

Phenome-wide burden of copy number variation in UK Biobank

Matthew Aguirre^{1,2}, Manuel Rivas¹, James Priest^{2,3}

Abstract:

Copy number variations (CNV) represent a significant proportion of the genetic differences between individuals and many CNVs associate causally with syndromic disease and clinical outcomes. Here, we characterize the landscape of copy number variation and their phenome-wide effects in a sample of 472,228 array-genotyped individuals from the UK Biobank. In addition to population-level selection effects against genic loci conferring high-mortality, we describe genetic burden from syndromic and previously uncharacterized CNV loci across nearly 2,000 quantitative and dichotomous traits, with separate analyses for common and rare classes of variation. Specifically, we highlight the effects of CNVs at two well-known syndromic loci *16p11.2* and *22q11.2*, as well as novel associations at *9p23*, in the context of acute coronary artery disease and high body mass index. Our data constitute a deeply contextualized portrait of population-wide burden of copy number variation, as well as a series of known and novel dosage-mediated genic associations across the medical phenome.

Introduction:

Copy number variants (CNV) are a class of structural variation often defined as large deletions or duplications of at least 1 kilobase (kb) of genomic sequence. CNVs exhibit substantial variability in both size and frequency in the population and have been implicated across a variety of common and rare health outcomes¹. Regional deletion and duplication syndromes have also been described at many loci, clustering near microsatellite repeats or areas of segmental duplication which may potentiate non-allelic homologous recombination². For example, CNV-based architectures for neuropsychiatric (e.g. autism spectrum disorder), developmental (e.g. *16p11.2*)^{3,4}, and syndromic cardiac disease (e.g. *22q11.2*)⁵ phenotypes have been well established.

Despite a growing body research on CNV-related syndromes and disease etiologies, large-scale studies of CNV effects have been limited by their rarity in the general population. However, burden testing methods which address this rarity by pooling observed variation across gene regions have yielded reproducible associations in the context of congenital heart disease and various neurocognitive outcomes^{6,7}. Moreover, as studies which include either microarray or NGS-based genotype data have grown in size and scope, it has become possible to describe the distribution of CNVs at kilobase-level resolution in the general population^{8,9}. Furthermore, the aggregation of richly annotated phenotype data in biobanks has diversified the set of phenotypes available for well-powered association studies, and allows for more precise characterization of syndromic CNV-associated disease^{10,11,12}.

¹Department of Biomedical Data Science, School of Medicine, Stanford University, Stanford, CA, USA.

²Department of Pediatrics, School of Medicine, Stanford University, Stanford, CA, USA.

³Stanford Cardiovascular Institute, Stanford University, Stanford, CA, USA.

41

42 We here describe the landscape of copy number variation and their associations with 1,937
43 phenotypes in a cohort of 472,228 participants from the UK Biobank¹³. We replicate well-
44 established syndromic effects of common CNVs — namely *22q11.2* deletion (DiGeorge)
45 syndrome and two variants of *16p11.2* deletion syndrome — and highlight known and novel
46 associations for body mass index (BMI), acute coronary artery disease (CAD), and related
47 adipose and cardiovascular phenotypes. Summary statistics from traditional genome-wide
48 associations for common CNVs, as well as from gene-level aggregate burden tests of rare
49 variants across all phenotypes are available for download on the Global Biobank Engine¹⁴.

50

51 **Results:**

52

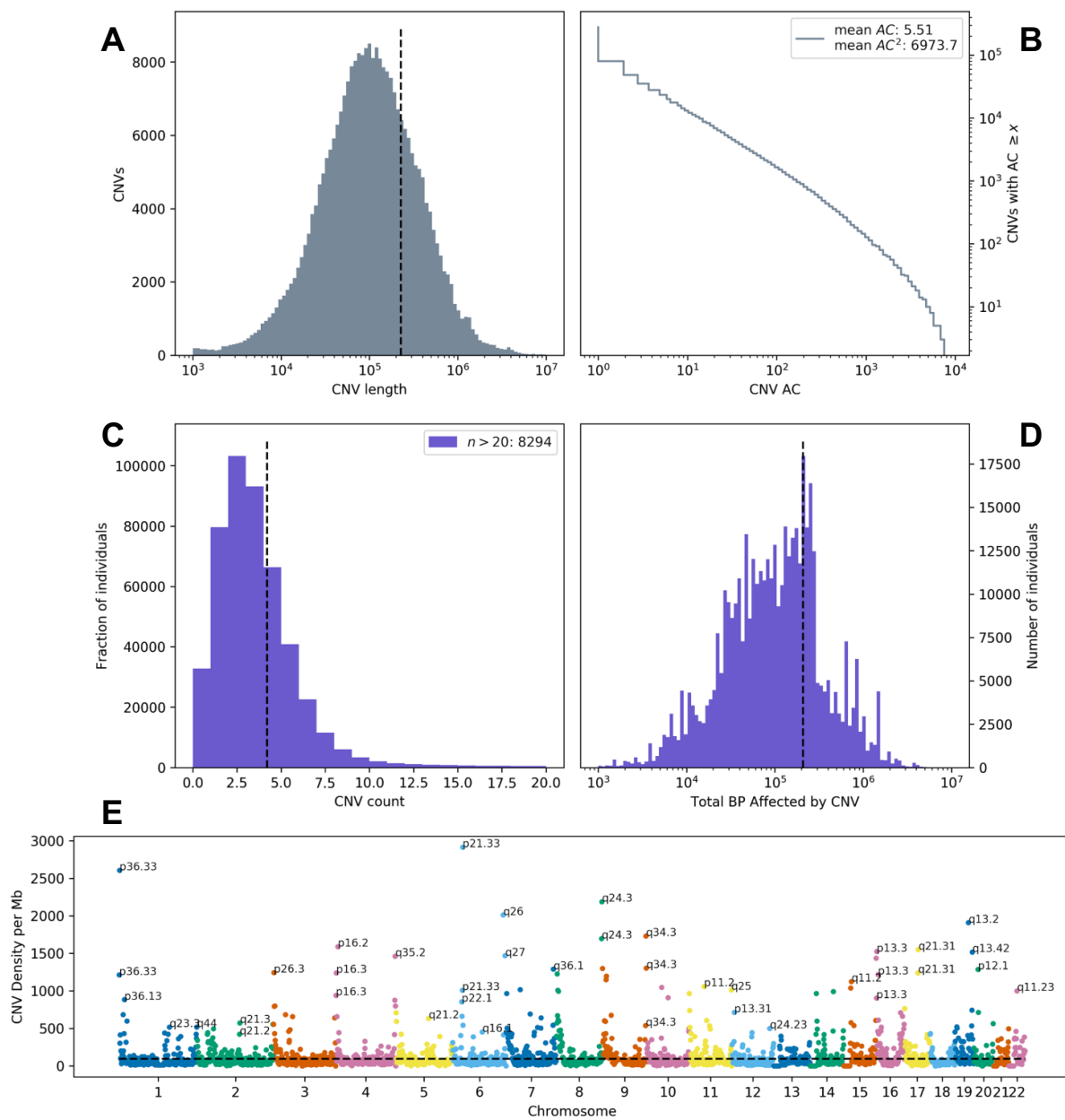
53 **Landscape of common and rare CNVs in a large volunteer cohort**

54

55 To call copy number variants in UK Biobank, we apply PennCNV¹⁵ separately within each
56 genotyping batch, resulting in 278,455 unique CNVs among 472,228 individuals after sample
57 quality control. We also observe heavy-tailed distributions in size and allele count of CNVs, with
58 average CNV length ~226kb and the majority of called variants singleton in the cohort (Figure
59 1a,b). This translates to notable burden of variation for nearly all individuals, with 439,464
60 (93.1%) of the individuals possessing at least one CNV detectable at kilobase- resolution
61 (Figure 1c,d). Among individuals with at least one CNV, we estimate an average burden of 5.5
62 variants covering >200kb of genomic sequence (median 3 variants affecting ~100kb, Figure
63 1c,d). While in-line with previous reports⁸, these estimates of individual-level burden are likely
64 conservative, as regions where array markers are sparse or missing limit the accuracy of variant
65 calling. Furthermore, we are unable to call smaller (<1kb) variants due to inconsistent marker
66 density across all chromosomal regions on the Axiom and BiLEVE UK Biobank genotyping
67 arrays. This limitation is visible in the histogram of called CNV lengths (Figure 1a); we call
68 substantially fewer variants on the order of hundreds of base-pairs than on the order of
69 thousands.

70

71 We also observe a highly non-uniform burden of variation across genomic position, with
72 breakpoints most common near the ends of chromosomes, and at known regions of segmental
73 duplication (Figure 1e). Among them are *1p36*, *8q24.3*, *9q34.3*, and *19q13*, all of which have
74 associated microdeletion syndromes causing developmental delay with uniquely characteristic
75 growth patterns^{16–19}. Other CNV-hotspots like *6p21.33*, which contains the major
76 histocompatibility complex protein gene family, may be influenced by high marker density (in
77 this case for HLA allelotyping) in addition to these biological features which underlie structural
78 mutagenesis. However, these loci do not categorically correspond to areas where structural
79 variation is commonly observed in the population (Figure S1). For example, *1p36* and *19q13* are
80 also the respective sites of common CNVs overlapping *RHD* and *FUT2* (Rhesus and Lewis
81 blood groups), but there are no such common variants within the telomeric *16p13* cytoband.



82

Figure 1: Burden and distribution of copy number variation in UK Biobank. (A) Log-scale histogram of CNV lengths. Mean length (dashed line) is 226.5kb. **(B)** Cumulative density of CNV allele count (AC), displayed in log-log axes. Average AC is 5.5, but average frequency as experienced by the population (weighted by count, hence AC^2) is $\sim 1.6\%$. **(C)** Histogram of CNV counts and **(D)** log-scale base-pairs affected by CNV per individual. Sample-level burden is heavy-tailed, with the average individual carrying 4.2 variants (dashed line), affecting mean ~ 207.6 kb of genomic sequence. **(E)** Genome-wide density of CNV, defined as the number of unique CNVs overlapping 10 megabase (Mb) windows tiling each chromosome. Hotspots of structural variation are labeled by cytogenic band.

83

84

85 Survivorship bias due to genetic selection against early-onset diseases

86
87 We estimate gene-level intolerance to structural variation by adapting a method for estimating
88 regional selective constraint⁸. Relative to the general population, the volunteers within the UK
89 Biobank are described to have a “healthy-cohort” enrollment bias²⁰ and were enrolled between
90 the ages of 40 to 69, which informs our findings. Within the tail of positive constraint z-scores,
91 which indicate the strongest intolerance to structural variation, we observe enrichment for genes
92 which cause early onset diseases, particularly cancer. Among the top fifteen constrained genes
93 (Table 1) are *BRCA1* and *BRCA2*, which are associated with early-onset breast cancer^{21,22};
94 *MLH1*, *MSH2*, *MSH6*, which cause early onset colorectal cancer (Lynch syndrome)^{23–25}; and
95 *ATM* and *APC*, which are involved with mismatch repair cancers^{26,27}.

	Constraint z	CNV pLI
BRCA2	3.262	0.9905
BRCA1	2.470	0.9830
APC	1.997	0.9421
ATM	1.196	0.9886
MSH2	1.192	0.9875
MLH1	1.176	0.9959
MSH6	0.918	0.9928
SBDS	0.875	0.9739
RB1	0.872	0.9552
SPATA31D1	0.869	0.9979
CYP3A4	0.863	0.9979
OTOP1	0.848	0.9930
PABPC3	0.847	0.9923
KRT16	0.844	0.9979
ZNF302	0.844	0.9979

Table 1: (Left) 15 genes most intolerant to copy number variation. Columns are gene label, constraint z-score, and probability of CNV intolerance (pLI, see Methods for definitions).

Table 2: (Below): 15 pathways most enriched for constrained genes (t-test, gene set members versus all others). Columns are GO pathway ID/name, change in z-score between set and non- set members, indicating mean strength of selective effect in the pathway, and *p*-value.

	Pathway	Δz	P
GO:0000137	Golgi cis cisterna	0.42	1.16×10^{-31}
GO:0045095	keratin filament	0.26	7.29×10^{-31}
GO:0052697	xenobiotic glucuronidation	0.49	1.03×10^{-21}
GO:0005515	protein binding	0.07	1.24×10^{-21}
GO:0031424	keratinization	0.18	1.51×10^{-21}
GO:0000800	lateral element	0.46	3.46×10^{-21}
GO:0008194	UDP-glycosyltransferase activity	0.36	6.13×10^{-21}
GO:0015020	glucuronosyltransferase activity	0.35	1.32×10^{-17}
GO:0008202	steroid metabolic process	0.28	9.31×10^{-17}
GO:0005132	type I interferon receptor binding	0.35	1.78×10^{-16}
GO:0008274	gamma-tubulin ring complex	0.39	2.66×10^{-16}
GO:0002323	natural killer cell activation involved in immune response	0.36	1.54×10^{-14}
GO:0005131	growth hormone receptor binding	0.54	4.33×10^{-14}
GO:0052696	flavonoid glucuronidation	0.45	5.40×10^{-14}
GO:0042954	lipoprotein transporter activity	0.40	5.78×10^{-14}

96
97
98 Selections from the most highly constrained pathways from Gene Ontology Consortium²⁸
99 resources (Table 2) also suggest strong intolerance to CNV for genes involved with core
100 biological processes like protein binding, cellular structural integrity (keratinization),
101 development (growth hormone receptor binding), and immune regulation (natural killer cell
102 activation). Similar results at the gene- and pathway-level are observed for deletion-specific

103 constraint (Table S1,S2), whereas duplication-specific analysis suggests autoimmune-related
104 genes and pathways are most strongly intolerant to dosage effects (Table S1,S2). These results
105 indicate strong selective effects occurring prior to enrollment in the UK Biobank during childhood
106 and early adulthood against loss of function variation in core developmental, metabolic, and
107 tumor-suppressing genes, and against dosage-altering variation in immune-related genes.

108

109 **Association testing identifies CNVs at novel and syndromic loci**

110

111 We compute genome-wide associations across 1,893 phenotypes for 8,274 common CNVs
112 observed at 0.005% allele frequency (1 in 20,000) in our cohort, using regression as
113 implemented in the analysis software PLINK²⁹. We also perform L1-regularized regression for
114 rare-variant burden tests, pooled by gene. For these tests, we measure the net effect of rare
115 CNVs (AF < 0.1%) overlapping within 10kb of the gene region as defined by HGNC³⁰ for 7,614
116 protein coding genes with at least 5 individuals affected by such a variant. A complete list of
117 phenotypes analyzed is available in Table S3. Here, we describe representative results for one
118 common disease and one quantitative measure with established genetic risk factors and large
119 sample sizes in UK Biobank: acute coronary artery disease (CAD) and body mass index (BMI).

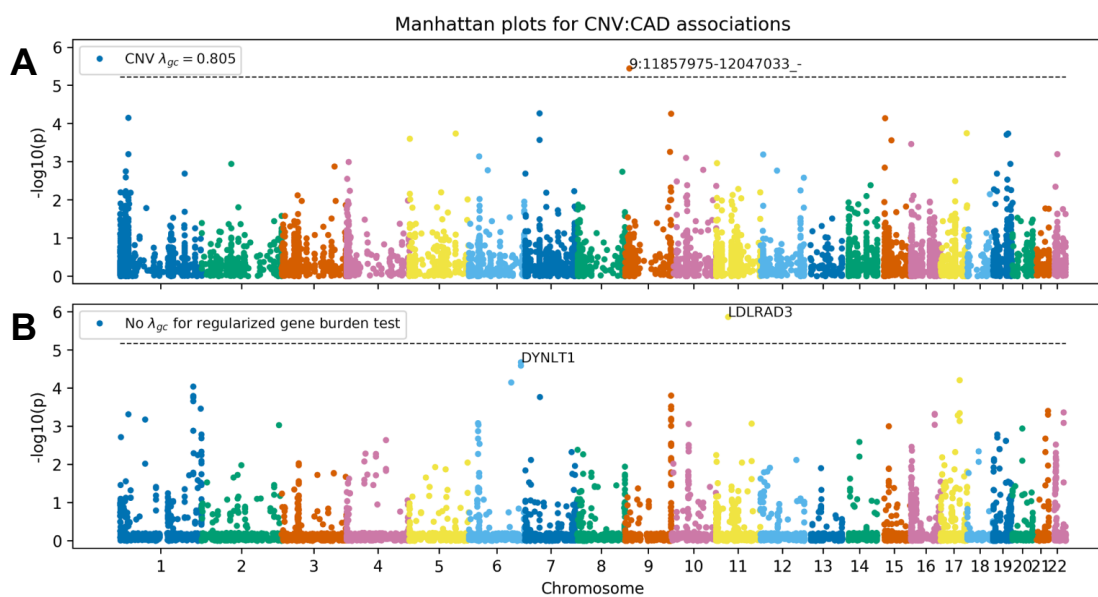
120

121 For Acute CAD, we identify one statistically significant ($p < 6 \times 10^{-6}$) association after Bonferroni
122 correction for the common CNV GWAS: an intergenic deletion at chromosome 9p23. Intergenic
123 variants at the 9p21 locus have been implicated in previous association studies of blood-based
124 biomarkers relevant to cardiac outcomes, specifically, decreases in hematocrit and hemoglobin
125 concentration³¹, as well as carotid plaque burden³². A recent meta-analysis³³ using data from
126 UK Biobank and CARDIoGRAMplusC4D identified a lead variant in the vicinity of this locus
127 (rs2891168) associated with 6% unit increase in risk for similarly defined coronary artery
128 disease. However, the CNV we here identify confers an estimated 12.4-fold increased risk
129 (95%CI: 4.3-35.9, $p=3.6 \times 10^{-6}$) and is at least 2Mb distant from the nearest SNPs (rs10961206)
130 at genome-wide significance near the 9p21/9p23 locus in the meta-analysis. This and the
131 absence of linkage between the 9p23 CNVs and rs10961206 ($r = 0.013$) are suggestive of
132 independent effects.

133

134 Gene-level burden testing of rare CNVs in individuals with CAD implicates *LDLRAD3*, a member
135 of the low density lipoprotein (LDL) receptor family. CNVs called in this gene are predominantly
136 deletions affecting the coding sequence — in aggregate ($n=27$), these variants confer an
137 estimated 10-fold increase in risk of Acute CAD (95% CI: 3.9-25.6, $p=1.4 \times 10^{-6}$). Though the role
138 of lipoprotein receptors in cholesterol metabolism is a well established mechanism of risk for
139 cardiovascular disease, *LDLRAD3* is not known to participate in cholesterol metabolism. It is,
140 however, a receptor widely expressed throughout adult tissues which may participate in
141 proteolysis in the central nervous system^{34,35}. We therefore sought to replicate these findings
142 using two-sample mendelian randomization³⁶ on expression quantitative trait loci (eQTLs) from
143 CAD summary statistics from a CARDIoGRAMplusC4D meta-analysis³⁷. We identify a nominally
144 significant protective effect between an eQTL increasing expression of *LDLRAD3* and CAD
145 (OR=0.85 [95%CI: 0.62-0.97], $p=0.012$), the direction of which is consistent with a dosage
146 model of *LDLRAD3*-mediated risk for CAD.

147



148

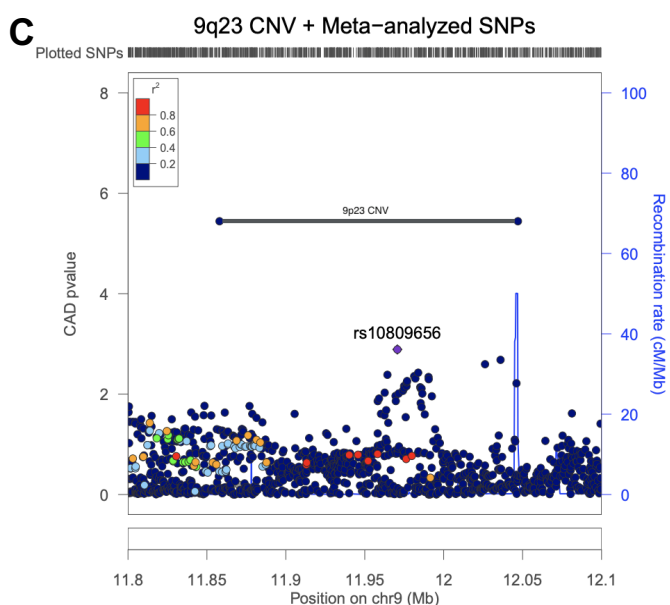
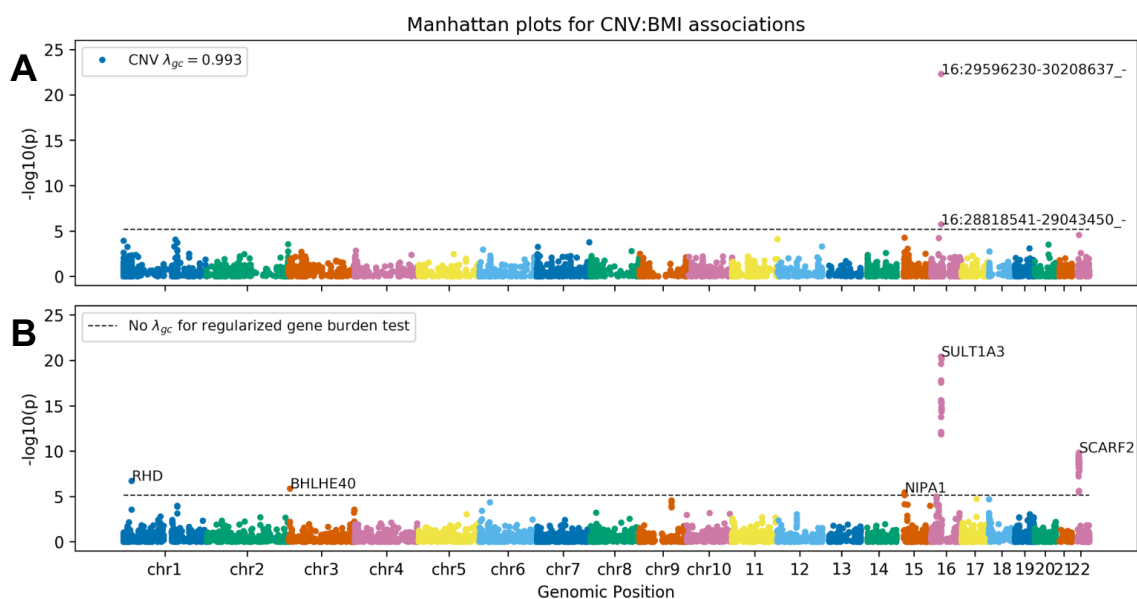


Figure 2: Genome-wide CNV associations for acute coronary artery disease (CAD). Manhattan plots for (A) genome-wide association of common copy number variants, and (B) genome-wide burden test of rare variants. (C) LocusZoom³⁸ of 9p23 CNV and summary statistics from meta-analysis of CAD³³ colored by marker LD with lead regional GWAS SNP (rs10809656) in HapMap³⁹ European samples.

149

150 We also find two significant associations for BMI, both deletions at chromosome *16p11.2*, a
 151 locus implicated in syndromic early onset obesity and developmental delay. Each of these
 152 CNVs appears to correspond to a distinct form of *16p11.2* deletion syndrome. The smaller
 153 $\sim 220\text{kb}$ deletion ($\beta = 4.5 \text{ kg/m}^2$ [95%CI: 2.7-6.3 kg/m^2], $p = 1.8 \times 10^{-6}$, AC=35) has been
 154 associated with early onset obesity, and spans *ATXN2L*, *TUFM*, *SH2B1*, *ATP2A1*, *RABEP2*,
 155 *CD19*, *NFATC2IP*, *SPNS1*, and *LAT*, with *SH2B1* the suspected causal obesity gene³. Obesity
 156 is also a phenotypic consequence of a larger $\sim 593\text{kb}$ deletion ($\beta = 7.8 \text{ kg/m}^2$ [95%CI: 6.2-9.4
 157 kg/m^2], $p = 5.0 \times 10^{-23}$, AC=58), which is further associated with neurodevelopmental delay and
 158 related conditions⁴. However, this deletion spans a wholly distinct set of genes which are
 159 suspected to play complex dosage-dependent roles in the phenotypic consequences of the

160 syndrome⁴⁰. As both subtypes of *16p11.2* deletion syndrome may present in early childhood, it
 161 is noteworthy that the effect we measure on BMI is in a cohort comprised entirely of older
 162 individuals, indicating burden of adult disease associated with the CNV locus.
 163



164

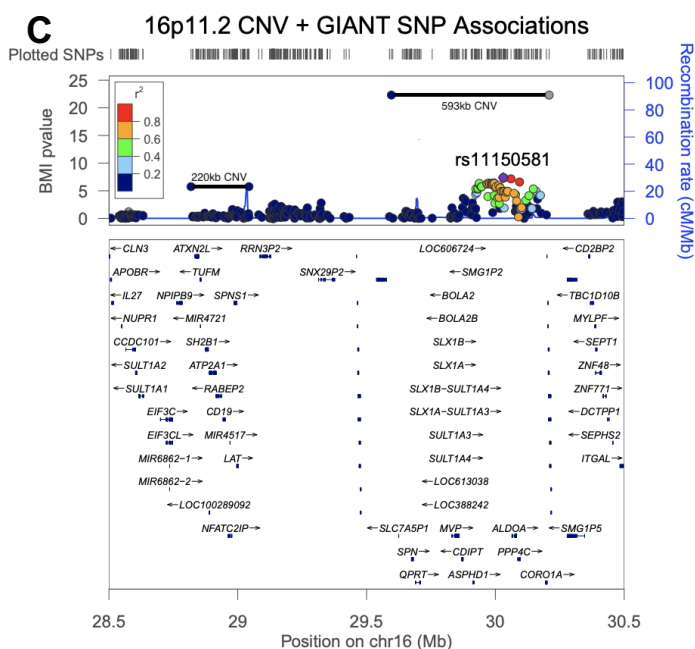


Figure 3: Genome-wide CNV associations for body mass index (BMI). Manhattan plots for (A) genome-wide association of common copy number variants, and (B) genome-wide burden test of rare variants. (C) LocusZoom of *16p11.2* CNVs and BMI summary statistics from the GIANT study⁴¹, colored by marker LD with the lead SNP at the locus (rs11150581), computed from HapMap European samples.

165 After controlling for multiple comparisons, burden testing for BMI identifies a group of genes at
 166 chromosome *22q11.2* and recapitulates the list of genes affected by each of the *16p11.2*
 167 deletions. Variation at the *22q11.2* locus also constitutes one of the first named microdeletion
 168 syndromes, DiGeorge syndrome, which has variable phenotypic consequences including
 169 craniofacial dysmorphisms and conotruncal congenital heart disease, along with increased risk
 170 for an adverse cardiovascular outcomes and neuropsychiatric disease later in life⁵. Among
 171

172 individuals affected with 22q11.2 deletion syndrome obesity is a recognized manifestation of
 173 disease⁴², and we estimate a 0.3-0.4 point increase in BMI for genic CNVs near 22q11.2. as
 174 well as 0.55-0.7 for genic CNV at 16p11.2. The presence of these associations in a large
 175 volunteer cohort offers further evidence that syndromic alleles may contribute to the risk of
 176 common diseases in the general population.
 177

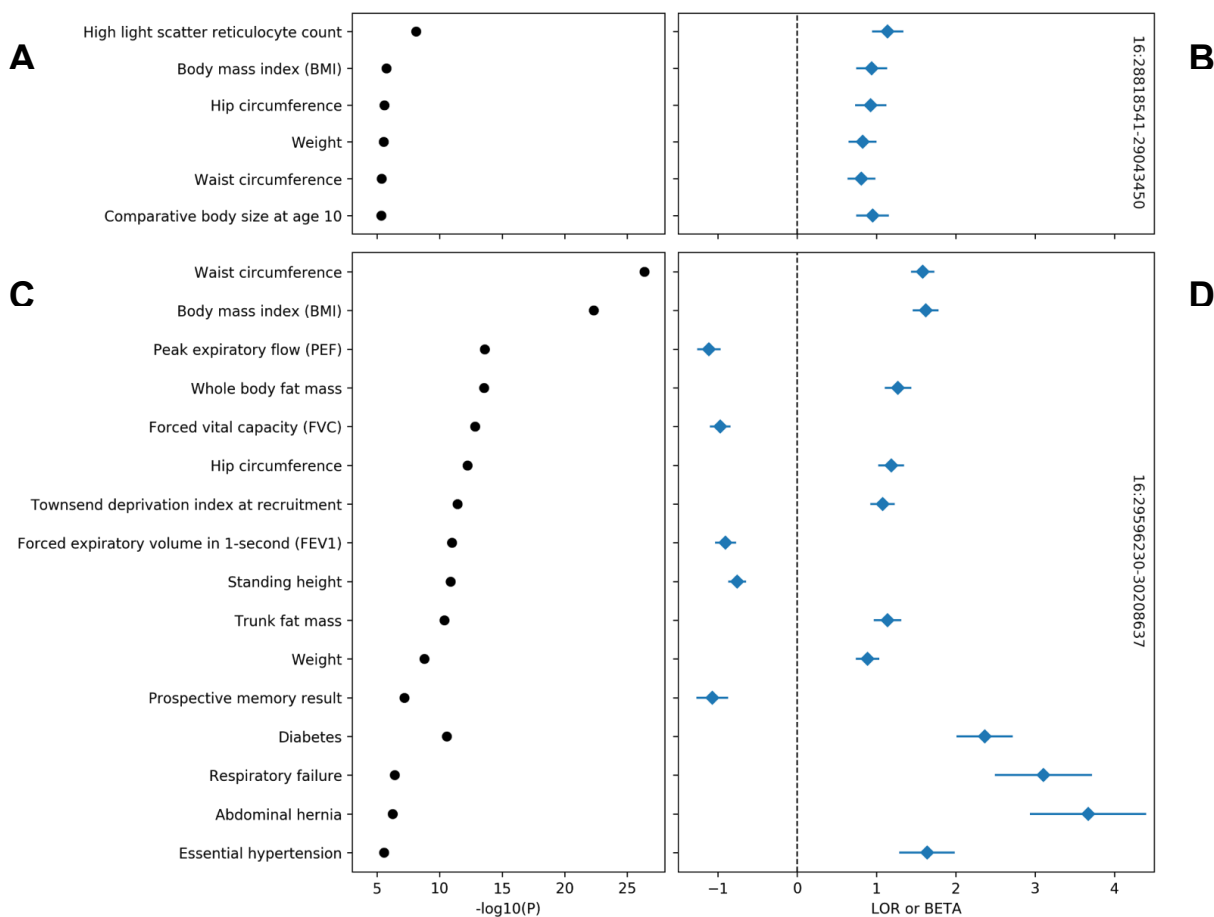


Figure 4: PheWAS of 16p11.2 CNVs. Selected genome-wide significant ($p < 6 \times 10^{-6}$) associations for 220kb (top panels) and 593kb (bottom panels) 16p11.2 CNVs. Traits are grouped by type (binary/quantitative) then sorted by p -value (left panels). Log-odds ratio and standardized betas (right panels) align with trait names on the y-axis, with the horizontal dashed line separating positive and negative direction of association.

179
 180 Phenome-wide associations for each of the CNVs at 16p11.2 further highlight changes in
 181 biomarkers, biomeasures, and increased risk of common disease, consistent with high BMI over
 182 the course of a lifetime (Figure 4). Genome-wide significant phenotypes for the 220kb CNV
 183 recapitulate the established syndromic effects from early onset obesity. We observe significant
 184 increases, on the order of one standard deviation, in weight, BMI, hip and waist circumference,
 185 reticulocyte count, and comparative body size at age 10 for these individuals. The larger 593kb
 186 CNV associates with similar measures of body size and fat, as well as hypertension, diabetes,
 187 and abdominal hernia. These results are also indicative of effects due to developmental delay;

188 namely, decreased measures of memory, higher Townsend deprivation, and lower lung capacity
189 (FEV, FVC) with higher associated risk of respiratory failure. Taken together these results
190 highlight the variable expressivity of CNV-related disease, as well as its long-term effects across
191 the medical phenome.

192

193 **Discussion:**

194

195 In calling copy number variants and performing genetic association studies at scale from a large
196 cohort of array-genotyped individuals with richly annotated phenotype data, we provide a
197 portrait of the phenome-wide burden of genomic copy number variation. Our estimates of the
198 individual-level burden of CNV and population-wide allele frequencies are consistent with
199 previous reports, the deep phenotypic information available in the UK Biobank permits more
200 finely tuned measures of the genic intolerance to CNV which include estimates of variation
201 absent from our cohort of predominantly healthy, middle-aged individuals. We consider both
202 rare and common CNVs in our association studies and identify effects of previously known
203 (*22q11.2*, *16p11.2*) syndromes and potentially novel (*9p23*) loci which have not yet been
204 characterized or associated with disease.

205

206 Our study has significant limitations which inform our analysis. While arrays are an inexpensive
207 way to genotype large cohorts, permitting straightforward algorithms to infer the presence of
208 structural variation, the resulting CNV calls are limited by the density and placement of markers
209 across chromosomes. For UK Biobank genotyping arrays in particular, there are large portions
210 of genomic sequence with low marker density (in particular near centromeric regions) which
211 bias our resulting genotype calls away from such regions. Array-derived CNV calls also suffer
212 from limitations, in their inability to differentiate other structural events like inversions or
213 translocations, or to determine breakpoint position at base-pair resolution. Complicating these
214 barriers is the fact that our sample was genotyped on two distinct arrays, which may cause
215 identical CNVs to present with different breakpoints across individuals in the call set. Our choice
216 to present gene-level burden tests which include the vast majority of variants included in our
217 CNV GWAS was informed by this realization.

218

219 Our associations are also heavily impacted by a known “healthy-cohort” bias, which results in
220 null results for several phenotypes with known genetic contributions; notably, there are no hits in
221 our burden tests for cancers other than leukemia and lymphomas. With this in mind, our
222 constraint scores constitute a sobering observation of genetic survivorship bias. Estimates of
223 gene-level intolerance to structural variation are derived from people who did not enroll in UK
224 Biobank; the absence of individuals who were not healthy enough to participate or did not
225 survive until age 40 constitutes an enrollment bias against severe early-onset disease. We take
226 this opportunity to honor these non-participating individuals and their implicit contribution to our
227 understanding of genetic disease. The observation of selection bias colors the interpretation of
228 genetic findings from UK Biobank in general, as the cohort is relatively depleted of disease of
229 early-onset morbidity and mortality and any genetic variation associated with these diseases will
230 likewise be difficult to detect. While UK Biobank is unprecedented in size, scope, and scientific

231 yield, our data illustrate that the anticipated findings from the proliferation of large biobank
232 studies around the world will be influenced by implicit and explicit barriers to participation.

233
234 Despite selection against high-penetrance alleles causing early-onset disease, we detect a
235 novel and strong association for coronary artery disease at *LDLRAD3*. While this locus has prior
236 putative association with bone mineral density⁴³, existing large-scale GWAS do not detect a
237 strong association with coronary artery disease or established cardiometabolic risk-factors.
238 However, the absence of gene level intolerance to truncating and missense variation in
239 *LDLRAD3* does not suggest a compelling rationale for important biological function or role in
240 disease. In our study, CNVs at this locus are associated with some established cardiometabolic
241 risk factors, such as diabetes onset, smoking status, and arterial stiffness, but not obesity or
242 other fat-related phenotypes (Figure S6). Consistent with our findings that a decrease in
243 *LDLRAD3* dosage increases the risk of disease, a strong eQTL increasing *LDLRAD3*
244 expression decreases the risk of disease when used as an instrument in a two-sample
245 mendelian randomization in a large-scale study of coronary artery disease. Thus, our findings of
246 an mRNA dosage effect are replicated at the gene level. These results highlight the utility of
247 analyzing genic CNV which, when directly impacting mRNA dosage, offer a more easily
248 interpretable mechanism distinct from alterations of protein structure or small changes in
249 transcriptional regulation.

250
251 The observation of variation at the *16p11.2* and *22q11.2* loci sheds further light on the
252 penetrance of syndromic loci in the general population. The *16p11.2* recurrent microdeletion
253 syndrome was first described in individuals with autism and neuropsychiatric disease and may
254 include seizures, brain and other anatomic abnormalities. Even accounting for the enrollment
255 bias inherent to the UK Biobank our PheWAS detects a modest relationship to neurocognitive
256 measurements via secondary markers of intellectual differences such as prospective memory
257 for one variant at this locus. People carrying variation at the *22q11.2* locus within the general
258 population are known to be at increased risk of neuropsychiatric diseases⁴⁴ for which variable
259 phenotypic penetrance is well recognized³⁴⁵. To wit, individuals with genetic variation at both
260 syndromic loci were by and large sufficiently healthy and capable of volunteering to participate
261 in the Biobank. Our findings support a growing recognition that the penetrance and effect sizes
262 of syndromic alleles will likely require revision in the context of broad population-based surveys
263 of rare genetic variation^{46,47}.

264
265 Our findings add to the growing body of literature measuring global burden of structural variation
266 across healthy and diseased individuals. Our estimates of the effects of common CNV suggest
267 a notable role of structural variation in population-wide burden of common disease, and suggest
268 genomic loci where novel CNV-derived syndromic disease may exist. For rare variants, our data
269 offer broad phenotypic characterization of the effects of gene-specific knockouts, which may
270 inform development of pharmacological and genetic therapies. While the functional
271 consequence and pathogenicity of missense, synonymous, and noncoding single nucleotide
272 variation within a gene may be difficult to classify, the mechanism of most genic CNV are clear:
273 a dosage effect upon mRNA transcription. This population-scale catalog of variation and the
274 described associations with a multiplicity of diseases should be of immediate use by genetic

275 clinicians in classification of novel and rare CNV detected in clinical testing. Full support for
276 gene- and variant-level browsing is forthcoming in a future version of the Global Biobank Engine
277 (biobankengine.stanford.edu). In the interim, summary statistics from association studies
278 described here, as well as for all phenotypes present on the engine, are freely available for
279 download on the site. We hope that these data will be leveraged to empower future analyses of
280 the phenome-wide effects of structural variation and gene-level dosage effects.

281

282

283

284 Acknowledgements:

285

286 This research has been conducted using the UK Biobank Resource under application numbers
287 24983, 16698, 13721, and 15860. We thank all the participants in the study. The primary and
288 processed data used to generate the analyses presented here are available in the UK Biobank
289 access management system (<https://amsportal.ukbiobank.ac.uk/>) for application 24983,
290 "Generating effective therapeutic hypotheses from genomic and hospital linkage data"
291 (<http://www.ukbiobank.ac.uk/wp-content/uploads/2017/06/24983-Dr-Manuel-Rivas.pdf>), and the
292 results are displayed in the Global Biobank Engine (<https://biobankengine.stanford.edu>).

293

294 M.A.R. is supported by Stanford University and a National Institute of Health center for Multi-
295 and Trans-ethnic Mapping of Mendelian and Complex Diseases grant (5U01 HG009080). This
296 work was supported by National Human Genome Research Institute (NHGRI) of the National
297 Institutes of Health (NIH) under awards R01HG010140. The content is solely the responsibility
298 of the authors and does not necessarily represent the official views of the National Institutes of
299 Health.

300 Methods

301
302 CNVs were called using PennCNV v1.0.4 on raw signal intensity data from each genotyping
303 array. Phenotype data was derived from data-fields collected for UK Biobank corresponding to
304 various body measurements, biomarkers, disease diagnoses and medical procedures from
305 medical records, as well as a questionnaire about lifestyle and medical history. Summary-level
306 data from all statistical tests described here, as well as more thorough documentation on
307 phenotyping, will be available on the Global Biobank Engine¹⁴ (biobankengine.stanford.edu) and
308 can be found in related publications⁴⁸.

309 310 CNV calling in UK Biobank:

311
312 Methods for genetic data acquisition and quality control as performed by the UK Biobank have
313 been previously described¹³. In brief, two similar arrays were used for targeted genotyping
314 within the study population: the UK BiLEVE Axiom Array ($n=49,950$) by Affymetrix and the UK
315 Biobank Axiom Array ($n=438,427$), which was custom-designed by Applied Biosystems.
316 Samples and array markers were subject to threshold-based filtration and quality control prior to
317 public release. Specifically, markers were tested for discordance across control replicates,
318 departures from Hardy-Weinberg equilibrium, as well as effects due to batch, plate, array, and
319 sex; affected markers were set as missing in affected batches or removed. Similarly, samples
320 were tested for missingness ($>5\%$) and heterozygosity across a set of high-quality markers, but
321 samples identified as low quality ($n=968$) were not excluded. We also chose to include these
322 samples in this analysis, considering that large structural variants may have been responsible
323 for their poor quality with respect to metrics used for filtration.

324
325 We used PennCNV v1.0.4¹⁵ to call CNVs within each of the 106 genotyping batches from UK
326 Biobank. We first estimate genomic runs of heterozygosity (RoH) for each sample using a
327 previously developed pipeline in PLINK^{49,50} using the `--homozyg` option. We then select $n=100$
328 samples with total RoH covering less than 20MB to train a hidden markov model (HMM) of copy
329 state on each chromosome. HMM training was initialized with conditions optimized for
330 Affymetrix arrays (`affygw6.hmm`), provided in PennCNV resources. We used the general calling
331 mode, which performs likelihood-based testing for copy-number state ($CN=0,1,2,3,4$) at each
332 input marker using its log-normalized signal intensity and allele balance in a given sample. We
333 also apply adjustment for GC content across sites using waviness factor correction⁵¹. After CNV
334 calling, we exclude 1,360 samples with over 30 called CNVs from downstream analysis,
335 resulting in a cohort of 472,228 individuals with 278,455 unique variants.

336 337 Gene-level constraint estimation:

338
339 Regional selective constraint to CNV was estimated for all protein-coding genes, with genic
340 CNV defined as any variant overlapping within 10kb of the HGNC gene region. We estimate a
341 null model of structural mutation empirically, and model burden of genic CNV as a linear
342 function of gene size, fraction of genic sequence covered by regions of segmental duplication as
343 extracted from the UCSC Genome Browser^{52,53}. We also account for biased observations due to

344 array genotyping by including the number of genic markers as a covariate. The formula for this
345 null model can be written as:

346

$$347 \quad n_{cnv} = \beta_1 \cdot len(gene) + \beta_2 \cdot frac(segdup) + \beta_3 \cdot n_{markers} + \epsilon$$

348

349 From this model, we compute constraint z-scores for each gene using its negated standardized
350 residual for each gene, winsorizing the negative tail at the lowest 5% of values. We also
351 compute the probability of loss of function intolerance (pLI) as the non-normalized residual over
352 the number of expected CNV, with negative values rounded to zero.

353

354 Genetic associations:

355

356 Variant-level associations in UK Biobank were estimated with PLINK v2.00a (5 Jan 2018). We
357 used the --glm firth-fallback option for computation. This option is a hybrid algorithm for logistic
358 regression which defaults to a standard regression solver for computation, falling back to Firth's
359 regression (<https://cran.r-project.org/web/packages/logistf/index.html>) in cases where one of the
360 cells of the 2x2 contingency table is zero, or where the traditional method fails to converge in a
361 pre-specified number of iterations. These analyses were performed in a subset of 337,538
362 unrelated individuals of self-reported white British ancestry, and were controlled for age, sex,
363 and 4 marker-based genomic principal components from the UK Biobank PCA calculation. To
364 ensure adequate power for estimating genetic effects, we perform these tests on 8,274 CNVs
365 observed at a frequency of 0.005% (1 in 20,000, or 18 individuals) in the whole sample of
366 individuals with called CNVs.

367

368 Gene-level burden tests were conducted across all gene:phenotype pairs for genes with at least
369 5 individuals with overlapping CNV. Genic burden was encoded as a binary variable which
370 indicates whether an individual has a CNV which contains any overlap within 10,000 base pairs
371 of the HGNC gene region. CNVs which overlapped several gene regions were used for analysis
372 in each gene. We treat deletions and duplications identically, with the assumption that any CNV
373 which overlaps a gene in this fashion will disrupt its normal function. Effects of genic CNV
374 burden were estimated by linear and regularized logistic regression with the python package
375 Statsmodels, respectively computed using the .fit() and .fit_regularized() methods. We included
376 the following as covariates in both models: age, sex, four marker-based genomic principal
377 components from UK Biobank's PCA calculation, and the number and combined length of CNVs
378 in each individual.

379

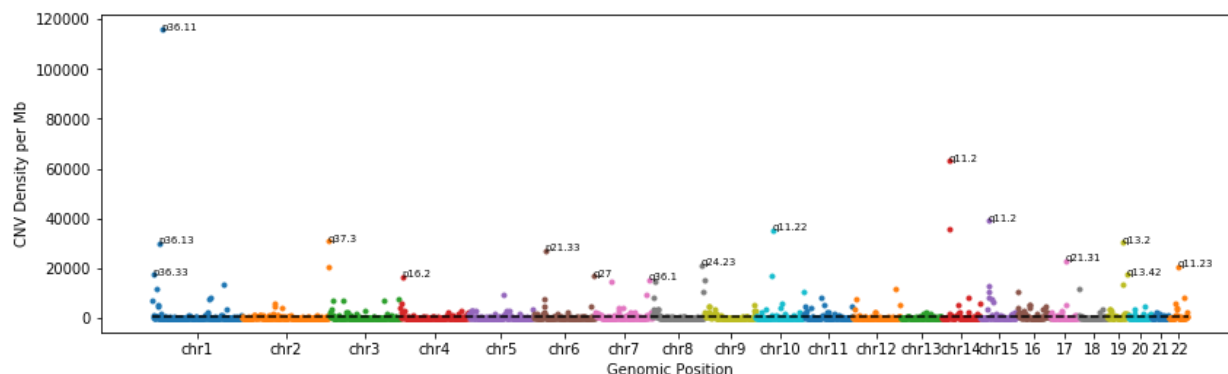
380 Two-sample mendelian randomization was performed via the MR Base web app using GWAS
381 summary statistics for *LDLRAD3* expression QTLs from a CARDIoGRAMplusC4D meta-
382 analysis³⁷. We report Wald summary statistics from inverse-variance weighted Egger
383 regression; these are the default analysis options for the web interface.

384

385

386
387

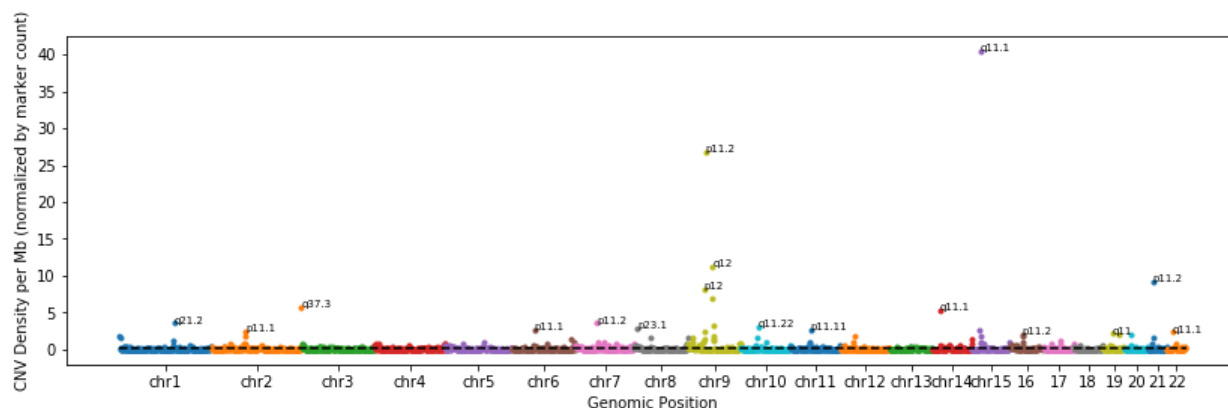
Supplementary Figures and Tables:



388

Figure S1: CNV density weighted by allele count in UK Biobank. Per-megabase genomic density of CNV, weighted by number of observations across all samples in UK Biobank. Variants are counted by whether the CNV has any overlap with 10 megabase (Mb) windows tiling each chromosome. Selected hotspots of structural variation are labeled by the region's corresponding cytogenetic band.

389
390
391



392

Figure S2: CNV density normalized by array marker density in UK Biobank. Variants are counted by whether the CNV has any overlap with 10 megabase (Mb) windows tiling each chromosome, then divided by the number of markers in the window. Regions with no array markers are defined to have density of zero. Selected hotspots of structural variation are labeled by the region's corresponding cytogenetic band.

393

	Deletion z		Duplication z
BRCA2	2.974	HLA-DRB1	0.581
BRCA1	2.212	FRG2B	0.581
APC	1.857	SPATA31D1	0.581
ATM	1.104	SLC35G6	0.580
MSH2	1.087	NAT8	0.580
MLH1	1.072	TUBB8	0.580
MYH7	0.934	CSH2	0.580
PMS2	0.880	ZNF302	0.580
TTN	0.833	CSHL1	0.580
PVRL2	0.823	GH1	0.580
MSH6	0.821	CGB2	0.579
SBDS	0.816	FAM215A	0.579
CYP3A4	0.816	OR4F17	0.579
SPATA31D1	0.815	CGB5	0.579
OTOP1	0.809	CGB7	0.579

394

Table S1: 15 genes most intolerant to overlapping deletion (left), and whole-gene duplication (right), with respective constraint z-scores.

395

396

Deletion-intolerant Pathway		Δz	P
GO:0045095	keratin filament	0.23	2.22×10^{-28}
GO:0000137	Golgi cis cisterna	0.36	2.10×10^{-26}
GO:0052697	xenobiotic glucuronidation	0.44	5.29×10^{-20}
GO:0008194	UDP-glycosyltransferase activity	0.32	4.79×10^{-19}
GO:0005515	protein binding	0.05	1.19×10^{-17}
GO:0000800	lateral element	0.39	1.36×10^{-17}
GO:0031424	keratinization	0.15	1.33×10^{-16}
GO:0015020	glucuronosyltransferase activity	0.31	2.56×10^{-16}
GO:0042954	lipoprotein transporter activity	0.37	3.33×10^{-14}
GO:0005131	growth hormone receptor binding	0.50	3.48×10^{-13}
GO:0008202	steroid metabolic process	0.24	4.49×10^{-13}
GO:0046703	natural killer cell lectin-like receptor binding	0.39	7.08×10^{-13}
GO:0008274	gamma-tubulin ring complex	0.32	1.08×10^{-12}
GO:0005132	type I interferon receptor binding	0.28	1.40×10^{-12}
GO:0052696	flavonoid glucuronidation	0.40	1.59×10^{-12}

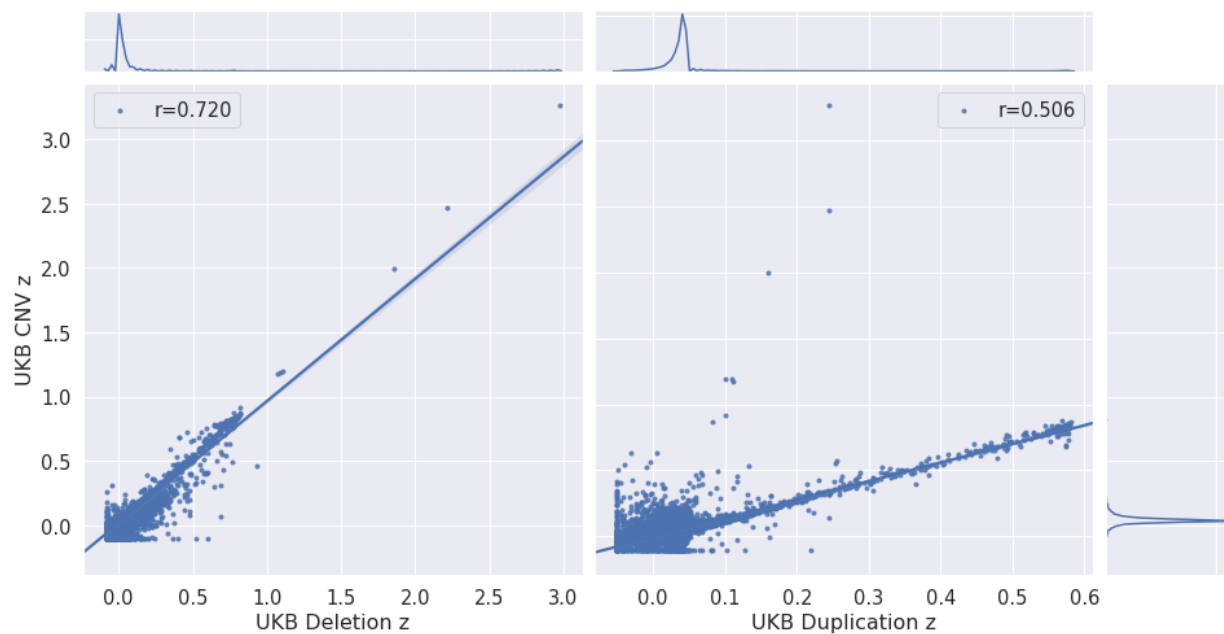
397

398

Duplication-intolerant Pathway		Δz	P
GO:0000137	Golgi cis cisterna	0.34	2.57×10^{-46}
GO:0045095	keratin filament	0.19	4.45×10^{-34}
GO:0005515	protein binding	0.05	1.56×10^{-30}
GO:0031424	keratinization	0.14	2.12×10^{-23}
GO:0008202	steroid metabolic process	0.22	1.60×10^{-22}
GO:0005132	type I interferon receptor binding	0.27	8.04×10^{-20}
GO:0008194	UDP-glycosyltransferase activity	0.24	5.47×10^{-19}
GO:0005801	cis-Golgi network	0.17	6.16×10^{-19}
GO:0046703	natural killer cell lectin-like receptor binding	0.37	6.81×10^{-19}
GO:0052697	xenobiotic glucuronidation	0.30	9.68×10^{-18}
GO:0002323	natural killer cell activation involved in immune response	0.27	1.49×10^{-17}
GO:0033141	positive regulation of peptidyl-serine phosphorylation of STAT protein	0.24	7.73×10^{-17}
GO:0006805	xenobiotic metabolic process	0.16	1.73×10^{-15}
GO:0005131	growth hormone receptor binding	0.39	2.39×10^{-15}
GO:0042271	susceptibility to natural killer cell mediated cytotoxicity	0.24	4.69×10^{-15}

399

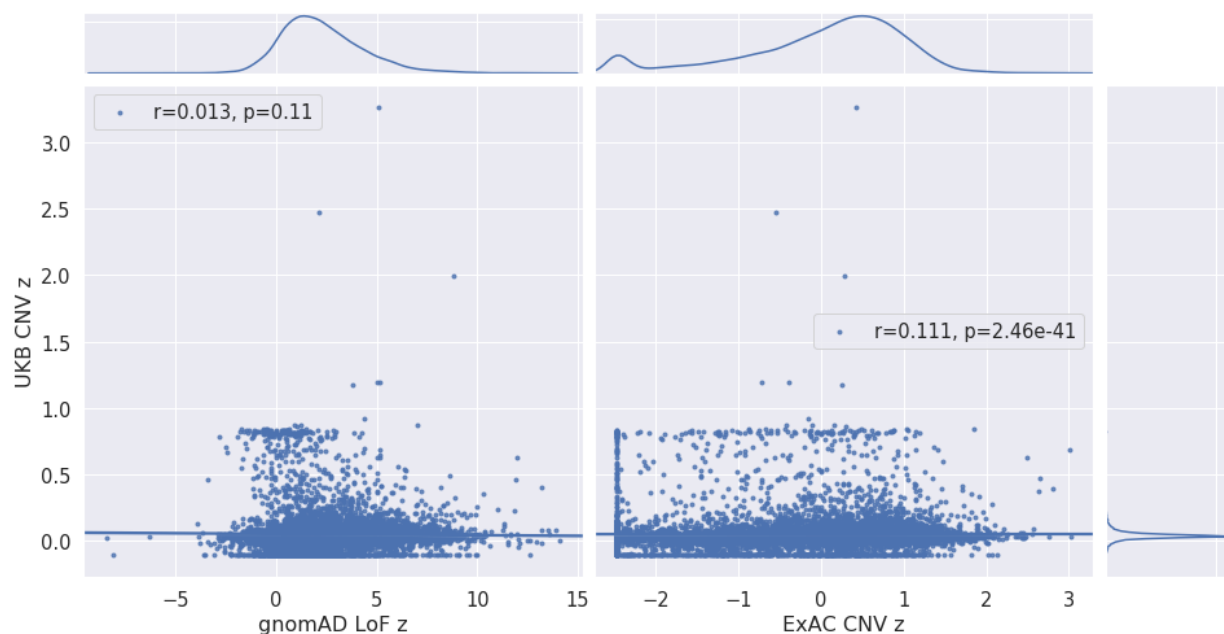
Table S2: GO pathways most intolerant to overlapping deletion (top), and whole-gene duplication (bottom), with change in constraint z-scores and significance thereof (t-test) relative to other pathways.



400

Figure S3: Correlation between intolerance measures for partial-gene deletion, whole-gene duplication, and CNV burden. The legend for each panel denotes correlation (Spearman's r) between burden-constraint and each other measure. Kernel density estimates for each distribution of constraint scores are in the panels opposite their corresponding axis labels.

401



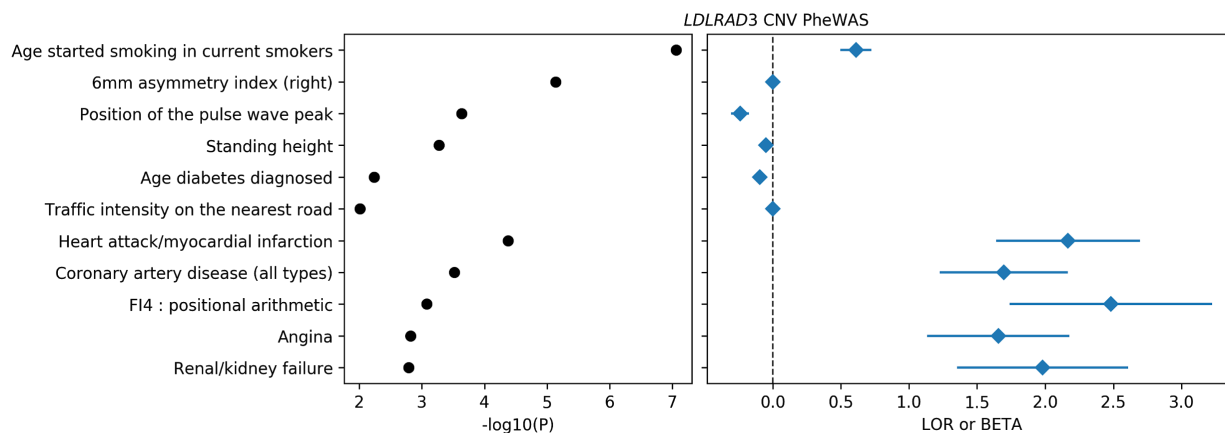
402

Figure S4: Distribution of constraint z-scores from UK Biobank and ExAC/gnomAD. Our measures of gene-level intolerance to structural variation show nominal correlation with gnomAD loss of function constraint z-scores (Spearman's $r = 0.013$, left), and modest correlation with CNV-intolerance in ExAC (Spearman's $r = 0.11$, right panel). Gaussian kernel density estimates for each distribution of z-scores are opposite their corresponding axes.

While correlation between constraint measures across datasets is non-random, we suspect cohort-specific effects and varying definitions of genic burden of variation drive these departures. As a cohort of predominantly healthy adults, intolerance to variation in UK Biobank constraint is driven by severe early onset disease, while the same measures in ExAC/gnomAD, whose samples have a more diverse age range and relatively higher of burden of disease, highlight genes involved with fundamental biological processes whose loss of function likely confer phenotypic consequences causing embryonic lethality.

403

404



409

Figure S6: *LDLRAD3* burden test PheWas. Significant ($p < 10^{-3}$) associations between regularized burden tests for *LDLRAD3* CNV and phenotypes. We highlight quantitative traits with $n > 15,000$ observations and binary traits with $n > 100$ cases. Traits are grouped by data type then sorted by p -value (left). Log-odds ratio and standardized betas (right; for binary and quantitative traits, respectively) align with trait names on the y-axis, with the horizontal dashed line separating positive and negative direction of association.

410

411

412

413

414 **Citations:**

- 415 1. Mikhail, F. M. Copy number variations and human genetic disease. *Curr. Opin. Pediatr.* **26**,
416 646–652 (2014).
- 417 2. Carvill, G. L. & Mefford, H. C. Microdeletion syndromes. *Curr. Opin. Genet. Dev.* **23**, 232–
418 239 (2013).
- 419 3. Bachmann-Gagescu, R. *et al.* Recurrent 200-kb deletions of 16p11.2 that include the
420 SH2B1 gene are associated with developmental delay and obesity. *Genet. Med.* **12**, 641–
421 647 (2010).
- 422 4. Zufferey, F. *et al.* A 600 kb deletion syndrome at 16p11.2 leads to energy imbalance and
423 neuropsychiatric disorders. *J. Med. Genet.* **49**, 660–668 (2012).
- 424 5. McDonald-McGinn, D. M., Emanuel, B. S. & Zackai, E. H. 22q11.2 Deletion Syndrome. in
425 *GeneReviews* (eds. Adam, M. P. *et al.*) (University of Washington, Seattle, 1999).
- 426 6. Raychaudhuri, S. *et al.* Accurately assessing the risk of schizophrenia conferred by rare
427 copy-number variation affecting genes with brain function. *PLoS Genet.* **6**, e1001097
428 (2010).
- 429 7. Priest, J. R. *et al.* De Novo and Rare Variants at Multiple Loci Support the Oligogenic
430 Origins of Atrioventricular Septal Heart Defects. *PLoS Genet.* **12**, e1005963 (2016).
- 431 8. Ruderfer, D. M. *et al.* Patterns of genic intolerance of rare copy number variation in 59,898
432 human exomes. *Nat. Genet.* **48**, 1107–1111 (2016).
- 433 9. Kirov, G. *et al.* The Uk Biobank: A Resource For Cnv Analysis. *Eur.*
434 *Neuropsychopharmacol.* **27**, S491 (2017).
- 435 10. Crawford, K. *et al.* Medical consequences of pathogenic CNVs in adults: analysis of the UK
436 Biobank. *J. Med. Genet.* (2018). doi:10.1136/jmedgenet-2018-105477
- 437 11. Owen, D. *et al.* Effects of pathogenic CNVs on physical traits in participants of the UK
438 Biobank. *BMC Genomics* **19**, 867 (2018).
- 439 12. Kendall, K. M. *et al.* Cognitive Performance Among Carriers of Pathogenic Copy Number

- 440 Variants: Analysis of 152,000 UK Biobank Subjects. *Biol. Psychiatry* **82**, 103–110 (2017).
- 441 13. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
442 *Nature* **562**, 203–209 (2018).
- 443 14. McInnes, G. *et al.* Global Biobank Engine: enabling genotype-phenotype browsing for
444 biobank summary statistics. *Bioinformatics* (2018). doi:10.1093/bioinformatics/bty999
- 445 15. Wang, K. *et al.* PennCNV: An integrated hidden Markov model designed for high-resolution
446 copy number variation detection in whole-genome SNP genotyping data. *Genome Res.* **17**,
447 1665–1674 (2007).
- 448 16. Jordan, V. K., Zaveri, H. P. & Scott, D. A. 1p36 deletion syndrome: an update. *Appl. Clin.*
449 *Genet.* **8**, 189–200 (2015).
- 450 17. Akbaroghli, S., Tonekaboni, S. H., Kariminejad, R., Liehr, T. & Coci, E. G. De-novo
451 interstitial 2.33 Mb deletion in 8q24.3: new insights on a very rare partial monosomy
452 syndrome. *Clin. Dysmorphol.* **27**, 97–100 (2018).
- 453 18. Iwakoshi, M. *et al.* 9q34.3 deletion syndrome in three unrelated children. *Am. J. Med.*
454 *Genet. A* **126A**, 278–283 (2004).
- 455 19. Cario, H., Bode, H., Gustavsson, P., Dahl, N. & Kohne, E. A microdeletion syndrome due to
456 a 3-Mb deletion on 19q13.2--Diamond-Blackfan anemia associated with macrocephaly,
457 hypotonia, and psychomotor retardation. *Clin. Genet.* **55**, 487–492 (1999).
- 458 20. Fry, A. *et al.* Comparison of Sociodemographic and Health-Related Characteristics of UK
459 Biobank Participants With Those of the General Population. *Am. J. Epidemiol.* **186**, 1026–
460 1034 (2017).
- 461 21. Hall, J. M. *et al.* Closing in on a breast cancer gene on chromosome 17q. *Am. J. Hum.*
462 *Genet.* **50**, 1235–1242 (1992).
- 463 22. Wooster, R. *et al.* Identification of the breast cancer susceptibility gene BRCA2. *Nature*
464 **378**, 789–792 (1995).
- 465 23. Papadopoulos, N. *et al.* Mutation of a mutL homolog in hereditary colon cancer. *Science*

- 466 **263**, 1625–1629 (1994).
- 467 24. Papadopoulos, N. *et al.* Mutations of GTBP in genetically unstable cells. *Science* **268**,
- 468 1915–1917 (1995).
- 469 25. Fishel, R. *et al.* The human mutator gene homolog MSH2 and its association with
- 470 hereditary nonpolyposis colon cancer. *Cell* **75**, 1027–1038 (1993).
- 471 26. Savitsky, K. *et al.* A single ataxia telangiectasia gene with a product similar to PI-3 kinase.
- 472 *Science* **268**, 1749–1753 (1995).
- 473 27. Horii, A. *et al.* The APC gene, responsible for familial adenomatous polyposis, is mutated in
- 474 human gastric cancer. *Cancer Res.* **52**, 3231–3233 (1992).
- 475 28. Gene Ontology Consortium & Gene Ontology Consortium. The Gene Ontology (GO)
- 476 database and informatics resource. *Nucleic Acids Res.* **32**, 258D–261 (2004).
- 477 29. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
- 478 datasets. *Gigascience* **4**, 7 (2015).
- 479 30. Povey, S. *et al.* The HUGO Gene Nomenclature Committee (HGNC). *Hum. Genet.* **109**,
- 480 678–680 (2001).
- 481 31. Astle, W. J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation and Links to
- 482 Common Complex Disease. *Cell* **167**, 1415–1429.e19 (2016).
- 483 32. Pott, J. *et al.* Genome-wide meta-analysis identifies novel loci of plaque burden in carotid
- 484 artery. *Atherosclerosis* **259**, 32–40 (2017).
- 485 33. van der Harst, P. & Verweij, N. Identification of 64 Novel Genetic Loci Provides an
- 486 Expanded View on the Genetic Architecture of Coronary Artery Disease. *Circ. Res.* **122**,
- 487 433–443 (2018).
- 488 34. Ranganathan, S. *et al.* LRAD3, a novel low-density lipoprotein receptor family member that
- 489 modulates amyloid precursor protein trafficking. *J. Neurosci.* **31**, 10836–10846 (2011).
- 490 35. Noyes, N. C., Hampton, B., Migliorini, M. & Strickland, D. K. Regulation of Itch and Nedd4
- 491 E3 Ligase Activity and Degradation by LRAD3. *Biochemistry* **55**, 1204–1213 (2016).

- 492 36. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the
493 human phenome. *Elife* **7**, (2018).
- 494 37. Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association meta-
495 analysis of coronary artery disease. *Nat. Genet.* **47**, 1121–1130 (2015).
- 496 38. Pruim, R. J. *et al.* LocusZoom: regional visualization of genome-wide association scan
497 results. *Bioinformatics* **26**, 2336–2337 (2010).
- 498 39. International HapMap Consortium. The International HapMap Project. *Nature* **426**, 789–796
499 (2003).
- 500 40. Qiu, Y. *et al.* Oligogenic Effects of 16p11.2 Copy Number Variation on Craniofacial
501 Development. (2019).
- 502 41. Wojciechowski, P. *et al.* Impact of FTO genotypes on BMI and weight in polycystic ovary
503 syndrome: a systematic review and meta-analysis. *Diabetologia* **55**, 2636–2645 (2012).
- 504 42. Voll, S. L. *et al.* Obesity in adults with 22q11.2 deletion syndrome. *Genet. Med.* **19**, 204–
505 208 (2017).
- 506 43. Medina-Gomez, C. *et al.* Life-Course Genome-wide Association Study Meta-analysis of
507 Total Body BMD and Assessment of Age-Specific Effects. *Am. J. Hum. Genet.* **102**, 88–102
508 (2018).
- 509 44. Olsen, L. *et al.* Prevalence of rearrangements in the 22q11.2 region and population-based
510 risk of neuropsychiatric and developmental disorders in a Danish population: a case-cohort
511 study. *Lancet Psychiatry* **5**, 573–580 (2018).
- 512 45. Klaassen, P. *et al.* Explaining the variable penetrance of CNVs: Parental intelligence
513 modulates expression of intellectual impairment caused by the 22q11.2 deletion. *Am. J.*
514 *Med. Genet. B Neuropsychiatr. Genet.* **171**, 790–796 (2016).
- 515 46. Castel, S. E. *et al.* Modified penetrance of coding variants by cis-regulatory variation
516 contributes to disease risk. *Nat. Genet.* **50**, 1327–1334 (2018).
- 517 47. Wang, N. K. & Chiang, J. P. W. Increasing evidence of combinatory variant effects calls for

- 518 revised classification of low-penetrance alleles. *Genet. Med.* (2018). doi:10.1038/s41436-
519 018-0347-3
- 520 48. DeBoever, C. *et al.* Medical relevance of protein-truncating variants across 337,205
521 individuals in the UK Biobank study. *Nat. Commun.* **9**, 1612 (2018).
- 522 49. Howrigan, D. P. *et al.* Genome-wide autozygosity is associated with lower general cognitive
523 ability. *Mol. Psychiatry* **21**, 837–843 (2016).
- 524 50. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based
525 linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 526 51. Diskin, S. J. *et al.* Adjustment of genomic waves in signal intensities from whole-genome
527 SNP genotyping platforms. *Nucleic Acids Res.* **36**, e126 (2008).
- 528 52. Bailey, J. A., Yavor, A. M., Massa, H. F., Trask, B. J. & Eichler, E. E. Segmental
529 duplications: organization and impact within the current human genome project assembly.
530 *Genome Res.* **11**, 1005–1017 (2001).
- 531 53. Haeussler, M. *et al.* The UCSC Genome Browser database: 2019 update. *Nucleic Acids*
532 *Res.* **47**, D853–D858 (2019).