

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

Siwei Chen<sup>1,2,†</sup>, Laurent C. Francioli<sup>1,2,†</sup>, Julia K. Goodrich<sup>1,2</sup>, Ryan L. Collins<sup>1,3,4</sup>, Qingbo Wang<sup>1,5</sup>, Jessica Alföldi<sup>1,2</sup>, Nicholas A. Watts<sup>1,2</sup>, Christopher Vittal<sup>1,2</sup>, Laura D. Gauthier<sup>6</sup>, Timothy Poterba<sup>1,2,7</sup>, Michael W. Wilson<sup>1,2</sup>, Yekaterina Tarasova<sup>1</sup>, William Phu<sup>1,2</sup>, Mary T. Yohannes<sup>1</sup>, Zan Koenig<sup>1</sup>, Yossi Farjoun<sup>6</sup>, Eric Banks<sup>6</sup>, Stacey Donnelly<sup>7</sup>, Stacey Gabriel<sup>1,7</sup>, Namrata Gupta<sup>1,7</sup>, Steven Ferriera<sup>7</sup>, Charlotte Tolonen<sup>6</sup>, Sam Novod<sup>6</sup>, Louis Bergelson<sup>6</sup>, David Roazen<sup>6</sup>, Valentin Ruano-Rubio<sup>6</sup>, Miguel Covarrubias<sup>6</sup>, Christopher Llanwarne<sup>6</sup>, Nikelle Petrillo<sup>6</sup>, Gordon Wade<sup>6</sup>, Thibault Jeandet<sup>6</sup>, Ruchi Munshi<sup>6</sup>, Kathleen Tibbetts<sup>6</sup>, gnomAD Project Consortium<sup>\*</sup>, Anne O'Donnell-Luria<sup>1,2,8</sup>, Matthew Solomonson<sup>1,2</sup>, Cotton Seed<sup>2,9</sup>, Alicia R. Martin<sup>1,2</sup>, Michael E. Talkowski<sup>1,3</sup>, Heidi L. Rehm<sup>1,3</sup>, Mark J. Daly<sup>1,2,10</sup>, Grace Tiao<sup>1,2</sup>, Benjamin M. Neale<sup>1,2,9,‡</sup>, Daniel G. MacArthur<sup>1,2,11,12,‡</sup>, Konrad J. Karczewski<sup>1,2,9</sup>

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

<sup>2</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

<sup>3</sup>Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA

<sup>4</sup>Division of Medical Sciences, Harvard Medical School, Boston, MA, USA

<sup>5</sup>Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita, Japan

<sup>6</sup>Data Science Platform, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>7</sup>Broad Genomics, Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA

<sup>8</sup>Division of Genetics and Genomics, Boston Children's Hospital, Boston, MA

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

<sup>10</sup>Institute for Molecular Medicine Finland, (FIMM) Helsinki, Finland

<sup>11</sup>Centre for Population Genomics, Garvan Institute of Medical Research and UNSW Sydney, Sydney, Australia

<sup>12</sup>Centre 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 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 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 characterized 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 reduction of variants in a gene compared to an expectation conditioned on that gene’s mutability. Using coding variant data from sequencing of more than 125K humans<sup>1</sup>, we have previously developed a constraint metric that classifies each protein-coding gene along a spectrum of LoF intolerance<sup>1</sup>, which provides a valuable resource for studying the functional significance of human genes<sup>2-5</sup>. 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 recognized<sup>6-10</sup>.

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, 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, which 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<sup>11,12</sup>.

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., Orion<sup>13</sup>, CDTs<sup>14</sup>); 2) machine-learning classifiers that are trained to differentiate between disease-associated variants and benign variants (e.g., CADD<sup>15</sup>, GWAVA<sup>16</sup>, gwRVIS<sup>17</sup>); 3) evolutionary methods that employ phylogenetic conservation scores to predict the consequence of non-coding variants on fitness (e.g., LINSIGHT<sup>18</sup>). 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) strong dependence on the availability of well-characterized functional mutations as training data; 3) compromised ability to detect regions that have only been under selection recently in the human lineage, which may have a functional impact on human phenotypic 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 quantify depletion of variation in a tiled window across the entire genome. By building a fuller picture of genic constraint than solely focusing on coding variation, we demonstrate that it facilitates the functional interpretation of non-coding regions and improves 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 include 36,811 (48.3%) of non-European ancestry, including 20,744 individuals with African ancestries and 7,647 individuals with admixed Amerindigenous ancestries. After stringent quality control (see Supplemental Information), we were left with a set of 644,267,978 high-confidence short nuclear variants (gnomAD v3.1.2), of which 390,393,900 rare (allele frequency [AF]≤1%) single nucleotide variants were used for building the genome-wide constraint map. These correspond to approximately one variant every 4.9 bp (one rare 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 -> ATG), with adjustment for the effect of methylation on mutation rate at CpG sites, which become saturated for mutation at sample sizes above ~10K genomes<sup>19</sup> (Extended Fig. 1a,b; Methods). We then jointly analyzed these local-context-specific estimates with regional genomic features (e.g., replication timing and recombination rate) to establish an expected number of variants. We trained our model on 413,304 *de novo* mutations previously detected in family-based whole-genome sequencing studies<sup>20,21</sup> (Extended Fig. 1c) and applied it to assess constraint across the entire genome (Methods).

We quantified the deviation from expectation for each 1kb window using a Z score<sup>22</sup> (Methods; Extended Fig. 1d,e), which was centered slightly below zero for non-coding regions (median=-0.50), and was significantly higher (more constrained) for windows containing any protein-coding sequences (median=0.17, Wilcoxon  $P<10^{-200}$ ; **Fig. 1a**). The constraint Z score positively correlated with the percentage of coding bases in a window and presented a substantial shift towards higher constraint for exonic sequences calculated from directly concatenating coding exons into 1kb windows (median=1.48; Extended Fig. 2a-c). About 13.5% and 0.5% of the non-coding windows exhibited constraint as strong as the 50<sup>th</sup> and 90<sup>th</sup> 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)<sup>23</sup>, we found a strong correlation (linear regression  $P=1.54\cdot10^{-40}$ , **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 Z score and several annotations of putatively functional non-coding sequences (**Fig. 2a**). First, we found that candidate cis-regulatory elements (cCREs, derived from ENCODE<sup>24</sup> integrated DNase- and ChIP-seq data) are significantly enriched in the most constrained percentile of the genome ( $Z\ge4$ , OR=1.28 compared to the genome-wide average, Fisher's exact  $P=2.72\cdot10^{-44}$ ), indicating that a large fraction of the constrained non-coding regions may serve a regulatory role. Similarly, a significant enrichment was found for an independent set of active, *in vivo*-

transcribed enhancers (identified by FANTOM CAGE analyses<sup>25</sup>; OR=1.36,  $P=1.03 \cdot 10^{-12}$ ). Moreover, we found that super enhancers<sup>26</sup> – groups of enhancers in close genomic proximity – exhibited a stronger enrichment (OR=1.44,  $P=4.98 \cdot 10^{-10}$ ), which is in line with their prominent roles in regulating genes important for cell type specification<sup>27</sup>. We also detected a modest but significant enrichment for long non-coding RNAs (lncRNAs, implicated by the FANTOM CAT catalog<sup>28</sup>; OR=1.21,  $P=2.71 \cdot 10^{-25}$ ), lending support to their potential functionality in the pervasive transcription in the human genome<sup>29</sup>.

A pronounced level of constraint was observed for the 9p21 locus (OR=7.10, Fisher's exact  $P=2.90 \cdot 10^{-4}$ ; **Fig. 2b**), a widely recognized and replicated risk factor for coronary artery disease (CAD) and type 2 diabetes (T2D)<sup>30</sup>. Although devoid of genes, the locus tended to be as constrained as regions encompassing coding sequences (median Z=0.38 versus 0.17, Wilcoxon  $P=0.19$ ;  $P=2.11 \cdot 10^{-4}$  when compared to the rest of the genome). The signal of constraint was correlated with the presence of genetic markers identified by CAD/T2D genome-wide association studies<sup>31</sup> (GWAS; median Z=2.04 versus 0.003, Wilcoxon  $P=0.017$ ) and coincided with the high density of regulatory elements<sup>24</sup> within the locus (median Z=0.74 versus -0.70, Wilcoxon  $P=0.011$ ; **Fig. 2c**). A well-replicated CAD marker rs10738607<sup>32-34</sup>, for example, was found to reside in a region under strong selection (Z=4.48, chr9:22088000-22089000) ~80kb upstream of the cell cycle suppressor gene *CDKN2B*, where three consecutive enhancer signatures occupied >80% of its sequence<sup>24</sup>. The risk allele of rs10738607 was predicted to disrupt a binding site for STAT5A<sup>30</sup>, a known transcriptional activator of *CDKN2B* that negatively regulates cell proliferation<sup>35,36</sup>; reduced expression of *CDKN2B* and enhanced proliferation have been previously-recognized phenotypes linked to the 9p21 CAD risk locus<sup>37-41</sup>. Therefore, our constraint analysis adds weight to the hypothesis that loss of *STAT5A-CDKN2B* antiproliferative effect may present one mechanism that explains the genotype-phenotype association for 9p21 risk alleles.

### Prioritizing non-coding variants through the constraint map

The genome-wide constraint map allows us to systematically evaluate each genetic variant in the genome, particularly expanding our ability to study variants in the vast and under-characterized non-coding regions. Examining the distribution of non-coding variants identified by GWAS on the constraint spectrum, we found a significant enrichment for non-coding GWAS hits in the constrained end of the genome (621/14,808 constrained windows [ $Z \geq 4$ ] overlapped with GWAS Catalog<sup>31</sup> annotations, OR=1.47 compared to the genome-wide average of 48,045/1,665,599, Fisher's exact  $P=3.76 \cdot 10^{-19}$ , **Fig. 3a**; Methods). The enrichment appeared to become increasingly stronger for hits that had a replication experiment (internally replicated within the same study: OR=1.54,  $P=1.39 \cdot 10^{-9}$ ; externally replicated by a different study: OR=1.76,  $P=1.51 \cdot 10^{-6}$ ). Furthermore, substantial selection signals were found for likely causal GWAS variants fine-mapped from 94 complex diseases and traits in UK Biobank (UKB)<sup>42,43</sup> (**Fig. 3b**; Methods). Leading enrichments were found for pathological conditions including CAD, inguinal hernia, insomnia, fibroblastic disorders, and hypothyroidism (OR>2.5), as well as quantitative molecular traits such as hemoglobin, uric acid, and bilirubin levels. These results revealed a high positive correlation between the level of constraint and the occurrence of candidate functional variants in the non-coding genome. Thus, our constraint metric could effectively serve to prioritize non-coding variants discovered in large-scale human disease or trait association studies.

For example, in the constraint analysis of CAD GWAS, we found seven variants from a 95% credible set located in highly constrained non-coding regions related to the gene *PLG* – four were mapped to its antisense transcript ENST00000659713 (rs144679016, rs116480834, rs10945688, rs148517610), one resided in its 9<sup>th</sup> intron (rs143556831), and two colocalized with its enhancer marks ~55-57kb downstream of the gene (rs12529023, rs139095347). *PLG* encodes the plasminogen protein that circulates in blood

plasma and is converted to the active protease, plasmin, which dissolves the fibrin of blood clots. Dysregulation of the PLG-plasmin system has been frequently implicated in the pathogenesis of CAD<sup>44-49</sup>. Therefore, together with previous evidence on *PLG*, our findings suggest a plausible role of the prioritized non-coding variants in increasing CAD risk, likely acting through changes in the constrained regulatory elements of *PLG*.

### Exploring non-coding dosage sensitivity in the constrained genome

Copy number variants (CNVs) or dosage alterations (deletions/loss or duplications/gain) of DNA are well-established risk factors for human developmental disorders<sup>50-55</sup>. CNVs can exert pathogenic effects in more than one mechanism<sup>56,57</sup>, such as direct gene dosage effect, positional effect, and transvection effect. Still, the most common is alteration of dosage-sensitive genes within the loci, whereas the functional relevance of the large proportion of non-coding sequences being affected are much less understood. With 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 ~90K CNVs from a genome-wide CNV morbidity map of developmental delay and congenital birth defects<sup>58,59</sup>. We observed a substantial excess of CNVs that affect constrained non-coding regions ( $Z \geq 4$ ) among individuals with developmental disorders (DD cases) in comparison to healthy controls (42.2% versus 10.9%, OR=6.0, Fisher's exact  $P < 10^{-200}$ , **Fig. 4a**; Methods). More strikingly, of the 18 loci that had been previously classified as pathogenic<sup>58</sup>, 16 (88.9%) coincided with constrained non-coding regions. The high incidence was recapitulated in a curated set of ~4K putatively pathogenic CNVs (ClinVar<sup>60</sup>), of which 82.8% exhibited high non-coding constraint (**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 subtype of CNV deletions and duplications (**Fig. 4b**; Methods). Non-coding constraint manifested a comparable level of association as gene constraint when modeled simultaneously to predict DD CNVs ( $\log[OR] = 1.13$  and  $1.48$ , Logit  $P < 2 \cdot 10^{-16}$  for both), suggesting that dosage imbalance of constrained non-coding regions could confer risk for developmental disorders in addition to gene dosage sensitivity.

One classic example of non-coding dosage imbalance is a set of duplications involving the regulatory domain of *IHH* associated with synpolydactyly and craniosynostosis<sup>61-63</sup>. The four implicated duplications covered a stretch of ~102kb sequence upstream of *IHH*, with a smallest ~10kb critical overlapping region (**Fig. 4c**). The critical region contained no exons but showed high non-coding constraint (median  $Z = 1.22$ , Wilcoxon  $P = 0.011$  compared to the rest of the genome). Interestingly, the highest constraint signal colocalized 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 skull<sup>63</sup>. Collectively, we suggest that non-coding constraint can be a useful indicator of dosage-sensitive regulatory CNVs and improves understanding of potential non-coding mechanisms underlying CNV disease associations.

### Leveraging non-coding constraint to improve gene function characterization

Given the significant roles 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 surveyed enhancer-gene links from the Roadmap Epigenomics Enhancer-Gene Linking database<sup>64</sup> (Methods), where we found overall, more constrained non-coding regions were more frequently linked to a gene (**Fig. 5a**), and more constrained enhancers tend to be associated with more constrained genes (e.g., haploinsufficient genes; median  $Z = 2.21$  versus  $1.43$ , Wilcoxon  $P = 4.24 \cdot 10^{-18}$ , **Fig. 5b**; Methods).

On the other hand, a particular interesting set of associations are the links between constrained enhancers and the “unconstrained” genes classified by gene constraint metrics, because these links may reflect functional significance of the “unconstrained” genes that had been previously underrecognized. More than 40% of the least constrained genes (last decile scored by LOEUF [loss-of-function observed/expected upper bound fraction]<sup>1</sup>) were linked with a constrained enhancer (last decile with constraint  $Z \geq 2.23$ ); the lack of predicted gene constraint can be explained by the intrinsic design of LOEUF as a measure of intolerance to rare LoF variation, where small genes with few expected LoF variants are likely underpowered. Indeed, when stratifying genes by the number of expected LoF variants, we found a significantly higher enhancer constraint for genes that were underpowered (with  $\leq 5$  expected LoF variants by LOEUF) compared to those that were sufficiently powered and scored as unconstrained (median  $Z=2.06$  versus  $1.68$ , Wilcoxon  $P=1.68 \cdot 10^{-3}$ , **Fig. 5b**). This pattern suggests that a portion of these underpowered genes may play an important functional role 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 exomes<sup>1</sup>; 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 ( $Z=6.26$ ), located ~16kb upstream of the gene, where >40% of the expected variants were depleted (200.6 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* expression<sup>65</sup>; elevated *ASCL2* expression has been implicated in the development and progression of several human cancers<sup>66-68</sup>. This example highlights the value of non-coding constraint – as a complementary metric to gene constraint – for identifying functionally important genes.

Gene-set enrichment analysis of the underpowered genes with constrained enhancers identified strong enrichment for genes encoding histones, secreted proteins involved in cell signaling and communication (e.g., cytokines, hormones, growth factors; **Fig. 5c**). These genes are generally small and are thus likely to be underrecognized by LoF-intolerance-based metrics. Applying an alternative, phylogeny-based conservation score<sup>69</sup> to evaluate these underpowered genes, we found a significantly higher evolutionary constraint on these genes compared to those with less constrained enhancers (Wilcoxon  $P=4.28 \cdot 10^{-12}$ ; **Fig. 5d**), providing additional evidence for their functional significance. Moreover, by dissecting enhancers in a tissue-specific manner, we examined how tissue-specific enhancer constraint correlates with tissue-specific gene regulation. We found that even conditioning on gene constraint, enhancer constraint remained a significant predictor for the expression level of target genes in matched tissue types (**Fig. 5e**; Methods). These results further support the application of our constraint metric to improve the characterization of human gene function.

A practical implementation of this notion would be refining gene constraint models to incorporate contributions from constrained regulatory elements. This essentially borrows power from extending the functional unit of a gene in the context of the gene regulatory network. Such an approach would allow constraint modeling for specific tissue/cell types and conditions given the diverse range of epigenomics data. Taking enhancer-gene links from brain tissue as a proof-of-principle, a simple logistic regression test conditioning on gene constraint indicated significant contribution from enhancer constraint for predicting functionally relevant genes (e.g., targets of Fragile X Mental Retardation Protein (FMRP)<sup>70</sup>, Logit  $P=2.84 \cdot 10^{-8}$ ; Methods). While we acknowledge that the biological consequences of mutations in enhancers are not clearly understood and thus natural selection may differ in its interest depending on mechanistic consequence, an extended model to incorporate non-coding variation information in a biologically-informed way hold the promise to provide a finer characterization of gene function and facilitate a better 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 LoF variation for each protein-coding gene<sup>1,19</sup>. 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 mutation rate for all substitution and trinucleotide context 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 to distinguish between selection and neutral variation. To this end, we intended not to include features that are directly reflective of regulatory functions in our model, such as histone modifications as commonly used in classifiers predicting functional non-coding variants. Instead, we employed such features as part of our validation, where we showed that our constraint Z score is positively correlated with the regulatory elements derived from epigenomic signatures. Meanwhile, to demonstrate the ability of our metric in prioritizing non-coding variants independent of established regulatory elements, we re-examined the enrichment analyses with previously annotated regulatory sequences excluded – all results persisted (Extended Fig. 3a,b). Likewise, we demonstrated that our constraint metric captures additional information to phylogeny-based conservation scores, which are likely to be blind to functional non-coding regions that have high evolutionary turnover (e.g., enhancers<sup>71–76</sup>; Extended Fig. 4). We also note that we restricted all our validation analyses to non-coding regions to explicitly evaluate the metric for characterizing the non-coding genome, and we further eliminated potential bias from nearby genes by recapitulating the results within regions >10kb away from any protein-coding genes (Extended Fig. 5). Further validating our metric on experimental data from multiplex assays of variant effect (MAVE)<sup>77</sup>, we showed that our constraint Z score correlates well with the experimental measurements, along with comparison to other genome-wide predictive scores (Extended Fig. 6; Methods). Altogether, our constraint metric presented reliable and consistent performance in identifying important non-coding regions in the human genome.

Our analyses revealed considerable correlation between constraint of the non-coding regulatory elements and the functional importance of their target genes. The implication is twofold. First, the constraint metric can be applied to further prioritize existing regulatory annotations. For instance, categorizing ENCODE cCREs by constraint Z score identified an increasingly stronger GWAS signals in the more constrained cCREs (Extended Fig. 3c). This also suggests the integration of multiple lines of evidence indicative of function. Second, the prioritization of regulatory elements can be applied to improve identification of important genes. The example of *ASCL2* demonstrates how leveraging the non-coding functional variation of a gene can build out a clearer picture of its importance from a regulatory point of view. The analyses of tissue-specific enhancer constraint further demonstrate how non-coding constraint could be applied to refine gene constraint modeling.

Despite the clear constraint signal identified for non-coding regions, many limitations exist. First, the mutational model explains ~75% of the variation ( $R^2=0.746$ ) in *de novo* mutation rate, indicating

underrecognized sources or properties of spontaneous human mutation. This situation can be partially alleviated by an increased sample size of *de novo* mutation data, which would allow inclusion of additional genomic features at smaller scales as well as more accurate estimation of feature contributions. Even with a sufficiently large dataset, however, comprehensive modeling of the heterogeneity in the mutation rate requires advanced knowledge about the underlying mutational mechanisms and ideally a more interpretable, biologically-informed statistical framework. Further, the practical 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 directly built in current constraint scores nor in the mutational dataset, thus downstream analysis of functional genomics data (e.g., incorporating tissue-specific enhancer signatures as in this study) is often necessary for justifying specific biological implications. It should also be noted that, since the detection of depletion of variation is immune to negative selection after reproductive age, genomic regions involved in late-onset phenotypes are therefore 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. Using constraint seen in coding regions as benchmarks (Extended Fig. 7), we are currently well-powered to detect non-coding regions with constraint as strong as the 90<sup>th</sup> percentile of coding exons of similar size, and we estimate a sample size of ~700K to detect constraint that is about the average of coding regions. Moreover, as functional non-coding elements are generally small, to further increase our resolution, for example at a 100bp scale, we would need ~4.8M individuals to detect individual constrained non-coding elements.

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

### Genome data aggregation, variant-calling, and quality control

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 efforts<sup>1</sup>. We harmonized these data using the GATK Best Practices pipeline and joint-called all samples using Hail<sup>78</sup>, and developed and utilized an updated pipeline of sample, variant, and genotype QC 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).

### Mutation rate model

To estimate the baseline mutation rate for each substitution and context, we tallied each trinucleotide context in the human genome. As previously shown<sup>19</sup>, the methylated CpG variants begin to saturate at a sample size of ~10K genomes, and therefore we downsampled the gnomAD dataset to 1,000 genomes for use in calculating the mutation rate. Sites were further excluded if there were variants observed in gnomAD but flagged as low quality, or found in greater than 5 copies in the downsampled set. Using this dataset, we computed the proportion of possible variants observed for each trinucleotide change ( $XY_1Z \rightarrow XY_2Z$ ), with CpG sites stratified by their methylation levels (see DNA methylation analysis). These proportions represent the relative probability of a given nucleotide mutating into one of the three other possible bases in a trinucleotide context, and we scaled this factor so that the weighted genome-wide average is the human per-base, per-generation mutation rate ( $1.2 \cdot 10^{-8}$ ) to obtain the absolute mutation rates  $\hat{f}$ . The  $\hat{f}$  estimates were well correlated with the proportion of possible variants observed in the 76,156 gnomAD genomes ( $R^2=0.999$ , Extended Fig. 1a,b; Supplementary Data 1), a dataset of 6,079,733,538 possible variants at 2,026,577,846 autosomal sites with 30-32X coverage, where 390,393,900 high-quality rare ( $AF \leq 1\%$ ) variants were observed. These fitted proportion observed (mutabilities) were then used for establishing the context-dependent expected variation in construction of constraint Z scores.

### DNA methylation analysis

To correct for the effect of DNA methylation on 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 (from averaging different tissues<sup>1</sup>), 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)<sup>79</sup> and six from primordial germ cells (7Wk, 10Wk, 11Wk, 13Wk, 17Wk, and 19Wk)<sup>80</sup>. 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. In order 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], (0.3,0.5], (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.

### 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:

$$\chi^2 = (Obs - Exp)^2 / Exp$$
$$Z = \{\sqrt{\chi^2} \text{ if } Obs < Exp - \sqrt{\chi^2} \text{ if } Obs \geq Exp\}$$

Here, the observed variant count (*Obs*) is the number of unique rare (AF<0.1%) variants in a 1kb window identified in the gnomAD dataset. To establish the expected number of variants (*Exp*), we first created a context-dependent mutability for each window by summing all trinucleotide mutabilities (proportion observed modeled by  $\hat{}$ ) within the 1kb sequence. At the same time, we computed for each window 17 genomic features that have been shown to influence mutation rate (e.g., replication timing and recombination rate; see Collection of genomic features). To determine the relative contribution of the 18 features (trinucleotide context plus 17 genomic features), we trained these features to predict the occurrence of *de novo* mutations (DNMs), as a proxy of spontaneous mutations, using a random forest regression model. A set of 413,304 unique DNMs were compiled from two large-scale family-based whole-genome sequencing studies{Halldorsson, 2019 #17;An, 2018 #18}. Given the data sparsity, we counted the incidence of DNMs per 1Mb and regressed it against the 18 features computed correspondingly ( $R^2=0.746$  by a 90/10 train/test split, Extended Fig. 1c). The fitted model was then applied to establish the expected variation (*Exp*) in the gnomAD dataset.

Constraint Z scores were created for 2,442,347 non-overlapping 1kb windows across the human genome (2,330,252 on autosomes and 112,095 on chromosome X). Due to the lack of DNM data on chromosome X, a model was trained using autosomal regions and was extrapolated to chromosome X for computing constraint scores. We performed downstream analyses separately for autosomes and chromosome X and presented the former as primary, with the latter provided in Extended Fig. 8. In 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, and 3) the mean coverage in the gnomAD genomes was between 25-35X (or 20-25X for chromosome X). This resulted in 1,797,153 windows (73.6% of initial) for the primary analyses (Supplementary Data 2), with 131,554 encompassing coding sequences and 1,665,599 exclusively non-coding.

### Collection of genomic features

The 17 regional genomic features used for predicting the expected variation are 1) replication timing<sup>81</sup>, 2) male and 3) female recombination rate<sup>20</sup>, 4) GC content<sup>82</sup>, 5) low-complexity region<sup>83</sup>, 6) short and 7) long interspersed nuclear element<sup>82</sup>, distance from the 8) telomere and the 9) centromere<sup>82</sup>, 10) CpG island<sup>82</sup>, 11) nucleosome density<sup>81</sup>, 12) maternal and 13) paternal DNM cluster<sup>84</sup>, DNA methylation in 14) sperm<sup>79</sup>, 15) oocyte<sup>79</sup>, 16) preimplantation embryo<sup>79</sup>, and 17) primordial germ cell<sup>80</sup>. Data were downloaded from the referenced resources, lifted over to GRCh38 coordinates when needed using CrossMap<sup>85</sup>, and files in .bed or .BigWig format were processed using bedtools<sup>86</sup> and bigWigAverageOverBed<sup>81</sup> to obtain feature values within specific genomic windows.

### Correlation between constraint Z and APS

As an internal validation, we compared our constraint Z score against the SV constraint score APS<sup>23</sup>. For each SV from the original study (gnomAD-SV), 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 205,699 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 bins, with 95% confidence intervals computed from 100-fold bootstrapping.

### Correlation of constraint Z with candidate functional non-coding elements

We validated the constraint metric using a number of external functional annotations, including 926,535 ENCODE cCREs<sup>24</sup>, 63,285 FANTOM5<sup>25</sup> enhancers, 19,175 FANTOM6<sup>28</sup> lncRNAs, 331,601 super enhancers (SEdb<sup>26</sup>), 111,308 GWAS Catalog<sup>31</sup> variants (with an association  $P$ -value  $\leq 5.0 \cdot 10^{-8}$ ; 32,708 with an internal replication and 6,748 with an external replication), 601,191 UKB GWAS variants fine-mapped from 94 heritable traits (95% credible set)<sup>42</sup>, and 90,423 unique CNVs from a CNV morbidity map of developmental delay<sup>58,59</sup>.

To assess the correlation between constraint Z scores and the collected candidate functional elements, we intersected each annotation with the scored 1kb windows binned by constraint Z ( $<-4$ ,  $[-4,-3]$ ,  $[-3,-2]$ ,  $[-2,-1]$ ,  $[-1,-0]$ ,  $[0,1]$ ,  $[1,2]$ ,  $[2,3]$ ,  $[3,4]$ ,  $\geq 4$ ), 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 constraint<sup>1</sup> 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.

### **Coordination between non-coding constraint and gene constraint**

To examine whether constraint of non-coding regulatory elements implicates 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 ClinGen<sup>87</sup> haploinsufficient genes, 2,454 MGI<sup>88</sup> essential genes mapped to human orthologs, 1,771 OMIM<sup>89</sup> autosomal dominant genes, and 1,920 LOEUF<sup>1</sup> 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 decile). The LOEUF underpowered list included 1,117 genes with  $\leq 5$  expected LoF variants. Enhancers associated with 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<sup>90,91</sup>. For each gene, we aggregated and merged enhancers predicted from all 127 reference epigenomes and assigned the most constrained enhancer to each gene for analysis of enhancer-gene constraint coordination (Supplementary Data 3).

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 metadata<sup>64</sup>). For each gene and tissue type, we searched for tissue-specific gene expression in the Genotype-Tissue Expression (GTEx<sup>65</sup>) database (RNASeQCv1.1.9) and computed a normalized median expression for each gene ( $\log_2(\text{TPM}+1)$ ). Enhancer constraint and gene expression levels were calculated for 11 matched tissue types and their correlation within each tissue type was evaluated by regressing gene expression on the enhancer constraint of corresponding gene. Importantly, we included gene constraint (LOEUF score) as a covariate to test the additional contribution from enhancer constraint. Similarly, we used a logistic regression model conditioning on LOEUF to examine the value of brain-tissue enhancer constraint in identifying genes functional in brain. As a proof-of-principle, we regressed genes encoding the mRNA targets of Fragile X Mental Retardation Protein (FMRP<sup>70</sup>; 0 or 1) against the constraint of enhancers derived from brain tissue, including LOEUF score as a covariate.

### Correlation of constraint Z and other predictive scores with MAVE

We compared our metric with other four genome-wide predictive scores – Orion<sup>13</sup>, CDTs<sup>14</sup>, gwRVIS<sup>17</sup>, and JARVIS<sup>17</sup>) – in their correlation with experimental measurements on 11 enhancers tested by MAVE<sup>77</sup>. Each predictive score was downloaded from the original study, lifted over to GRCh38 (for Orion), and applied to score enhancers by taking the average over corresponding genomic regions. Enhancers were ranked by the proportion of mutations identified by MAVE as causing significant changes in gene expression (high to low: *SORT1*, *IRF4*, *IRF6*, *ZRS*, *ZFAND3*, *RET*, *TCF7L2*, *MYC* rs11986220 [*MYCs2*], *BCL11A*, *UC88*, *MYC* rs11986220 [*MYCs1*]), and the rank correlation between MAVE and each predictive score was evaluated using a Spearman's rank correlation test. The correlation was computed on all scorable enhancers for each score, as well as enhancers that are scorable by all metrics ( $n=7$ ; *RET* and *IRF4* were not scored by constraint Z due to feature missingness near the telomere/centromere and *MYCs1* and *SORT1* were excluded due to quality filtering yet their scores are provided in the “Unfiltered Z” penal. In addition, we excluded *UC88* in favor of other scores as it appeared to be an “outlier” deflecting their ranks (Extended Fig. 6).

## Figure legends

**Fig. 1:** Distribution of constraint Z scores across the genome. **a**, Histograms of constraint Z scores for 1,797,153 1kb windows across the human autosomes. Windows encompassing coding regions ( $n=131,554$  with  $\geq 1$ bp coding sequence; red) overall exhibit a higher constraint Z (stronger negative selection) than windows that are exclusively non-coding ( $n=1,665,599$ ; blue). **b**, The correlation between constraint Z score and the adjusted proportion of singletons (APS) score developed for structural variation (SV) constraint. A collection of 205,699 autosomal SVs were assessed using constraint Z score by assigning each SV the highest Z among all overlapping 1kb windows, which shows a strong positive correlation with the SV constraint metric APS. Error bars indicate 100-fold bootstrapped 95% confidence intervals of the mean values.

**Fig. 2:** Evaluation of constraint in putatively functional non-coding regions. **a**, Distributions of candidate regulatory elements 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. **b,c**, The level of constraint of the 9p21 locus, in comparison to coding regions (**b**) and in concordance with external functional annotations (**c**). For **b**, the kernel density estimate (KDE) plots display the distribution of constraint Z scores of the 9p21 locus (blue) and regions encompassing coding sequences (red); the scatter plot shows the enrichment of 9p21 locus across constraint Z bins, with error bars indicating 95% confidence intervals of the odds ratios. For **c**, each bar shows the constraint Z score of a 1kb window within the locus; gaps indicate windows removed by quality filters.

**Fig. 3:** Prioritization of non-coding variants for genome-wide association studies (GWAS). **a**, The distribution of GWAS Catalog variants on the non-coding constraint spectrum. Variants were further partitioned by the degree of experimental replication (Int/Ext repl: internal/external replication), with the latter exhibiting the highest constraint. **b**, Enrichment of UKB GWAS variants in constrained non-coding regions ( $Z \geq 4$ ). Analysis was performed separately for 94 complex diseases and traits, and results with nominal significance are shown, ordered by the lower bound of 95% confidence interval. For both panels, enrichment was evaluated by comparing the proportion of non-coding 1kb windows (at a given constraint level) that contain at least one GWAS variant to the genome-wide average. Error bars indicate 95% confidence intervals of the odds ratios; for **b**, only lower bounds are shown for presentation purposes.

**Fig. 4:** Contribution of non-coding constraint in evaluating copy number variants (CNVs). **a**, The proportion of constrained CNVs ( $Z \geq 4$ ) identified in individuals with developmental delay (DD cases) in comparison to healthy controls. The non-coding constraint of a given CNV was assessed by the highest constraint Z score among all non-coding windows being affected. CNVs were categorized and the frequency of constrained CNVs ( $Z \geq 4$ ) was compared based on pathogenic potential: constrained CNVs are more common in DD cases than controls ( $6,631/15,717=42.2\%$  versus  $8,107/74,706=10.9\%$ ) and are most frequent for CNVs previously implicated as pathogenic ( $16/18=88.9\%$  by DD and  $3,219/3,886=82.8\%$  by ClinVar). **b**, The contribution of non-coding constraint in predicting CNVs in DD cases versus controls. Non-coding constraint remains a strong predictor for the case/control status of CNVs after adjusting for gene constraint (LOEUF score), gene number, and size of CNVs, using a logistic regression test. 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 bars) overlap in a  $\sim 10$ kb region that exhibit high non-coding constraint (blue), with the highest Z score coinciding with the major

*IHH* enhancers (dark blue). Each blue bar shows the constraint Z score of a 1kb window within the locus; gaps indicate windows removed by quality filters.

**Fig. 5:** Relationship between non-coding constraint and gene constraint. **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 ("LOUEF underpowered") present a higher constraint than those powered yet unconstrained genes ("LOUEF unconstrained"). **c**, Gene sets enriched for the LOUEF underpowered genes. Enrichment analysis was performed for a set of 506 genes underpowered by LOUEF while being linked with a constrained enhancer. Red dashed line indicates FDR=0.05. **d**, Comparison of evolutionary constraint for the LOUEF underpowered genes with and without a constrained enhancer. PhastCons, a phylogeny-based conservation metric, supports the functional importance of genes with constrained enhancers. **e**, Correlations between enhancer constraint and gene expression in specific tissue types. Enhancer constraint was reprocessed in a tissue-specific manner and was modeled to predict the expression level of target genes in matched tissue types; a linear regression test suggests significant contribution from enhancer constraint even conditioning on gene constraint (LOUEF score). Error bars indicate 95% confidence intervals of the beta coefficient estimates.

**Extended Data Fig. 1:** Construction of mutational model and constraint Z score. **a,b**, Estimating the relative mutability for each trinucleotide context across the genome. The proportion of possible variants observed in 76,156 gnomAD genomes (mutability; y-axis) is exponentially correlated with the absolute mutation rate estimated from 1,000 downsampled genomes ( $\mu$ ; x-axis). Fit lines were modeled separately for human autosomes (**a**) and chromosome X (**b**). **c**, Estimating the relative contribution of local sequence context and regional genomic features in predicting the expected variation. Mutational model incorporating trinucleotide context and 17 genomic features was trained on DNMAs using a random forest regression model (with a 90/10 train/test split). Trinucleotide mutability appears the most important feature in predicting the expected variation. **d,e**, The distribution of constraint Z score as a function of expected and observed variation. Each point represents Z score for a 1kb window on the genome ( $n=1,797,153$  on autosomes (**d**) and  $n=49,936$  on chromosome X (**e**)), which quantifies the deviation of observed variation from expectation; positive Z (red) indicates deletion of variation ( $\text{obs} < \text{exp}$ ) and the higher the Z the stronger the depletion (i.e., more constrained).

**Extended Data Fig. 2:** Comparison of constraint Z score between coding and non-coding regions. **a**, Exonic regions (1kb windows created by directly concatenating coding exons,  $n=26,987$ ; purple) exhibit a significantly higher constraint Z than windows that are exclusively non-coding ( $n=1,665,599$ ; blue). **b**, The proportion of highly constrained ( $Z \geq 4$ ) windows is positively correlated with the percentage of coding sequences in a window and is substantially higher for the exonic windows. **c**, The proportion of highly constrained ( $Z \geq 4$ ) windows increases linearly as more exonic windows are included in the analysis through random sampling. **d**, Constraint Z score percentiles of non-coding versus exonic windows. About 0.5% (100-99.5) and 13.5% (100-86.5) of the non-coding windows exhibit similar constraint to top 10% (90<sup>th</sup> percentile) and top 50% of exonic regions.

**Extended Data Fig. 3:** Distributions of GWAS variants in non-coding regions with regard to non-coding constraint and ENCODE cCRE annotations. **a,b**, The enrichment of GWAS catalog (**a**) and UKB GWAS (**b**) variants in constrained non-coding regions persists after excluding candidate cis-regulatory elements

(cCREs) annotated by ENCODE. **c**, GWAS variants occur more frequently in cCREs under higher constraint, suggesting the value of non-coding constraint in prioritizing existing functional annotations.

**Extended Data Fig. 4:** Enrichment of enhancers (**a**) and GWAS variants (**b**) across the spectrum of constraint Z and conservation (PhastCons) score. Non-coding 1kb windows were binned by their constraint Z and PhastCons conservation scores, from the least constrained/conserved (1<sup>st</sup> decile) to the most constrained/conserved (10<sup>th</sup> decile), and enrichment within each bin was evaluated by comparing the proportion of 1kb windows that overlap with an enhancer annotation (**a**) or a GWAS hit (**b**) to the genome-wide average. Red color indicates enriched and blue color indicates depleted; odds ratio is shown for nominally significant ( $p < 0.05$ ) enrichment/depletion. As expected, enrichment increases as both scores increase (upper right), yet each score captures independent signals – for instance, regions with high constraint Z scores present significant enrichment across the deciles of conservation score (rightmost column).

**Extended Data Fig. 5:** The enrichment of candidate regulatory elements (**a**) and GWAS variants (**b**, GWAS Catalog; **c**, UKB GWAS) in constrained non-coding regions persists when restricting to regions farther away from protein-coding genes ( $\pm 10\text{kb}$ ).

**Extended Data Fig. 6:** Correlation of constraint Z and other predictive scores with experimental data on enhancers tested by multiplex assays of variant effect (MAVE). **a**, Predictive scores on MAVE-tested enhancers, the importance of which was ranked by the proportion of mutations MAVE identified as causing significant changes in gene expression. Predictive scores of enhancers were aligned below in the same order and were modified when necessary (multiplied by -1 for CDTS and gwRVis) such that a higher value represents higher importance for all scores. An increasing trend (dotted line) indicates positive correlation between the predictive score and the experimental measurement from MAVE (Spearman's rank correlations are shown in 1b). Gaps indicate unavailable values; *RET* and *IRF4* were not scored by constraint Z due to feature missingness near the telomere/centromere and *MYCs1* and *SORT1* were excluded due to quality filtering (yet they are provided in the “Unfiltered Z” panel). **b**, Spearman's rank correlation between MAVE and each predictive score. Constraint Z score exhibits the highest correlation with the experimental measurements, either on all scorable enhancers (All) or the subset of seven enhancers scored by constraint Z (Z scored); the high performance persists even when excluding *UC88* (Exl. UC88) in favor of other scores as it appeared to be an outlier deflecting their ranks. \*Of note, JARVIS and LINSIGHT employed enhancers (as well as other functional annotations) in constructing their scores, which may introduce circularity in favor of their performance in scoring enhancers.

**Extended Data Fig. 7:** The sample size required for well-powered non-coding constraint detection. The proportion of non-coding regions powered to detect constraint ( $Z>4$ ) 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. For a given level of depletion of variation, the minimum number of expected variants to achieve a  $Z>4$  was determined, and the number of individuals required to achieve a given expected number of variants was extrapolated using a linear model of  $\log(\text{number of expected variants}) \sim \log(\text{number of individuals})$ , computed from downsamplings of the gnomAD datasets. Lighter color indicates milder deletion of variation (weaker selection), which requires a larger sample size to detect constraint; the grey dashed vertical line indicates current sample size. Dotted curves (left to right) benchmark the 95<sup>th</sup>, 90<sup>th</sup>, and 50<sup>th</sup> percentile of depletion of variation observed in coding exons of similar size. The number of samples required to obtain 80% detection power is labeled at corresponding benchmarks.

**Extended Data Fig. 8:** Analyses of constraint for chromosome X. **a**, Histograms of constraint Z scores for 49,936 1kb windows on chromosome X. Windows encompassing coding sequences ( $n=2,346$ ; red) overall exhibit a higher constraint Z (stronger negative selection) than windows that are exclusively non-coding ( $n=47,590$ ; blue). **b**, Distributions of candidate regulatory elements on the non-coding constraint spectrum of chromosome X. Enrichment was evaluated by comparing the proportion of non-coding 1kb windows, binned by constraint Z ( $<-4$ ,  $[-4,-3]$ ,  $[-3,-2]$ ,  $[-2,-1]$ ,  $[-1,0]$ ,  $[0,1]$ ,  $[1,2]$ ,  $[2,3]$ ,  $[3,4]$ ,  $\geq 4$ ), that overlap with a given functional annotation to the chromosome-wide average. Error bars indicate 95% confidence intervals of the odds ratios.

## 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 Septyarskiy, 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 Septyarskiy, 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 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).
- 16 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).
- 17 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).
- 18 Huang, Y. F., Gulko, B. & Siepel, A. Fast, scalable prediction of deleterious noncoding variants from functional and population genomic data. *Nat Genet* **49**, 618-624, doi:10.1038/ng.3810 (2017).
- 19 Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285-291, doi:10.1038/nature19057 (2016).
- 20 Halldorsson, B. V. *et al.* Characterizing mutagenic effects of recombination through a sequence-level genetic map. *Science* **363**, doi:10.1126/science.aau1043 (2019).
- 21 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).
- 22 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).
- 23 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).

- 24 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).
- 25 Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455-461, doi:10.1038/nature12787 (2014).
- 26 Jiang, Y. *et al.* SEdb: a comprehensive human super-enhancer database. *Nucleic Acids Res* **47**, D235-D243, doi:10.1093/nar/gky1025 (2019).
- 27 Pott, S. & Lieb, J. D. What are super-enhancers? *Nat Genet* **47**, 8-12, doi:10.1038/ng.3167 (2015).
- 28 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).
- 29 Kung, J. T., Colognori, D. & Lee, J. T. Long noncoding RNAs: past, present, and future. *Genetics* **193**, 651-669, doi:10.1534/genetics.112.146704 (2013).
- 30 Harismendy, O. *et al.* 9p21 DNA variants associated with coronary artery disease impair interferon-gamma signalling response. *Nature* **470**, 264-268, doi:10.1038/nature09753 (2011).
- 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 Wakil, S. M. *et al.* A genome-wide association study reveals susceptibility loci for myocardial infarction/coronary artery disease in Saudi Arabs. *Atherosclerosis* **245**, 62-70, doi:10.1016/j.atherosclerosis.2015.11.019 (2016).
- 33 AlRasheed, M. M. *et al.* The role of CDKN2B in cardiovascular risk in ethnic Saudi Arabs: A validation study. *Gene* **673**, 206-210, doi:10.1016/j.gene.2018.06.024 (2018).
- 34 Silander, K. *et al.* Worldwide patterns of haplotype diversity at 9p21.3, a locus associated with type 2 diabetes and coronary heart disease. *Genome Med* **1**, 51, doi:10.1186/gm51 (2009).
- 35 Yu, J. H. *et al.* The transcription factors signal transducer and activator of transcription 5A (STAT5A) and STAT5B negatively regulate cell proliferation through the activation of cyclin-dependent kinase inhibitor 2b (Cdkn2b) and Cdkn1a expression. *Hepatology* **52**, 1808-1818, doi:10.1002/hep.23882 (2010).
- 36 Grange, M. *et al.* Control of CD8 T cell proliferation and terminal differentiation by active STAT5 and CDKN2A/CDKN2B. *Immunology* **145**, 543-557, doi:10.1111/imm.12471 (2015).
- 37 Almontashiri, N. A. M. The 9p21.3 risk locus for coronary artery disease: A 10-year search for its mechanism. *J Taibah Univ Med Sci* **12**, 199-204, doi:10.1016/j.jtumed.2017.03.001 (2017).
- 38 Jarinova, O. *et al.* Functional analysis of the chromosome 9p21.3 coronary artery disease risk locus. *Arterioscler Thromb Vasc Biol* **29**, 1671-1677, doi:10.1161/ATVBAHA.109.189522 (2009).
- 39 Motterle, A. *et al.* Functional analyses of coronary artery disease associated variation on chromosome 9p21 in vascular smooth muscle cells. *Hum Mol Genet* **21**, 4021-4029, doi:10.1093/hmg/dds224 (2012).
- 40 Visel, A. *et al.* Targeted deletion of the 9p21 non-coding coronary artery disease risk interval in mice. *Nature* **464**, 409-412, doi:10.1038/nature08801 (2010).
- 41 Almontashiri, N. A. *et al.* 9p21.3 Coronary Artery Disease Risk Variants Disrupt TEAD Transcription Factor-Dependent Transforming Growth Factor beta Regulation of p16 Expression in Human Aortic Smooth Muscle Cells. *Circulation* **132**, 1969-1978, doi:10.1161/CIRCULATIONAHA.114.015023 (2015).
- 42 Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med* **12**, e1001779, doi:10.1371/journal.pmed.1001779 (2015).
- 43 Kanai, M. *et al.* Insights from complex trait fine-mapping across diverse populations. *medRxiv*, 2021.09.2003.21262975, doi:10.1101/2021.09.03.21262975 (2021).
- 44 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).
- 45 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).
- 46 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).
- 47 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).

- 48 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).
- 49 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).
- 50 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).
- 51 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).
- 52 Sebat, J. *et al.* Strong association of de novo copy number mutations with autism. *Science* **316**, 445-449, doi:10.1126/science.1138659 (2007).
- 53 Stefansson, H. *et al.* Large recurrent microdeletions associated with schizophrenia. *Nature* **455**, 232-236, doi:10.1038/nature07229 (2008).
- 54 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).
- 55 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).
- 56 Rice, A. M. & McLysaght, A. Dosage sensitivity is a major determinant of human copy number variant pathogenicity. *Nat Commun* **8**, 14366, doi:10.1038/ncomms14366 (2017).
- 57 Zhang, F., Gu, W., Hurles, M. E. & Lupski, J. R. Copy number variation in human health, disease, and evolution. *Annu Rev Genomics Hum Genet* **10**, 451-481, doi:10.1146/annurev.genom.9.081307.164217 (2009).
- 58 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).
- 59 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).
- 60 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).
- 61 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).
- 62 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).
- 63 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).
- 64 Roadmap Epigenomics, C. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317-330, doi:10.1038/nature14248 (2015).
- 65 Consortium, G. T. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580-585, doi:10.1038/ng.2653 (2013).
- 66 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).
- 67 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).
- 68 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).
- 69 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).
- 70 Darnell, J. C. *et al.* FMRP stalls ribosomal translocation on mRNAs linked to synaptic function and autism. *Cell* **146**, 247-261, doi:10.1016/j.cell.2011.06.013 (2011).
- 71 Cotney, J. *et al.* The evolution of lineage-specific regulatory activities in the human embryonic limb. *Cell* **154**, 185-196, doi:10.1016/j.cell.2013.05.056 (2013).

- 72 Xiao, S. *et al.* Comparative epigenomic annotation of regulatory DNA. *Cell* **149**, 1381-1392, doi:10.1016/j.cell.2012.04.029 (2012).
- 73 Vierstra, J. *et al.* Mouse regulatory DNA landscapes reveal global principles of cis-regulatory evolution. *Science* **346**, 1007-1012, doi:10.1126/science.1246426 (2014).
- 74 Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554-566, doi:10.1016/j.cell.2015.01.006 (2015).
- 75 Reilly, S. K. *et al.* Evolutionary genomics. Evolutionary changes in promoter and enhancer activity during human corticogenesis. *Science* **347**, 1155-1159, doi:10.1126/science.1260943 (2015).
- 76 Young, R. S. *et al.* The frequent evolutionary birth and death of functional promoters in mouse and human. *Genome Res* **25**, 1546-1557, doi:10.1101/gr.190546.115 (2015).
- 77 Kircher, M. *et al.* Saturation mutagenesis of twenty disease-associated regulatory elements at single base-pair resolution. *Nat Commun* **10**, 3583, doi:10.1038/s41467-019-11526-w (2019).
- 78 Hail Team. Hail 0.2.62-84fa81b9ea3d. <https://github.com/hail-is/hail/commit/84fa81b9ea3d>.
- 79 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).
- 80 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).
- 81 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).
- 82 Karolchik, D. *et al.* The UCSC Table Browser data retrieval tool. *Nucleic Acids Res* **32**, D493-496, doi:10.1093/nar/gkh103 (2004).
- 83 Li, H. Toward better understanding of artifacts in variant calling from high-coverage samples. *Bioinformatics* **30**, 2843-2851, doi:10.1093/bioinformatics/btu356 (2014).
- 84 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).
- 85 Zhao, H. *et al.* CrossMap: a versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* **30**, 1006-1007, doi:10.1093/bioinformatics/btt730 (2014).
- 86 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).
- 87 Rehm, H. L. *et al.* ClinGen--the Clinical Genome Resource. *N Engl J Med* **372**, 2235-2242, doi:10.1056/NEJMsr1406261 (2015).
- 88 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).
- 89 McKusick, V. A. Mendelian Inheritance in Man and its online version, OMIM. *Am J Hum Genet* **80**, 588-604, doi:10.1086/514346 (2007).
- 90 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).
- 91 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).

## Genome Aggregation Database Consortium

Maria Abreu<sup>1</sup>, Carlos A. Aguilar Salinas<sup>2</sup>, Tariq Ahmad<sup>3</sup>, Christine M. Albert<sup>4,5</sup>, Jessica Alföldi<sup>6,7</sup>, Diego Ardiissino<sup>8</sup>, Irina M. Armean<sup>6,7,9</sup>, Gil Atzmon<sup>10,11</sup>, Eric Banks<sup>12</sup>, John Barnard<sup>13</sup>, Samantha M. Baxter<sup>6</sup>, Laurent Beaugerie<sup>14</sup>, Emelia J. Benjamin<sup>15,16,17</sup>, David Benjamin<sup>12</sup>, Louis Bergelson<sup>12</sup>, Michael Boehnke<sup>18</sup>, Lori L. Bonnycastle<sup>19</sup>, Erwin P. Bottinger<sup>20</sup>, Donald W. Bowden<sup>21,22,23</sup>, Matthew J. Bown<sup>24,25</sup>, Steven Brant<sup>26</sup>, Sarah E. Calvo<sup>6,27</sup>, Hannia Campos<sup>28,29</sup>, John C. Chambers<sup>30,31,32</sup>, Juliana C. Chan<sup>33</sup>, Katherine R. Chao<sup>6</sup>, Sinéad Chapman<sup>6,7,34</sup>, Daniel Chasman<sup>4,35</sup>, Siwei Chen<sup>6,7</sup>, Rex Chisholm<sup>36</sup>, Judy Cho<sup>20</sup>, Rajiv Chowdhury<sup>37</sup>, Mina K. Chung<sup>38</sup>, Wendy Chung<sup>39,40,41</sup>, Kristian Cibulskis<sup>12</sup>, Bruce Cohen<sup>35,42</sup>, Ryan L. Collins<sup>6,27,43</sup>, Kristen M. Connolly<sup>44</sup>, Adolfo Correa<sup>45</sup>, Miguel Covarrubias<sup>12</sup>, Beryl Cummings<sup>6,43</sup>, Dana Dabelea<sup>46</sup>, Mark J. Daly<sup>6,7,47</sup>, John Danesh<sup>37</sup>, Dawood Darbar<sup>48</sup>, Joshua Denny<sup>49</sup>, Stacey Donnelly<sup>6</sup>, Ravindranath Duggirala<sup>50</sup>, Josée Dupuis<sup>51,52</sup>, Patrick T. Ellinor<sup>6,53</sup>, Roberto Elosua<sup>54,55,56</sup>, James Emery<sup>12</sup>, Eleina England<sup>6</sup>, Jeanette Erdmann<sup>57,58,59</sup>, Tõnu Esko<sup>6,60</sup>, Emily Evangelista<sup>6</sup>, Yossi Farjoun<sup>12</sup>, Diane Fatkin<sup>61,62,63</sup>, Steven Ferriera<sup>64</sup>, Jose Florez<sup>35,65,66</sup>, Laurent C. Franciolini<sup>6,7</sup>, Andre Franke<sup>67</sup>, Martti Färkkilä<sup>68</sup>, Stacey Gabriel<sup>64</sup>, Kiran Garimella<sup>12</sup>, Laura D. Gauthier<sup>12</sup>, Jeff Gentry<sup>12</sup>, Gad Getz<sup>35,69,70</sup>, David C. Glahn<sup>71,72</sup>, Benjamin Glaser<sup>73</sup>, Stephen J. Glatt<sup>74</sup>, David Goldstein<sup>75,76</sup>, Cicerio Gonzalez<sup>77</sup>, Julia K. Goodrich<sup>6</sup>, Leif Groop<sup>78,79</sup>, Sanna Gudmundsson<sup>6,7,80</sup>, Namrata Gupta<sup>6,64</sup>, Andrea Haessly<sup>12</sup>, Christopher Haiman<sup>81</sup>, Ira Hall<sup>82</sup>, Craig Hanis<sup>83</sup>, Matthew Harms<sup>84,85</sup>, Mikko Hiltunen<sup>86</sup>, Matti M. Holi<sup>87</sup>, Christina M. Hultman<sup>88,89</sup>, Chaim Jalas<sup>90</sup>, Thibault Jeandet<sup>12</sup>, Mikko Kallela<sup>91</sup>, Diane Kaplan<sup>12</sup>, Jaakkko Kaprio<sup>79</sup>, Konrad J. Karczewski<sup>6,7,34</sup>, Sekar Kathiresan<sup>27,35,92</sup>, Eimear Kenny<sup>89,93</sup>, Bong-Jo Kim<sup>94</sup>, Young Jin Kim<sup>94</sup>, George Kirov<sup>95</sup>, Zan Koenig<sup>6</sup>, Jaspal Kooner<sup>31,96,97</sup>, Seppo Koskinen<sup>98</sup>, Harlan M. Krumholz<sup>99</sup>, Subra Kugathasan<sup>100</sup>, Soo Heon Kwak<sup>101</sup>, Markku Laakso<sup>102,103</sup>, Nicole Lake<sup>104</sup>, Trevyn Langford<sup>12</sup>, Kristen M. Laricchia<sup>6,7</sup>, Terho Lehtimäki<sup>105</sup>, Monkol Lek<sup>104</sup>, Emily Lipscomb<sup>6</sup>, Christopher Llanwarne<sup>12</sup>, Ruth J.F. Loos<sup>20,106</sup>, Steven A. Lubitz<sup>6,53</sup>, Teresa Tusie Luna<sup>107,108</sup>, Ronald C.W. Ma<sup>33,109,110</sup>, Daniel G. MacArthur<sup>6,111,112</sup>, Gregory M. Marcus<sup>113</sup>, Jaume Marrugat<sup>55,114</sup>, Alicia R. Martin<sup>6</sup>, Kari M. Mattila<sup>105</sup>, Steven McCarroll<sup>34,115</sup>, Mark I. McCarthy<sup>116,117,118</sup>, Jacob McCauley<sup>119,120</sup>, Dermot McGovern<sup>121</sup>, Ruth McPherson<sup>122</sup>, James B. Meigs<sup>6,35,123</sup>, Olle Melander<sup>124</sup>, Andres Metspalu<sup>125</sup>, Deborah Meyers<sup>126</sup>, Eric V. Minikel<sup>6</sup>, Braxton Mitchell<sup>127</sup>, Vamsi K. Motha<sup>6,128</sup>, Ruchi Munshi<sup>12</sup>, Aliya Naheed<sup>129</sup>, Saman Nazarian<sup>130,131</sup>, Benjamin M. Neale<sup>6,7</sup>, Peter M. Nilsson<sup>132</sup>, Sam Novod<sup>12</sup>, Anne H. O'Donnell-Luria<sup>6,7,80</sup>, Michael C. O'Donovan<sup>95</sup>, Yukinori Okada<sup>133,134,135</sup>, Dost Ongur<sup>35,42</sup>, Lorena Orozco<sup>136</sup>, Michael J. Owen<sup>95</sup>, Colin Palmer<sup>137</sup>, Nicholette D. Palmer<sup>138</sup>, Aarno Palotie<sup>7,34,79</sup>, Kyong Soo Park<sup>101,139</sup>, Carlos Pato<sup>140</sup>, Nikelle Petrillo<sup>12</sup>, William Phu<sup>6,80</sup>, Timothy Poterba<sup>6,7,34</sup>, Ann E. Pulver<sup>141</sup>, Dan Rader<sup>130,142</sup>, Nazneen Rahman<sup>143</sup>, Heidi L. Rehm<sup>6,27</sup>, Alex Reiner<sup>144,145</sup>, Anne M. Remes<sup>146</sup>, Dan Rhodes<sup>6</sup>, Stephen Rich<sup>147,148</sup>, John D. Rioux<sup>149,150</sup>, Samuli Ripatti<sup>79,151,152</sup>, David Roazen<sup>12</sup>, Dan M. Roden<sup>153,154</sup>, Jerome I. Rotter<sup>155</sup>, Valentin Ruano-Rubio<sup>12</sup>, Nareh Sahakian<sup>12</sup>, Danish Saleheen<sup>156,157,158</sup>, Veikko Salomaa<sup>159</sup>, Andrea Saltzman<sup>6</sup>, Nilesh J. Samani<sup>24,25</sup>, Kaitlin E. Samocha<sup>6,27</sup>, Jeremiah Scharf<sup>6,27,34</sup>, Molly Schleicher<sup>6</sup>, Heribert Schunkert<sup>160,161</sup>, Sebastian Schönherr<sup>162</sup>, Eleanor Seaby<sup>6</sup>, Cotton Seed<sup>7,34</sup>, Svati H. Shah<sup>163</sup>, Megan Shand<sup>12</sup>, Moore B. Shoemaker<sup>164</sup>, Tai Shyong<sup>165,166</sup>, Edwin K. Silverman<sup>167,168</sup>, Moriel Singer-Berk<sup>6</sup>, Pamela Sklar<sup>169,170,171</sup>, J. Gustav Smith<sup>152,172,173</sup>, Jonathan T. Smith<sup>12</sup>, Hilkka Soininen<sup>174</sup>, Harry Sokol<sup>175,176,177</sup>, Matthew Solomonson<sup>6,7</sup>, Rachel G. Son<sup>6</sup>, Jose Soto<sup>12</sup>, Tim Spector<sup>178</sup>, Christine Stevens<sup>6,7,34</sup>, Nathan Stitziel<sup>82,179</sup>, Patrick F. Sullivan<sup>88,180</sup>, Jaana Suvisaari<sup>159</sup>, E. Shyong Tai<sup>181,182,183</sup>, Michael E. Talkowski<sup>6,27,34</sup>, Yekaterina Tarasova<sup>6</sup>, Kent D. Taylor<sup>155</sup>, Yik Ying Teo<sup>181,184,185</sup>, Grace Tiao<sup>6,7</sup>, Kathleen Tibbetts<sup>12</sup>, Charlotte Tolonen<sup>12</sup>, Ming Tsuang<sup>186,187</sup>, Tiinamaija Tuomi<sup>79,188,189</sup>, Dan Turner<sup>190</sup>, Teresa Tusie-Luna<sup>191,192</sup>, Erkki Vartiainen<sup>193</sup>, Marquis Vawter<sup>194</sup>, Christopher Vittal<sup>6,7</sup>, Gordon Wade<sup>12</sup>, Arcturus Wang<sup>6,7,34</sup>, Qingbo Wang<sup>6,133</sup>, James S. Ware<sup>6,195,196</sup>, Hugh Watkins<sup>197</sup>, Nicholas A. Watts<sup>6,7</sup>, Rinse K. Weersma<sup>198</sup>, Ben Weisburd<sup>12</sup>, Maija Wessman<sup>79,199</sup>, Nicola Whiffin<sup>6,200,201</sup>, Michael W. Wilson<sup>6,7</sup>, James G. Wilson<sup>202</sup>, Ramnik J. Xavier<sup>203,204</sup>, Mary T. Yohannes<sup>6</sup>

<sup>1</sup>University of Miami Miller School of Medicine, Gastroenterology, Miami, USA

<sup>2</sup>Unidad de Investigacion de Enfermedades Metabolicas, Instituto Nacional de Ciencias Medicas y Nutricion, Mexico City, Mexico

<sup>3</sup>Peninsula College of Medicine and Dentistry, Exeter, UK

<sup>4</sup>Division of Preventive Medicine, Brigham and Women's Hospital, Boston, MA, USA

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

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

<sup>7</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

<sup>8</sup>Department of Cardiology University Hospital, Parma, Italy

<sup>9</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK

<sup>10</sup>Department of Biology Faculty of Natural Sciences, University of Haifa, Haifa, Israel

<sup>11</sup>Departments of Medicine and Genetics, Albert Einstein College of Medicine, Bronx, NY, USA

<sup>12</sup>Data Science Platform, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>13</sup>Department of Quantitative Health Sciences, Lerner Research Institute Cleveland Clinic, Cleveland, OH, USA

<sup>14</sup>Sorbonne Université, APHP, Gastroenterology Department Saint Antoine Hospital, Paris, France

<sup>15</sup>NHLBI and Boston University's Framingham Heart Study, Framingham, MA, USA

<sup>16</sup>Department of Medicine, Boston University School of Medicine, Boston, MA, USA

<sup>17</sup>Department of Epidemiology, Boston University School of Public Health, Boston, MA, USA

<sup>18</sup>Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA

<sup>19</sup>National Human Genome Research Institute, National Institutes of Health Bethesda, MD, USA

<sup>20</sup>The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA

<sup>21</sup>Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC, USA

<sup>22</sup>Center for Genomics and Personalized Medicine Research, Wake Forest School of Medicine, Winston-Salem, NC, USA

<sup>23</sup>Center for Diabetes Research, Wake Forest School of Medicine, Winston-Salem, NC, USA

<sup>24</sup>Department of Cardiovascular Sciences, University of Leicester, Leicester, UK

<sup>25</sup>NIHR Leicester Biomedical Research Centre, Glenfield Hospital, Leicester, UK

<sup>26</sup>John Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

<sup>27</sup>Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA

<sup>28</sup>Harvard School of Public Health, Boston, MA, USA

<sup>29</sup>Central American Population Center, San Pedro, Costa Rica

<sup>30</sup>Department of Epidemiology and Biostatistics, Imperial College London, London, UK

<sup>31</sup>Department of Cardiology, Ealing Hospital, NHS Trust, Southall, UK

<sup>32</sup>Imperial College, Healthcare NHS Trust Imperial College London, London, UK

<sup>33</sup>Department of Medicine and Therapeutics, The Chinese University of Hong Kong, Hong Kong, China

<sup>34</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>35</sup>Department of Medicine, Harvard Medical School, Boston, MA, USA

<sup>36</sup>Northwestern University, Evanston, IL, USA

<sup>37</sup>University of Cambridge, Cambridge, England

<sup>38</sup>Departments of Cardiovascular, Medicine Cellular and Molecular Medicine Molecular Cardiology, Quantitative Health Sciences, Cleveland Clinic, Cleveland, OH, USA

<sup>39</sup>Department of Pediatrics, Columbia University Irving Medical Center, New York, NY, USA

<sup>40</sup>Herbert Irving Comprehensive Cancer Center, Columbia University Medical Center, New York, NY, USA

<sup>41</sup>Department of Medicine, Columbia University Medical Center, New York, NY, USA

<sup>42</sup>McLean Hospital, Belmont, MA, USA

<sup>43</sup>Division of Medical Sciences, Harvard Medical School, Boston, MA, USA

<sup>44</sup>Genomics Platform, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>45</sup>Department of Medicine, University of Mississippi Medical Center, Jackson, MI, USA

<sup>46</sup>Department of Epidemiology Colorado School of Public Health Aurora, CO, USA

<sup>47</sup>Institute for Molecular Medicine Finland, (FIMM) Helsinki, Finland

<sup>48</sup>Department of Medicine and Pharmacology, University of Illinois at Chicago, Chicago, IL, USA

<sup>49</sup>Vanderbilt University Medical Center, Nashville, TN, USA

<sup>50</sup>Department of Genetics, Texas Biomedical Research Institute, San Antonio, TX, USA

<sup>51</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA, USA

<sup>52</sup>National Heart Lung and Blood Institute's Framingham Heart Study, Framingham, MA, USA

<sup>53</sup>Cardiac Arrhythmia Service and Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA

<sup>54</sup>Cardiovascular Epidemiology and Genetics Hospital del Mar Medical Research Institute, (IMIM) Barcelona Catalonia, Spain

<sup>55</sup>CIBER CV Barcelona, Catalonia, Spain

<sup>56</sup>Departament of Medicine, Medical School University of Vic-Central, University of Catalonia, Vic Catalonia, Spain

<sup>57</sup>Institute for Cardiogenetics, University of Lübeck, Lübeck, Germany

<sup>58</sup>German Research Centre for Cardiovascular Research, Hamburg/Lübeck/Kiel, Lübeck, Germany

<sup>59</sup>University Heart Center Lübeck, Lübeck, Germany

<sup>60</sup>Estonian Genome Center, Institute of Genomics University of Tartu, Tartu, Estonia

<sup>61</sup>Victor Chang Cardiac Research Institute, Darlinghurst, NSW, Australia

<sup>62</sup>Faculty of Medicine, UNSW Sydney, Kensington, NSW, Australia

<sup>63</sup>Cardiology Department, St Vincent's Hospital, Darlinghurst, NSW, Australia

<sup>64</sup>Broad Genomics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>65</sup>Diabetes Unit and Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA

<sup>66</sup>Programs in Metabolism and Medical & Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>67</sup>Institute of Clinical Molecular Biology, (IKMB) Christian-Albrechts-University of Kiel, Kiel, Germany

<sup>68</sup>Helsinki University and Helsinki University Hospital Clinic of Gastroenterology, Helsinki, Finland

<sup>69</sup>Bioinformatics Program MGH Cancer Center and Department of Pathology, Boston, MA, USA

<sup>70</sup>Cancer Genome Computational Analysis, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>71</sup>Department of Psychiatry and Behavioral Sciences, Boston Children's Hospital and Harvard Medical School, Boston, MA, USA

<sup>72</sup>Harvard Medical School Teaching Hospital, Boston, MA, USA

<sup>73</sup>Department of Endocrinology and Metabolism, Hadassah Medical Center and Faculty of Medicine, Hebrew University of Jerusalem, Israel

<sup>74</sup>Department of Psychiatry and Behavioral Sciences, SUNY Upstate Medical University, Syracuse, NY, USA

<sup>75</sup>Institute for Genomic Medicine, Columbia University Medical Center Hammer Health Sciences, New York, NY, USA

<sup>76</sup>Department of Genetics & Development Columbia University Medical Center, Hammer Health Sciences, New York, NY, USA

<sup>77</sup>Centro de Investigacion en Salud Poblacional, Instituto Nacional de Salud Publica, Mexico

<sup>78</sup>Lund University Sweden, Sweden

<sup>79</sup>Institute for Molecular Medicine Finland, (FIMM) HiLIFE University of Helsinki, Helsinki, Finland

<sup>80</sup>Division of Genetics and Genomics, Boston Children's Hospital, Boston, MA, USA

<sup>81</sup>Lund University Diabetes Centre, Malmö, Skåne County, Sweden

<sup>82</sup>Washington School of Medicine, St Louis, MI, USA

<sup>83</sup>Human Genetics Center University of Texas Health Science Center at Houston, Houston, TX, USA

<sup>84</sup>Department of Neurology Columbia University, New York City, NY, USA

<sup>85</sup>Institute of Genomic Medicine, Columbia University, New York City, NY, USA

<sup>86</sup>Institute of Biomedicine, University of Eastern Finland, Kuopio, Finland

<sup>87</sup>Department of Psychiatry, Helsinki University Central Hospital Lapinlahdentie, Helsinki, Finland

<sup>88</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

<sup>89</sup>Icahn School of Medicine at Mount Sinai, New York, NY, USA

<sup>90</sup>Bonei Olam, Center for Rare Jewish Genetic Diseases, Brooklyn, NY, USA

<sup>91</sup>Department of Neurology, Helsinki University, Central Hospital, Helsinki, Finland

<sup>92</sup>Cardiovascular Disease Initiative and Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>93</sup>Charles Bronfman Institute for Personalized Medicine, New York, NY, USA

<sup>94</sup>Division of Genome Science, Department of Precision Medicine, National Institute of Health, Republic of Korea

<sup>95</sup>MRC Centre for Neuropsychiatric Genetics & Genomics, Cardiff University School of Medicine, Cardiff, Wales

<sup>96</sup>Imperial College, Healthcare NHS Trust, London, UK

<sup>97</sup>National Heart and Lung Institute Cardiovascular Sciences, Hammersmith Campus, Imperial College London, London, UK

<sup>98</sup>Department of Health THL-National Institute for Health and Welfare, Helsinki, Finland

<sup>99</sup>Section 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

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

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

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

Teresa Tusie-Luna: CONACyT Project 312688

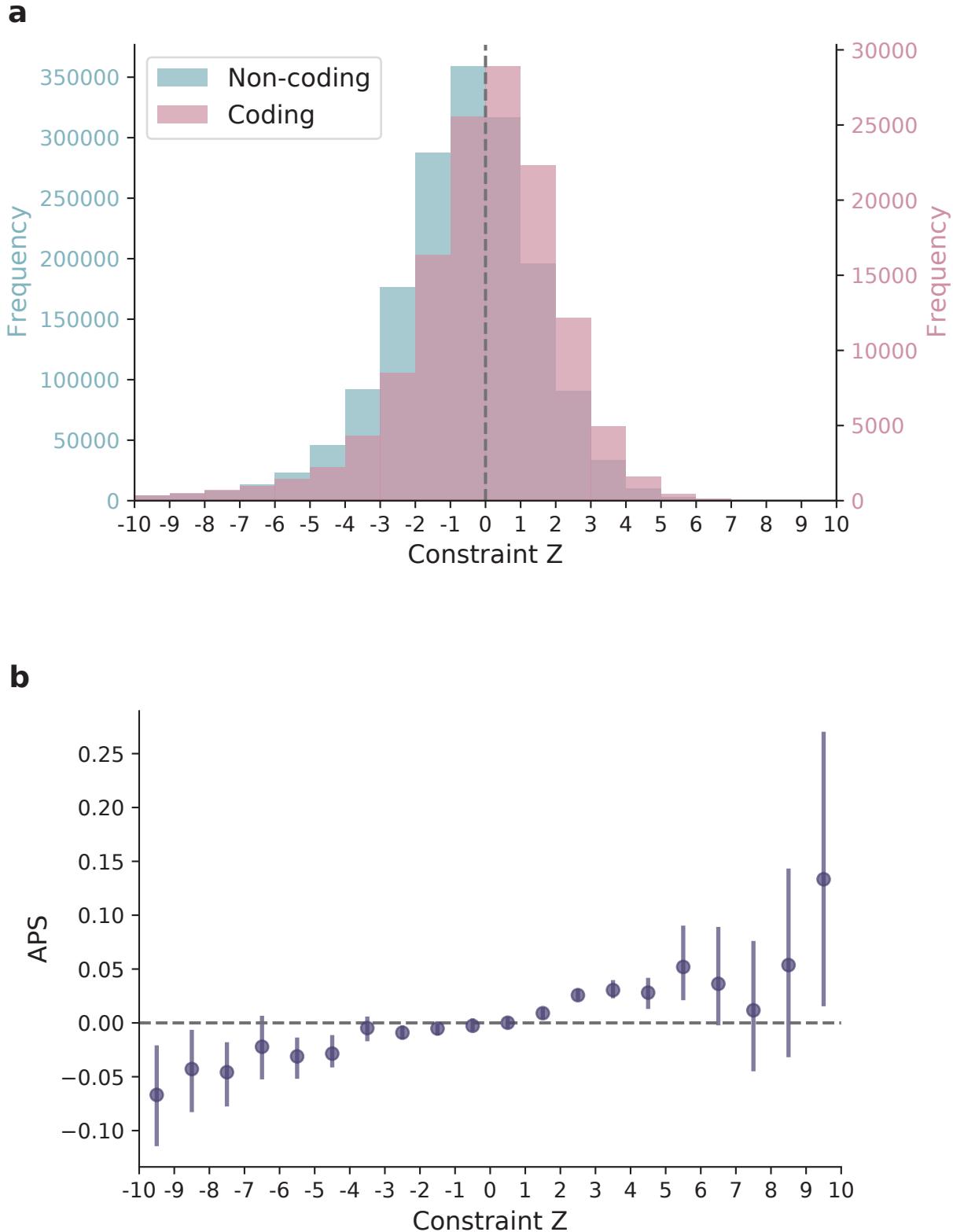
James S. Ware: Wellcome Trust [107469/Z/15/Z], Medical Research Council (UK), NIHR Imperial College Biomedical Research Centre

Rinse K. Weersma: The Lifelines Biobank initiative has been made possible by subsidy from the Dutch Ministry of

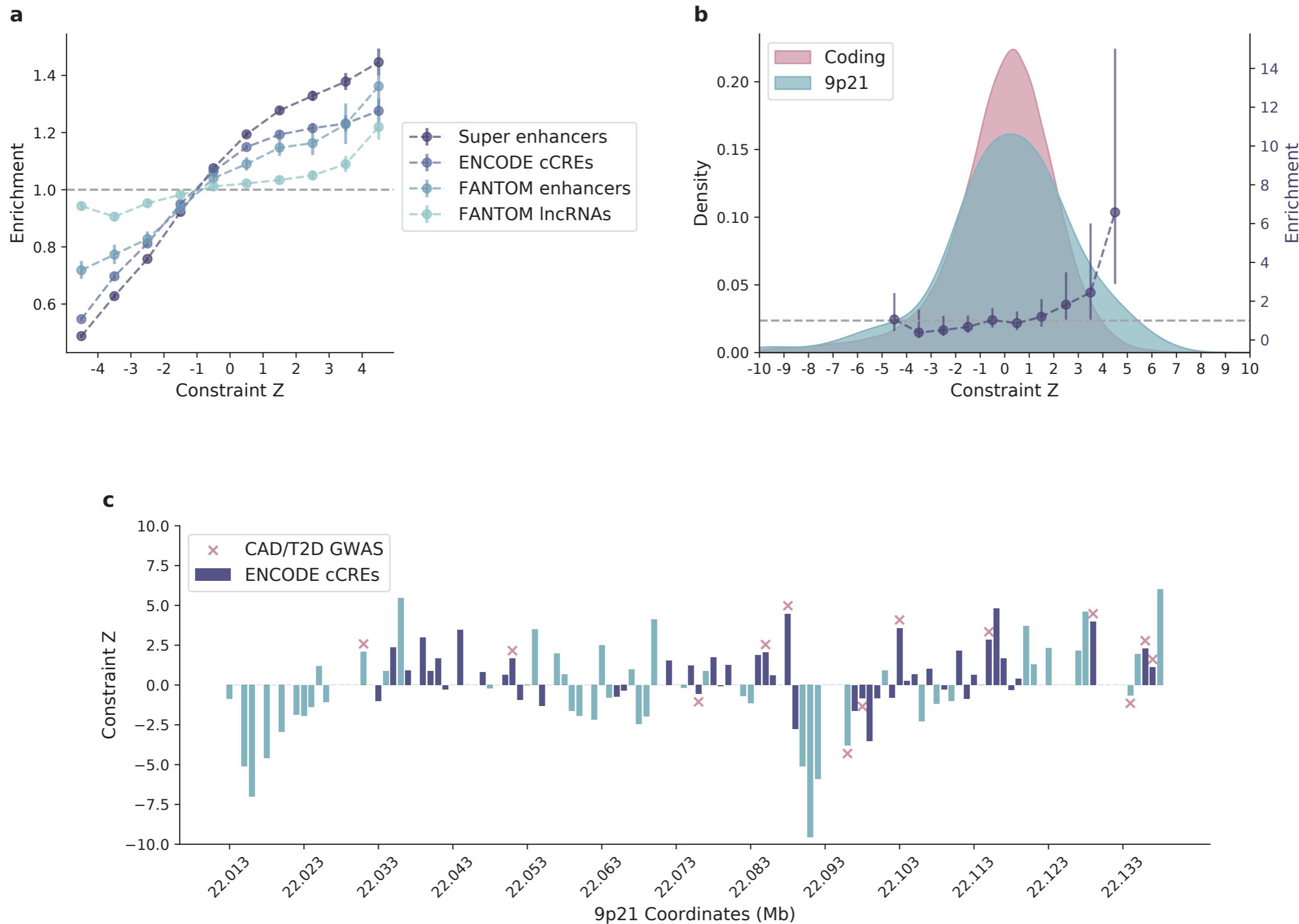
Health Welfare and Sport the Dutch Ministry of Economic Affairs the University Medical Centre Groningen (UMCG the Netherlands ) the University of Groningen and the Northern Provinces of the Netherlands

No conflicts of interest to declare

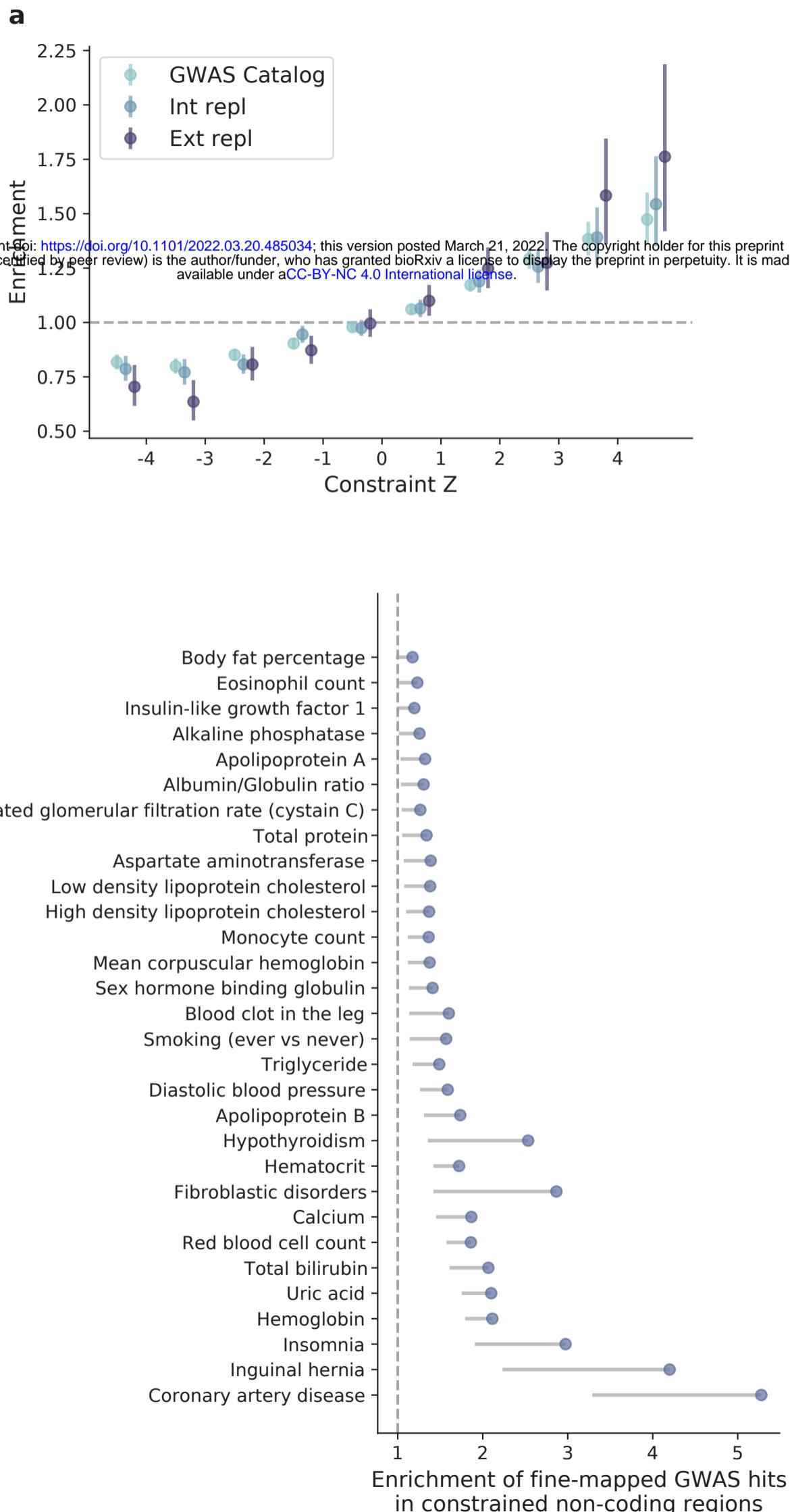
## Figure 1



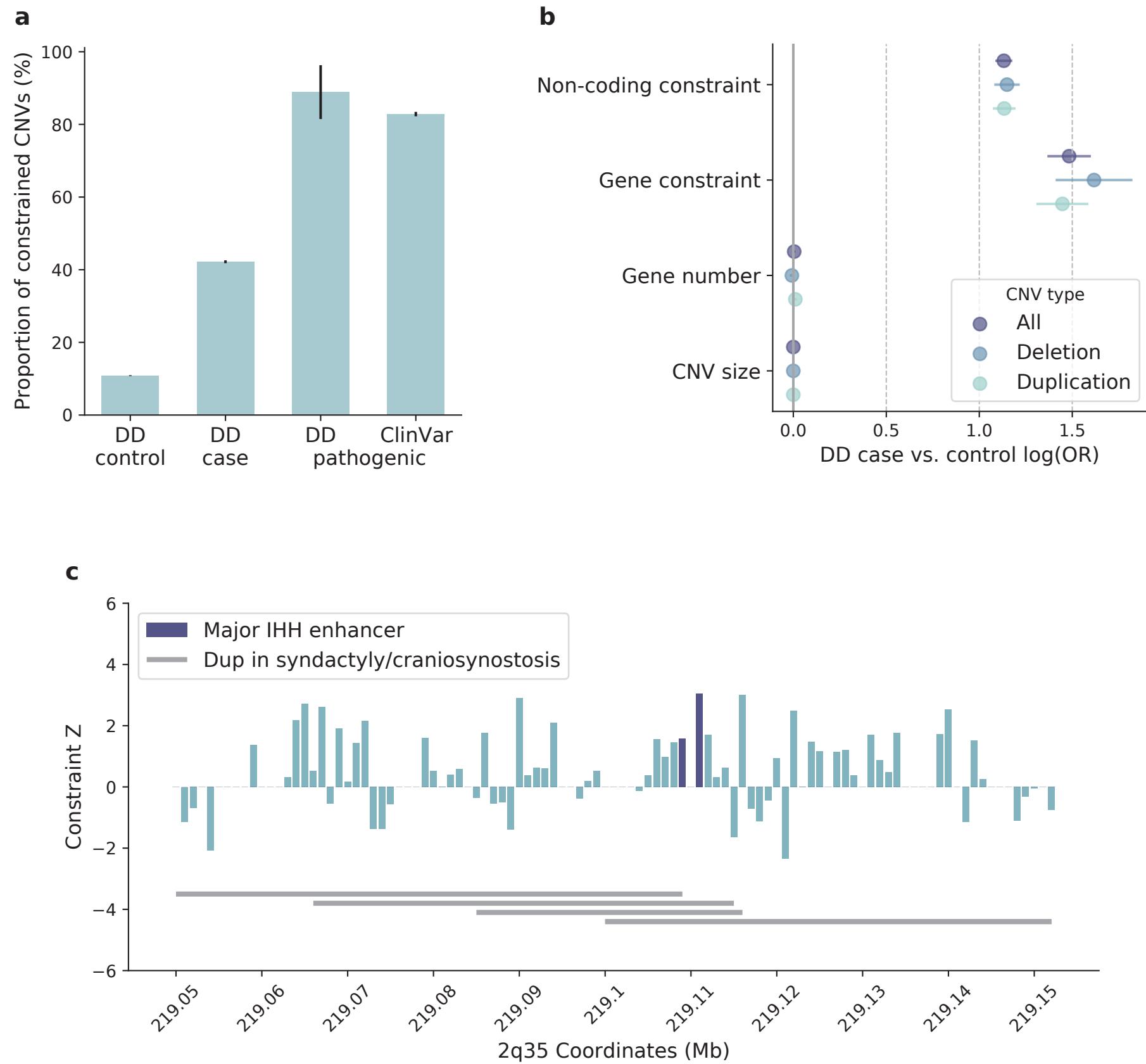
**Figure 2**



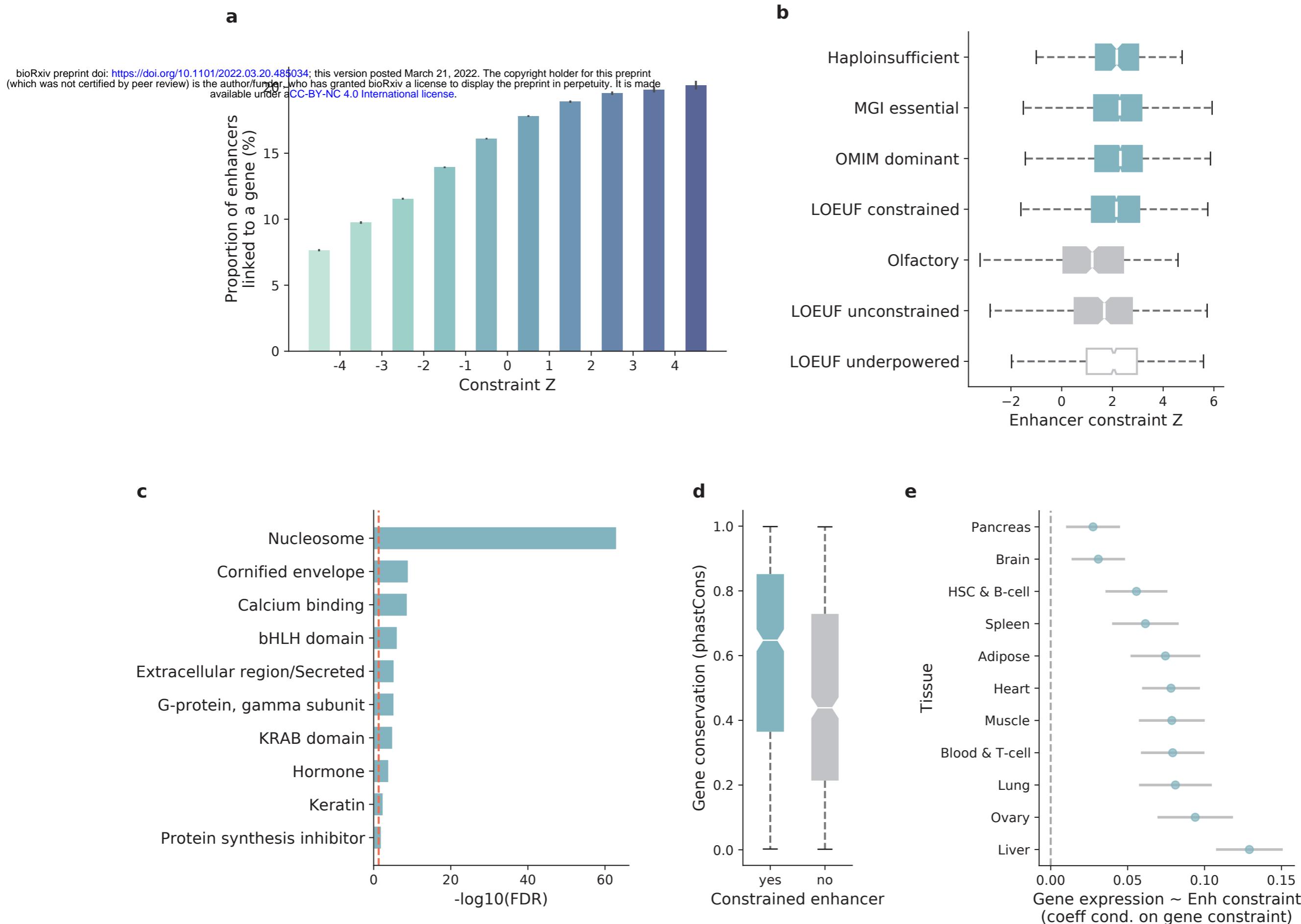
**Figure 3**



**Figure 4**

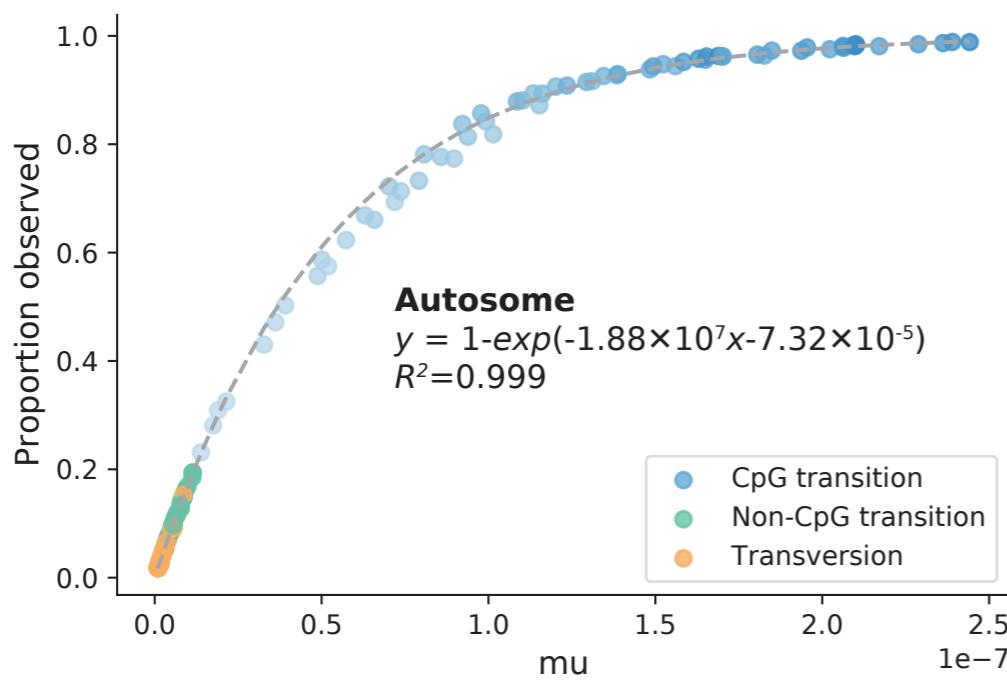


**Figure 5**

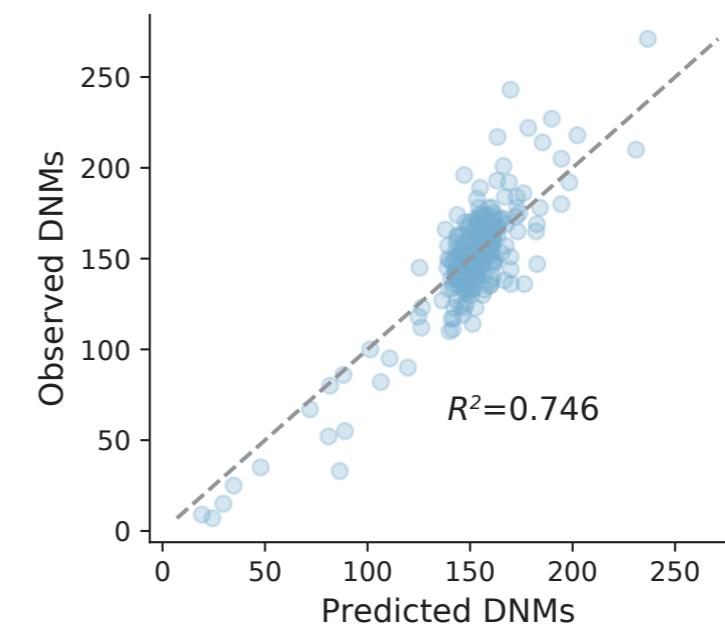


## Extended Data Figure 1

a

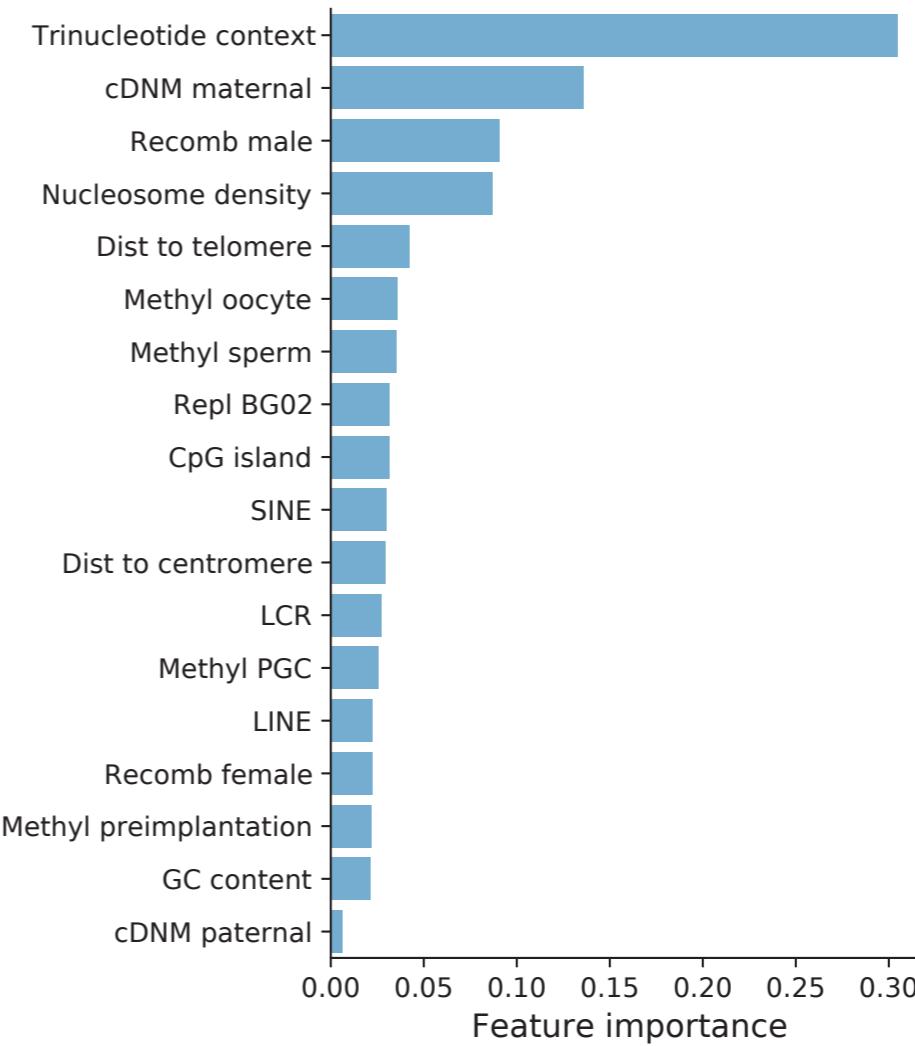
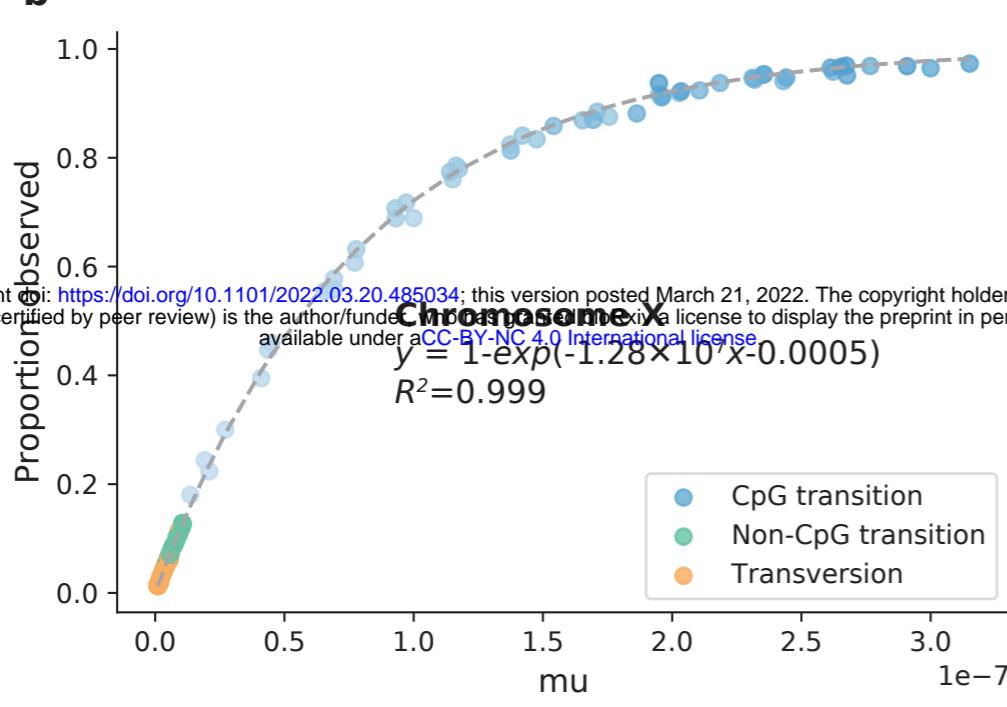


c

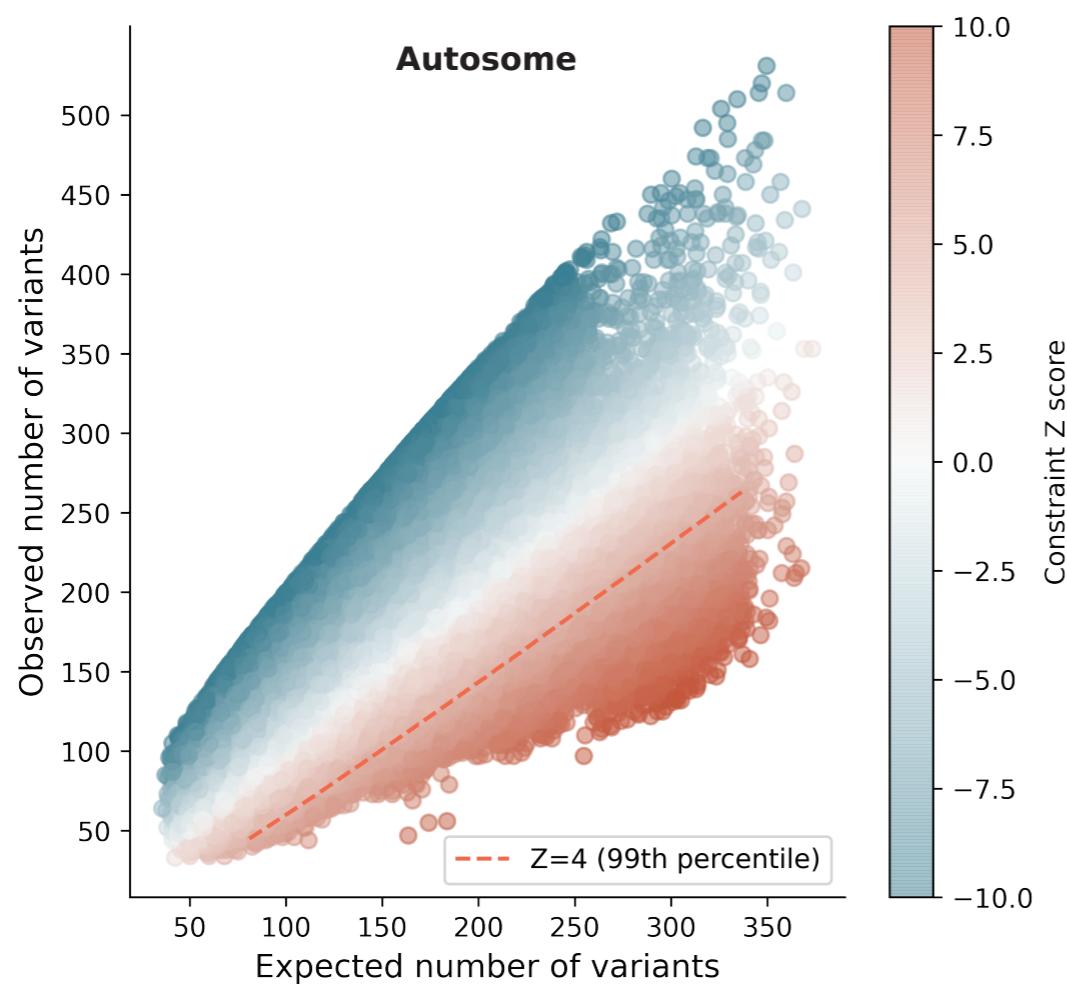


b

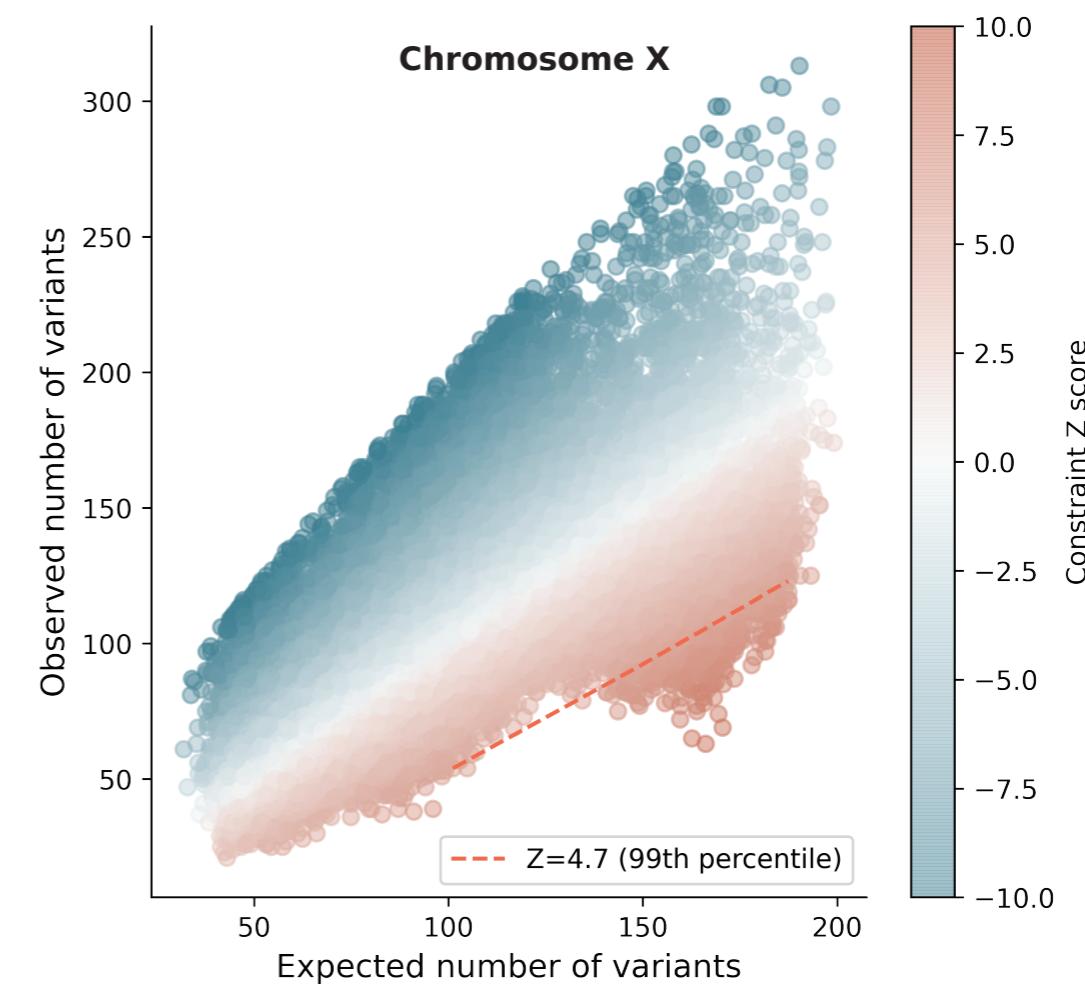
bioRxiv preprint doi: <https://doi.org/10.1101/2022.03.20.485034>; this version posted March 21, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



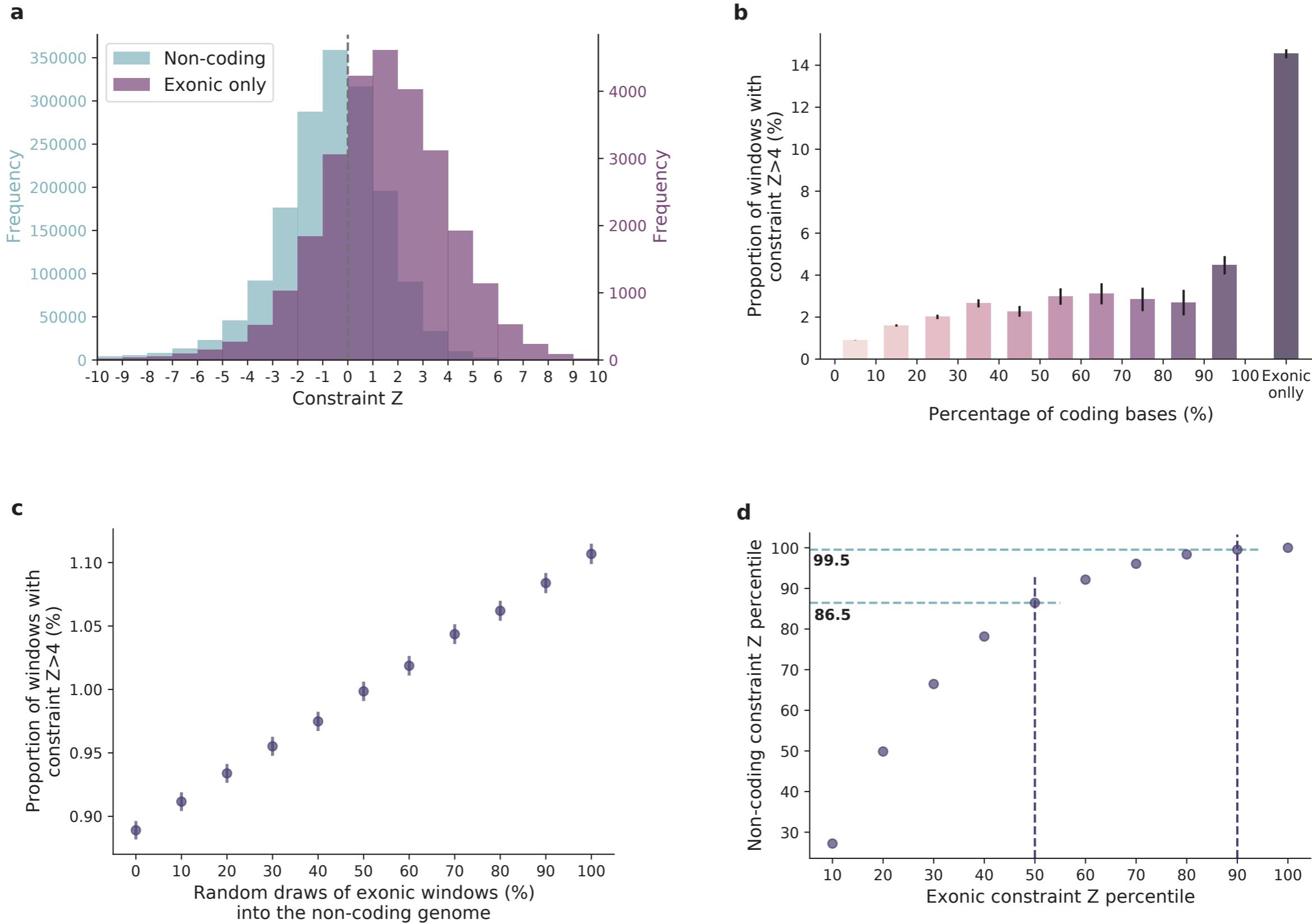
d



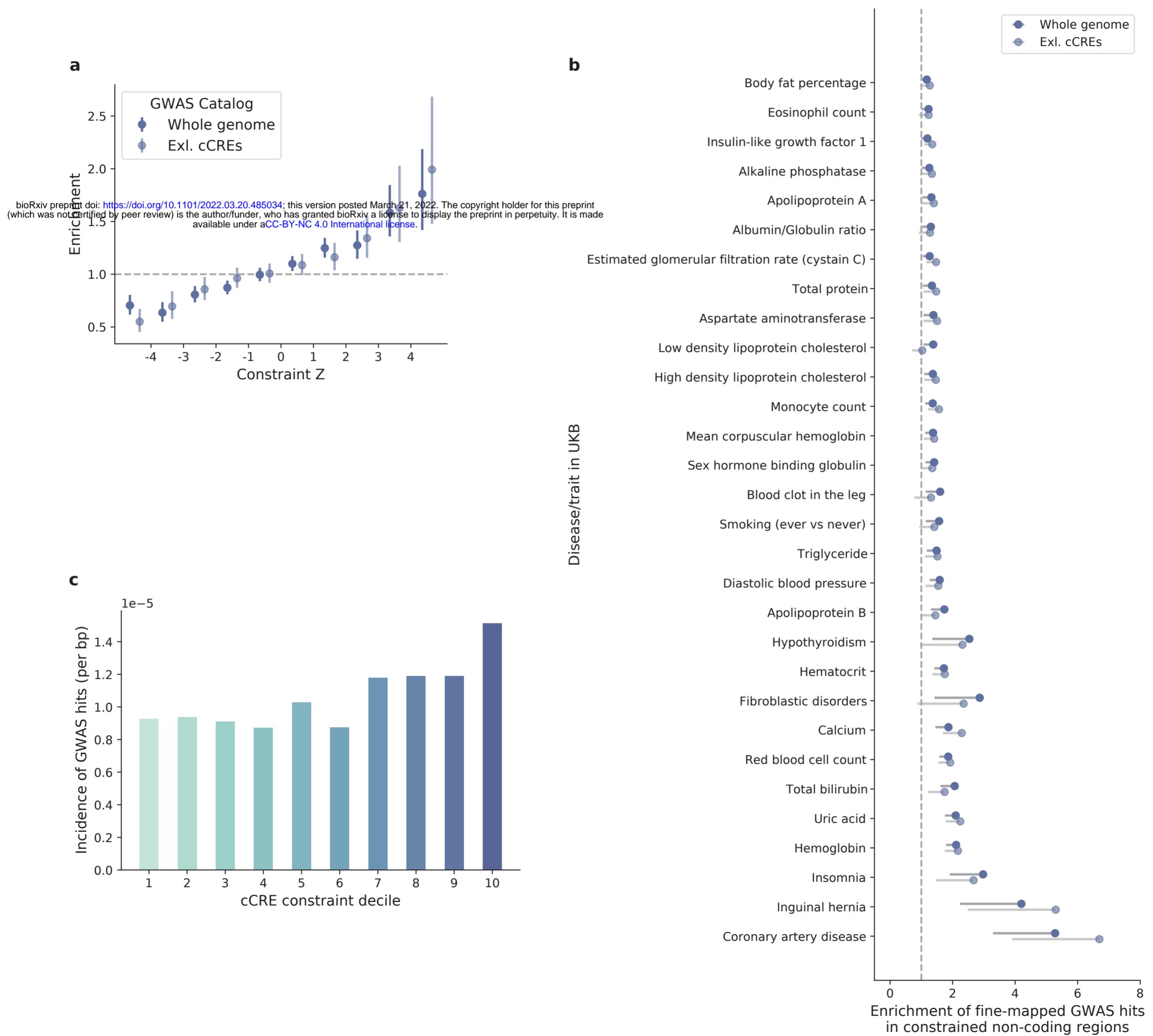
e



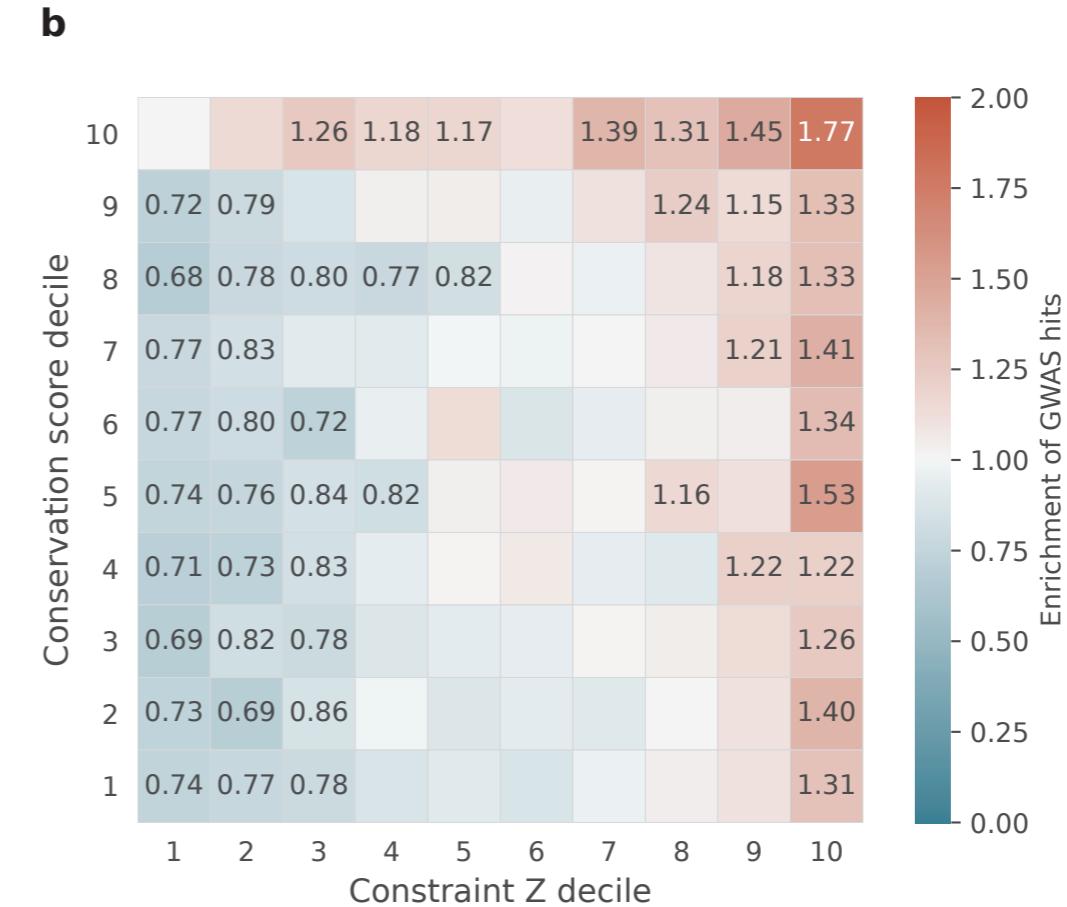
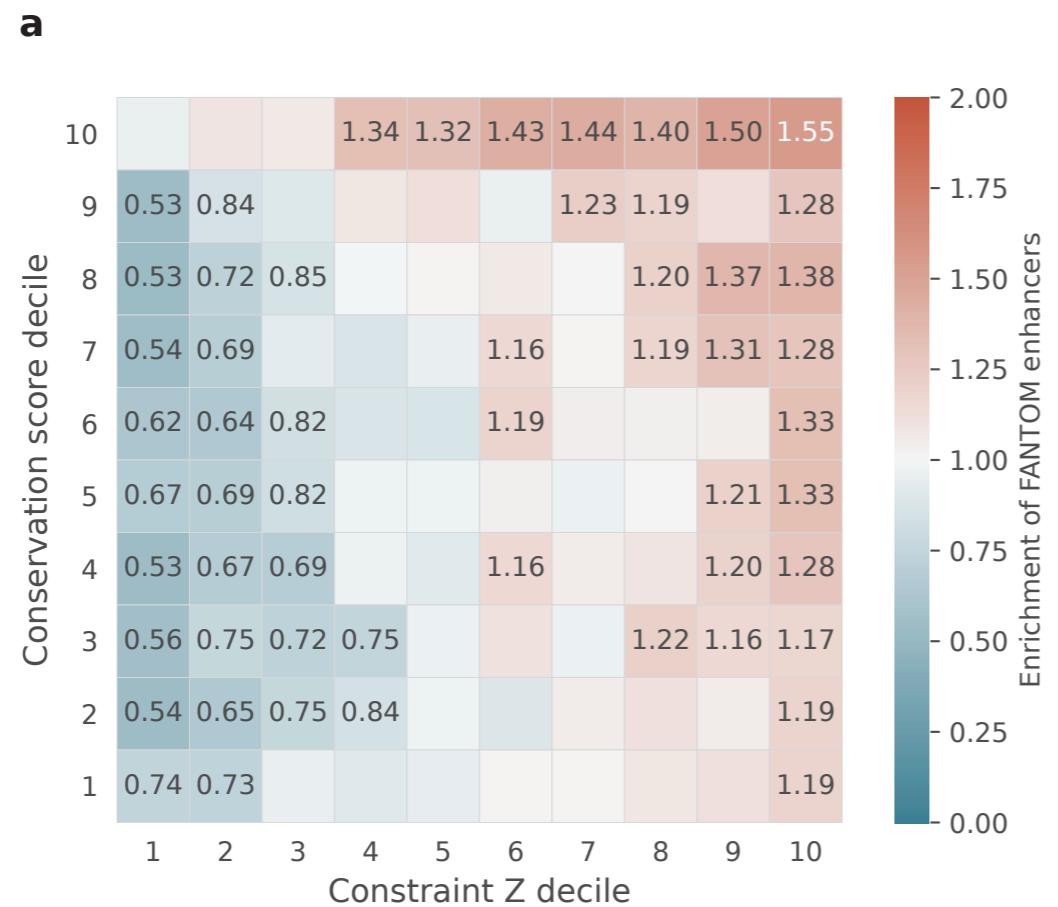
## Extended Data Figure 2



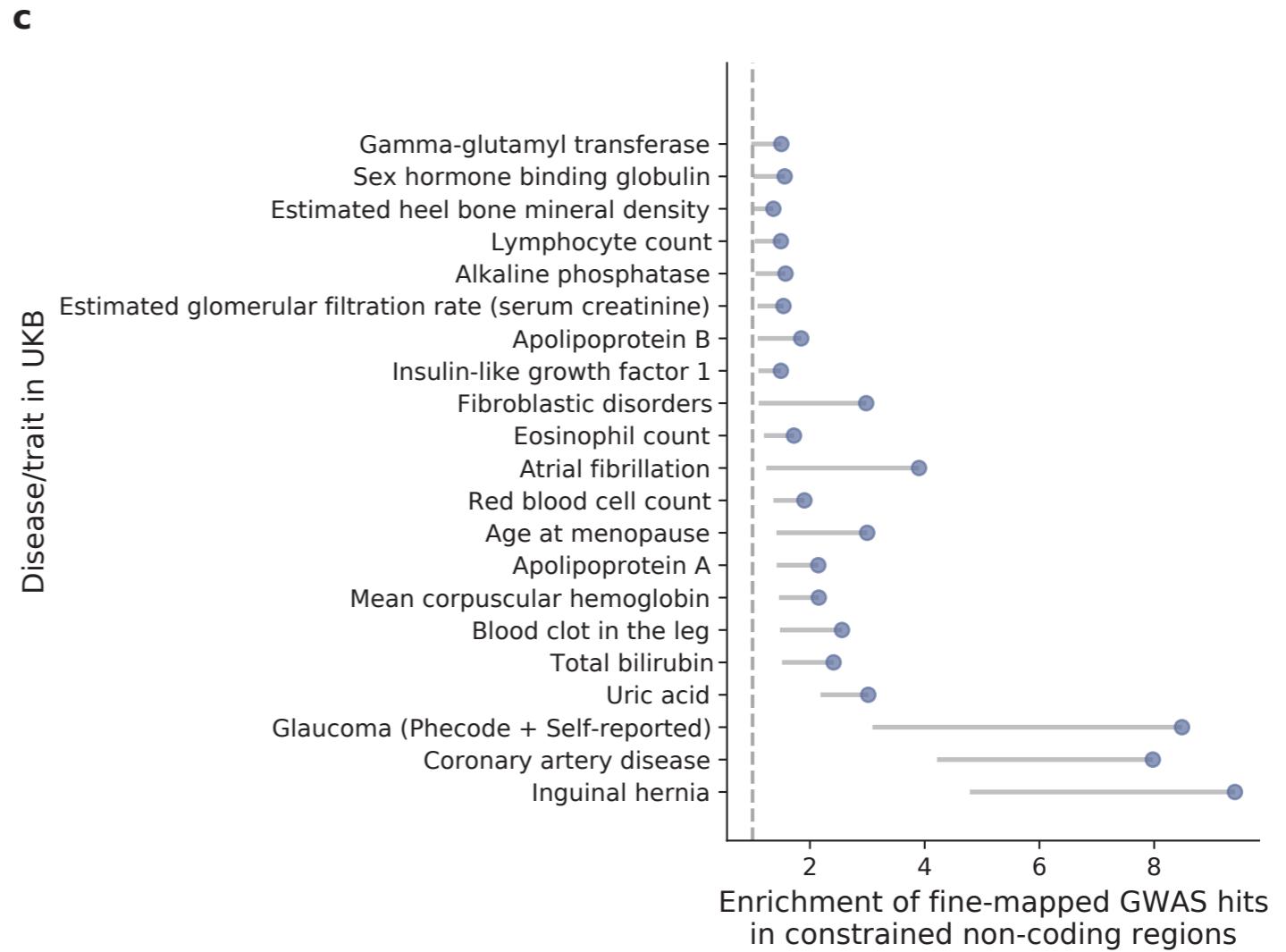
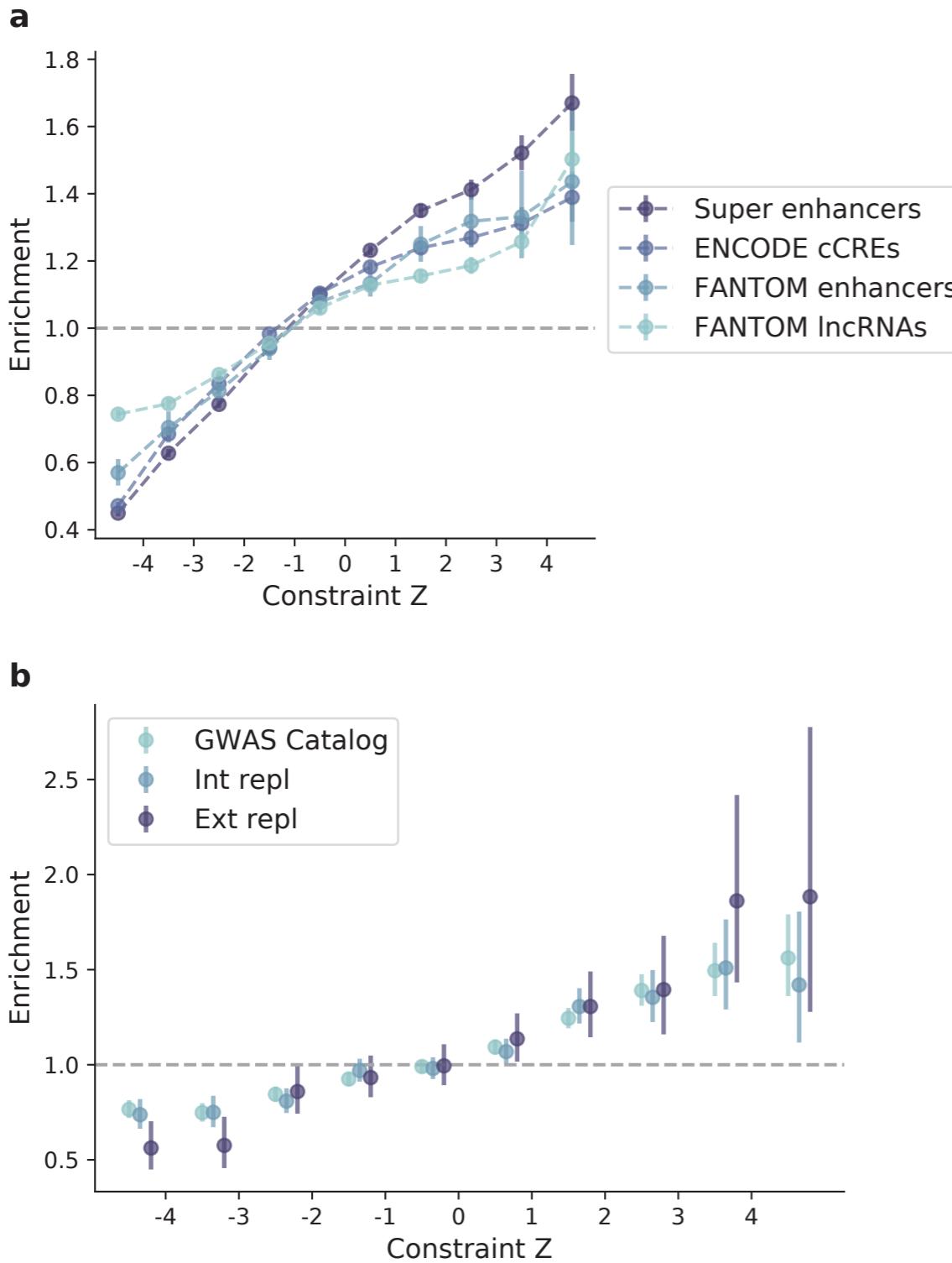
### Extended Data Figure 3



## Extended Data Figure 4

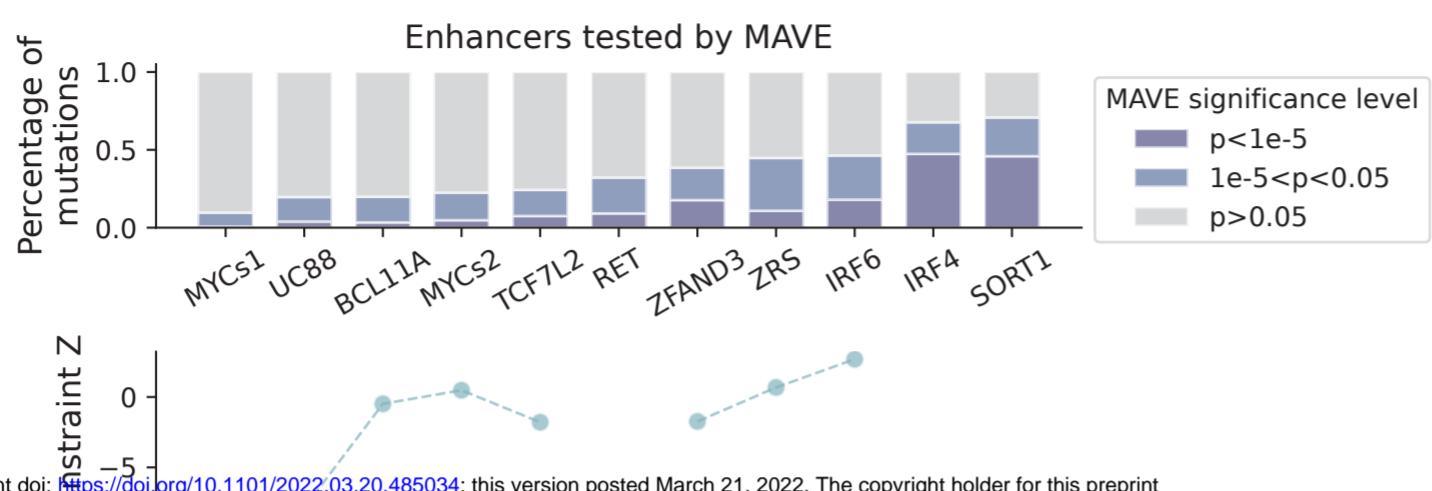


## Extended Figure Data 5

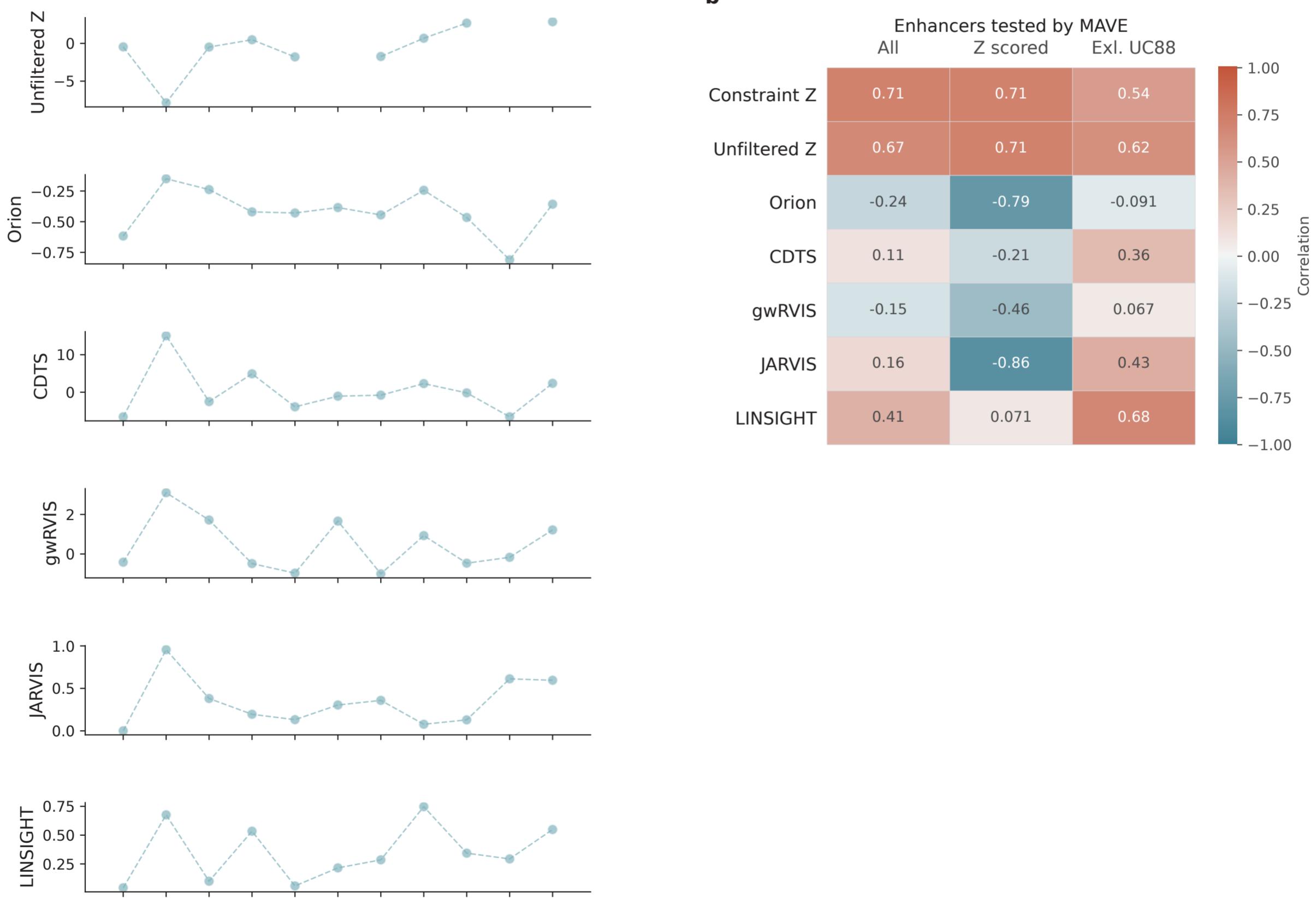


## Extended Data Figure 6

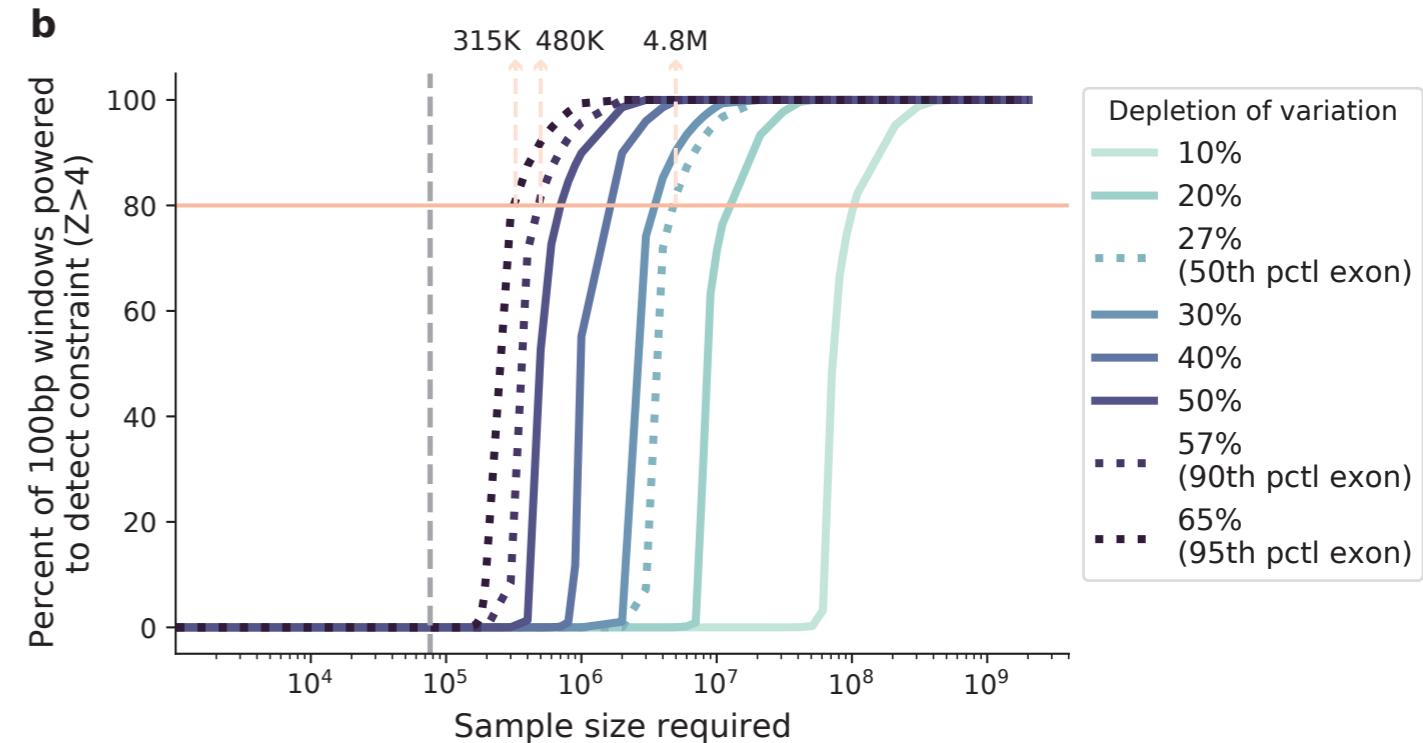
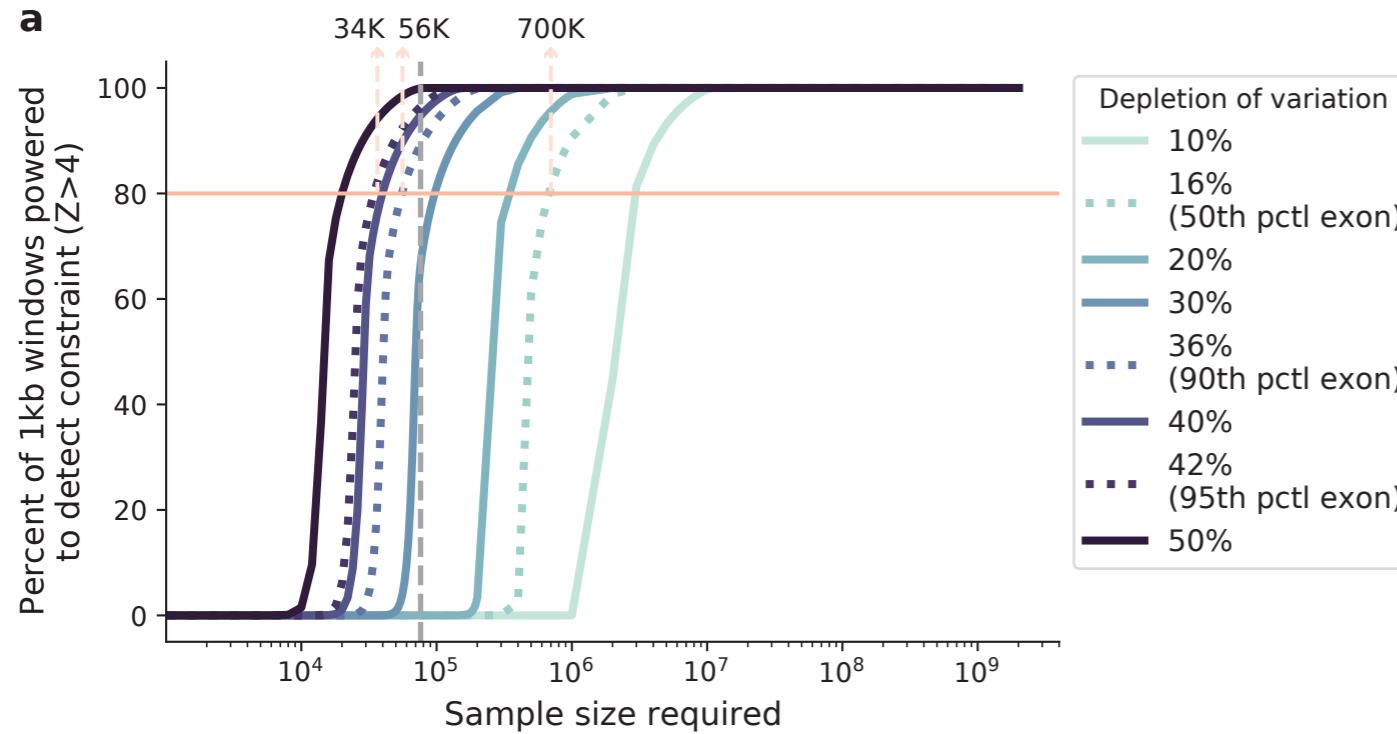
**a**



**b**



## Extended Data Figure 7



## Extended Data Figure 8

