets the Human Germline Epigenome for Development. *Cell* **161**, 1453-1467, doi:10.1016/j.cell.2015.04.053 (2015).
- 66 Ross, D. A., Lim, J., Lin, R.-S. & Yang, M.-H. Incremental learning for robust visual tracking. *International journal of computer vision* **77**, 125-141 (2008).
- 67 Karolchik, D. *et al.* The UCSC Table Browser data retrieval tool. *Nucleic Acids Res* **32**, D493-496, doi:10.1093/nar/gkh103 (2004).
- 68 Li, H. Toward better understanding of artifacts in variant calling from high-coverage samples. *Bioinformatics* **30**, 2843-2851, doi:10.1093/bioinformatics/btu356 (2014).
- 69 Davis, C. A. *et al.* The Encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic Acids Res* **46**, D794-D801, doi:10.1093/nar/gkx1081 (2018).
- 70 Goldmann, J. M. *et al.* Germline de novo mutation clusters arise during oocyte aging in genomic regions with high double-strand-break incidence. *Nat Genet* **50**, 487-492, doi:10.1038/s41588-018-0071-6 (2018).
- 71 Zhao, H. *et al.* CrossMap: a versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* **30**, 1006-1007, doi:10.1093/bioinformatics/btt730 (2014).
- 72 Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841-842, doi:10.1093/bioinformatics/btq033 (2010).
- 73 Harrow, J. *et al.* GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res* **22**, 1760-1774, doi:10.1101/gr.135350.111 (2012).
- 74 Ramilowski, J. A. *et al.* Functional annotation of human long noncoding RNAs via molecular phenotyping. *Genome Res* **30**, 1060-1072, doi:10.1101/gr.254219.119 (2020).
- 75 Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in regression, with application to genetic fine-mapping. *BioRxiv*, 501114 (2020).

- 76 Budescu, D. V. Dominance analysis: a new approach to the problem of relative importance of predictors in multiple regression. *Psychological bulletin* **114**, 542 (1993).
- 77 Azen, R. & Budescu, D. V. The dominance analysis approach for comparing predictors in multiple regression. *Psychological methods* **8**, 129 (2003).
- 78 Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature* **473**, 43-49, doi:10.1038/nature09906 (2011).
- 79 Liu, Y., Sarkar, A., Kheradpour, P., Ernst, J. & Kellis, M. Evidence of reduced recombination rate in human regulatory domains. *Genome Biol* **18**, 193, doi:10.1186/s13059-017-1308-x (2017).
- 80 Robin, X. *et al.* pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC bioinformatics* **12**, 1-8 (2011).

Genome Aggregation Database Consortium

Maria Abreu¹, Carlos A. Aguilar Salinas², Tariq Ahmad³, Christine M. Albert^{4,5}, Jessica Alföldi^{6,7}, Diego Ardissino⁸, Irina M. Armean^{6,7,9}, Gil Atzmon^{10,11}, Eric Banks¹², John Barnard¹³, Samantha M. Baxter⁶, Laurent Beaugerie¹⁴, Emelia J. Benjamin^{15,16,17}, David Benjamin¹², Louis Bergelson¹², Michael Boehnke¹⁸, Lori L. Bonnycastle¹⁹, Erwin P. Bottinger²⁰, Donald W. Bowden^{21,22,23}, Matthew J. Bown^{24,25}, Steven Brant²⁶, Sarah E. Calvo^{6,27}, Hannia Campos^{28,29}, John C. Chambers^{30,31,32}, Juliana C.N. Chan^{33,109,110}, Katherine R. Chao⁶, Sinéad Chapman^{6,7,34}, Daniel Chasman^{4,35}, Siwei Chen^{6,7}, Rex L. Chisholm³⁶, Judy Cho²⁰, Rajiv Chowdhury³⁷, Mina K. Chung³⁸, Wendy K. Chung^{39,40,41}, Kristian Cibulskis¹², Bruce Cohen^{35,42}, Ryan L. Collins^{6,27,43}, Kristen M. Connolly⁴⁴, Adolfo Correa⁴⁵, Miguel Covarrubias¹², Beryl Cummings^{6,43}, Dana Dabelea⁴⁶, Mark J. Daly^{6,7,47}, John Danesh³⁷, Dawood Darbar⁴⁸, Joshua Denny⁴⁹, Stacey Donnelly⁶, Ravindranath Duggirala⁵⁰, Josée Dupuis^{51,52}, Patrick T. Ellinor^{6,53}, Roberto Elosua^{54,55,56}, James Emery¹², Eleina England⁶, Jeanette Erdmann^{57,58,59}, Tõnu Esko^{6,60}, Emily Evangelista⁶, Yossi Farjoun¹², Diane Fatkin^{61,62,63}, Steven Ferriera⁶⁴, Jose Florez^{35,65,66}, Laurent C. Francioli^{6,7}, Andre Franke⁶⁷, Martti Färkkilä⁶⁸, Stacey Gabriel⁶⁴, Kiran Garimella¹², Laura D. Gauthier¹², Jeff Gentry¹², Gad Getz^{35,69,70}, David C. Glahn^{71,72}, Benjamin Glaser⁷³, Stephen J. Glatt⁷⁴, David Goldstein^{75,76}, Clicerio Gonzalez⁷⁷, Julia K. Goodrich⁶, Leif Groop^{78,79}, Sanna Gudmundsson^{6,7,80}, Namrata Gupta^{6,64}, Andrea Haessler¹², Christopher Haiman²⁰⁵, Ira Hall⁸², Craig Hanis⁸³, Matthew Harms^{84,85}, Mikko Hiltunen⁸⁶, Matti M. Holi⁸⁷, Christina M. Hultman^{88,89}, Chaim Jalas⁹⁰, Thibault Jeandet¹², Mikko Kallela⁹¹, Diane Kaplan¹², Jaakko Kaprio⁷⁹, Konrad J. Karczewski^{6,7,34}, Sekar Kathiresan^{27,35,92}, Eimear Kenny^{89,93}, Bong-Jo Kim⁹⁴, Young Jin Kim⁹⁴, George Kirov⁹⁵, Zan Koenig⁶, Jaspal Kooner^{31,96,97}, Seppo Koskinen⁹⁸, Harlan M. Krumholz⁹⁹, Subra Kugathasan¹⁰⁰, Soo Heon Kwak¹⁰¹, Markku Laakso^{102,103}, Nicole Lake¹⁰⁴, Trevyn Langsford¹², Kristen M. Laricchia^{6,7}, Terho Lehtimäki¹⁰⁵, Monkol Lek¹⁰⁴, Emily Lipscomb⁶, Christopher Llanwarne¹², Ruth J.F. Loos^{20,106}, Steven A. Lubitz^{6,53}, Teresa Tusie Luna^{107,108}, Ronald C.W. Ma^{33,109,110}, Daniel G. MacArthur^{6,111,112}, Gregory M. Marcus¹¹³, Jaume Marrugat^{55,114}, Alicia R. Martin⁶, Kari M. Mattila¹⁰⁵, Steven McCarroll^{34,115}, Mark I. McCarthy^{116,117,118}, Jacob McCauley^{119,120}, Dermot McGovern¹²¹, Ruth McPherson¹²², James B. Meigs^{6,35,123}, Olle Melander¹²⁴, Andres Metspalu¹²⁵, Deborah Meyers¹²⁶, Eric V. Minikel⁶, Braxton D. Mitchell¹²⁷, Vamsi K. Mootha^{6,128}, Ruchi Munshi¹², Aliya Naheed¹²⁹, Saman Nazarian^{130,131}, Benjamin M. Neale^{6,7}, Peter M. Nilsson¹³², Sam Novod¹², Anne H. O'Donnell-Luria^{6,7,80}, Michael C. O'Donovan⁹⁵, Yukinori Okada^{133,134,135}, Dost Ongur^{35,42}, Lorena Orozco¹³⁶, Michael J. Owen⁹⁵, Colin Palmer¹³⁷, Nicholette D. Palmer¹³⁸, Aarno Palotie^{7,34,79}, Kyong Soo Park^{101,139}, Carlos Pato¹⁴⁰, Nikelle Petrillo¹², William Phu^{6,80}, Timothy Poterba^{6,7,34}, Ann E. Pulver¹⁴¹, Dan Rader^{130,142}, Nazneen Rahman¹⁴³, Heidi L. Rehm^{6,27}, Alex Reiner^{144,145}, Anne M. Remes¹⁴⁶, Dan Rhodes⁶, Stephen S. Rich^{147,148}, John D. Rioux^{149,150}, Samuli Ripatti^{79,151,152}, David Roazen¹², Dan M. Roden^{153,154}, Jerome I. Rotter¹⁵⁵, Valentin Ruano-Rubio¹², Nareh Sahakian¹², Danish Saleheen^{156,157,158}, Veikko Salomaa¹⁵⁹, Andrea Saltzman⁶, Nilesh J. Samani^{24,25}, Kaitlin E. Samocha^{6,27}, Jeremiah Scharf^{6,27,34}, Molly Schleicher⁶, Heribert Schunkert^{160,161}, Sebastian Schönherr¹⁶², Eleanor Seaby⁶, Cotton Seed^{7,34}, Svati H. Shah¹⁶³, Megan Shand¹², Moore B. Shoemaker¹⁶⁴, Tai Shyong^{165,166}, Edwin K. Silverman^{167,168}, Moriel Singer-Berk⁶, Pamela Sklar^{169,170,171}, J. Gustav Smith^{152,172,173}, Jonathan T. Smith¹², Hilka Soininen¹⁷⁴, Harry Sokol^{175,176,177}, Matthew Solomonson^{6,7}, Rachel G. Son⁶, Jose Soto¹², Tim Spector¹⁷⁸, Christine Stevens^{6,7,34}, Nathan Stitzel^{82,179}, Patrick F. Sullivan^{88,180}, Jaana Suvisaari¹⁵⁹, E. Shyong Tai^{181,182,183}, Michael E. Talkowski^{6,27,34}, Yekaterina Tarasova⁶, Kent D. Taylor¹⁵⁵, Yik Ying Teo^{181,184,185}, Grace Tiao^{6,7}, Kathleen Tibbetts¹², Charlotte Tolonen¹², Ming Tsuang^{186,187}, Tiinamaija Tuomi^{79,188,189}, Dan Turner¹⁹⁰, Teresa Tusie-Luna^{191,192}, Erkki Vartiainen¹⁹³, Marquis P. Vawter¹⁹⁴, Christopher Vittal^{6,7}, Gordon Wade¹², Arcturus Wang^{6,7,34}, Qingbo Wang^{6,133}, James S. Ware^{6,195,196}, Hugh Watkins¹⁹⁷, Nicholas A. Watts^{6,7}, Rinse K. Weersma¹⁹⁸, Ben Weisburd¹², Maija Wessman^{79,199}, Nicola Whiffin^{6,200,201}, Michael W. Wilson^{6,7}, James G. Wilson²⁰², Ramnik J. Xavier^{203,204}, Mary T. Yohannes⁶

¹University of Miami Miller School of Medicine, Gastroenterology, Miami, USA

²Unidad de Investigacion de Enfermedades Metabolicas, Instituto Nacional de Ciencias Medicas y Nutricion, Mexico City, Mexico

³Peninsula College of Medicine and Dentistry, Exeter, UK

⁴Division of Preventive Medicine, Brigham and Women's Hospital, Boston, MA, USA

⁵Division of Cardiovascular Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA

⁶Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

⁷Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

- ⁸Department of Cardiology University Hospital, Parma, Italy
- ⁹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK
- ¹⁰Department of Biology Faculty of Natural Sciences, University of Haifa, Haifa, Israel
- ¹¹Departments of Medicine and Genetics, Albert Einstein College of Medicine, Bronx, NY, USA
- ¹²Data Science Platform, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ¹³Department of Quantitative Health Sciences, Lerner Research Institute Cleveland Clinic, Cleveland, OH, USA
- ¹⁴Sorbonne Université, APHP, Gastroenterology Department Saint Antoine Hospital, Paris, France
- ¹⁵NHLBI and Boston University's Framingham Heart Study, Framingham, MA, USA
- ¹⁶Department of Medicine, Boston University School of Medicine, Boston, MA, USA
- ¹⁷Department of Epidemiology, Boston University School of Public Health, Boston, MA, USA
- ¹⁸Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA
- ¹⁹National Human Genome Research Institute, National Institutes of Health Bethesda, MD, USA
- ²⁰The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA
- ²¹Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC, USA
- ²²Center for Genomics and Personalized Medicine Research, Wake Forest School of Medicine, Winston-Salem, NC, USA
- ²³Center for Diabetes Research, Wake Forest School of Medicine, Winston-Salem, NC, USA
- ²⁴Department of Cardiovascular Sciences, University of Leicester, Leicester, UK
- ²⁵NIHR Leicester Biomedical Research Centre, Glenfield Hospital, Leicester, UK
- ²⁶John Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
- ²⁷Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA
- ²⁸Harvard School of Public Health, Boston, MA, USA
- ²⁹Central American Population Center, San Pedro, Costa Rica
- ³⁰Department of Epidemiology and Biostatistics, Imperial College London, London, UK
- ³¹Department of Cardiology, Ealing Hospital, NHS Trust, Southall, UK
- ³²Imperial College, Healthcare NHS Trust Imperial College London, London, UK
- ³³Department of Medicine and Therapeutics, The Chinese University of Hong Kong, Hong Kong, China
- ³⁴Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ³⁵Department of Medicine, Harvard Medical School, Boston, MA, USA
- ³⁶Northwestern University Feinberg School of Medicine, Chicago, IL, USA
- ³⁷University of Cambridge, Cambridge, England
- ³⁸Departments of Cardiovascular, Medicine Cellular and Molecular Medicine Molecular Cardiology, Quantitative Health Sciences, Cleveland Clinic, Cleveland, OH, USA
- ³⁹Department of Pediatrics, Columbia University Irving Medical Center, New York, NY, USA
- ⁴⁰Herbert Irving Comprehensive Cancer Center, Columbia University Medical Center, New York, NY, USA
- ⁴¹Department of Medicine, Columbia University Medical Center, New York, NY, USA
- ⁴²McLean Hospital, Belmont, MA, USA
- ⁴³Division of Medical Sciences, Harvard Medical School, Boston, MA, USA
- ⁴⁴Genomics Platform, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ⁴⁵Department of Medicine, University of Mississippi Medical Center, Jackson, MI, USA
- ⁴⁶Department of Epidemiology Colorado School of Public Health Aurora, CO, USA
- ⁴⁷Institute for Molecular Medicine Finland, (FIMM) Helsinki, Finland
- ⁴⁸Department of Medicine and Pharmacology, University of Illinois at Chicago, Chicago, IL, USA
- ⁴⁹Vanderbilt University Medical Center, Nashville, TN, USA
- ⁵⁰Department of Genetics, Texas Biomedical Research Institute, San Antonio, TX, USA
- ⁵¹Department of Biostatistics, Boston University School of Public Health, Boston, MA, USA
- ⁵²National Heart Lung and Blood Institute's Framingham Heart Study, Framingham, MA, USA
- ⁵³Cardiac Arrhythmia Service and Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA
- ⁵⁴Cardiovascular Epidemiology and Genetics Hospital del Mar Medical Research Institute, (IMIM) Barcelona Catalonia, Spain

- ⁵⁵CIBER CV Barcelona, Catalonia, Spain
- ⁵⁶Department of Medicine, Medical School University of Vic-Central, University of Catalonia, Vic Catalonia, Spain
- ⁵⁷Institute for Cardiogenetics, University of Lübeck, Lübeck, Germany
- ⁵⁸German Research Centre for Cardiovascular Research, Hamburg/Lübeck/Kiel, Lübeck, Germany
- ⁵⁹University Heart Center Lübeck, Lübeck, Germany
- ⁶⁰Estonian Genome Center, Institute of Genomics University of Tartu, Tartu, Estonia
- ⁶¹Victor Chang Cardiac Research Institute, Darlinghurst, NSW, Australia
- ⁶²Faculty of Medicine, UNSW Sydney, Kensington, NSW, Australia
- ⁶³Cardiology Department, St Vincent's Hospital, Darlinghurst, NSW, Australia
- ⁶⁴Broad Genomics, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ⁶⁵Diabetes Unit and Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA
- ⁶⁶Programs in Metabolism and Medical & Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ⁶⁷Institute of Clinical Molecular Biology, (IKMB) Christian-Albrechts-University of Kiel, Kiel, Germany
- ⁶⁸Helsinki University and Helsinki University Hospital Clinic of Gastroenterology, Helsinki, Finland
- ⁶⁹Bioinformatics Program MGH Cancer Center and Department of Pathology, Boston, MA, USA
- ⁷⁰Cancer Genome Computational Analysis, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ⁷¹Department of Psychiatry and Behavioral Sciences, Boston Children's Hospital and Harvard Medical School, Boston, MA, USA
- ⁷²Harvard Medical School Teaching Hospital, Boston, MA, USA
- ⁷³Department of Endocrinology and Metabolism, Hadassah Medical Center and Faculty of Medicine, Hebrew University of Jerusalem, Israel
- ⁷⁴Department of Psychiatry and Behavioral Sciences, SUNY Upstate Medical University, Syracuse, NY, USA
- ⁷⁵Institute for Genomic Medicine, Columbia University Medical Center Hammer Health Sciences, New York, NY, USA
- ⁷⁶Department of Genetics & Development Columbia University Medical Center, Hammer Health Sciences, New York, NY, USA
- ⁷⁷Centro de Investigacion en Salud Poblacional, Instituto Nacional de Salud Publica, Mexico
- ⁷⁸Lund University Sweden, Sweden
- ⁷⁹Institute for Molecular Medicine Finland, (FIMM) HiLIFE University of Helsinki, Helsinki, Finland
- ⁸⁰Division of Genetics and Genomics, Boston Children's Hospital, Boston, MA, USA
- ⁸¹Lund University Diabetes Centre, Malmö, Skåne County, Sweden
- ⁸²Washington School of Medicine, St Louis, MI, USA
- ⁸³Human Genetics Center University of Texas Health Science Center at Houston, Houston, TX, USA
- ⁸⁴Department of Neurology Columbia University, New York City, NY, USA
- ⁸⁵Institute of Genomic Medicine, Columbia University, New York City, NY, USA
- ⁸⁶Institute of Biomedicine, University of Eastern Finland, Kuopio, Finland
- ⁸⁷Department of Psychiatry, Helsinki University Central Hospital Lapinlahdentie, Helsinki, Finland
- ⁸⁸Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden
- ⁸⁹Icahn School of Medicine at Mount Sinai, New York, NY, USA
- ⁹⁰Bonei Olam, Center for Rare Jewish Genetic Diseases, Brooklyn, NY, USA
- ⁹¹Department of Neurology, Helsinki University, Central Hospital, Helsinki, Finland
- ⁹²Cardiovascular Disease Initiative and Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ⁹³Charles Bronfman Institute for Personalized Medicine, New York, NY, USA
- ⁹⁴Division of Genome Science, Department of Precision Medicine, National Institute of Health, Republic of Korea
- ⁹⁵MRC Centre for Neuropsychiatric Genetics & Genomics, Cardiff University School of Medicine, Cardiff, Wales
- ⁹⁶Imperial College, Healthcare NHS Trust, London, UK
- ⁹⁷National Heart and Lung Institute Cardiovascular Sciences, Hammersmith Campus, Imperial College London, London, UK
- ⁹⁸Department of Health THL-National Institute for Health and Welfare, Helsinki, Finland
- ⁹⁹Section of Cardiovascular Medicine, Department of Internal Medicine, Yale School of Medicine, Center for Outcomes Research and Evaluation Yale-New Haven Hospital, New Haven, CT, USA

- ¹⁰⁰Division of Pediatric Gastroenterology, Emory University School of Medicine, Atlanta, GA, USA
- ¹⁰¹Department of Internal Medicine, Seoul National University Hospital, Seoul, Republic of Korea
- ¹⁰²The University of Eastern Finland, Institute of Clinical Medicine, Kuopio, Finland
- ¹⁰³Kuopio University Hospital, Kuopio, Finland
- ¹⁰⁴Department of Genetics, Yale School of Medicine, New Haven, CT, USA
- ¹⁰⁵Department of Clinical Chemistry Fimlab Laboratories and Finnish Cardiovascular Research Center-Tampere Faculty of Medicine and Health Technology, Tampere University, Finland
- ¹⁰⁶The Mindich Child Health and Development, Institute Icahn School of Medicine at Mount Sinai, New York, NY, USA
- ¹⁰⁷National Autonomous University of Mexico, Mexico City, Mexico
- ¹⁰⁸Salvador Zubirán National Institute of Health Sciences and Nutrition, Mexico City, Mexico
- ¹⁰⁹Li Ka Shing Institute of Health Sciences, The Chinese University of Hong Kong, Hong Kong, China
- ¹¹⁰Hong Kong Institute of Diabetes and Obesity, The Chinese University of Hong Kong, Hong Kong, China
- ¹¹¹Centre for Population Genomics, Garvan Institute of Medical Research and UNSW Sydney, Sydney, Australia
- ¹¹²Centre for Population Genomics, Murdoch Children's Research Institute, Melbourne, Australia
- ¹¹³University of California San Francisco Parnassus Campus, San Francisco, CA, USA
- ¹¹⁴Cardiovascular Research REGICOR Group, Hospital del Mar Medical Research Institute, (IMIM) Barcelona, Catalonia, Spain
- ¹¹⁵Department of Genetics, Harvard Medical School, Boston, MA, USA
- ¹¹⁶Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford, Churchill Hospital Old Road Headington, Oxford, OX, UK
- ¹¹⁷Welcome Centre for Human Genetics, University of Oxford, Oxford, OX, UK
- ¹¹⁸Oxford NIHR Biomedical Research Centre, Oxford University Hospitals, NHS Foundation Trust, John Radcliffe Hospital, Oxford, OX, UK
- ¹¹⁹John P. Hussman Institute for Human Genomics, Leonard M. Miller School of Medicine, University of Miami, Miami, FL, USA
- ¹²⁰The Dr. John T. Macdonald Foundation Department of Human Genetics, Leonard M. Miller School of Medicine, University of Miami, Miami, FL, USA
- ¹²¹F. Widjaja Foundation Inflammatory Bowel and Immunobiology Research Institute Cedars-Sinai Medical Center, Los Angeles, CA, USA
- ¹²²Atherogenomics Laboratory University of Ottawa, Heart Institute, Ottawa, Canada
- ¹²³Division of General Internal Medicine, Massachusetts General Hospital, Boston, MA, USA
- ¹²⁴Department of Clinical Sciences University, Hospital Malmö Clinical Research Center, Lund University, Malmö, Sweden
- ¹²⁵Estonian Genome Center, Institute of Genomics, University of Tartu, Tartu, Estonia
- ¹²⁶University of Arizona Health Science, Tucson, AZ, USA
- ¹²⁷Department of Medicine, University of Maryland School of Medicine, Baltimore, MD, USA
- ¹²⁸Howard Hughes Medical Institute and Department of Molecular Biology, Massachusetts General Hospital, Boston, MA, USA
- ¹²⁹International Centre for Diarrhoeal Disease Research, Bangladesh
- ¹³⁰Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA
- ¹³¹Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
- ¹³²Lund University, Dept. Clinical Sciences, Skåne University Hospital, Malmö, Sweden
- ¹³³Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita, Japan
- ¹³⁴Laboratory of Statistical Immunology, Immunology Frontier Research Center (WPI-IFReC), Osaka University, Suita, Japan
- ¹³⁵Integrated Frontier Research for Medical Science Division, Institute for Open and Transdisciplinary Research Initiatives, Osaka University, Suita, Japan
- ¹³⁶Instituto Nacional de Medicina Genómica, (INMEGEN) Mexico City, Mexico
- ¹³⁷Medical Research Institute, Ninewells Hospital and Medical School University of Dundee, Dundee, UK
- ¹³⁸Wake Forest School of Medicine, Winston-Salem, NC, USA
- ¹³⁹Department of Molecular Medicine and Biopharmaceutical Sciences, Graduate School of Convergence Science and Technology, Seoul National University, Seoul, Republic of Korea

- ¹⁴⁰Department of Psychiatry Keck School of Medicine at the University of Southern California, Los Angeles, CA, USA
- ¹⁴¹Department of Psychiatry and Behavioral Sciences, Johns Hopkins University School of Medicine, Baltimore, MD, USA
- ¹⁴²Children's Hospital of Philadelphia, Philadelphia, PA, USA
- ¹⁴³Division of Genetics and Epidemiology, Institute of Cancer Research, London, SM, NG
- ¹⁴⁴University of Washington, Seattle, WA, USA
- ¹⁴⁵Fred Hutchinson Cancer Research Center, Seattle, WA, USA
- ¹⁴⁶Medical Research Center, Oulu University Hospital, Oulu Finland and Research Unit of Clinical Neuroscience Neurology University of Oulu, Oulu, Finland
- ¹⁴⁷Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA
- ¹⁴⁸Department of Public Health Sciences, University of Virginia, Charlottesville, VA, USA
- ¹⁴⁹Research Center Montreal Heart Institute, Montreal, Quebec, Canada
- ¹⁵⁰Department of Medicine, Faculty of Medicine Université de Montréal, Québec, Canada
- ¹⁵¹Department of Public Health Faculty of Medicine, University of Helsinki, Helsinki, Finland
- ¹⁵²Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ¹⁵³Department of Biomedical Informatics Vanderbilt, University Medical Center, Nashville, TN, USA
- ¹⁵⁴Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA
- ¹⁵⁵The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA, USA
- ¹⁵⁶Department of Biostatistics and Epidemiology, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA
- ¹⁵⁷Department of Medicine, Perelman School of Medicine at the University of Pennsylvania, Philadelphia, PA, USA
- ¹⁵⁸Center for Non-Communicable Diseases, Karachi, Pakistan
- ¹⁵⁹National Institute for Health and Welfare, Helsinki, Finland
- ¹⁶⁰Deutsches Herzzentrum, München, Germany
- ¹⁶¹Technische Universität München, Germany
- ¹⁶²Institute of Genetic Epidemiology, Department of Genetics and Pharmacology, Medical University of Innsbruck, 6020 Innsbruck, Austria
- ¹⁶³Duke Molecular Physiology Institute, Durham, NC
- ¹⁶⁴Division of Cardiovascular Medicine, Nashville VA Medical Center, Vanderbilt University School of Medicine, Nashville, TN, USA
- ¹⁶⁵Division of Endocrinology, National University Hospital, Singapore
- ¹⁶⁶NUS Saw Swee Hock School of Public Health, Singapore
- ¹⁶⁷Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA, USA
- ¹⁶⁸Harvard Medical School, Boston, MA, USA
- ¹⁶⁹Department of Psychiatry, Icahn School of Medicine at Mount Sinai, New York, NY, USA
- ¹⁷⁰Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA
- ¹⁷¹Institute for Genomics and Multiscale Biology, Icahn School of Medicine at Mount Sinai, New York, NY, USA
- ¹⁷²The Wallenberg Laboratory/Department of Molecular and Clinical Medicine, Institute of Medicine, Gothenburg University and the Department of Cardiology, Sahlgrenska University Hospital, Gothenburg, Sweden
- ¹⁷³Department of Cardiology, Wallenberg Center for Molecular Medicine and Lund University Diabetes Center, Clinical Sciences, Lund University and Skåne University Hospital, Lund, Sweden
- ¹⁷⁴Institute of Clinical Medicine Neurology, University of Eastern Finland, Kuopio, Finland
- ¹⁷⁵Sorbonne Université, INSERM, Centre de Recherche Saint-Antoine, CRSA, AP-HP, Saint Antoine Hospital, Gastroenterology department, F-75012 Paris, France
- ¹⁷⁶INRA, UMR1319 Micalis & AgroParisTech, Jouy en Josas, France
- ¹⁷⁷Paris Center for Microbiome Medicine, (PaCeMM) FHU, Paris, France
- ¹⁷⁸Department of Twin Research and Genetic Epidemiology King's College London, London, UK
- ¹⁷⁹The McDonnell Genome Institute at Washington University, Seattle, WA, USA
- ¹⁸⁰Departments of Genetics and Psychiatry, University of North Carolina, Chapel Hill, NC, USA
- ¹⁸¹Saw Swee Hock School of Public Health National University of Singapore, National University Health System, Singapore
- ¹⁸²Department of Medicine, Yong Loo Lin School of Medicine National University of Singapore, Singapore

- ¹⁸³Duke-NUS Graduate Medical School, Singapore
- ¹⁸⁴Life Sciences Institute, National University of Singapore, Singapore
- ¹⁸⁵Department of Statistics and Applied Probability, National University of Singapore, Singapore
- ¹⁸⁶Center for Behavioral Genomics, Department of Psychiatry, University of California, San Diego, CA, USA
- ¹⁸⁷Institute of Genomic Medicine, University of California San Diego, San Diego, CA, USA
- ¹⁸⁸Endocrinology, Abdominal Center, Helsinki University Hospital, Helsinki, Finland
- ¹⁸⁹Institute of Genetics, Folkhalsan Research Center, Helsinki, Finland
- ¹⁹⁰Juliet Keidan Institute of Pediatric Gastroenterology Shaare Zedek Medical Center, The Hebrew University of Jerusalem, Jerusalem, Israel
- ¹⁹¹Instituto de Investigaciones Biomédicas, UNAM, Mexico City, Mexico
- ¹⁹²Instituto Nacional de Ciencias Médicas y Nutrición Salvador Zubirán, Mexico City, Mexico
- ¹⁹³Department of Public Health Faculty of Medicine University of Helsinki, Helsinki, Finland
- ¹⁹⁴Department of Psychiatry and Human Behavior, University of California Irvine, Irvine, CA, USA
- ¹⁹⁵National Heart & Lung Institute & MRC London Institute of Medical Sciences, Imperial College, London, UK
- ¹⁹⁶Cardiovascular Research Centre Royal Brompton & Harefield Hospitals, London, UK
- ¹⁹⁷Radcliffe Department of Medicine, University of Oxford, Oxford, UK
- ¹⁹⁸Department of Gastroenterology and Hepatology, University of Groningen and University Medical Center Groningen, Groningen, Netherlands
- ¹⁹⁹Folkhälsan Institute of Genetics, Folkhälsan Research Center, Helsinki, Finland
- ²⁰⁰National Heart & Lung Institute and MRC London Institute of Medical Sciences, Imperial College London, London, UK
- ²⁰¹Cardiovascular Research Centre, Royal Brompton & Harefield Hospitals NHS Trust, London, UK
- ²⁰²Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson, MS, USA
- ²⁰³Program in Infectious Disease and Microbiome, Broad Institute of MIT and Harvard, Cambridge, MA, USA
- ²⁰⁴Center for Computational and Integrative Biology, Massachusetts General Hospital, Boston, MA, USA
- ²⁰⁵Center for Genetic Epidemiology, Department of Population and Public Health Sciences, University of Southern California, Los Angeles, CA, USA

Authors received funding as follows:

Matthew J. Bown: British Heart Foundation awards CS/14/2/30841 and RG/18/10/33842
 Josée Dupuis: National Heart Lung and Blood Institute's Framingham Heart Study Contract (HHSNI); National Institute for Diabetes and Digestive and Kidney Diseases (NIDDK) R DK
 Martti Färkkilä: State funding for university level health research
 Laura D. Gauthier: Intel, Illumina
 Stephen J. Glatt: U.S. NIMH Grant R MH
 Leif Groop: The Academy of Finland and University of Helsinki: Center of Excellence for Complex Disease Genetics (grant number 312063 and 336822), Sigrid Jusélius Foundation; IMI 2 (grant No 115974 and 15881)
 Mikko Hiltunen: Academy of Finland (grant 338182) Sigrid Jusélius Foundation the Strategic Neuroscience Funding of the University of Eastern Finland
 Chaim Jalas: Bonei Olam
 Ronald Ma: Research Grants Council Theme-based Research Scheme (T12-402/13N), RGC Research Impact Fund (CU R4012-18) and a Croucher Foundation Senior Medical Research Fellowship.
 Jaakko Kaprio: Academy of Finland (grants 312073 and 336823)
 Jacob McCauley: National Institute of Diabetes and Digestive and Kidney Disease Grant R01DK104844
 Yukinori Okada: JSPS KAKENHI (19H01021, 20K21834), AMED (JP21km0405211, JP21ek0109413, JP21gm4010006, JP21km0405217, JP21ek0410075), JST Moonshot R&D (JPMJMS2021)
 Michael J. Owen: Medical Research Council UK: Centre Grant No. MR/L010305/1, Program Grant No. G0800509
 Aarno Palotie: the Academy of Finland Center of Excellence for Complex Disease Genetics (grant numbers 312074 and 336824) and Sigrid Jusélius Foundation
 John D. Rioux: National Institute of Diabetes and Digestive and Kidney Diseases (NIDDK; DK062432), from the Canadian Institutes of Health (CIHR GPG 102170), from Genome Canada/Génomique Québec (GPH-129341), and a Canada Research Chair (#230625)
 Samuli Ripatti: the Academy of Finland Center of Excellence for Complex Disease Genetics (grant number) Sigrid

Jusélius Foundation

Jerome I. Rotter: Trans-Omics in Precision Medicine (TOPMed) program was supported by the National Heart, Lung and Blood Institute (NHLBI). WGS for “NHLBI TOPMed: Multi-Ethnic Study of Atherosclerosis (MESA)” (phs001416.v1.p1) was performed at the Broad Institute of MIT and Harvard (3U54HG003067-13S1). Core support including centralized genomic read mapping and genotype calling, along with variant quality metrics and filtering were provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I). Core support including phenotype harmonization, data management, sample-identity QC, and general program coordination were provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract HHSN268201800001I). We gratefully acknowledge the studies and participants who provided biological samples and data for MESA and TOPMed. JSK was supported by the Pulmonary Fibrosis Foundation Scholars Award and grant K23-HL-150301 from the NHLBI. MRA was supported by grant K23-HL-150280, AJP was supported by grant K23-HL-140199, and AM was supported by R01-HL131565 from the NHLBI. EJB was supported by grant K23-AR-075112 from the National Institute of Arthritis and Musculoskeletal and Skin Diseases. The MESA project is conducted and supported by the National Heart, Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators. Support for MESA is provided by contracts 75N92020D00001, HHSN268201500003I, N01-HC-95159, 75N92020D00005, N01-HC-95160, 75N92020D00002, N01-HC-95161, 75N92020D00003, N01-HC-95162, 75N92020D00006, N01-HC-95163, 75N92020D00004, N01-HC-95164, 75N92020D00007, N01-HC-95165, N01-HC-95166, N01-HC-95167, N01-HC-95168, N01-HC-95169, UL1-TR-000040, UL1-TR-001079, and UL1-TR-001420. Also supported in part by the National Center for Advancing Translational Sciences, CTSI grant UL1TR001881, and the National Institute of Diabetes and Digestive and Kidney Disease Diabetes Research Center (DRC) grant DK063491 to the Southern California Diabetes Endocrinology Research Center

Edwin K. Silverman: NIH Grants U01 HL089856 and U01 HL089897

J. Gustav Smith: The Swedish Heart-Lung Foundation (2019-0526), the Swedish Research Council (2017-02554), the European Research Council (ERC-STG-2015-679242), Skåne University Hospital, governmental funding of clinical research within the Swedish National Health Service, a generous donation from the Knut and Alice Wallenberg foundation to the Wallenberg Center for Molecular Medicine in Lund, and funding from the Swedish Research Council (Linnaeus grant Dnr 349-2006-237, Strategic Research Area Exodiab Dnr 2009-1039) and Swedish Foundation for Strategic Research (Dnr IRC15-0067) to the Lund University Diabetes Center

Kent D. Taylor: Trans-Omics in Precision Medicine (TOPMed) program was supported by the National Heart, Lung and Blood Institute (NHLBI). WGS for “NHLBI TOPMed: Multi-Ethnic Study of Atherosclerosis (MESA)” (phs001416.v1.p1) was performed at the Broad Institute of MIT and Harvard (3U54HG003067-13S1). Core support including centralized genomic read mapping and genotype calling, along with variant quality metrics and filtering were provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I). Core support including phenotype harmonization, data management, sample-identity QC, and general program coordination were provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract HHSN268201800001I). We gratefully acknowledge the studies and participants who provided biological samples and data for MESA and TOPMed. JSK was supported by the Pulmonary Fibrosis Foundation Scholars Award and grant K23-HL-150301 from the NHLBI. MRA was supported by grant K23-HL-150280, AJP was supported by grant K23-HL-140199, and AM was supported by R01-HL131565 from the NHLBI. EJB was supported by grant K23-AR-075112 from the National Institute of Arthritis and Musculoskeletal and Skin Diseases. The MESA project is conducted and supported by the National Heart, Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators. Support for MESA is provided by contracts 75N92020D00001, HHSN268201500003I, N01-HC-95159, 75N92020D00005, N01-HC-95160, 75N92020D00002, N01-HC-95161, 75N92020D00003, N01-HC-95162, 75N92020D00006, N01-HC-95163, 75N92020D00004, N01-HC-95164, 75N92020D00007, N01-HC-95165, N01-HC-95166, N01-HC-95167, N01-HC-95168, N01-HC-95169, UL1-TR-000040, UL1-TR-001079, and UL1-TR-001420. Also supported in part by the National Center for Advancing Translational Sciences, CTSI grant UL1TR001881, and the National Institute of Diabetes and Digestive and Kidney Disease Diabetes Research Center (DRC) grant DK063491 to the Southern California Diabetes Endocrinology Research Center

Tiinamaija Tuomi: The Academy of Finland and University of Helsinki: Center of Excellence for Complex Disease Genetics (grant number 312072 and 336826), Folkhalsan Research Foundation, Helsinki University Hospital, Ollqvist Foundation, Liv och Halsä foundation; NovoNordisk Foundation

Teresa Tusie-Luna: CONACyT Project 312688

James S. Ware: Sir Jules Thorn Charitable Trust [21JTA], Wellcome Trust [107469/Z/15/Z], Medical Research Council (UK), NIHR Imperial College Biomedical Research Centre

Rinse K. Weersma: The Lifelines Biobank initiative has been made possible by subsidy from the Dutch Ministry of Health Welfare and Sport the Dutch Ministry of Economic Affairs the University Medical Centre Groningen (UMCG the Netherlands) the University of Groningen and the Northern Provinces of the Netherlands

No conflicts of interest to declare.

Figure 1

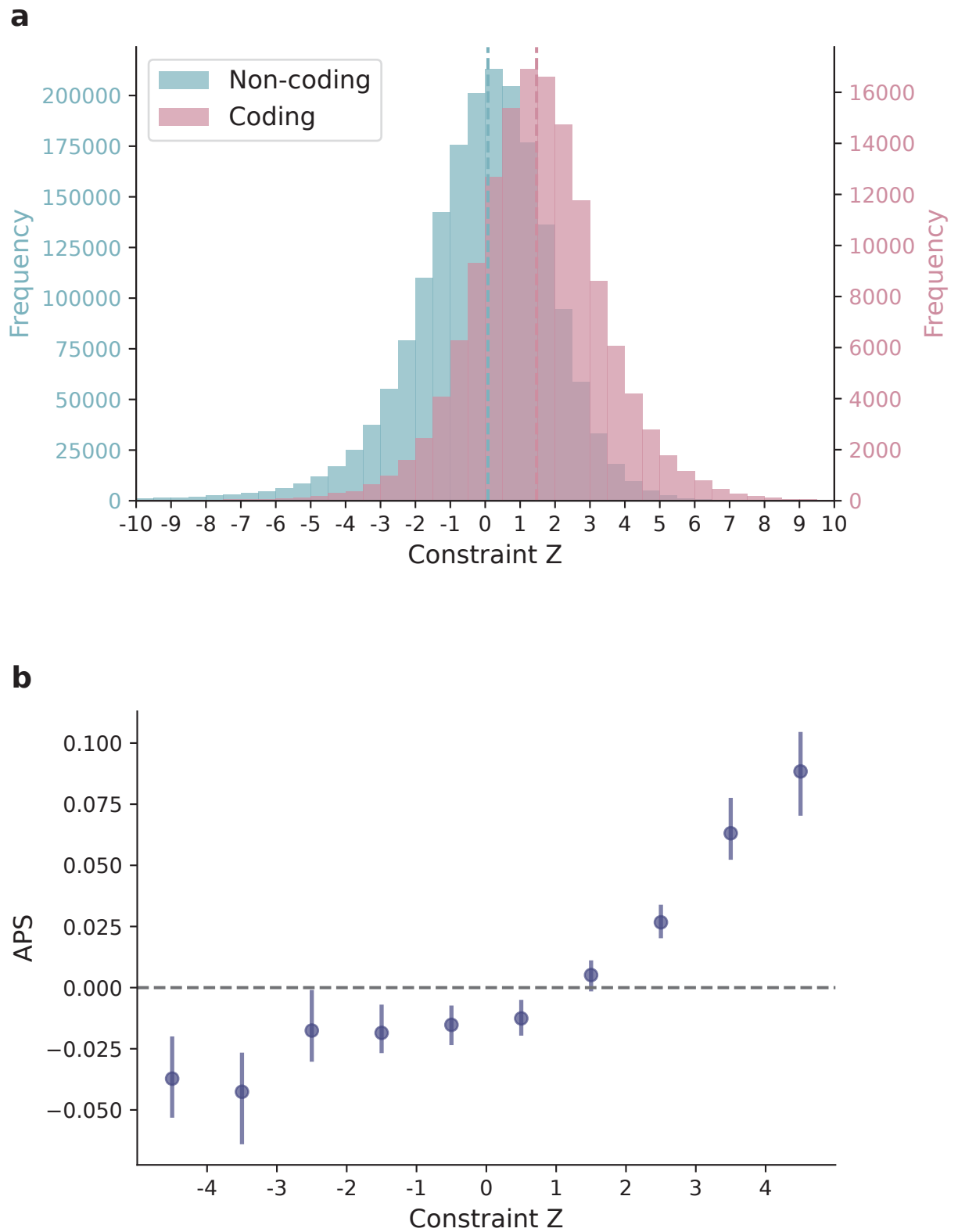


Figure 2

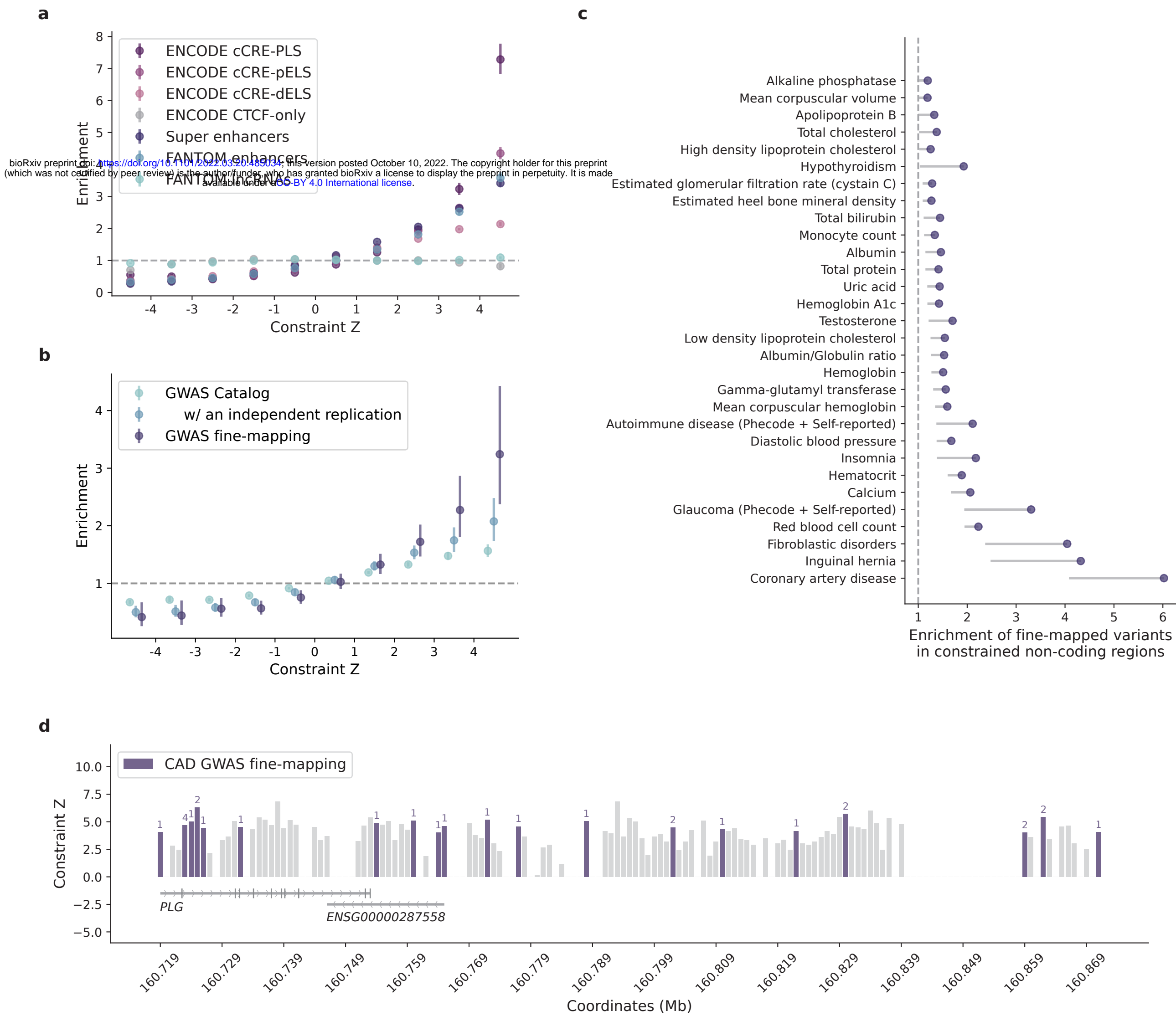


Figure 3

bioRxiv preprint doi: <https://doi.org/10.1101/2022.03.20.485034>; this version posted October 10, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

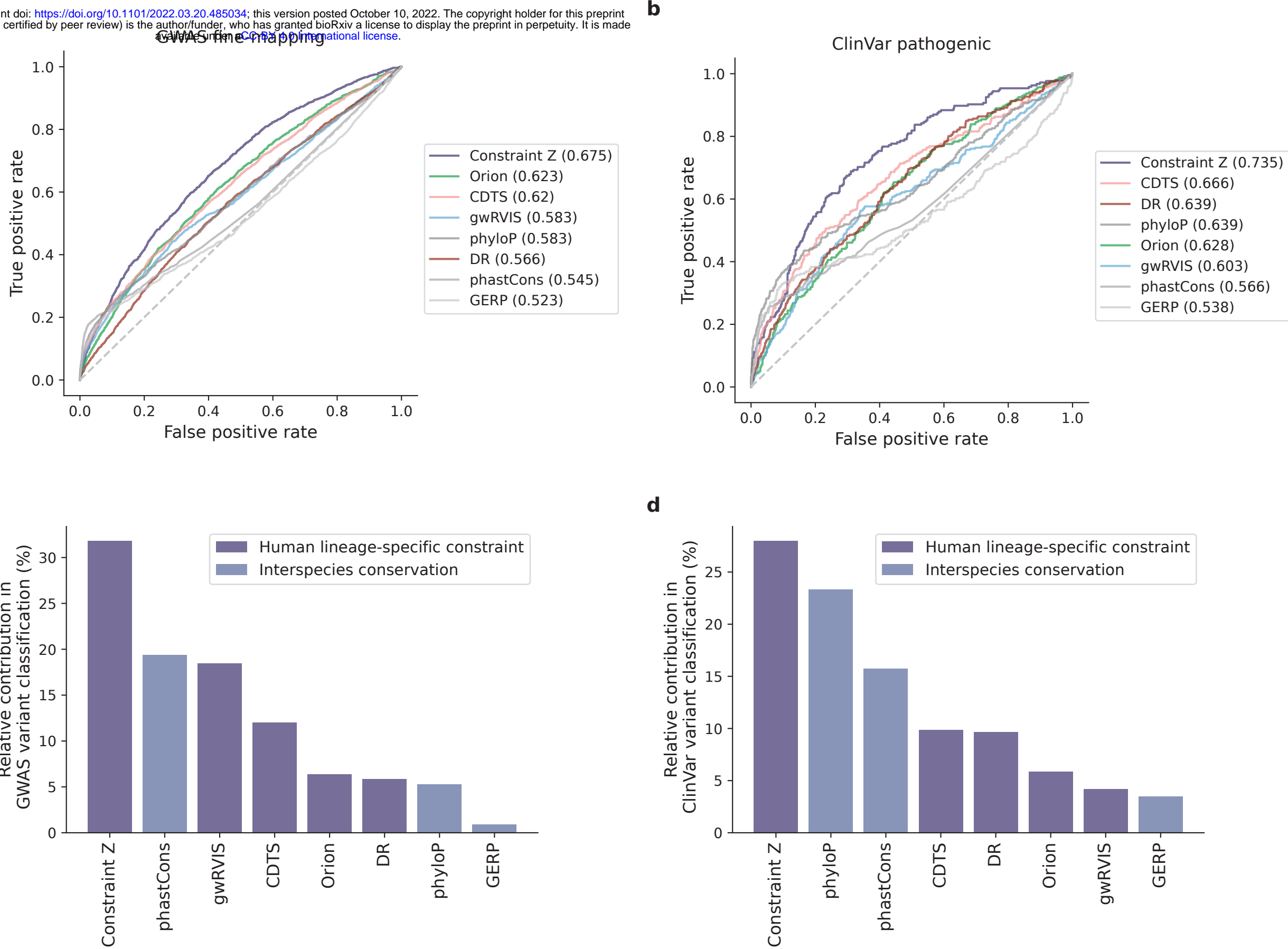


Figure 4

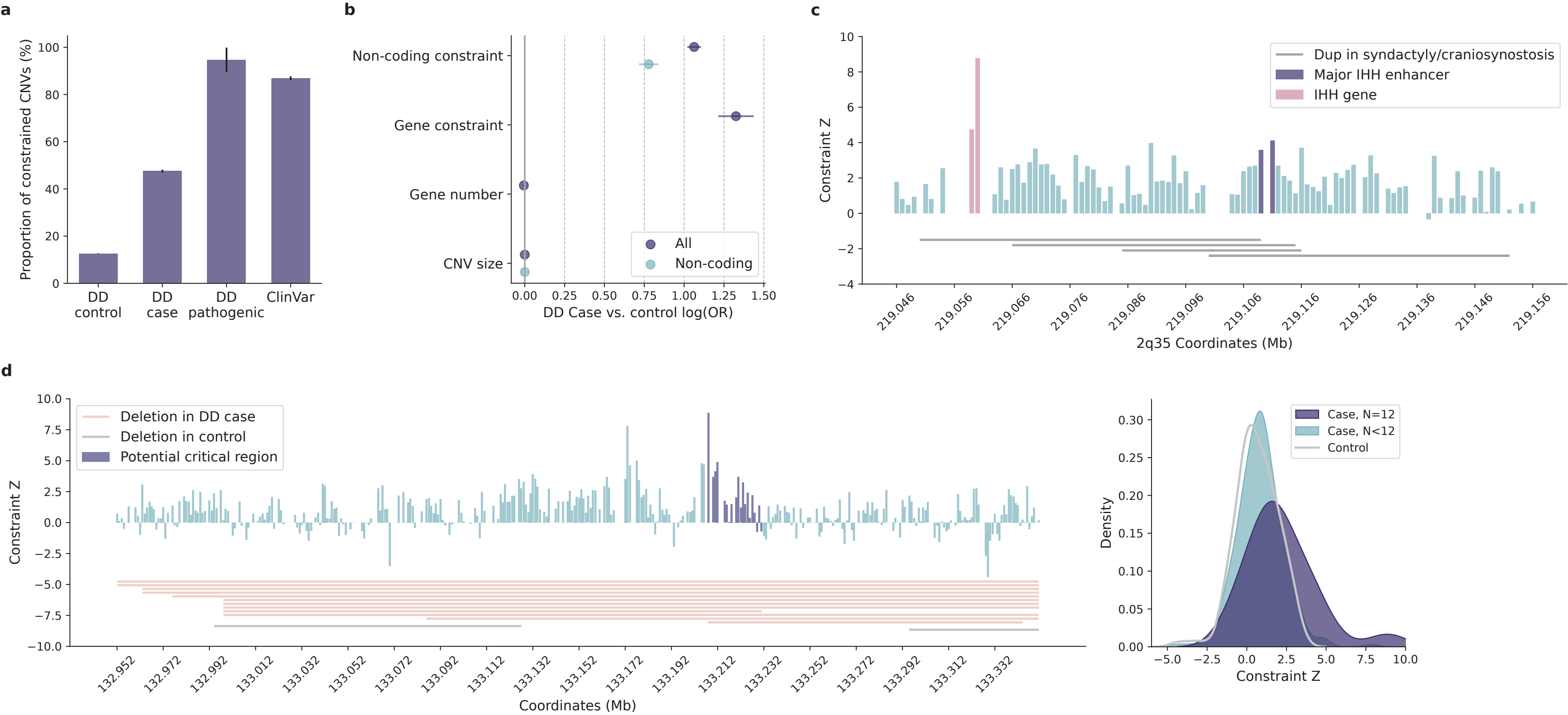
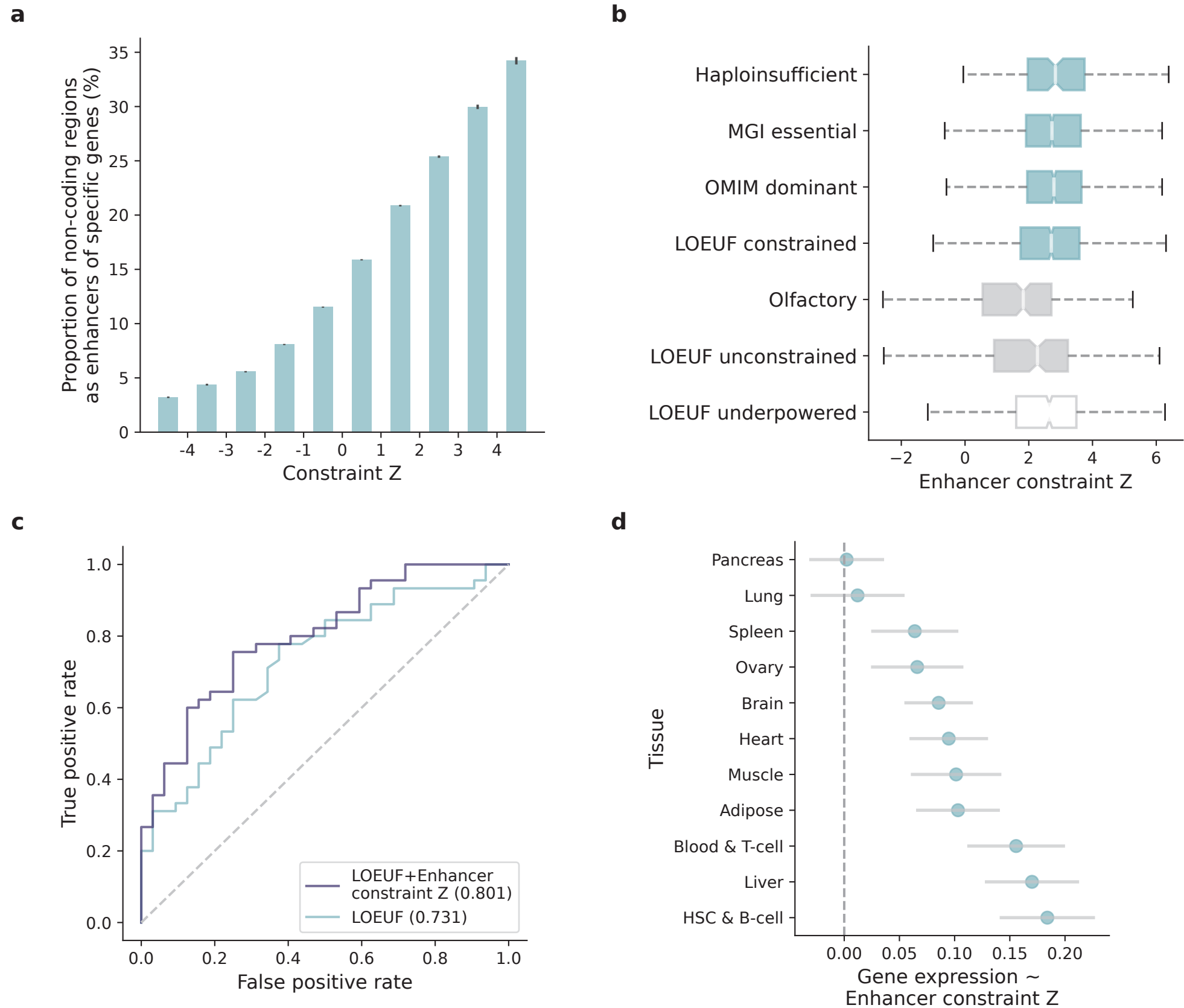
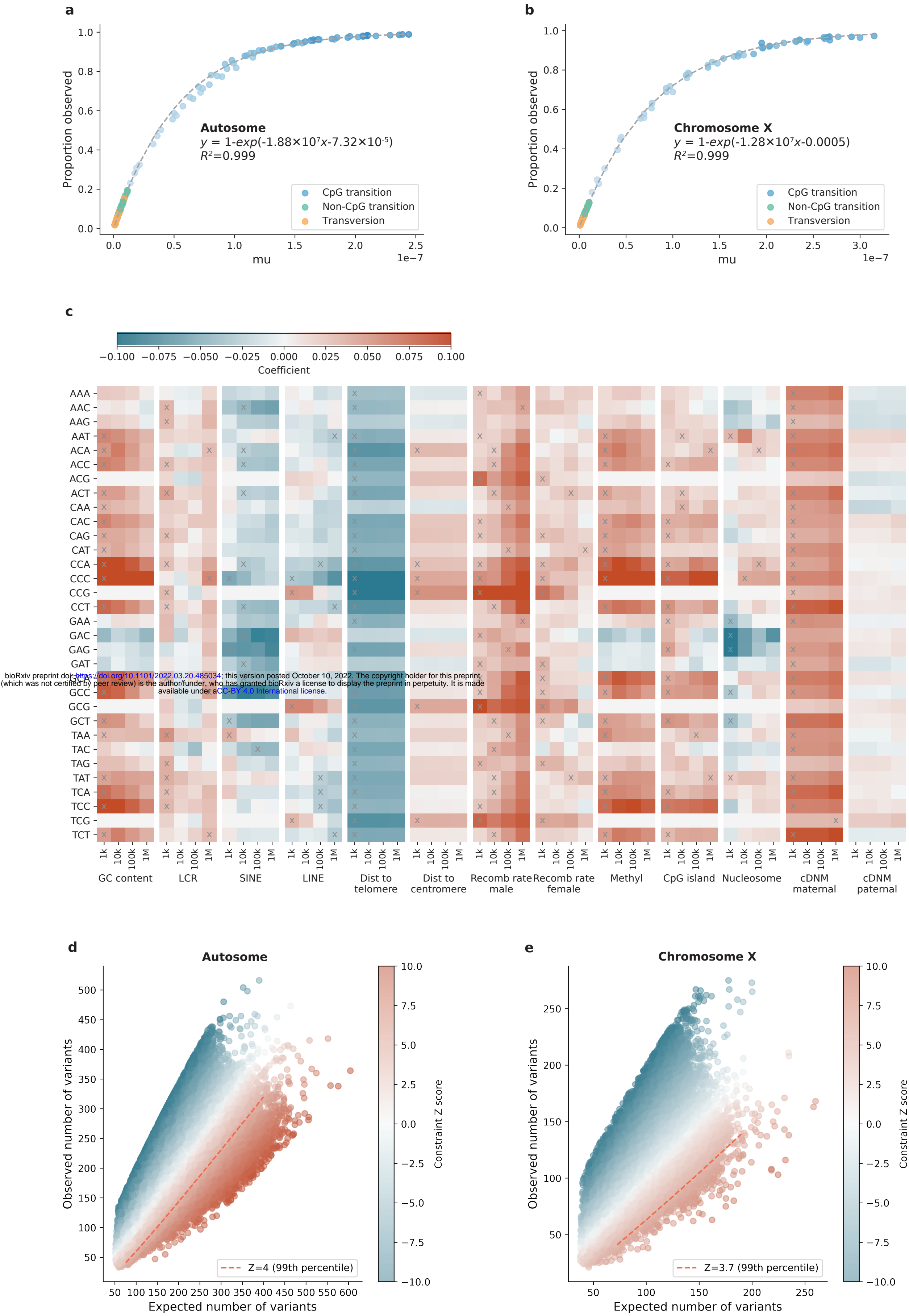


Figure 5

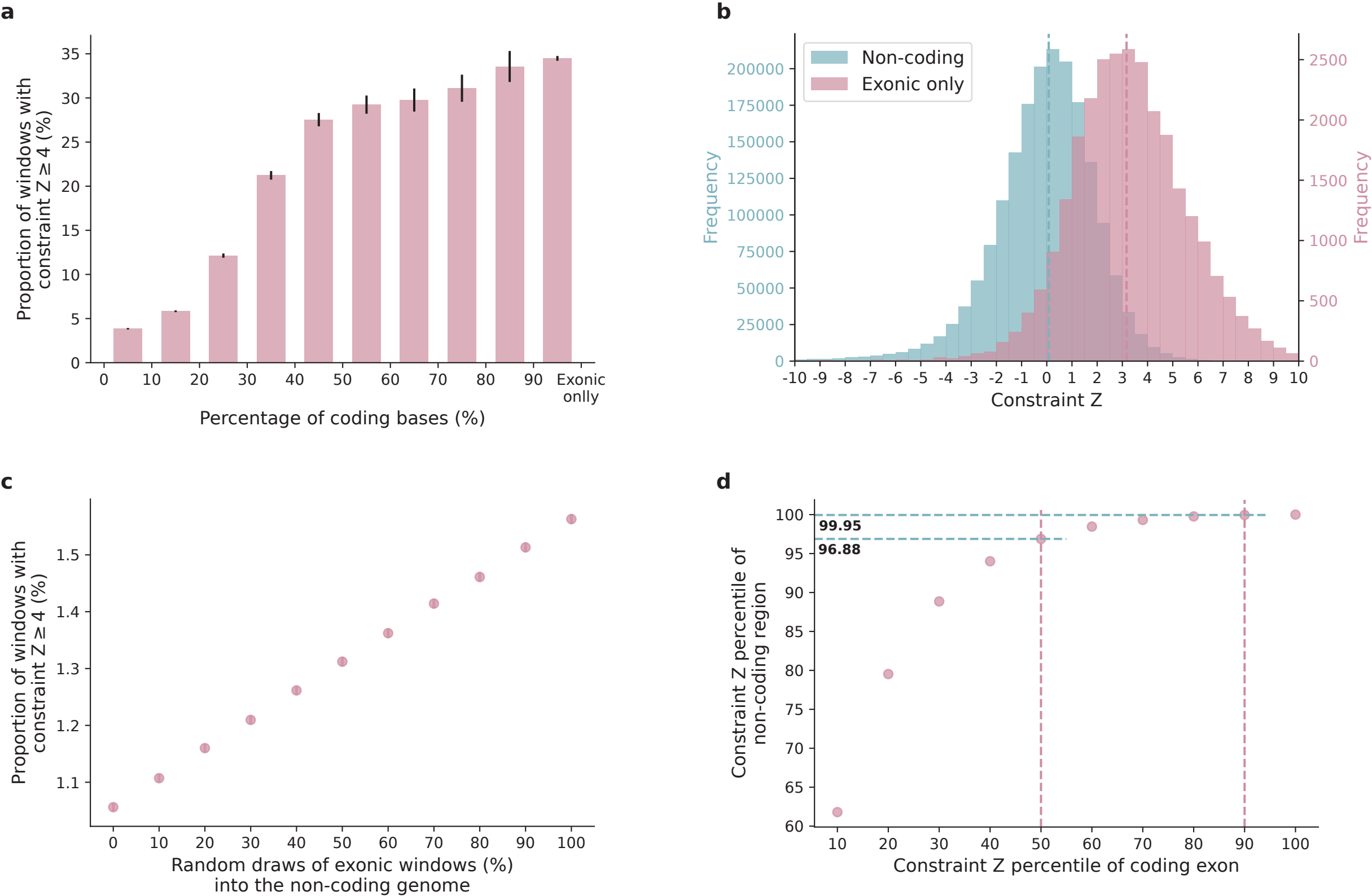
bioRxiv preprint doi: <https://doi.org/10.1101/2022.03.20.485034>; this version posted October 10, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



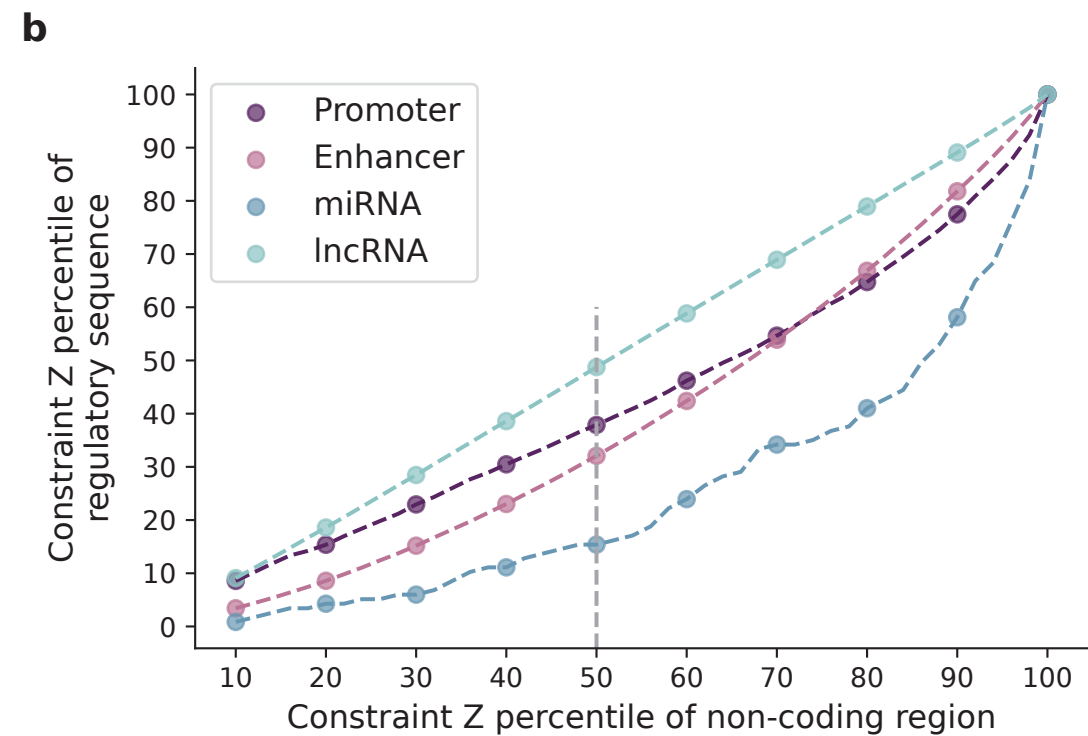
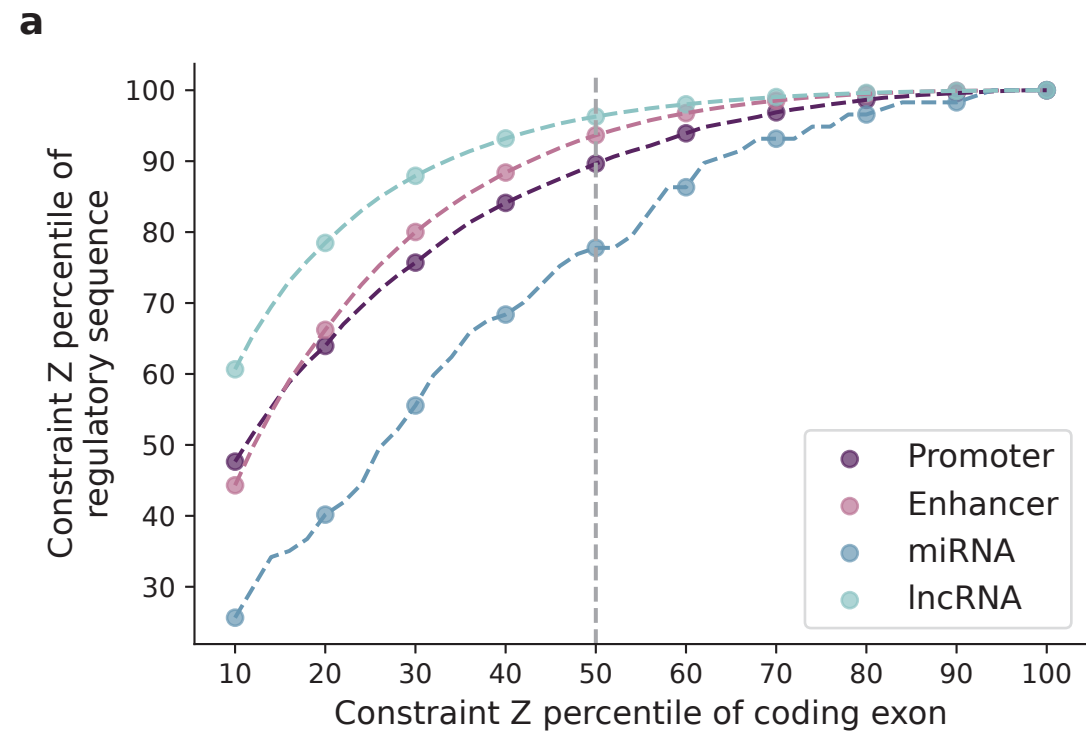
Extended Data Figure 1



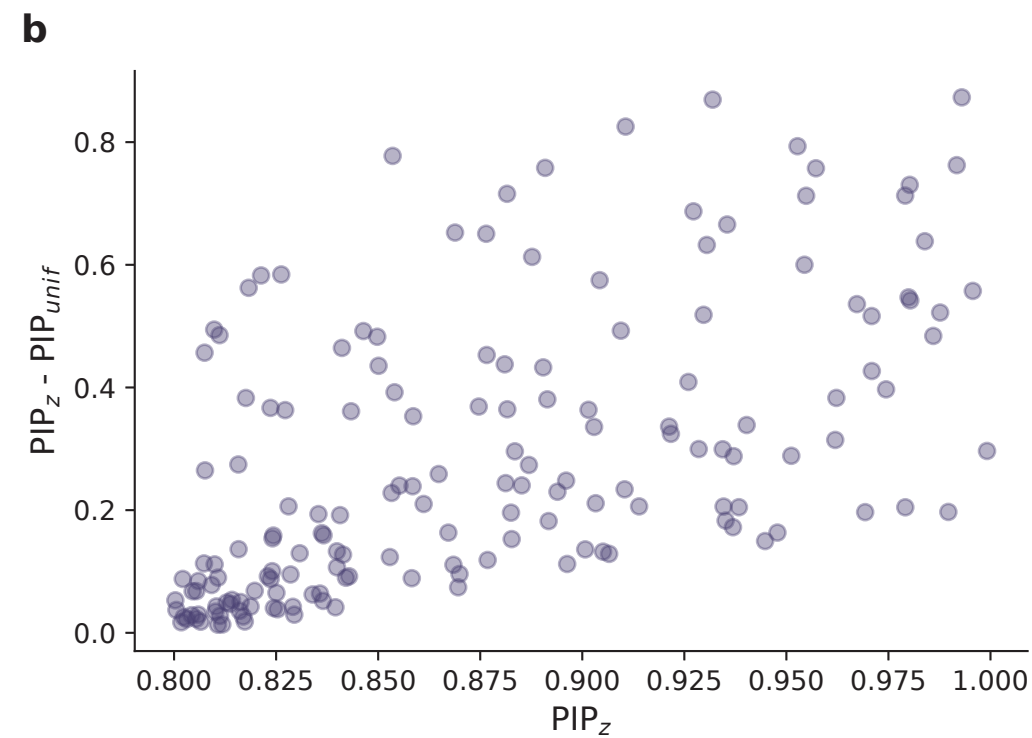
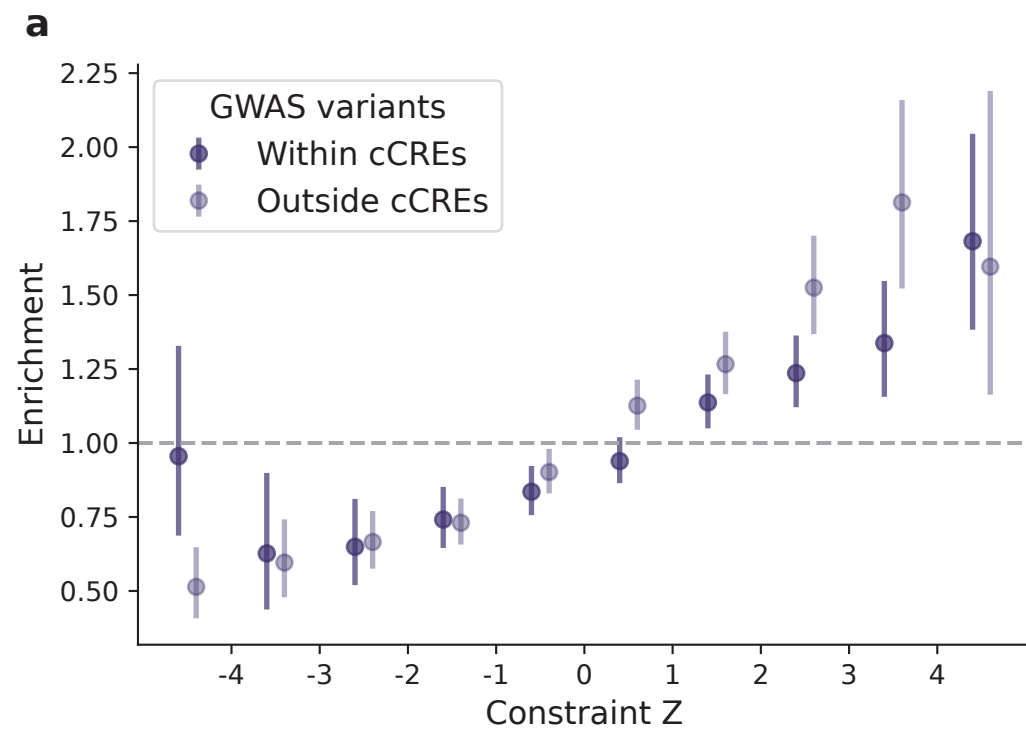
Extended Data Figure 2



Extended Data Figure 3

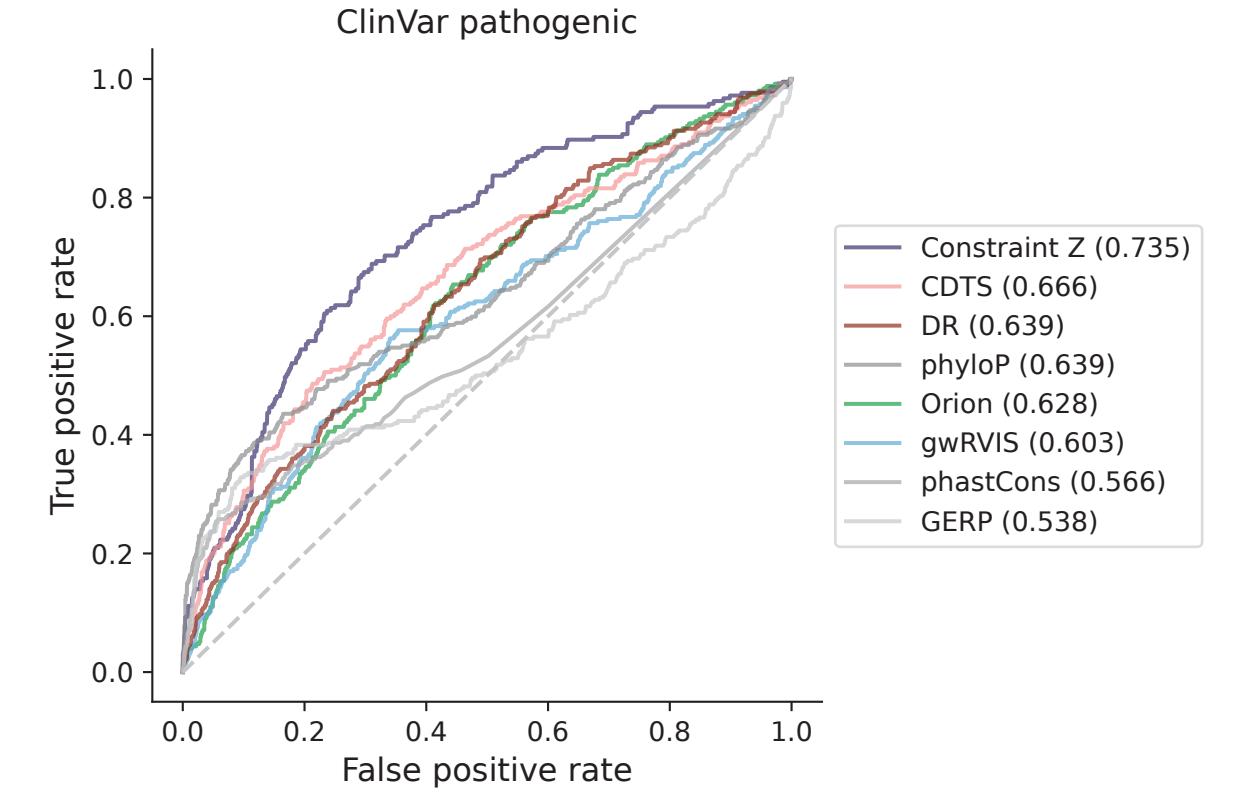
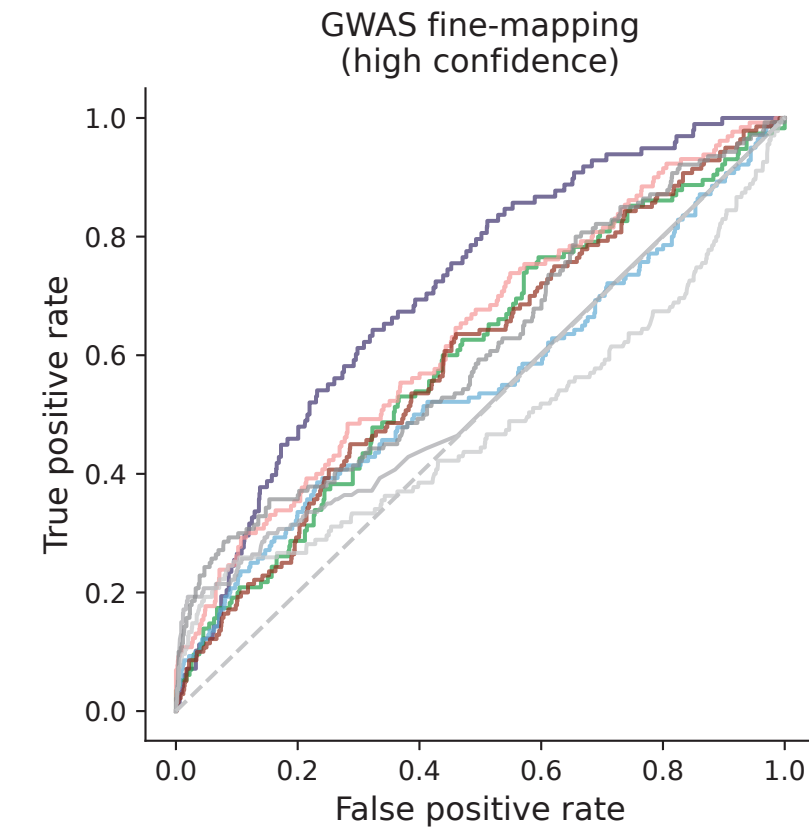
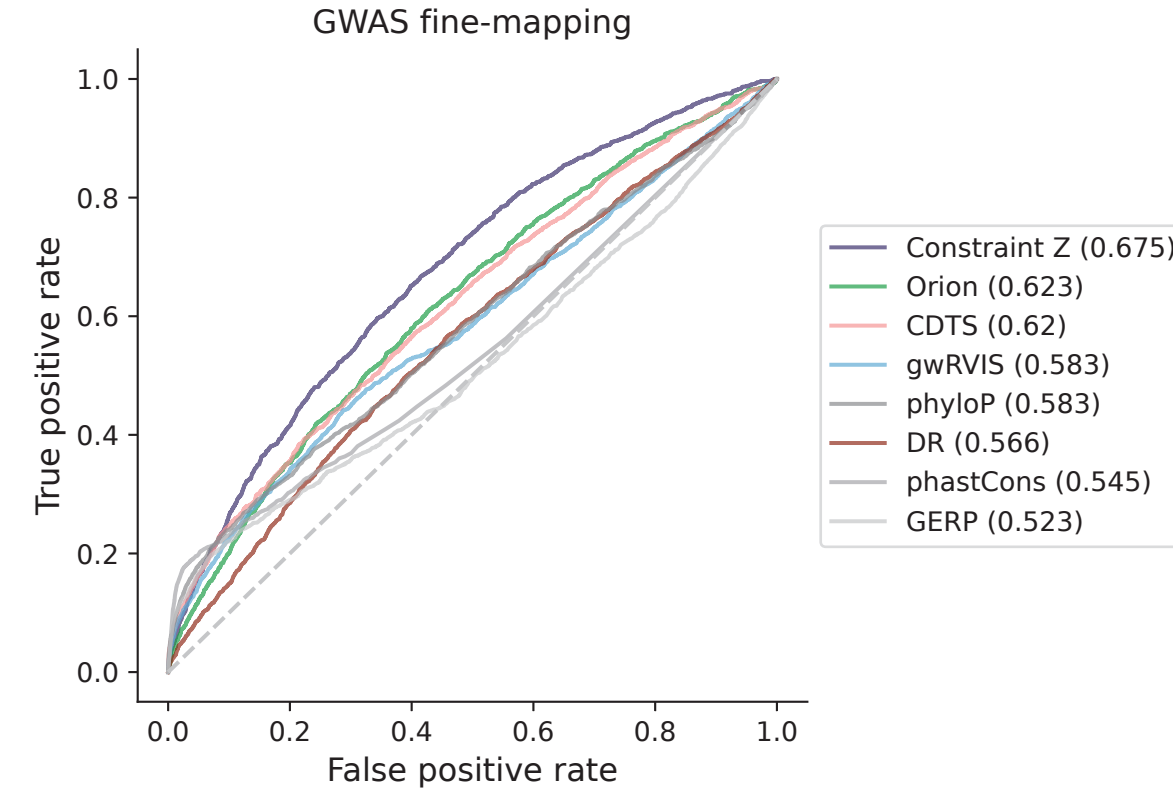
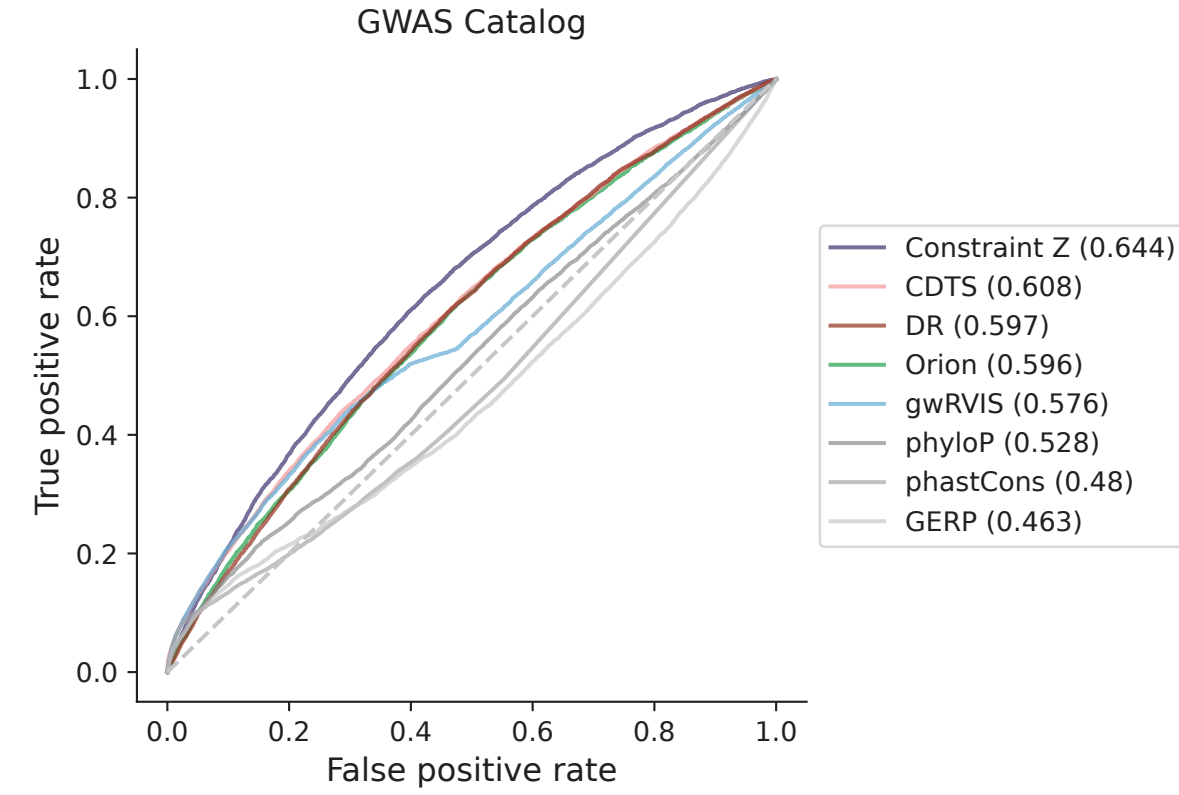
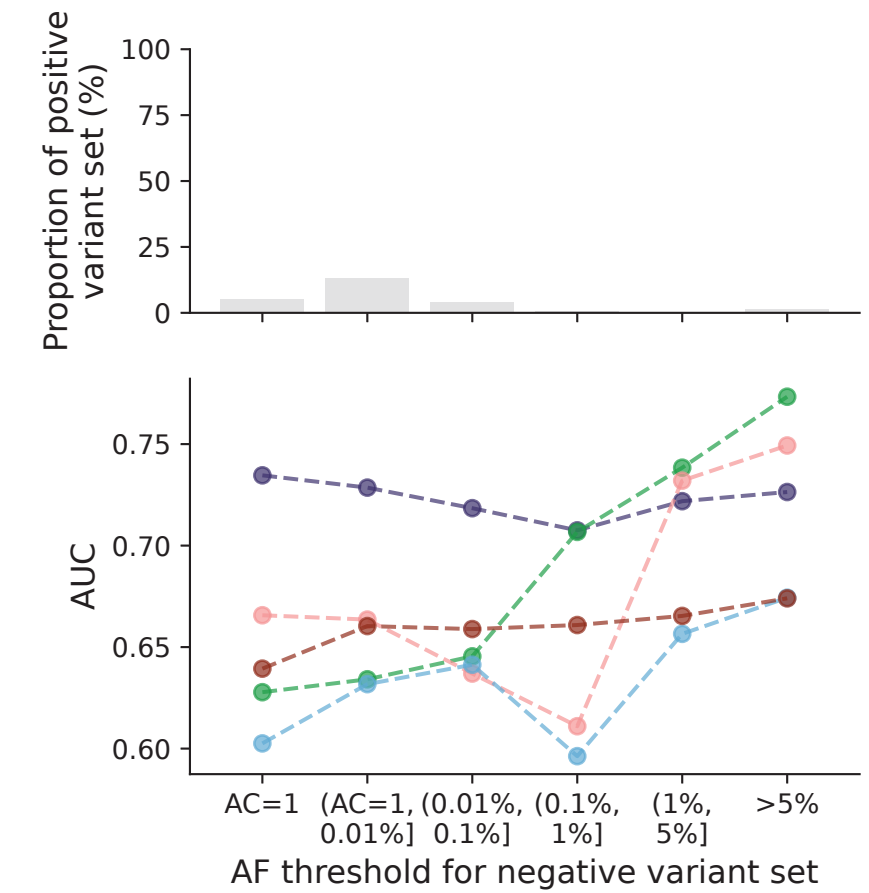
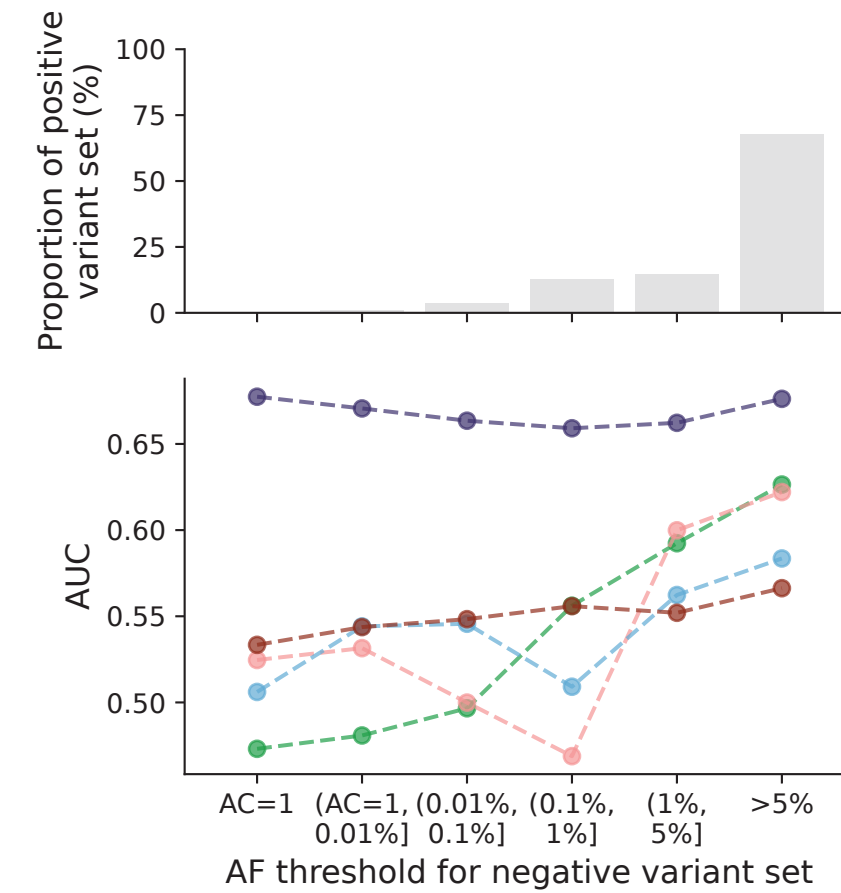
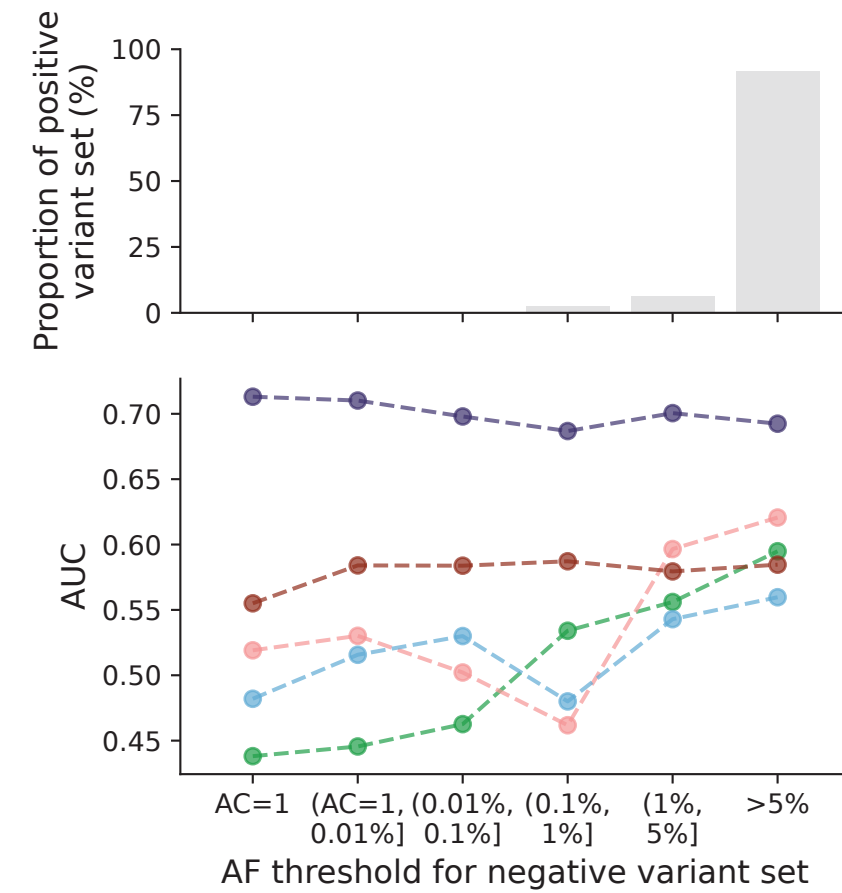
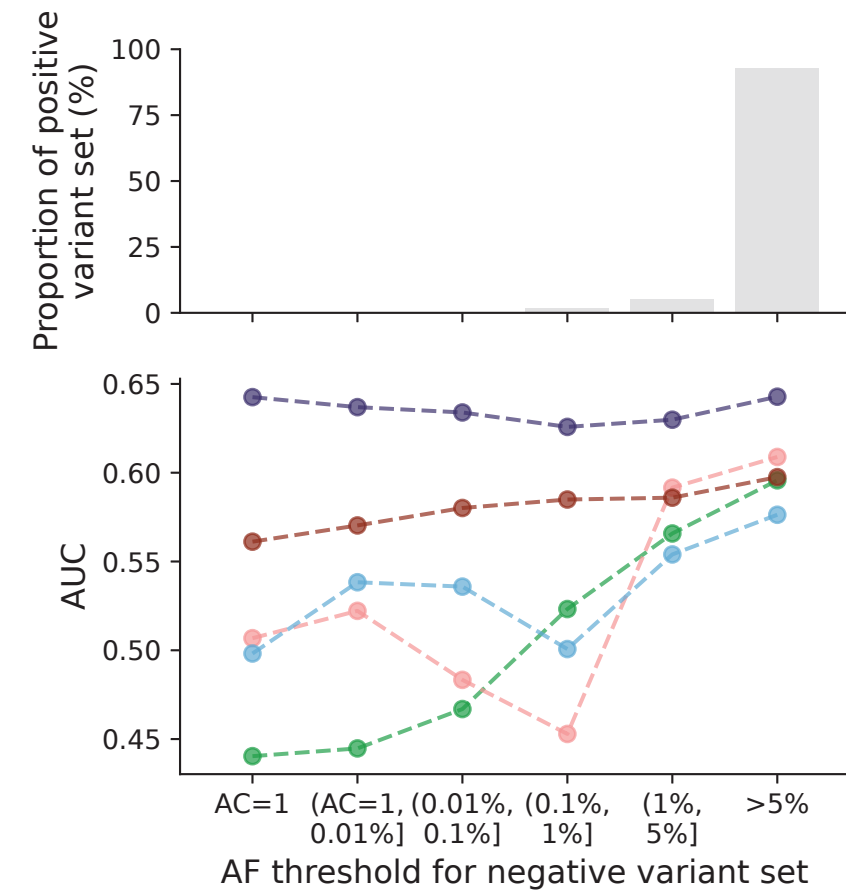


Extended Data Figure 4

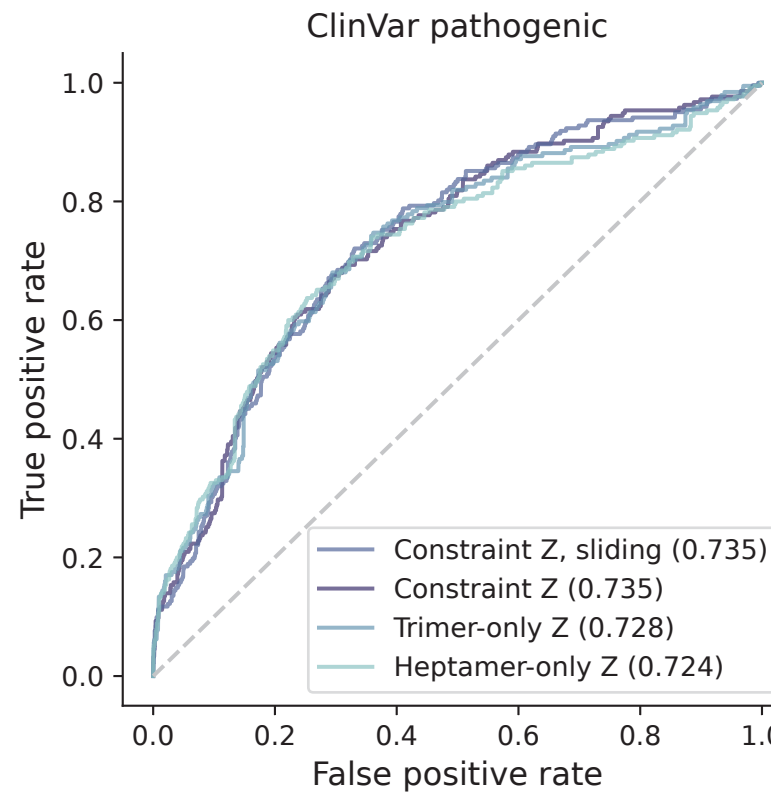
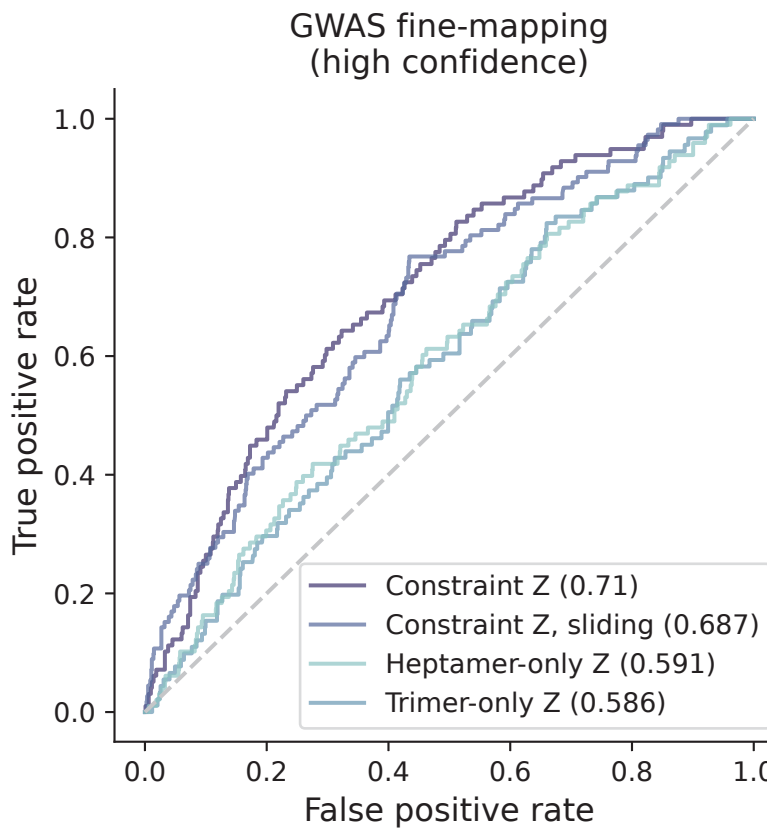
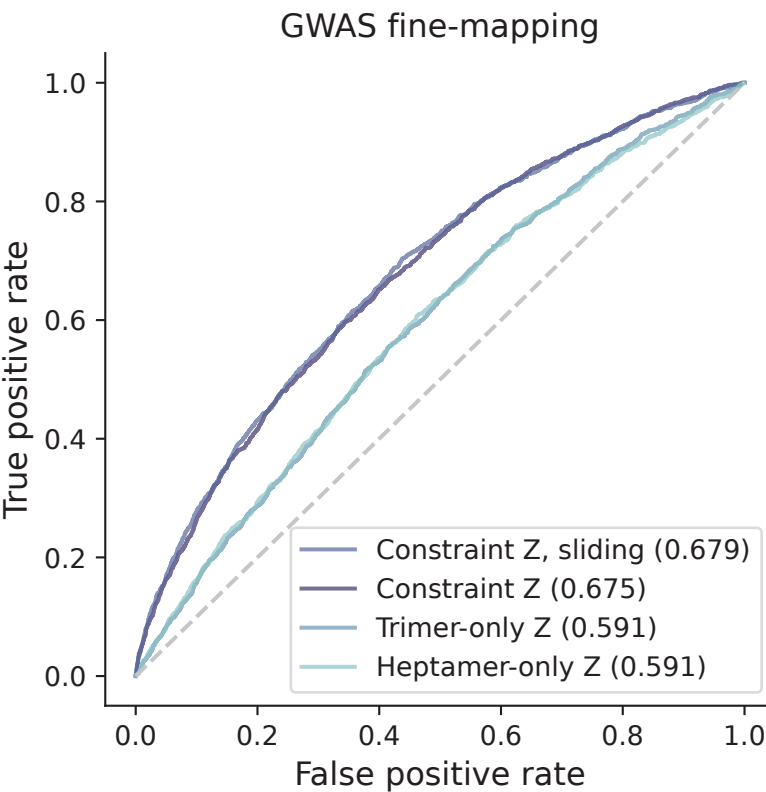
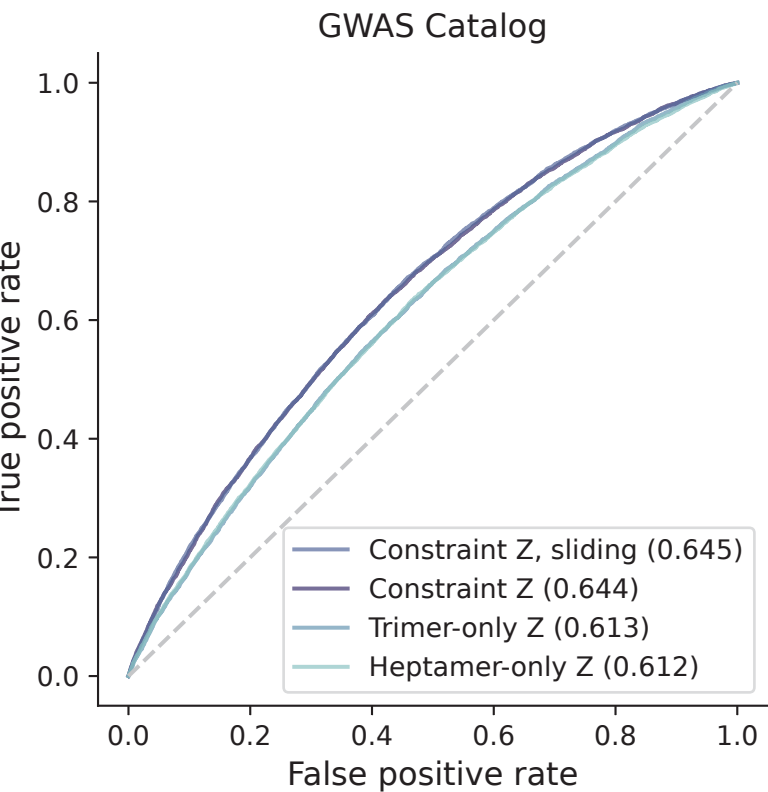


Extended Data Figure 5

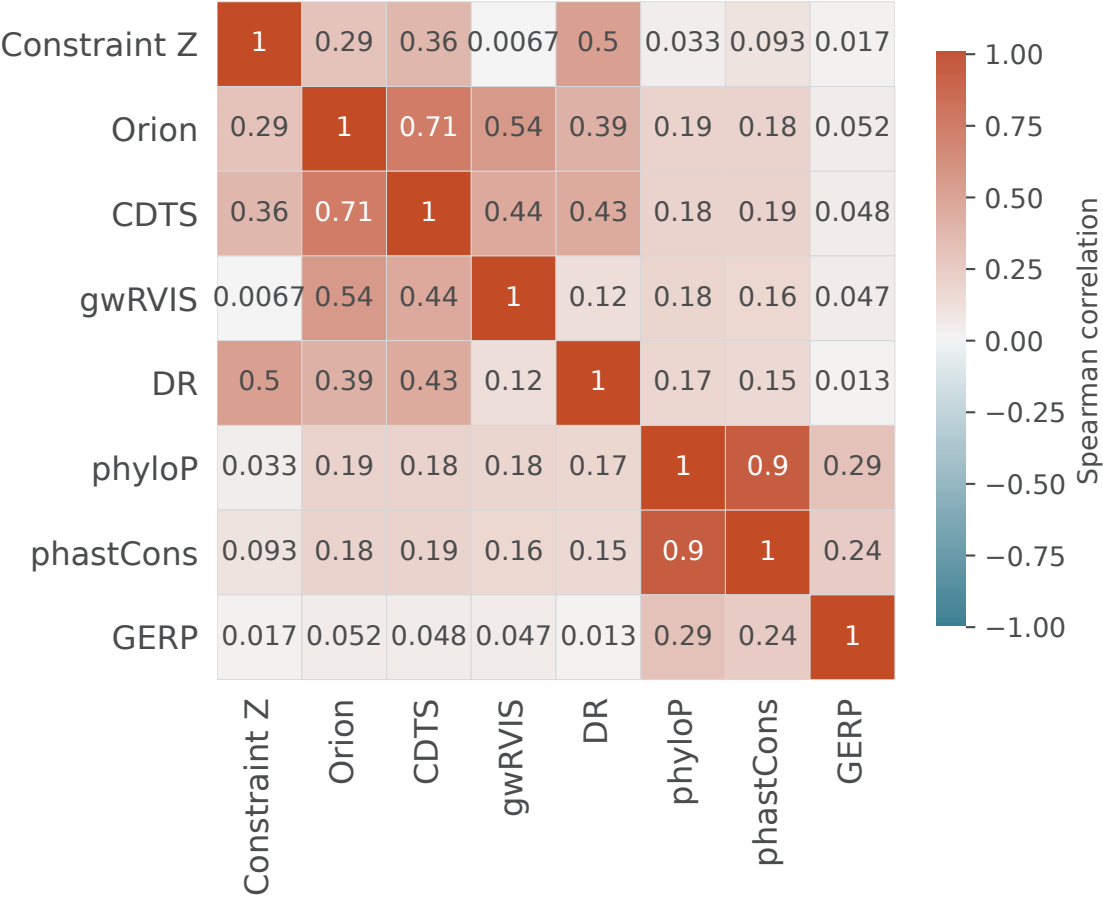
bioRxiv preprint doi: <https://doi.org/10.1101/2022.03.20.485034>; this version posted October 10, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

**b**

Extended Data Figure 6



Extended Data Figure 7



Extended Data Figure 3

