# A global overview of pleiotropy and genetic architecture in complex traits

**Authors:** Kyoko Watanabe[1], Sven Stringer[1], Oleksandr Frei[2], Maša Umićević Mirkov[1], Tinca J.C. Polderman[1], Sophie van der Sluis[1,3], Ole A. Andreassen[2,4], Benjamin M. Neale[5-7], Danielle Posthuma[1,3]*

**Affiliations:**
1. Department of Complex Trait Genetics, Center for Neurogenomics and Cognitive Research, Neuroscience Campus Amsterdam, VU University Amsterdam, The Netherlands.
2. NORMENT, KG Jebsen Centre for Psychosis Research, Institute of Criminal Medicine, University of Oslo, Oslo, Norway
3. Department of Clinical Genetics, Section of Complex Trait Genetics, Neuroscience Campus Amsterdam, VU Medical Center, Amsterdam, the Netherlands.
4. Division of Mental health and addiction Oslo University hospital, Oslo, Norway
5. Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA
6. Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, MA, USA
7. Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA

*Correspondence to: Danielle Posthuma, Department of Complex Trait Genetics, VU University, De Boelelaan 1085, 1081 HV, Amsterdam, The Netherlands. Phone: +31 20 5982823, Fax: +31 20 5986926, Email: d.posthuma@vu.nl

**Word count**: Abstract 181 words, Main text 5,762 words and Methods 4,572 words
**References**: 40
**Display items**: 4 figures and 2 tables
**Extended Data**: 11 figures
**Supplementary Information**: Text 3,401 words and 25 tables

1

## 2 ABSTRACT

3 After a decade of genome-wide association studies (GWASs), fundamental questions in

4 human genetics are still unanswered, such as the extent of pleiotropy across the genome, the

5 nature of trait-associated genetic variants and the disparate genetic architecture across human

6 traits. The current availability of hundreds of GWAS results provide the unique opportunity

7 to gain insight into these questions. In this study, we harmonized and systematically analysed

8 4,155 publicly available GWASs. For a subset of well-powered GWAS on 558 unique traits,

9 we provide an extensive overview of pleiotropy and genetic architecture. We show that trait

10 associated loci cover more than half of the genome, and 90% of those loci are associated with

11 multiple trait domains. We further show that potential causal genetic variants are enriched in

12 coding and flanking regions, as well as in regulatory elements, and how trait-polygenicity is

13 related to an estimate of the required sample size to detect 90% of causal genetic variants.

14 Our results provide novel insights into how genetic variation contributes to trait variation. All

15 GWAS results can be queried and visualized at the GWAS ATLAS resource

16 (http://atlas.ctglab.nl).

17    **MAIN TEXT**

18    Since the first genome-wide association study (GWAS) on macular degeneration in 2005[1],

19    over 3,000 GWASs have been published, for more than 1,000 traits, reporting on over tens of

20    thousands of significantly associated genetic variants[2]. Results from GWASs have increased

21    our insight into the genetic architectures of investigated traits, and for some traits, GWAS

22    results have led to further insight into disease mechanisms[3,4], such as autophagy for Crohn's

23    disease[5], immunodeficiency for Rheumatoid arthritis[6] and transcriptome regulation through

24    *FOXA2* in the pancreatic islet and liver for Type 2 diabetes[7]. The emerging picture after over

25    a decade of GWASs is that the majority of studied traits are highly polygenic and thus

26    influenced by many genetic variants each of small effect[4,8], with disparate genetic

27    architectures across traits[9]. Fundamental questions, such as whether all genetic variants or all

28    genes in the human genome are associated with at least one trait, with many or even all traits,

29    and whether the polygenic effects for specific traits are functionally clustered or whether they

30    are randomly spread across the genome, are however still unanswered[4,10,11]. Answers to these

31    questions would greatly enhance our understanding of how genetic variation leads to trait

32    variation and trait correlation. Whereas GWAS primarily aims to discover genetic variants

33    associated with specific traits, the current availability of a vast amount of GWAS results can

34    be used to investigate some of these fundamental questions.

35    To this end, we compiled a catalogue of 4,155 GWAS results across 2,965 unique traits from

36    295 studies, including publicly available GWASs and new results for 600 traits from the UK

37    Biobank (http://atlas.ctglab.nl). These GWAS results were used in the current study to

38    achieve the following aims; *i*) charting the extent of pleiotropy at trait-associated locus, gene,

39    SNP and gene-set levels, *ii*) characterizing the nature of trait-associated variants (i.e. the

40    distribution of effect size, minor allele frequency and biological functionality of trait-

41    associated or credible SNPs), and *iii*) understanding the nature of the genetic architecture

3

42     across a variety of traits and domains in terms of SNP heritability and trait polygenicity (see

43     **Extended Data Fig. 1)**.

44

45     **Catalogue of 4,155 GWAS summary statistics for 2,965 unique traits**

46     We collected publicly available full GWAS summary statistics (last update 23rd October

47     2018; see **Methods**). This resulted in 3,555 GWAS summary statistics from 294 studies. We

48     additionally performed GWAS on 600 traits available from the UK Biobank release 2 cohort

49     (UKB2; release May 2017)[12], by selecting non-binary traits with >50,000 European

50     individuals with non-missing phenotypes, and binary traits for which the number of available

51     cases and controls were each >10,000 and total sample size was >50,000 (see **Methods**,

52     **Supplementary Information 1** and **Supplementary Table 1-2**). In total, we collected 4,155

53     GWASs from 295 unique studies and 2,965 unique traits (see **Supplementary Table 3** for a

54     full list of collected GWASs). Traits were manually classified into 27 standard domains

55     based on previous studies[13,14]. The average sample size across curated GWASs was 56,250

56     subjects. The maximum sample size was 898,130 subjects for a Type 2 Diabetes meta-

57     analysis[15]. The 4,155 GWAS results are made available in an online database

58     (http://atlas.ctglab.nl). The database provides a variety of information per trait, including

59     SNP-based and gene-based Manhattan plots, gene-set analyses[16], SNP heritability

60     estimates[17], genetic correlations, cross GWAS comparisons and phenome-wide plots.

61     For the present study, we restricted our analyses to reasonably powered GWASs (i.e. sample

62     size >50,000), to avoid including SNP effect estimates with relatively large standard errors

63     (see **Methods**). By selecting a GWAS with the largest sample size per trait, it resulted in 558

64     GWASs for 558 unique traits across 24 trait domains. The average sample size of these 558

65     GWASs was 256,276, and 478 GWASs (85.7%) were based on the UKB2 including 11 meta-

66     analyses with UKB2, 46 (8.2%) on the UK Biobank release 1 cohort (UKB1) including 8

4

67    meta-analyses with UKB1, and the remaining were non-UKB cohorts. All results presented

68    hereafter concern these selected 558 GWASs unless specified otherwise. The online database,

69    however, allows researchers to reproduce similar analyses with custom selections of GWASs.

70

71    **The extent of pleiotropy**

72    Results of previous GWASs have shown significant associations of thousands of genomic

73    loci with a large number of traits[2,4]. Given a finite number of segregating variants on the

74    human genome, this suggests the presence of widespread pleiotropy. Pleiotropy may be

75    informative to the reasons of co-morbidity between traits, as it may indicate an underlying

76    shared genetic mechanism, and may aid in resolving questions regarding causal effects of one

77    trait on another. However, the exact extent of pleiotropy across the genome is currently

78    unknown[4]. We therefore investigated pleiotropy at locus, gene, SNP and gene-set levels. We

79    defined pleiotropy as the presence of statistically significant associations with more than one

80    trait domain as traits within domain tend to show stronger phenotypic correlations than

81    between domains (see **Supplementary Information 2** and **Extended Data Fig. 2**). Our

82    definition thus refers to 'statistical pleiotropy', and includes situations of true pleiotropy (e.g.

83    one SNP directly influences multiple traits), or situations where statistical associations to

84    multiple traits are induced via causational effects of one trait on another, via phenotypic

85    correlations between traits, or via a third common factor[18]. We defined the level of pleiotropy

86    by the number of associated domains, and further grouped into four categories; multi-domain

87    (associated with traits from multiple domains), domain-specific (associated with multiple

88    traits from a single domain), trait-specific (associated with a single trait) and non-associated

89    (**Methods**). We then assessed whether pleiotropic associations at the locus, gene, SNP or

90    gene set level are structurally or functionally different from trait- or domain-specific

91    associations or non-associated sites.

5

92

*Pleiotropic genomic loci*

93

94    The 558 GWASs yielded 41,511 trait-associated loci (from 470 traits, as 88 traits did not

95    yield any genome-wide significant association after QC; see **Methods**). After grouping

96    physically overlapping trait-associated loci, we obtained 3,362 grouped loci (**Methods**,

97    **Extended Data Fig. 3,** and **Supplementary Table 4**). The total summed length of these loci

98    (1706.0 Mb) covered 61.0% of the genome. Of these, 93.3% were associated with more than

99    one trait and 90.0% were multi-domain loci (**Table 1** and **Extended Data Fig. 4a, b**). The

100    multi-domain and domain-specific loci showed a significantly higher density of protein

101    coding genes compared to non-associated genomic regions ($p$=5.3e-16 and $p$=2.6e-4; **Fig. 1a**

102    and **Supplementary Table 5**).

103    The locus associated with the largest number of traits and domains (i.e. the most pleiotropic

104    locus) was the MHC region (chr 6:25Mb-37Mb), which contained 441 trait-associated loci

105    from 213 traits across 23 trait domains. The MHC region is well-known for its complex

106    structure of linkage disequilibrium, spanning over 300 genes. The extremely pleiotropic

107    nature of this region might, therefore, be explained by its long-ranged LD block due to

108    overlap of multiple independent signals from multiple traits. Similarly, high locus pleiotropy,

109    not limited to the MHC region, can occur purely due to the overlap of the LD blocks of the

110    loci in the grouped locus, and they may not share the same causal SNPs. By performing

111    colocalization (i.e. statistically identifying loci sharing the same causal SNP) for all possible

112    pairs of physically overlapping trait-associated loci (see **Methods, Supplementary**

113    **Information 3** and **Extended Data Fig. 3**), we indeed observed a decrease in the number of

114    associated traits and trait domains per group of colocalized loci compared to loci defined by

115    physical overlap (**Extended Data Fig. 4** and **Supplementary Table 6**). In addition, loci

116    grouped based on physical overlap often contained multiple independent groups of

6

117     colonized loci (**Supplementary Table 6**). Therefore, physical overlap of trait-associated

118     loci does not necessary mean that the same causal SNPs are involved in the traits associated

119     with such a grouped locus. Examination of pleiotropy at the gene or SNP level will provide

120     further insight into the nature of the pleiotropy observed at the locus level.

121

122     *Pleiotropic genes*

123     We next investigated the extent of pleiotropy at the gene level. For this, we conducted a

124     gene-based analysis on 17,444 protein-coding genes using MAGMA for each trait[16]

125     (**Methods**). Of the 558 traits, 516 yielded at least one significantly associated gene and

126     11,443 (65.6%) genes were significantly associated to at least one trait (**Supplementary**

127     **Table 7**). Of these, 81.0% were associated with more than one trait and 66.9% were

128     associated with traits from multiple domains (**Table 1** and **Extended Data Fig. 5a, b**). We

129     found that genes associated with at least one trait are significantly longer than genes that are

130     not associated with any of the 558 tested traits ($p$=2.1e-194, $p$=8.7e-12 and $p$=3.8e-29 for

131     multi-domain, domain-specific and trait-specific genes, respectively; **Fig. 1b** and

132     **Supplementary Table 8**). As the MAGMA algorithm is insensitive to bias caused by gene-

133     length, these findings are unlikely to be due to larger genes having an increased statistical

134     probability to be significantly associated (**Supplementary Information 4, Extended Data**

135     **Fig. 5c** and **Supplementary Table 9**). The multi-domain genes showed a significantly higher

136     probability of being intolerant to loss of function mutations (pLI score)[19] compared to trait-

137     /domain-specific and non-associated genes ($p$=1.2e-79, $p$=4.8e-22 and $p$=2.8e-19,

138     respectively; **Fig. 1c** and **Supplementary Table 10**), suggesting that more pleiotropic genes

139     are on average less tolerant to loss of function variants. The most pleiotropic genes are

140     located in the MHC region, yet a region on chromosome 3 also spanned multiple genes with

141  high levels of pleiotropy (**Extended Data Fig. 5a**). In this region, *BSN* was associated with

142  the largest number of trait domains (94 traits across 17 domains).

143  We next tested whether tissue specificity of genes was related to the level of pleiotropy by

144  counting the number of active tissues per gene based on gene expression profiles for 53 tissue

145  types obtained from GTEx[20] (see **Methods**). The results showed that the proportion of genes

146  expressed in all 53 tissue types increases along with the level of pleiotropy (*p*=9.7e-05, **Fig.**

147  **1d** and **Supplementary Table 11**). This indicates that more pleiotropic genes tend to be

148  active in multiple tissue types, suggesting that those genes are involved in general biological

149  functions across the human body.

150

151  *Pleiotropic SNPs*

152  The level of pleiotropy at a locus or gene level does not necessarily translate to pleiotropy at

153  the level of the SNP. For example, within the same locus or gene, multiple SNPs may be

154  significantly associated with different traits. A locus or gene can thus show a higher level of

155  pleiotropy compared to individual SNPs. We, therefore, investigated the extent of pleiotropy

156  at the level of the SNP. To do so, we extracted 1,740,179 SNPs that were present in all 558

157  GWAS results. We first confirmed that this selection of SNPs had the same distribution of

158  their location across the genome and their functional consequences as all known SNPs on the

159  genome (**Methods** and **Extended Data Fig. 6a, b**). We note that some of the observed SNP-

160  pleiotropy may still be induced by LD, e.g. a SNP could reach genome-wide significance

161  because of its strong LD with a causal SNP. However, the purpose of this analysis is to

162  identify individual SNPs (not loci) that are associated with multiple trait domains and their

163  functions. Of these, 237,120 (13.6%) were genome-wide significant (*p*<5e-8) in at least one

164  of the 558 traits (**Extended Data Fig. 6c** and **Supplementary Table 12**). Out of 237,120

165  SNPs that were associated with at least one trait, 60.2% were associated with more than one

8

166 trait and 32.4% were associated with more than one domain (**Table 1** and **Extended Data**

167 **Fig. 6d**).

168 These pleiotropic SNPs spread broadly across the genome but were not evenly distributed,

169 i.e. chromosome 1, 11, 12, 15, 17, 20 and 22 showed relative enrichment of pleiotropic SNPs

170 (**Supplementary Information 5** and **Supplementary Table 13**). Of all associated SNPs, the

171 most pleiotropic SNP, located in the MHC region (rs707939; an intronic SNP of *MSH5*) was

172 associated with 48 traits from 13 domains. There were 45 SNPs associated with 12 trait

173 domains, of which 35 were located on chromosome 3, 49.8Mb-50.1Mb overlapping with 5

174 protein coding genes, *TRAIP*, *CAMKV*, *MST1R*, *MON1A* and *RBM6*. These SNPs include two

175 exonic SNPs, rs2681781 (synonymous on *CAMKV*) and rs2230590 (nonsynonymous on

176 *MST1R*; **Supplementary Table 12**).

177 To investigate whether SNPs with a higher level of pleiotropy have different functional

178 annotations than less pleiotropic SNPs, we investigated how functional consequence and

179 tissue specificity in terms of expression quantitative trait loci (eQTLs) were represented

180 across different levels of SNP pleiotropy (**Methods**). We found that the proportion of intronic

181 and exonic SNPs increased as a function of the level of pleiotropy ($p$=2.2e-3 and $p$=1.7e-2,

182 respectively); the proportion of exonic SNPs increased from less than 1% to over 5%, and the

183 proportion of intronic SNPs increased from less than 40% to over 50% (**Fig. 1e** and

184 **Supplementary Table 14**) with increasing levels of pleiotropy. The proportion of SNPs

185 within flanking regions such as 5' and 3' untranslated regions (UTR) also increased with the

186 number of associated domains. At the same time, we observed a steep decrease of the

187 proportion of intergenic SNPs with increasing level of SNP pleiotropy ($p$=8.1e-4; **Fig. 1e** and

188 **Supplementary Table 14**). Based on active eQTLs, the proportion of SNPs being eQTLs in

189 a greater number of tissue types (>24 tissue types out of 48) increased along with the number

190 of associated domains ($p$=8.4e-3 and $p$=1.1e-2 for eQTLs in between 25 and 36 tissues, and

9

191     between 37 and 48 tissues, respectively) while SNPs in genes expressed in a single or less

192     than half of available tissue types showed decreasing proportion (**Fig. 1f** and **Supplementary**

193     **Table 15**). These results suggest that highly pleiotropic SNPs are more likely to be genic

194     (exonic and intronic) and less likely to be tissue specific.

195

196     *Pleiotropic gene-sets*

197     Pleiotropy at the level of trait-associated loci, genes or SNPs do not necessarily suggest the

198     presence of shared biological pathways across multiple traits. To assess the level of

199     pleiotropy at the level of gene-sets, reflecting a biological meaningful grouping of genes, we

200     performed MAGMA gene-set analyses for 558 traits using 10,650 gene-sets (Methods). In

201     total, 235 (42.1%) traits showed significant association with one of 1,106 (10.4%) gene-sets.

202     The most pleiotropic gene-set was 'Regulation of transcription from RNA polymerase II

203     promoter' (GO biological process) associated with 61 traits from 9 domains, followed by 7

204     other gene-sets associated with 7 domains, of which 5 of them were also involved in

205     regulation of transcription (**Supplementary Table 16**). We observed that the number of

206     genes in a gene-set was significantly larger for highly pleiotropic gene-sets (associated with

207     more than one domain) compared to other gene-sets (domain-specific, trait-specific and non-

208     associated; $p$=4.1e-12, $p$=1.6e-13 and $p$=1.2e-29, respectively; **Extended Data Fig. 7a,** and

209     **Supplementary Table 17**). Since GO terms (55.6% of tested gene-sets) have a hierarchical

210     structure, the larger gene-sets are more likely to be located at the top of the hierarchy,

211     representing more general functional categories.

212     In contrast to the pleiotropy at gene level where 80.9% genes were associated with more than

213     one trait, we only found 54.8% of the associated gene-sets to be pleiotropic (**Table 1**). We

214     observed that the proportion of pleiotropic genes per gene-set is not uniformly distributed,

215     and pleiotropic genes tend to cluster into a subset of gene-sets, explaining the decreased

10

216 proportion of pleiotropic gene-sets compared to pleiotropic genes (**Extended Data Fig. 7b,**

217 **c**). At the same time, the higher proportion of trait-specific gene-sets (45.2%) compared to

218 trait-specific genes (19.2%) suggests that, given current definitions of gene-sets, the

219 combination of associated genes is rather unique to a trait and focusing on gene-sets to gain

220 insight into trait-specific biological mechanisms may be more informative than focusing on

221 single genes (**Supplementary Information 6**).

222

223 *Genetic correlations across traits*

224 Above we showed that of all trait-associated loci, genes and SNPs that are associated with at

225 least one trait, 90.0%, 66.9% and 32.6% are associated with more than one domain,

226 respectively. Such wide-spread pleiotropy indices non-zero genetic correlations between

227 traits. To test whether genetic correlations are evenly present across traits or cluster into trait

228 domains, we computed pairwise genetic correlations ($r_g$) across 558 traits using LDSC[17].

229 We calculated the proportion of trait pairs with an $r_g$ that is significantly different from zero

230 across all 558 traits, within domains and between domains. Out of 155,403 possible pairs

231 across 558 traits, 24,106 pairs (15.5%) showed significant genetic correlations after

232 Bonferroni correction ($p<0.05/155,403=3.2e-7$) with an average $|r_g|$ of 0.38.

233 In principle, if the trait domains contain traits that are biologically related, we would expect

234 that traits within the same domain have stronger genetic correlations than traits across

235 domains. The proportion of pairs with a significant genetic correlation within a domain was

236 especially high in cognitive, 'ear, nose, throat', metabolic and respiratory domains, and for

237 most of domains, average $|r_g|$ across significant trait pairs was higher than 0.38 (across all

238 traits). Note that the proportion of trait pairs with significant $r_g$ may be biased by sample size

239 and $h^2_{SNP}$ of traits within a domain; across 558 traits, the worst case scenarios with the

240 minimum observed $h^2_{SNP}$ (0.0045 with sample size 385,289) or the minimum sample size

11

241   (51,750 with $h^2_{SNP}$ =0.0704) required $r_g$ to be above 0.39 or 0.18, respectively, to gain a

242   power of 0.8 (**Methods**). Within domain, the majority of significant genetic correlations was

243   positive and the average $|r_g|$ was above 0.5 in most of the domains (**Fig. 2a** and

244   **Supplementary Table 18**). Between domains, the proportion of pairs with significant genetic

245   correlations was generally lower than within domains, and most of the domain pairs showed

246   average $|r_g|$<0.4 (**Fig. 2b** and **Supplementary Table 19**). Some trait domains showed a

247   predominance of negative genetic correlations with other domains, i.e. activity, cognitive,

248   reproduction and social interaction domains. We further clustered traits based on genetic

249   correlations, which resulted in the majority of clusters contained traits from multiple domains

250   (**Methods**, **Supplementary Information 7** and **Extended Data Fig. 8**). These results

251   suggest that although $|r_g|$ is higher within domain than across domains, the trait domains do

252   not necessary reflect genetic similarity across traits.

253

254   **The nature of trait-associated variants**

255   We now address the question whether trait-associated variants differ from genetic variants

256   that are not associated with any trait. For this purpose, we extracted all lead SNPs from each

257   of the 558 GWASs. Lead SNPs were defined per trait at the standard threshold for genome-

258   wide significance ($p$<5e-8) and using an $r^2$ of 0.1 to obtain near-independent lead SNPs,

259   based on the population-relevant reference panel (see **Methods**). Lead SNPs with minor

260   allele count (MAC) ≤100 (based on MAF and sample size of the SNP) were excluded due to

261   lower statistical power and a high false positive rate of effects of SNPs with extremely small

262   MAF. This resulted in 82,590 lead SNPs for 476 traits, reflecting 43,455 unique SNPs. Out of

263   558 traits, 82 traits did not yield any genome-wide significant lead SNP after QC.

264

265   *Distribution of MAF and effect sizes of lead SNPs*

12

266     12.3% of the 43,455 (unique) lead SNPs derived from the 558 GWASs had a MAF below

267     0.01 which is significantly less than expected given the proportion of rare variants in the

268     reference panels ($p$<1e-323; **Supplementary Information 8**), while the distribution of lead

269     SNPs with a MAF above 0.01 was nearly uniform (**Fig. 3a**).

270     To gain insight into the distribution of effect sizes across lead SNPs, we calculated the

271     standardized effect size ($\beta$) from Z-statistics as a function of MAF and sample size[21], and

272     inspected the distribution of the squared standardized effect sizes ($\beta^2$) for lead SNPs across

273     all traits (**Methods**). $\beta$ ranged between 0.01 and 1.70, and $\beta^2$ is proportional to the variance

274     explained. The median $\beta^2$ of the lead SNPs across all traits was 5.7e-4 (4.9e-4 and 6.0e-2 for

275     lead SNPs with MAF≥0.01 and <0.01, respectively), and 94.6% of lead SNPs had a $\beta^2$ below

276     0.05 (**Fig. 3b**). Thus, the vast majority of lead SNPs thus explained less than 0.05% of the

277     trait variance. We observed a relationship between MAF and standardized effect size, with

278     rare variants (MAF<0.01) showing larger effect sizes (**Fig. 3c**). This is in line with the notion

279     that rare variants are more likely to have large effects compared to common variants, as they

280     are less likely to be under strong selective pressure[22]. However, we also note that statistical

281     power for detecting the rare variants is un-stable[23]. Given that the proportion of rare lead

282     SNPs is larger than the proportions in other MAF bins, it is possible that the distribution of

283     the effect sizes has longer tails for SNPs with MAF<0.01. For most of the traits, a similar

284     relationship between MAF and standardized effect size was observed (**Extended Data Fig.**

285     **9**), but large variation across traits was seen in terms of the number of rare lead SNPs, with

286     e.g. a large proportion of rare variants influencing nutritional and connective tissue domains

287     (see **Supplementary Information 8, Extended Data Fig. 10** and **Supplementary Table 20-**

288     **21**).

289

290     *Characterization of trait-associated loci and lead SNPs*

291     Here we sought to characterise differences in the distribution of functional annotations when

292     comparing SNPs within trait-associated loci to all SNPs in the genome, and comparing lead

293     SNPs to SNPs in the trait-association loci (**Methods**). We first compared SNPs in the trait-

294     associated loci against the entire genome. The strongest enrichment of SNPs in trait

295     associated loci was seen in flanking regions (upstream, downstream, 5' and 3' UTR) with

296     average fold enrichment ($E$) 1.31 (**Fig. 3d** and **Table 2**). Non-coding SNPs, in total, covered

297     93.1% of SNPs in the trait-associated loci, while intergenic SNPs were significantly depleted

298     ($E$=0.83) and intronic SNPs significantly enriched compared to all SNPs in the genome

299     ($E$=1.17; **Table 2**). SNPs in trait-associated loci were also slightly enriched for being exonic

300     compared to the entire genome ($E$=1.07). Active chromatin states and eQTLs were also

301     significantly enriched with notably high enrichment of eQTLs ($E$=1.61 and 5.95,

302     respectively; **Table 2**).

303     We next compared lead SNPs with SNPs in the trait-associated loci. The strongest

304     enrichment for lead SNPs was seen in exonic SNPs ($E$=2.84) followed by flanking regions

305     ($E$=1.38), while intronic and intergenic regions were slightly depleted (average $E$=0.95; **Fig.**

306     **3d** and **Table 2**). These results clearly indicate that SNPs located in exonic and flanking

307     regions tend to show stronger effect sizes than other SNPs in the trait-associated loci. On the

308     other hand, active chromatin states showed slight enrichment ($E$=1.08) while eQTLs were

309     significantly depleted ($E$=0.80; **Fig. 3e-f** and **Table 2**). This suggests that SNPs within the

310     trait-associated loci largely overlap with regulatory elements but these elements do not

311     always have the strongest effect sizes within the loci.

312

313     *Characterization of credible set SNPs based on fine-mapping*

314     Owing to the small effect sizes of variants in complex traits and extensive LD throughout the

315     human genome, there is a reasonable chance that lead SNPs (i.e. defined based on LD and P-

14

316   values) are not the causal SNPs in the trait-associated loci[24], even when the causal SNPs are

317   actually measured or imputed. Statistical fine-mapping utilizes evidence of the associations at

318   each variant in the loci (effect sizes and LD structure) to assign posterior probability of each

319   specific model at particular locus, which are then used to infer the posterior probabilities of

320   each SNP being included in the model (posterior inclusion probability, PIP) and ascertain the

321   minimum set of SNPs required to capture the likely causal variant. We performed fine-

322   mapping using FINEMAP software[25] for each trait-associated locus, setting the maximum

323   number of SNPs in the causal configuration ($k$) to 10 and using randomly selected 100k

324   individuals from UKB2 as a reference panel (see **Methods**). From all of the loci associated

325   with at least one of the 558 traits, we obtained a list of credible SNPs with PIP>0.95 consists

326   of 196,542 SNPs (**Supplementary Information 9**).

327   Next we characterized credible SNPs in respect to their functional annotations, similar as

328   done above with lead SNPs. We thus compared SNPs in the fine-mapped regions to all SNPs

329   in the genome, and credible SNPs to SNPs in the fine-mapped regions. The enrichment

330   pattern of SNPs in the fine-mapped regions was similar to SNPs in the trait-associated loci;

331   i.e. significant enrichment of SNPs in intronic and flanking regions but the fold enrichment

332   was much smaller (**Fig. 3d** and **Table 2**). This is mainly because the fine-mapped regions are

333   often larger than the trait-associated loci by taking 50kb around the top SNPs of the trait-

334   associated loci. In contrast, fold enrichment of exonic SNPs was slightly higher than trait-

335   associated loci (**Table 2**). As we observed higher gene-density around the trait-associated

336   loci, expanding the loci resulted in larger proportion of exonic regions. Both active chromatin

337   state and eQTLs were significantly enriched, however, fold enrichment of eQTLs was

338   notably less than trait-associated loci (**Fig. 3e-f** and **Table 2**). Similar to the lead SNPs,

339   credible SNPs showed strong enrichment in exonic ($E$=1.40) and flanking regions ($E$=1.29),

340   as well as intronic regions ($E$=1.17; **Table 2**). Although an enrichment of active chromatin

15

341 state is consistent with the result observed in the lead SNPs ($E$=1.51), eQTLs were also

342 significantly enriched in credible SNPs with very strong fold increase ($E$=4.14; **Fig. 3e-f** and

343 **Table 2**).

344 In summary, the number of credible SNPs is 4.5 times larger than the number of lead SNPs,

345 since for determining lead SNPs, all SNPs that have high LD with lead SNPs are discarded

346 while the fine-mapping captures likely causal SNPs given the observed pattern of association

347 and LD structure. Lead SNPs and credible SNPs show different distributions of enrichment in

348 tested biological functions. We observed a decreased proportion of exonic SNPs and an

349 increased proportion of non-coding or regulatory SNPs within the credible SNPs compared to

350 the lead SNPs. These findings may be due to the fact that coding SNPs tend to have higher

351 effect sizes and are more often assigned as lead SNPs, while the fine-mapping in regions

352 containing some of these causal coding variants may disperse a proportion of probability to

353 adjacent variants. On the other hand, in loci where causal variants are acting through

354 regulatory mechanisms, the credible sets may be more likely to capture the actual, single or

355 multiple causal variants as compared to the lead SNPs.

356

357 **The nature of genetic architecture**

358 The genetic architecture of a trait reflects the characteristics of genetic variants that

359 contribute to the phenotypic variability, and is defined by e.g. the number of variants

360 affecting the trait, the distribution of effect sizes, the MAF and the level of interactions

361 between SNPs [9]. To gain insight into how the genetic architecture varies across multiple

362 complex traits, we assessed the SNP heritability ($h^2_{SNP}$) and the polygenicity of 558 traits.

363

364 *SNP heritability*

16

365     $h^2_{SNP}$ is an indication of the total amount of variance that is captured by the additive effects of

366     all variants included in a GWAS. $h^2_{SNP}$ depends on several factors, such as the number of

367     SNPs included in the analyses based on their MAF given the current sample size, the

368     polygenicity of the trait (i.e. how many SNPs have an effect) and the distribution of effect

369     sizes. We estimated $h^2_{SNP}$ for each trait using LDSC[17] and SumHer from LDAK[26,27]

370     (**Methods**). The estimates of $h^2_{SNP}$ using LDSC and SumHer showed strong positive

371     correlation ($r$=0.77 and $p$=3.8e-111; **Fig. 4a**). Therefore, we focus on estimates based on

372     LDSC, hereafter, however complete results are available in **Supplementary Table 22** and

373     discussed in **Supplementary Information 10** (**Extended Data Fig. 11**). The highest $h^2_{SNP}$

374     was observed for height ($h^2_{SNP}$=0.31) followed by bone mineral density ($h^2_{SNP}$=0.27). Of 558

375     traits, 214 traits, with an average sample size 292,267, showed $h^2_{SNP}$ less than 0.05. Most of

376     these traits are classically regarded as 'environmental' (e.g. current employment status,

377     illness of family members and transport types or activity traits including frequency and type

378     of physical activities and type of accommodation), and tend to have a low $H^2$[14]. For these

379     traits, the number of detected trait-associated loci is also very low with a median 3. Given the

380     combination of current sample size of > 200,000 and low $h^2_{SNP}$, this suggests that for these

381     traits increasing the sample size may not lead to a substantial increase in detected loci.

382

383     *Polygenicity and discoverability of complex traits*

384     The general observation from GWASs is that with increasing sample size, detected signals

385     become not only more reliable but also more numerous, as with increasing power, smaller

386     SNP effects may be detected. The total number of associated SNPs, the amount of variance

387     they collectively represent, the distribution of effect sizes across the associated SNPs and

388     how many additional individuals are expected to be needed for the detection of a fixed

389    number of novel SNPs, are indications of the polygenicity of a trait. Such polygenicity may

390    vary across traits, and can be informative for designing SNP-discovery studies.

391    To obtain an indication of trait-polygenicity, we applied the Causal Mixture Model for

392    GWAS summary statistics (MiXeR)[28] to estimate $\pi$ (fraction of independent causal SNPs,

393    polygenicity) and $\sigma_\beta^2$ (variance of effect sizes of the causal SNPs, discoverability; see

394    **Methods**). $\pi$ ranges between 0 and 1, and a high $\pi$ indicates a high level of polygenicity,

395    while a high $\sigma_\beta^2$ indicates a high level of discoverability of causal SNPs for the traits. Since

396    the standard error of the model estimates become larger for traits with very small $h^2_{SNP}$ due to

397    the small effect sizes, we only discuss the results of 197 out of 558 traits with $h^2_{SNP}>0.05$ and

398    standard error of $\pi$ less than 50% of the estimated value (as recommended by O. Frei; full

399    results are available in **Supplementary Table 23**). We observed, as expected, a negative

400    relationship between polygenicity and discoverability ($r$=-0.89 and $p$=4.93e-70), confirming

401    that highly polygenic traits tend to have less causal SNPs with larger effect sizes (**Fig. 4b**).

402    The majority of traits (i.e. 116 traits) showed high polygenicity with $\pi$>1e-3 (more than 0.1%

403    of all SNPs are causal). The highest polygenicity was observed in Major depressive disorder

404    with 0.6% of SNPs being causal, while some traits, such as fasting glucose and serum urate

405    level showed relatively low polygenicity (**Fig. 4b** and **Supplementary Table 23**). The traits

406    with polygenicity >0.1% showed, on average, 8 times less discoverability compared to other

407    traits with <0.1% of causal SNPs. The GWAS discoveries for traits with lower polygenicity

408    and high discoverability will saturate with a lower sample size compared to the traits with

409    higher polygenicity. Indeed, the estimated sample size, which is required to explain 90% of

410    SNP heritability by genome-wide significant SNPs, is positively correlated with polygenicity

411    ($r$=0.84 and $p$=6.30e-54), and extremely polygenic traits require tens of millions of subjects

412    to identify 90% of causal SNPs at a genome-wide significant level (**Fig. 4c**).

413

**Discussion**

414 The availability of hundreds of GWAS results provides the unique opportunity to gain insight

416 into currently understudied questions regarding the genetic architecture of human traits. To

417 facilitate such insight, we compiled a catalogue of 4,155 GWASs which can be queried

418 online (http://atlas.ctglab.nl). We selected 558 well-powered GWASs to answer fundamental

419 questions concerning the extent of pleiotropy of loci, genes, SNPs and gene-sets,

420 characteristics of trait-associated variants and the polygenicity of traits.

421 We found that the total summed length of trait-associated loci for the 558 analysed traits

422 covered more than half (60.1%) of the genome. 90% of the grouped loci contained

423 associations with multiple traits across multiple trait domains. High locus pleiotropy can

424 occur in two scenarios; *i*) when the same gene in a locus is associated with multiple traits or

425 *ii*) when different genes or SNPs in the same locus are associated with multiple traits but due

426 to LD the same locus is indicated. Our results showed that the proportion of pleiotropic

427 associations dropped from 90% at the locus level to 63% at the gene level, and to 31% at the

428 SNP level. These results show that although locus pleiotropy is widespread, pleiotropy at the

429 level of genes and SNPs is much less abundant. This suggests that a gene can be involved in

430 two distinct traits but how that gene is affected by the causal SNPs might differ. For instance,

431 the function of the gene can be disrupted through a coding SNP for one trait, but expression

432 of the same gene can be affected through a regulatory SNP for another trait.

433 Genes and SNPs that had a higher level of pleiotropy, were less tissue specific in terms of

434 gene expression and active eQTLs. This suggests that SNPs and genes associated with

435 multiple trait domains are more likely to be involved in general biological functions. Indeed,

436 the top highly pleiotropic gene-sets were mostly involved in regulation of transcription which

437 is an essential biological mechanism for any kind of cell to be functioning. Highly pleiotropic

438 genes, therefore, can explain general vulnerability to a wide variety of traits, yet they may be

19

439   less informative when the aim is to understand the causes of a specific trait. Although a large

440   proportion of trait-associated genes are pleiotropic, the majority of trait associated gene-sets

441   were trait-specific. Thus, the trait-specific combination of genes is highly informative, and

442   future studies aimed at improved annotation of gene-functions will be needed to understand

443   trait-specific gene association patterns.

444   It has been widely acknowledged that almost 90% of GWAS findings fall into non-coding

445   regions[2]. Our results indeed show that 89.1% of the lead SNPs are non-coding, including

446   intergenic (34.3%) and intronic (43.6%) SNPs. similarly, of the credible SNPs 92.4% were

447   non-coding (intergenic 33.4% and intronic 48.1%). However, we showed different patterns

448   when considering lead and credible SNPs; intergenic SNPs were depleted and the intronic

449   SNPs were enriched in both the lead and credible SNPs. We also observed strong enrichment

450   of the lead and credible SNPs in coding and flanking regions. These results indicate that both

451   SNPs with the largest effect size (the lead SNPs) and the most likely causal SNPs (credible

452   SNPs) within a locus tend to be located within or close to the genes. Although active

453   chromatin states were enriched in both lead and credible SNPs, eQTLs were only enriched in

454   credible SNPs but depleted in lead SNPs. This implies that likely causal regulatory SNPs do

455   not necessarily have the strongest effect sizes in a locus.

456   Our analyses showed that the majority of analysed traits are highly polygenic with more than

457   0.1% of SNPs being causal. For those highly polygenic traits, over 10s of millions of

458   individuals are required to identify all SNPs at genome wide significance ($p$<5e-8) that can

459   explain at least 90% of the phenotypic variance explained by additive genetic effects. In the

460   case of polygenic traits, individuals have almost unique combinations of risk/effect alleles for

461   a specific disease or trait. With higher levels of polygenicity, and thus larger quantities of

462   causal SNPs, the possible combinations of them also increase. This substantially increases the

463   degree of genetic heterogeneity of the trait, and complicates the detection of genetic effects as

20

464     the effect sizes of individual SNPs that are yet to be detected are even smaller than those

465     observed in current GWASs.

466     In conclusion, our analyses have provided novel insight into the extent of pleiotropy, the

467     nature of associated genetic regions and how traits differ in genetic architectures. This

468     knowledge can guide the design of future genetic studies.

469     **METHODS**

470     **Publicly available GWAS summary statistics**

471     GWAS summary statistics were curated from multiple resources and were included only

472     when the full set of SNPs were available. We excluded whole exome sequencing studies.

473     This yielded 2,288 GWASs from 33 consortia and any other resources where summary

474     statistics are available (last update 23rd October 2018). From dbGAP, we obtained 2,659

475     unique datasets (ftp://ftp.ncbi.nlm.nih.gov/dbgap/Analyses_Table_of_Contents.txt, last

476     accessed 4th July 2017) and extracted 896 GWAS summary statistics in which a matched

477     publication was available and sample size for a specific trait was explicitly mentioned in the

478     original study. We excluded non-GWAS studies (e.g. PAGE (Prenatal Assessment of

479     Genomes and Exomes) studies) and GWASs with immune-chip, whole exome sequencing

480     and replication cohorts (exact reasons of exclusion for each dataset is available in

481     **Supplementary Table 24**).

482     Together this resulted in a total of 3,555 GWAS summary statistics. The complete list and

483     detailed information for each GWAS with summary statistics is available in **Supplementary**

484     **Table 3** (atlas ID 1-3184, 3785-4155).

485

486     **UK Biobank GWAS summary statistics**

487     Additional to the summary statistics available from external studies, we performed GWASs

488     of traits from UK Biobank release 2 cohort (UKB2)[12] under application ID 16404. We only

489     used phenotype fields with first visit and first run (e.g. f.xxx.0.0) with exceptions for multi-

490     coded phenotypes, which allowed to assign more than one code for a single subject (see

491     **Supplementary Information 1, 2**). From the 1,940 unique field IDs to which we had access,

492     755 had >50,000 subjects with non-missing values. They are assigned to field name using

493     ukb_field.tsv obtained from http://biobank.ctsu.ox.ac.uk/crystal/download.cgi (last accessed

494 31st August 2017). Note that for newly available phenotypes for release 2, we annotated field

495 names manually based on the UK biobank data showcase. From these phenotypes, we

496 excluded baseline characteristics, phenotypes used as covariates, date and place phenotypes,

497 status phenotypes (i.e. completion status, answered a specific question), ethnicity, genomic

498 phenotypes and any other phenotypes that are not relevant for performing a GWAS. For each

499 phenotype, we provided reason of exclusions in **Supplementary Table 1**. This resulted in

500 434 unique fields including 49 multi-coded phenotypes. 385 phenotypes were considered

501 quantitative when the phenotype value was quantitative or categorical, and could be ordered.

502 Phenotypes coded by yes/no were considered as binary with a few exceptions

503 (**Supplementary Table 1**). For quantitative and binary phenotypes, subjects with phenotype

504 codes -1 for "Do not know" or -3 for "Prefer to not answer" were excluded and the original

505 phenotype code as described in the UK biobank data showcase was used unless specified in

506 Supplementary Text or **Supplementary Table 1, 2**. For 49 multi-coded phenotypes, we

507 dichotomized each code to dummy binary phenotypes (cases for 1 and controls for 0) and

508 included subjects with phenotype code -7 for "None of the above" as controls. Again,

509 subjects with phenotype codes -1 for "Do not know" or -3 for "Prefer to not answer" were

510 excluded. For example, field 670 based on UKB Data-Coding 100286 is coded from 1 to 5

511 and dichotomization results in five phenotypes such as 1 vs all others, 2 vs all others and so

512 on. Detailed definitions of multi-coded phenotypes are described in **Supplementary Table 2**.

513 After phenotyping, we selected phenotypes that had at least 50,000 European subjects. For

514 binary traits, we further restricted to traits with at least 10,000 cases and controls. This

515 resulted in a total of 600 traits (260 quantitative and 340 binary traits). Note that the final

516 total sample size encoded in the atlas database (http://atlas.ctglab.nl) might be less than

517 50,000 due to lack of genotype data or missing values in covariates.

23

518    GWAS was performed for up to 10,846,944 SNPs with MAF > 0.0001 using PLINK 2[29],

519    while correcting for array, age (f.54.0.0), sex (f.31.0.0), Townsend deprivation index

520    (f.189.0.0), assessment centre (f.21003.0.0) and 20 PCs. Linear or logistic models were used

521    for quantitative or binary traits, respectively.

522    The complete list of traits from UK biobank release 2 analysed in this study is available in

523    **Supplementary Table 3** (atlas ID 3185-3784).

524

525    **Pre-processing of GWAS summary statistics**

526    Curated summary statistics were pre-processed to standardize the format. SNPs with $p \leq 0$ or

527    $>1$, or non-numeric values such as "NA" were excluded. For summary statistics with non-

528    hg19 genome coordinates, liftOver software was used to align to hg19. When only rsID was

529    available in the summary statistics file without chromosome and position, genome

530    coordinates were extracted from dbSNP 146. When rsID was missing, it was assigned based

531    on dbSNP 146. When only the effect allele was reported, the other allele was extracted from

532    dbSNP 146.

533

534    **Definition of lead SNPs and trait-associated loci**

535    For each GWAS, we defined lead SNPs and genomic trait-associated loci as described before

536    [30]. First, we defined independent significant SNPs with $p<$5e-8 and independent at $r^2<0.6$,

537    and defined LD blocks for each of independent significant SNPs based on SNPs with $p<0.05$.

538    Of these SNPs, we further defined lead SNPs that are independent at $r^2<0.1$. We finally

539    defined genomic trait-associated loci by merging LD blocks closer than 250kb. Each trait-

540    associated locus was then represented by the top SNP (with the minimum P-value) and its

541    genomic region was defined by the minimum and maximum position of SNPs which are in

542 LD ($r^2 \geq 0.6$) with one of the independent significant SNPs within the (merged) locus. We

543 used 1000 genome phase 3 (1000G)[31] as a reference panel to compute LD for most of the

544 GWASs in the database. For each GWAS, the matched population (from AFR, AMR, EAS,

545 EUR, SAS) was used as the reference based on the information obtained from the original

546 study. For trans-ethnic GWASs, the population with the largest total sample size was used.

547 When the GWAS was based on the UKB release 1 cohort (UKB1), we used 10,000 randomly

548 sampled unrelated White British subjects from UKB1 as reference. For other GWASs

549 performed in this study or GWASs based on the UKB2, 10,000 randomly selected unrelated

550 EUR subjects were used as a reference. Non-bi-allelic SNPs were excluded from any

551 analyses.

552 The reference panel used for each GWAS is provided in the column "Population" of

553 **Supplementary Table 3**. For trans-ethnic GWASs, the first population was used as

554 reference, e.g. EUR+EAS+SAS means EUR had the largest sample. GWASs based on the

555 UKB cohort was encoded either "UKB1 (EUR)" for UKB release 1 or "UKB2 (EUR)" for

556 UKB release 2.

557

558 **MAGMA gene and gene-set analysis**

559 We performed MAGMA v1.06[16] gene and gene-set analyses for every GWAS in the

560 database. For gene-analysis, 20,260 protein-coding genes were obtained using the R package

561 BioMart (Ensembl build v92 GCRh37). SNPs were assigned to genes with 1kb window at

562 both sides. The reference panel of corresponding populations used for each GWAS was based

563 on either 1000G, UKB1 or UKB2 as described in the previous section. The gene-set analysis

564 was performed with default parameters (snp-wise mean model). Gene-set analysis was

565 performed for 4,737 curated gene-sets (C2) and 5,917 GO terms (C5; 4,436 biological

566    processes, 580 cellular components and 901 molecular functions) from MsigDB v6.1

567    (http://software.broadinstitute.org/gsea/msigdb, last accessed 20 Apr 2018)[32].

568

569    **SNP heritability and genetic correlation with LD score regression**

570    We performed LD score regression (LDSC)[17] for each GWAS to obtain SNP heritability and

571    pairwise genetic correlations. Pre-calculated LD scores for 1000G EUR and EAS populations

572    were obtained from https://data.broadinstitute.org/alkesgroup/LDSCORE/ (last accessed 26

573    Nov 2016) and LD score regression was only performed for GWASs with either an EUR or

574    EAS population and when the number of SNPs in the summary statistics file was > 450,000.

575    LDSR was performed only for HapMap3 SNPs excluding the MHC region (25Mb-34Mb).

576    When the signed effect size or odds ratio was not available in the summary statistics file, "--

577    a1-inc" flag was used. As recommended previously[33], we excluded SNPs with chi-square

578    >80. For binary traits, the population prevalence was curated from the literature (only for

579    diseases whose prevalence was available, **Supplementary Table 25**) to compute SNP

580    heritability at the liability scale with "--samp-prep" and "--pop-prep" flags. For most of the

581    personality/activity (binary) traits from UKB2 cohort, we assumed that the sample prevalence

582    is equal to the population prevalence since the UK Biobank is a population cohort and not

583    designed to study a certain disease/traits. Likewise, when population prevalence was not

584    available, sample prevalence was used as population prevalence for all other binary traits.

585    Genetic correlations were computed for pair-wise GWASs with the following criteria as

586    suggested previously[33]:

587    • GWASs of EUR population or more than 80% of samples are EUR.

588    • The number of SNPs >450,000

589    • Signed effect size or odds ratio is available

590    • Effect and non-effect alleles are explicitly mentioned in the header or elsewhere.

26

591      • SNP heritability Z score >2

592    In total, pairwise genetic correlations were computed for 1,090 GWASs in the database.
593
594
595    **Selection of GWASs for cross-phenotype analyses**

596    From the 4,155 curated GWASs in the database, we selected 558 GWASs with unique traits

597    for cross-phenotype analyses based on the following criteria.

598      • Minimum sample size 50,000 and both cases and controls are >10,000 for binary

599        phenotypes.

600      • The number of SNPs in the summary statistics is >450,000.

601      • GWAS is based on EUR population or >80% of the samples are EUR. If summary

602        statistics of both trans-ethnic and EUR-only are available, use EUR-only GWAS.

603      • Exclude sex-specific GWAS, unless the phenotype under study is only available for a

604        specific sex (e.g., age at menopause). If  sex-specific and sex-combined GWASs are

605        available, use sex-combined GWAS.

606      • Z-score of $h^2_{SNP}$ computed by LDSC is >2

607      • Signed effect size (beta or odds ratio) is available in the summary statistics.

608      • Effect and non-effect alleles are explicitly mentioned in the header or elsewhere.

609      • From GWASs that met the above criteria, we selected a GWAS per trait with the

610        maximum sample size.

611

612    UKB2 GWASs performed in this study are further filtered based on the following:

613      • Exclude cancer screening or test phenotypes.

614      • Exclude item level phenotypes (i.e., Neuroticism and Fluid intelligence tests)

615      • Exclude phenotypes of parents' age and parents' still alive.

616      • Exclude medication, treatment, supplements and vitamin traits.

27

617     • If exactly the same traits were diagnosed by an expert (e.g. doctor) and self-reported,

618      use the expert qualification.

619     • If exactly the same traits were present as main and secondary diagnoses, both are

620      included.

621     • Phenotypes with large extremes were excluded from the analyses when the difference

622      between the maximum value and 99 percentiles of the standardized phenotype value

623      is >50.

624   There was one exception for height GWAS, where a meta-analysis by Yengo et al.[34] (ID

625   4044) has the larger sample size, however the meta-analysis was limited to ~2.4 million

626   HapMap 2 SNPs. Since over 10 million SNPs are included in most of the selected GWASs,

627   this smaller number of SNPs can bias our analyses. Therefore, the second largest GWAS

628   (UKB2 GWAS performed in this study, ID 3187) was used instead. This resulted in total of

629   558 GWASs, across 24 domains, which were subsequently used in the cross-phenotype

630   analyses in this study. These 558 GWASs are specified in **Supplementary Table 3**.

631

632   **Pleiotropic trait-associated loci**

633   To define pleiotropic loci for the 558 traits (GWASs), we first extracted trait-associated loci

634   on autosomal chromosomes. We excluded any locus with a single SNP (no other SNPs have

635   $r^2 > 0.6$) as these loci are more likely to be false positives. We then grouped physically

636   overlapping loci across 558 traits. In a group of loci, it is not required that all individual trait-

637   associated loci are physically overlapping but merging them should result in a continuous

638   genomic region. For example, when trait-associated loci A and B physically overlap and trait-

639   associated loci B and C also physically overlap, but A and C do not, these three trait-

640   associated loci were grouped into a single group of loci (**Extended Data Fig. 3**). Therefore, a

641   grouped locus could contain more than one independent locus from a single trait when gaps

28

642    between them were filled by loci from other traits. The grouped loci were further assigned to

643    three categories, *i*) multi-domain locus when a loci group contained traits from more than one

644    domain, *ii*) domain specific locus when a loci group contained more than one trait from the

645    same domain, and *iii*) trait specific locus when a locus did not overlap with any other loci.

646    We compared the distribution of gene density across four association categories of the loci;

647    multi-domain, domain specific and trait specific loci, and non-associated genomic regions.

648    To define non-associated genomic regions, we extracted the minimum and maximum

649    positions that were covered by 1000G, and the gap regions of grouped trait-associated loci

650    were defined as non-associated regions. The gene density was computed as a proportion of a

651    region that was overlapping with one of 20,260 protein-coding genes obtained from Ensembl

652    v92 GRCh37. We then performed pairwise Wilcoxon rank sum test (two sided).

653
654    **Colocalization of trait-associated loci**

655    To evaluate if physically overlapping trait-associated loci also share the same causal SNPs,

656    we performed colocalization using the *coloc.abf* (Approximate Bayes Factor colocalization

657    analysis) function of the coloc package in R[35]. Colocalization analysis was performed for all

658    possible pairs of physically overlapping trait-associated loci across 558 traits. When two loci

659    from different traits were physically overlapping but there were no SNPs that were present in

660    both GWAS summary statistics in that overlapping region, colocalization was not performed.

661    The inputs of the *coloc.abf* function are P-value, MAF and sample size for each SNP. When

662    MAF was not available in the original summary statistics, it was extracted from the matched

663    reference panel. For binary traits, sample prevalence was additionally provided based on total

664    cases and controls of the study.

665    The *coloc.abf* function assumes a single causal SNP for each trait and estimates the posterior

666    probability of the following 5 scenarios for each testing region; $H_0$: neither trait has a genetic

29

667    association, $H_1$: only trait 1 has a genetic association, $H_2$: only trait 2 has a genetic

668    association, $H_3$: both trait 1 and 2 are associated but with different causal SNPs and $H_4$: both

669    trait 1 and 2 are associated with the same single causal SNP. In this study, as we pre-define

670    the trait-associated loci for each trait which already discard scenarios $H_0$ to $H_2$, we are only

671    interested whether $H_4$ is most likely. We therefore defined, a pair of loci as colocalised when

672    the posterior probability of $H_4$ is >0.9. We note that it is possible that genomic regions

673    outside of the pre-defined trait-associated loci can also colocalize with other traits. However,

674    we limited the analyses to the pre-defined trait-associated loci in this study, to be consistent

675    with the level of pleiotropy measured by physical overlap of the loci.

676    Within a grouped locus defined based on physical overlap (see above), we further grouped

677    loci based on a colocalization pattern. To do so, we considered colocalization pattern across

678    group of physically overlapping loci as a graph in which nodes represent trait-associated loci

679    and edges represent colocalization of the loci First, loci which did not colocalized with any

680    other loci were considered as independent loci. For the rest of the loci, we identified

681    connected components of the graph (**Extended Data Fig. 3**). This does not require all loci

682    within a component to be colocalized with each other. For example, when locus A is

683    colocalized with locus B, and locus B is colocalized with locus C, but locus A is not

684    colocalized with locus C, all loci A, B and C are grouped into a single connected component.

685    Detailed results are discussed in the **Supplementary Information 3**.

686

687    **Pleiotropic genes**

688    For gene level pleiotropy, we extracted MAGMA gene analysis results for the 558 traits

689    where 17,444 genes on autosomal chromosomes were tested in all GWASs. For each trait,

690    genes with $p<2.87e-6$ (0.05/17,444) were considered as significantly associated. We did not

691    correct the P-value for testing 558 traits since our purpose is not to identify genes associated

692    with one of the 558 traits but to evaluate the overlap of trait-associations (when GWAS was

693    performed for a single trait) across the 558 traits, and this applies to SNPs and gene-set level

694    pleiotropy. The trait associated genes were further categorized into three groups in a similar

695    way as for trait-associated loci, i.e. *i*) multi-domain genes that were significantly associated

696    with traits from more than one domain, *ii*) domain-specific genes that were significantly

697    associated with more than one trait from the same domain and *iii*) trait-specific genes that

698    were significantly associated with a single trait.

699    We compared gene length and pLI score across genes in three different association categories

700    and non-associated genes. Gene length was based on the start and end position of genes

701    extracted from the R package biomaRt and pLI score was obtained from

702    ftp://ftp.broadinstitute.org/pub/ExAC_release/release0.3.1/functional_gene_constraint (last

703    accessed 27 April 2017). We performed t-tests for gene length in log scale and Wilcoxon

704    rank sum tests for pLI scores (both two sided).

705    For each protein coding gene, we first assessed whether a gene is expressed or not in each of

706    53 tissue types based on expression profile obtained from GTEx v7[20]. We defined genes as

707    expressed in a given tissue type if the average TPM is >1. For each of 17,444 genes, we then

708    counted the number of tissue types where the gene is expressed and grouped them into six

709    categories, i.e. genes expressed in *i*) a single tissue type (tissue specific genes), *ii*) between 2

710    and 13, *iii*) between 14 and 26, *ix*) between 27 and 39, *x*) between 40 and 52, and *xi*) 53 (all)

711    tissue types. At each number of associated domains (from 1 to 10 or more domains), we re-

712    calculated the proportion of genes in each of the 6 categories, and performed the Fisher's

713    exact tests (one-sided) against baseline (the proportion relative to all 17,444 genes) to

714    evaluate if the proportion is higher than expected.

715

716    **Pleiotropic SNPs**

717     We extracted 1,740,179 SNPs that were present in all 558 GWASs. To evaluate if the select

718     ion of ~1.7 million SNPs biased the results, we compared distribution of these analysed SNPs

719     with the all known SNPs in the genome (SNPs exist in 1000G EUR population, UKB1 and

720     UKB2 reference panels) by computing the proportion of SNPs per chromosome. In addition,

721     distribution of functional consequences of SNPs annotated by ANNOVAR[36] was also

722     compared with the all SNPs in the genome. For each SNP, we counted the number of traits to

723     which a SNP was significantly associated at $p<$5e-8, and then grouped the associated SNPs

724     into multi-domain, domain-specific and trait-specific SNPs using the same definitions as at

725     the gene level.

726     Functional consequences of SNPs were annotated using ANNOVAR[36]. To test if a SNP from

727     a certain functional category is enriched at a given number of associated domains compared

728     to all analysed SNPs, a baseline proportion was calculated from the 1,740,179 SNPs for each

729     functional category. At each number of associated domains (from 1 to 10 or more domains),

730     we re-calculated the proportion of SNPs with each functional category and performed the

731     Fisher's exact test (one-sided) against the baseline (the proportion relative to all 1,740,179

732     SNPs), to test if the proportion if higher than expected.

733     eQTLs for 48 tissue types were obtained from GTEx v7 (https://www.gtexportal.org/home/;

734     last accessed 20 January 2018)[20] and we considered SNPs with gene q-value <0.05 with any

735     gene in any tissue as eQTLs. For each eQTL, we counted the number of tissue types of being

736     eQTL (regardless of associated genes) and categorized them into five groups, i.e. being

737     eQTLs in *i*) a single tissue type (tissue specific eQTLs), *ii*) between two and 12, *iii*) between

738     13 and 24, *ix*) between 25 and 36 and *x*) and being in more than 37 tissue types. At each

739     number of associated domains, we re-calculated the proportion of SNPs in each of the 5

740     categories, and performed the Fisher's exact test (one-sided) against baseline (the proportion

741     relative to all 1,740,179 SNPs), to test if the proportion if higher than expected.

742

**Pleiotropic gene-sets**

744 For gene-set level pleiotropy, we extracted 10,650 gene-sets tested in all 588 traits. We then

745 considered gene-sets with $p<4.69e-6$ (0.05/10,650) as significantly associated. The trait

746 associated gene-sets were grouped into multi-domain, domain-specific and trait-specific

747 gene-sets with the same definitions as at the gene level.

748 We compared the number of genes and average gene-length across gene-sets in different

749 association categories and non-associated genes. Gene length was based on the start and end

750 position of genes extracted from R package, biomaRt. We performed two-sided t-test in log

751 scale of the number of genes and average gene-length.

752

**Power calculation of genetic correlation**

754 Power calculations were performed using the bivariate analysis of GCTA-GRML power

755 calculator (http://cnsgenomics.com/shiny/gctaPower/)[37], to estimate the minimum $r_g$ that

756 obtain a power of 0.8 in the worst case scenario. From 558 traits, two traits with the worst

757 case scenarios were selected, one with the minimum $h^2_{SNP}$ estimated by LDSC and another

758 with the minimum sample size. For each case, we obtained the minimum $r_g$ to obtain power

759 of 0.8 by assuming both traits are quantitative with same sample size and $h^2_{SNP}$ and have

760 phenotypic correlation 0.1.

761

**Hierarchical clustering of trait based on genetic correlation**

763 Hierarchical clustering was performed on the matrix of pair-wise $r_g$'s as calculated between

764 the 558 traits. After Bonferroni correction for all possible trait pairs, non-significant genetic

765 correlations were replaced with 0. The number of clusters $k$ was optimized between 50 and

766 250 by maximizing the silhouette score with 30 iterations for each $k$.

767

**Estimated standardized effect size of lead SNPs**

768

769  To enable comparison of effect sizes across GWASs from different studies, we first

770  converted P-values into Z-statistics (two sided) and expressed the estimated effect size as a

771  function of MAF and sample size as described previously[21] using the following equations:

772
$$\hat{b} = \frac{z}{\sqrt{2p(1-p)(n+z^2)}}, \qquad SE = \frac{1}{\sqrt{2p(1-p)(n+z^2)}}$$

773  where $p$ is MAF and n is the total sample size. We used the MAF of a corresponding

774  European reference panel (either 1000G, UKB1 or UKB2) as described in the previous

775  section "Definition of lead SNPs and genomic trait-associated loci". Since we were not

776  interested in the direction of effect, we used squared standardized effect sizes for analyses in

777  this study.

778

**Fine-mapping of trait-associated loci**

779

780  We defined the region to fine-map by taking 50kb around the top SNPs of the trait-associated

781  loci. When trait-associated loci were larger than the 50kb window, the largest boundary was

782  taken. Due to the complex LD structure, loci overlapping with the MHC region (chr6:25Mb-

783  36Mb) were excluded. The fine-mapping was performed using the FINEMAP software

784  (http://www.christianbenner.com/#) with shotgun stochastic search algorithm[25]. Since the

785  coverage of true causal SNPs is affected by the sample size of the reference panel and

786  GWASs[38], we used randomly selected unrelated 100k EUR individuals from UKB2 cohort

787  for all 558 GWASs. We limited the number of maximum causal SNPs ($k$) per locus to 10.

788  When the number of SNPs within a locus is relatively small (around 30 or less), the algorithm

789  can fail to converge. In that case, k was decreased by 1 until FINEMAP was successfully run.

790  Loci with less than 10 SNPs were excluded from the fine-mapping.

791  FINEMAP outputs a set of models (all possible combination of $k$ causal SNPs in a locus)

792  with posterior probability (PP) of being a causal model. A 95% credible set was defined by

793  taking models from the highest PP until the cumulative sum of PP reached 0.95. Then 95%

794  credible set SNPs were defined as unique SNPs included in the 95% credible set of models.

795  For each SNP, a posterior inclusion probability (PIP) was calculated as the sum of PPs of all

796  models that contains that SNP. To select most likely causal SNPs, we further defined credible

797  SNPs consists of SNPs with PIP>0.95. Detailed results are discussed in **Supplementary**

798  **Information 9.**

799

800  **Annotation and characterization of lead SNPs and credible SNPs**

801  Functional consequences of SNPs were annotated using ANNOVAR[36] based on Ensembl

802  gene annotations on hg19. Prior to ANNOVAR, we aligned the ancestral allele with dbSNP

803  build 146. 15-core chromatin states of 127 cell/tissue types were obtained from Roadmap[39]

804  (http://egg2.wustl.edu/roadmap/data/byFileType/chromhmmSegmentations/ChmmModels/co

805  reMarks/jointModel/final/all.mnemonics.bedFiles.tgz; last accessed 16 Mar 2016) and we

806  annotated one of the 15-core chromatin states to each of the lead SNPs based on chromosome

807  coordinates. Subsequently, consequence state was assigned for each SNP by taking the most

808  common state across 127 cell/tissue types. SNPs with consequence state$\leq$7 were defined as

809  active. eQTLs in 48 tissue types were obtained from GTEx v7[20] and we only used the

810  significant eQTLs at gene q-value<0.05. eQTLs were assigned to SNPs by matching

811  chromosome coordinate and alleles.

812  As we showed that trait-associated loci have higher gene density compared to non-associated

813  regions, and GWAS signals are known to be enriched in regulatory elements[40], we first

814  identified background enrichment by comparing SNPs within trait-associated loci or fine-

815  mapped regions with the entire genome. For this all known SNPs were extracted by

816     combining all SNPs in 1000G, UKB1 and UKB2 reference panels (~28 million SNPs in

817     total). SNPs within the trait-associated loci were defined as the ones with P-value<0.05 and

818     $r^2$>0.6 with one of the independent significant SNPs as described above (see section

819     'Definition of lead SNPs and trait-associated loci'). Therefore, it does not necessary include

820     all SNPs physically located within the trait-associated loci. On the other hand, SNPs within

821     fine-mapped region include all SNPs physically located within 50kb window from the most

822     significant SNP of a locus. To characterize lead SNPs and credible SNPs given background

823     enrichments, we compared these SNPs against all SNPs within trait-associated loci or fine-

824     mapped regions, respectively.

825

826     **SNP heritability estimation with SumHer using LDAK model**

827     We estimated SNP heritability of 558 traits using the SumHer function from the LDAK

828     software v5.0 (http://dougspeed.com/ldak/) [27]. Since our purpose was to compare estimates

829     from LDSC and SumHer, we used the 1000G EUR reference panel and extracted HapMap3

830     SNPs as consistent with LDSC. We used unique ID's of SNPs (consisting of

831     chromosome:posision:allele 1:allele2) instead of rsID to maximize the match between

832     GWAS summary statistics and the reference panel. The MHC region (chr6:25Mb-34Mb) was

833     excluded. As recommended by the author, SNPs with large effects ($Z^2/(Z^2+n)$>100 where $Z^2$

834     is chi-squared statistics and $n$ is sample size of the SNP) were excluded.

835     To obtain SNP heritability in a liability scale, we provided population prevalence and sample

836     prevalence with flags '--prevalance' and '--ascertainment' for binary traits. The same

837     population prevalence was used as described in the section of SNP heritability estimate with

838     LDSC (**Supplementary Table 25**). Details results are discussed in **Supplementary**

839     **Information 10**.

840

**Estimation of polygenicity and discoverability with MiXeR**

In the causal mixture model for GWAS summary statistics (MiXeR) proposed by Holland et al., the distribution of SNP effect sizes is treated a mixture of two distributions for causal and non-causal SNPs as the following[28]:

$$\beta = \pi N\left(0, \sigma_\beta^2\right) + (1 - \pi)N(0,0)$$

where $\pi$ is the proportion of (independent) causal SNPs and $\sigma_\beta^2$ is the variance of the effect sizes of causal SNPs. Therefore, $\pi$ and $\sigma_\beta^2$ respectively represent polygenicity and discoverability of the trait. We estimated both parameters for the 558 traits using MiXeR software (https://github.com/precimed/mixer)[28]. As recommended in the original study, we used 1000G EUR as a reference panel and restricted to HapMap 3 SNPs. SNPs with $\chi^2>80$ and the MHC region (chr6:26Mb-34Mb) were excluded. To estimate the sample size required to explain 90% of the additive genetic variance of a phenotype, we used an output of GWAS power estimates calculated in the MiXeR software, which contains 51 data points of sample size and the proportion of chip heritability explained[28]. We then estimated the sample size required to reaches 90% by using the *interp1* function from the pracma package in R.

**Data and materials availability**

All publicly available GWAS summary statistics (original) files curated in this study are accessible from the original links provided at http://atlas.ctglab.nl. GWAS summary statistics for 600 traits from UK Biobank performed in this study are also provided at http://atlas.ctglab.nl and an archived file will be made available upon publication from https://ctg.cncr.nl/software/summary_statistics.

## REFERENCES

1. Edwards, A. O. *et al.* Complement factor H polymorphism and age-related macular degeneration. *Science (80-. ).* **308,** 421–425 (2005).

2. Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* **42,** D1001–D1006 (2014).

3. Lander, E. S. Initial impact of the sequencing of the human genome. *Nature* **470,** 187–197 (2011).

4. Visscher, P. M. *et al.* 10 Years of GWAS Discovery : Biology, Function, and Translation. *Am. J. Hum. Genet.* **101,** 5–22 (2017).

5. Henderson, P. & Stevens, C. The role of autophagy in Crohn's Disease. *Cells* **1,** 492–519 (2012).

6. Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature* **506,** 376–81 (2014).

7. Gaulton, K. J. *et al.* Genetic fine mapping and genomic annotation defines causal mechanisms at type 2 diabetes susceptibility loci. *Nat. Genet.* **47,** 1415–1425 (2015).

8. Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK Biobank. *Nat. Genet.* **50,** 1593–1599 (2018).

9. Timpson, N. J., Greenwood, C. M. T., Soranzo, N., Lawson, D. J. & Richards, J. B. Genetic architecture: The shape of the genetic contribution to human traits and disease. *Nat. Rev. Genet.* **19,** 110–124 (2018).

10. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An expanded view of complex traits: from polygenic to omnigenic. *Cell* **169,** 1177–1186 (2017).

11. Wray, N. R., Wijmenga, C., Sullivan, P. F., Yang, J. & Visscher, P. M. Common disease is more complex than implied by the core gene omnigenic model. *Cell* **173,** 1573–1580 (2018).

12. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562,** 203–209 (2018).

13. Goh, K. *et al.* The human disease network. *Proc. Natl. Acad. Sci.* **104,** 8685–8690 (2007).

14. Polderman, T. J. C. *et al.* Meta-analysis of the heritability of human traits based on fifty years of twin studies. *Nat. Publ. Gr.* **47,** 702–709 (2015).

15. Mahajan, A. *et al.* Fine-mapping type 2 diabetes loci to single-variant resolution using high-density imputation and islet-specific epigenome maps. *Nat. Genet.* (2018).

16. de Leeuw, C. A., Mooij, J. M., Heskes, T. & Posthuma, D. MAGMA: generalized gene-set analysis of GWAS data. *PLoS Comput. Biol.* **11,** e1004219 (2015).

17. Bulik-sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47,** 291–295 (2015).

18. Solovieff, N., Cotsapas, C., Lee, P. H., Purcell, S. M. & Smoller, J. W. Pleiotropy in complex traits: Challenges and strategies. *Nat. Rev. Genet.* **14,** 483–495 (2013).

19. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536,** 285–291 (2016).

20. The GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature* **550,** 204–213 (2017).

21. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* **48,** 481–487 (2016).

22. Manolio, T. A. *et al.* Finding the missing heritability of complex diseases. *Nature* **461,** 747–753 (2009).

23. Lee, S., Abecasis, G. R., Boehnke, M. & Lin, X. Rare-variant association analysis:

912        Study designs and statistical tests. *Am. J. Hum. Genet.* **95,** 5–23 (2014).

913   24.   van de Bunt, M., Cortes, A., Brown, M. A., Morris, A. P. & McCarthy, M. I.
914        Evaluating the performance of fine-mapping strategies at common variant GWAS loci.
915        *PLoS Genet.* **11,** e1005535 (2015).

916   25.   Benner, C. *et al.* FINEMAP : efficient variable selection using summary data from
917        genome-wide association studies. *Bioinformatics* **32,** 1493–1501 (2016).

918   26.   Speed, D. *et al.* Reevaluation of SNP heritability in complex human traits. *Nat. Genet.*
919        **49,** 986–992 (2017).

920   27.   Speed, D. & Balding, D. J. Better estimation of SNP heritability from summary
921        statistics provides a new understanding of the genetic architecture of complex traits.
922        *Nat. Genet.* (2018).

923   28.   Holland, D. *et al.* Beyond SNP heritability: polygenicity and discoverability estimated
924        for multiple phenotypes with a univariate gaussian mixture model. *bioRxiv* (2018).

925   29.   Purcell, S. *et al.* PLINK: A tool set for whole-genome association and population-
926        based linkage analyses. *Am. J. Hum. Genet.* **81,** 559–575 (2007).

927   30.   Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping
928        and annotation of genetic associations with FUMA. *Nat. Commun.* **8,** 1826 (2017).

929   31.   Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526,** 68–74
930        (2015).

931   32.   Liberzon, A. *et al.* Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27,**
932        1739–1740 (2011).

933   33.   ZHENG, J. *et al.* LD Hub: a centralized database and web interface to perform LD
934        score regression that maximizes the potential of summary level GWAS data for SNP
935        heritability and genetic correlation analysis. *Bioinformatics* **33,** 272–279 (2017).

936   34.   Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body
937        mass index in ∼700000 individuals of European ancestry. *Hum. Mol. Genet.* **27,** 3641–
938        3649 (2018).

939   35.   Giambartolomei, C. *et al.* Bayesian Test for Colocalisation between Pairs of Genetic
940        Association Studies Using Summary Statistics. *PLoS Genet.* **10,** e1004383 (2014).

941   36.   Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic
942        variants from high-throughput sequencing data. *Nucleic Acids Res.* **38,** e164 (2010).

943   37.   Visscher, P. M. *et al.* Statistical power to detect genetic (co)variance of complex traits
944        using SNP data in unrelated samples. *PLoS Genet.* **10,** e1004269 (2014).

945   38.   Benner, C. *et al.* Prospects of fine-mapping trait-associated genomic regions by using
946        summary statistics from genome-wide association studies. *Am. J. Hum. Genet.* **101,**
947        539–551 (2017).

948   39.   Roadmap Epigenomics Consortium. Integrative analysis of 111 reference human
949        epigenomes. *Nature* **518,** 317–330 (2015).

950   40.   Tak, Y. G. & Farnham, P. J. Making sense of GWAS: using epigenomics and genome
951        engineering to understand the functional relevance of SNPs in non-coding regions of
952        the human genome. *Epigenetics Chromatin* **8,** 57 (2015).

953

954     **END NOTES**

955     **Acknowledgement** We thank all consortiums and all other individual labs for making

956     GWAS summary statistics publicly available. We also thank Peter Visscher and Naomi Wray

957     for their thoughtful suggestions and discussions. We additionally thank Anders Dale for his

958     suggestions for the manuscript. This work was funded by Netherlands Organization for

959     Scientific Research (NWO VICI 453-14-005 and NWO VIDI 452-12-014).

960     **Author contribution** D.P. designed the study. K.W. curated the database and performed

961     analyses. T.J.C.P assisted with harmonization of phenotype labels of the database. S.S.

962     performed QC on the UK Biobank data and wrote the analysis pipeline for UKB analyses.

963     M.U.M assisted with the fine-mapping analyses. O.F. and O.A.A. developed software

964     UGMG and assisted with the analyses. S.v.d.S and B.M.N discussed and provided valuable

965     suggestions for analyses. K.W. and D.P. wrote the paper. All authors critically reviewed the

966     paper.

967     **Competing interests** The authors declare no competing financial interest.

968     **Corresponding author** Correspondence and requests for materials should be addressed to

969     D.P. (danielle.posthuma@vu.nl).

970

40

971 **Table 1. Count and proportion of pleiotropic trait-associated loci, genes, SNPs and**

972 **gene-sets.**

| | Loci | | Genes | | SNPs | | Gene-set | |
|---|---|---|---|---|---|---|---|---|
| | Length (Mb) | % | Count | % | Count | % | Count | % |
| **Total in genome** | 2796.10 | 100.00 | 17,444 | 100.00 | 1,740,179 | 100.00 | 10,650 | 100.00 |
| **Associated** | 1706.00 | 61.01 | 11,443 | 65.60 | 236,388 | 13.58 | 1,106 | 10.38 |
| Pleiotropic* | 1592.53 | 93.35 | 9,252 | 80.85 | 142,376 | 60.23 | 606 | 54.79 |
| Multi-domain | 1535.76 | 90.02 | 7,657 | 66.91 | 76,650 | 32.43 | 361 | 32.64 |
| Domain specific | 56.77 | 3.33 | 1,595 | 13.94 | 65,726 | 27.80 | 245 | 22.15 |
| Trait specific | 113.48 | 6.65 | 2,191 | 19.15 | 94,012 | 39.77 | 500 | 45.21 |
| **Non-associated** | 1090.10 | 38.99 | 6,001 | 34.40 | 1,503,791 | 86.42 | 9,544 | 89.61 |

973 *The count of pleiotropic loci, genes, SNPs and gene-sets is the sum of the multi-domain and

974 domain specific categories. Proportion of pleiotropic, multi-domain, domain specific and trait

975 specific categories are relative to the associated loci, SNPs, genes or gene-sets, respectively.

976

41

**Table 2. Characteristics of lead SNPs and credible SNPs with PIP>0.95 across 558 traits versus all SNPs in the genome.**

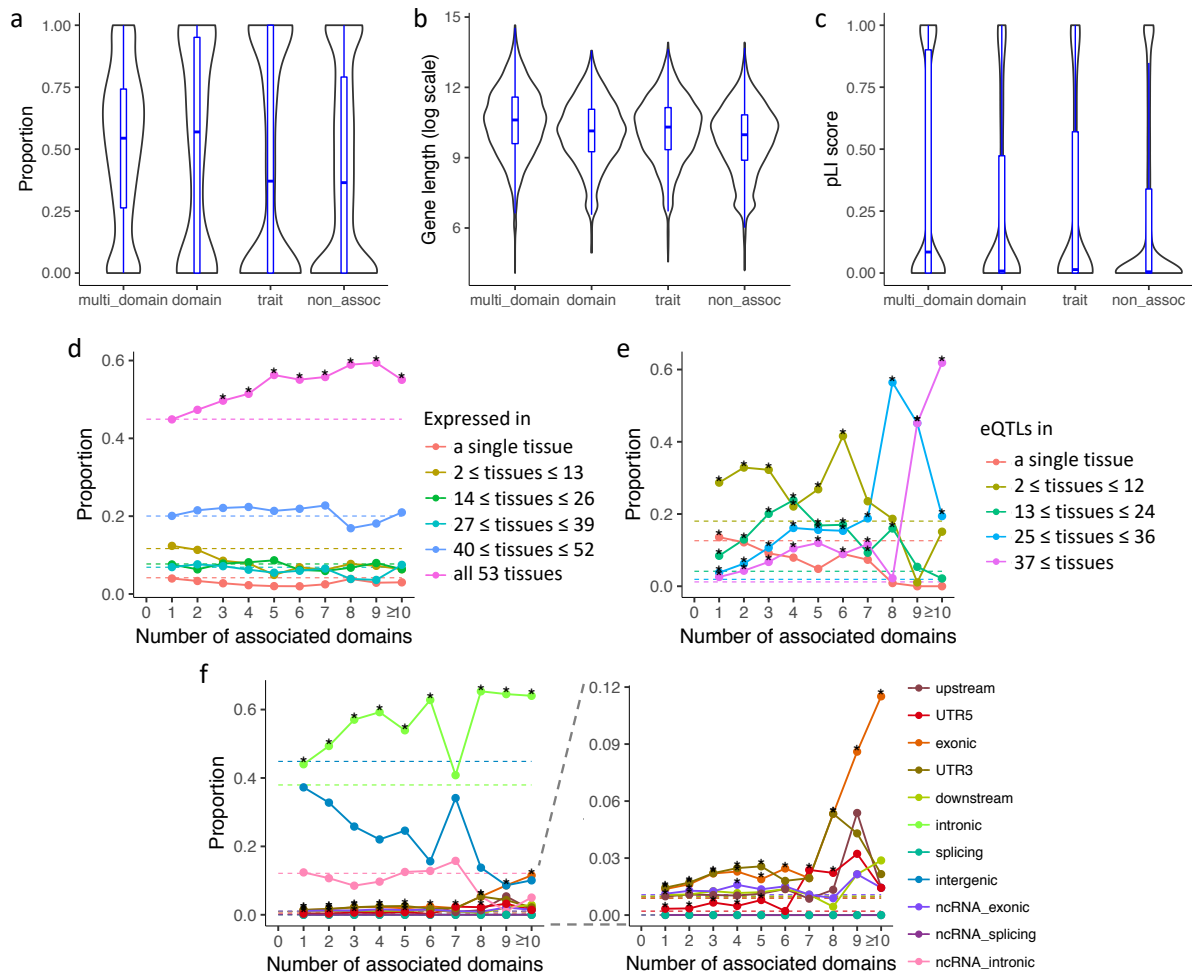| Annotation categories | Genome | Trait-associated loci | | | lead SNPs | | | 50kb around the top SNPs[a] | | | Credible SNPs (PIP>0.95)[b] | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | % | % | E | P[c] | % | E | P[d] | % | E | P[c] | % | E | P[e] |
| **Non-coding** | 94.37 | 93.06 | 0.99 | < 1e-323 | 89.13 | 0.96 | 1.14E-185 | 94.04 | 1.00 | < 1e-323 | 92.39 | 0.98 | 7.60E-192 |
| Intergenic | 44.11 | 36.88 | 0.84 | < 1e-323 | 34.31 | 0.93 | 1.20E-27 | 41.41 | 0.94 | < 1e-323 | 33.40 | 0.81 | < 1e-323 |
| Intronic | 38.29 | 44.88 | 1.17 | < 1e-323 | 43.85 | 0.98 | 2.38E-05 | 41.14 | 1.07 | < 1e-323 | 48.07 | 1.17 | < 1e-323 |
| scRNA intronic | 11.98 | 11.29 | 0.94 | 1.34E-115 | 10.98 | 0.97 | 0.044458 | 11.49 | 0.96 | < 1e-323 | 10.92 | 0.95 | 2.49E-15 |
| **Coding** | 2.15 | 2.40 | 1.12 | 7.42E-73 | 4.60 | 1.92 | 1.33E-147 | 2.27 | 1.06 | 4.02E-186 | 2.86 | 1.26 | 7.38E-63 |
| Exonic | 1.06 | 1.13 | 1.07 | 2.27E-14 | 3.22 | 2.84 | 1.30E-230 | 1.20 | 1.14 | < 1e-323 | 1.68 | 1.40 | 1.62E-73 |
| Splicing | 1.16E-02 | 1.13E-02 | 0.98 | 8.62E-01 | 2.11E-02 | 1.86 | 0.102234 | 1.29E-02 | 1.11 | 7.00E-05 | 1.95E-02 | 1.51 | 1.59E-02 |
| ncRNA exonic | 1.07 | 1.25 | 1.16 | 6.02E-71 | 1.36 | 1.09 | 0.04846 | 1.05 | 0.98 | 5.12E-11 | 1.16 | 1.10 | 4.14E-06 |
| ncRNA splicing | 5.40E-03 | 5.09E-03 | 0.94 | 7.03E-01 | 2.35E-03 | 0.46 | 0.72602 | 5.25E-03 | 0.97 | 5.06E-01 | 3.07E-03 | 0.59 | 2.66E-01 |
| **Flanking regions** | 3.48 | 4.54 | 1.31 | < 1e-323 | 6.27 | 1.38 | 4.60E-57 | 3.68 | 1.06 | 1.04E-299 | 4.75 | 1.29 | 1.48E-125 |
| Upstream | 1.09 | 1.33 | 1.22 | 9.09E-124 | 1.64 | 1.23 | 1.08E-07 | 1.09 | 1.00 | 7.59E-01 | 1.29 | 1.18 | 5.45E-16 |
| 5' UTR | 0.30 | 0.44 | 1.48 | 4.61E-151 | 0.78 | 1.76 | 1.64E-20 | 0.35 | 1.16 | 4.71E-183 | 0.57 | 1.66 | 8.75E-55 |
| 3' UTR | 0.98 | 1.32 | 1.34 | 2.41E-260 | 2.06 | 1.56 | 5.69E-34 | 1.13 | 1.15 | < 1e-323 | 1.67 | 1.48 | 2.47E-98 |
| Downstream | 1.10 | 1.45 | 1.32 | 4.18E-256 | 1.79 | 1.23 | 3.38E-08 | 1.11 | 1.01 | 9.73E-03 | 1.21 | 1.09 | 5.23E-05 |
| **Active chromatin** | 17.24 | 27.74 | 1.61 | < 1e-323 | 30.10 | 1.08 | 1.24E-27 | 20.63 | 1.20 | < 1e-323 | 31.06 | 1.51 | < 1e-323 |
| **eQTLs** | 9.66 | 57.41 | 5.95 | < 1e-323 | 46.15 | 0.80 | 7.54E-190 | 11.45 | 1.19 | < 1e-323 | 47.47 | 4.14 | < 1e-323 |

978 E: fold enrichment (proportion of SNPs with a certain annotation divided by the proportion of SNPs with the same annotation in background).

979 [a]Only including the fine-mapped regions (for loci larger than 50kb windows from the top SNPs, the largest boundaries were taken). [b]From 95%

980 credible set SNPs, only SNPs with posterior inclusion probability (PIP)>0.95 were selected. [c]P-value of Fisher's exact test (two-sided) against

981 the entire genome. [d]P-value of Fisher's exact test (two-sided) against trait-associated loci. [e]P-value of Fisher's exact test (two-sided) against
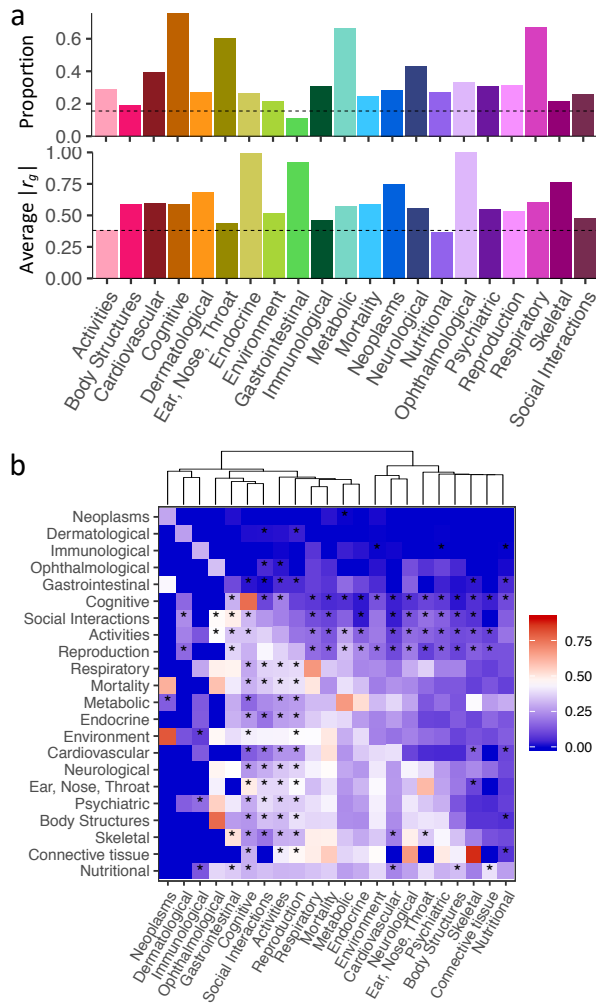
982 50kb around the top SNPs.

**Fig. 1. Trait-associated locus, gene and SNP pleiotropy across the genome. a.**
Distribution of gene density of loci with different association types. **b.** Distribution of gene
length in log scale with different association types. **c.** Distribution of pLI score of genes with
different association types. For **a-c**, multi_domain: associated with traits from >1 domain,
domain: associated with >1 traits from a single domain, trait: associated with a single trait,
non_assoc: not associated with any of 558 traits. **d.** Tissue specificity of genes at different
levels of pleiotropy. Each data point represents a proportion of genes expressed in a given
number of tissues for a specific number of associated domains. **e.** Proportion of SNPs with
different functional consequences at different levels of pleiotropy. Each data point represents
the proportion of SNPs with a given functional consequence for a specific number of

994      associated domains. **f.** Tissue specificity of SNPs based on active eQTLs at different levels of

995      pleiotropy. Each data point represents the proportion of SNPs being eQTLs in a given

996      number of tissues for a specific number of associated domains. For **d-f**, dashed lines refer to

997      the baseline proportions (relative to all 17,444 genes (d) or all 1,740,179 SNPs (e-f)), and

998      stars denote significant enrichment relative to the baseline (Fisher's exact test, one-sided).

999

**Fig. 2. Within and between domains genetic correlations. a.** Proportion of trait pairs with

significant $r_g$ (top) and average $|r_g|$ for significant trait pairs (bottom) within domains. Dashed

lines represent the proportion of trait pairs with significant $r_g$ (top) and average $|r_g|$ for

significant trait pairs (bottom) across all 558 traits, respectively. Connective tissue, muscular

and infection domains are excluded as these each contains less than 3 traits. **b.** Heatmap of

proportion of trait pairs with significant $r_g$ (upper right triangle) and average $|r_g|$ for

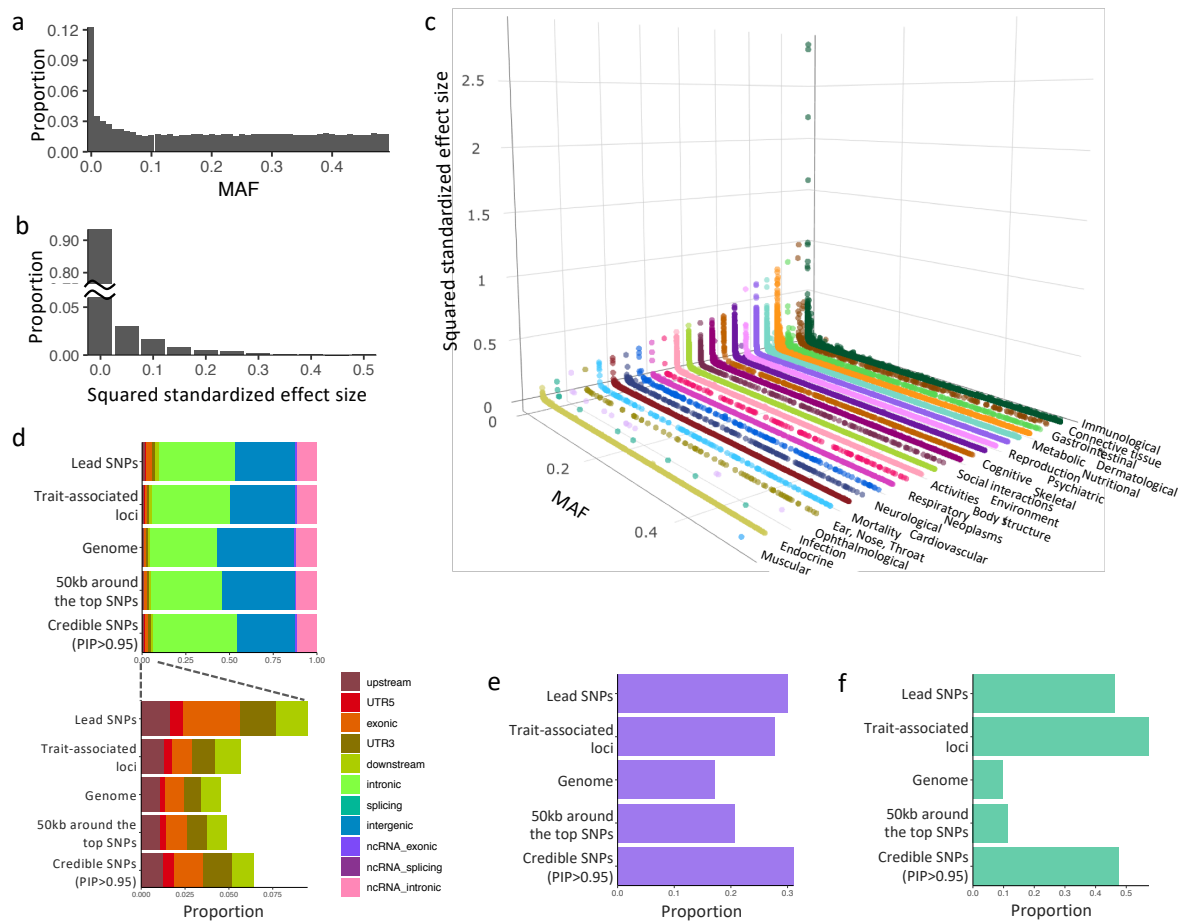significant trait pairs (lower left triangle) between domains. Connective tissue, muscular and

infection domains are excluded as each contains less than 3 traits. The diagonal represents the

proportion of trait pairs with significant $r_g$ within domains. Stars denote the pairs of domains

in which the majority (>50%) of significant $r_g$ are negative.

1011

**Fig. 3. Distribution and characterization of lead SNPs and credible SNPs of 558 traits. a.**
Histogram of MAF of the unique lead SNPs. **b.** Histogram of squared standardized effect size
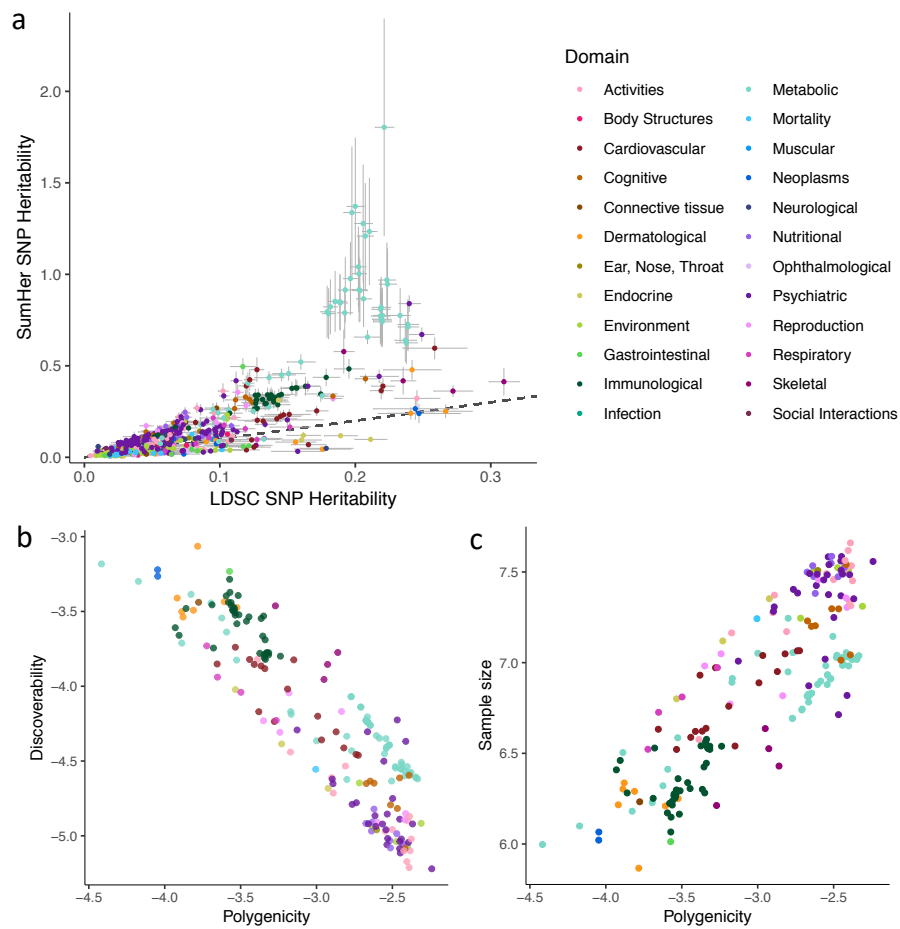of lead SNPs. **c.** Scatter plot of MAF and squared standardized effect sizes of lead SNPs
grouped by trait domains. **d.** Distribution of functional consequences of SNPs. **e.** Proportion
of SNPs that overlap with active consequence chromatin state (≤7) across 127 tissue/cell
types. **f.** Proportion of SNPs overlapping with significant eQTLs from any of 48 available
tissue types.

1019

**Fig. 4. SNP heritability and polygenicity of 558 traits. a.** Comparison of SNP heritability estimated by LDSC (x-axis) and SumHer (y-axis). Horizontal and vertical error bar represent standard errors of LDSC and SumHer estimates, respectively. **b.** Polygenicity and discoverability of traits, both on log 10 scale. Out of 558 traits, only 197 traits with reliable estimates (i.e. $h^2_{SNP}$>0.05 (estimated by MiXeR) and standard error of $\pi$ is less than 50% of the estimated value) are displayed. Traits are colored by domain. **c.** Polygenicity and estimated sample size required to reach 90% of total SNP heritability explained by genome-wide significant SNPs, both in log 10 scale. Traits are colored by domain. Full results are available in **Supplementary Table 22**, **23**.