



33 15. Division of Cancer Epidemiology and Genetics, National Cancer Institute, NIH,  
34 Department of Health and Human Services, Bethesda, MD, USA.

35

36 To whom correspondence should be addressed:

37 Dr William Reay, Medical Sciences Building, University Drive, The University of  
38 Newcastle, Callaghan, NSW, Australia, 2308. Email: [william.reay@newcastle.edu.au](mailto:william.reay@newcastle.edu.au)  
39 or Prof Murray Cairns, Medical Sciences Building, University Drive, The University of  
40 Newcastle, Callaghan, NSW, Australia, 2308, Email: [murray.cairns@newcastle.edu.au](mailto:murray.cairns@newcastle.edu.au)

41

## 42 **ABSTRACT**

43 Retinol is a fat-soluble vitamin that plays an essential role in many biological processes  
44 throughout the human lifespan. Previous work has characterised genetic influences on  
45 circulating retinol; however, small sample sizes have limited our ability to fully appreciate the  
46 genetic architecture of this trait. In this study, we performed the largest genome-wide  
47 association study (GWAS) of retinol to date in up to 22,274 participants. We identified eight  
48 common variant loci associated with retinol, as well as a rare-variant signal. An integrative  
49 gene prioritisation pipeline supported novel retinol-associated genes outside of the main retinol  
50 transport complex (*RBP4:TTR*) related to lipid biology, energy homeostasis, and endocrine  
51 signalling. Genetic proxies of circulating retinol were then used to estimate causal relationships  
52 with almost 20,000 clinical phenotypes via a phenome-wide Mendelian randomisation study  
53 (MR-pheWAS). The MR-pheWAS suggested that retinol may exert causal effects on  
54 inflammation, adiposity, ocular measures, the microbiome, and MRI-derived brain phenotypes,  
55 amongst several others. Conversely, circulating retinol may be causally influenced by factors  
56 including lipids and renal function. Finally, we demonstrated how a retinol polygenic score  
57 could identify individuals who are more likely to fall outside of the normative range of  
58 circulating retinol for a given age. In summary, this study provides a comprehensive evaluation  
59 of the genetics of circulating retinol, as well as revealing traits which should be prioritised for  
60 further clinical investigation with respect to retinol related therapies or nutritional intervention.

61

62

63

64

65

66

## 67 INTRODUCTION

68 Vitamin A is an essential micronutrient that is involved in a range of important biological  
69 processes, including, vision, immune function, cell division, and neurodevelopment<sup>1,2</sup>. Vitamin  
70 A does not refer to a single compound, but rather to a group of compounds that encompasses  
71 retinol (all-*trans* retinol), retinoids (metabolites of retinol, such as retinaldehyde and retinoic  
72 acid), and provitamin carotenoids (beta-carotene, alpha-carotene, and beta-cryptoxanthin).  
73 Retinol is the form of vitamin A dietarily consumed from animal products, along with retinyl  
74 ester<sup>3</sup>, while plant-based materials contain precursors termed carotenoids that can be converted  
75 to retinaldehyde<sup>4</sup>. Retinoic acid, an oxidised form of retinaldehyde, is a particularly potent  
76 signalling molecule that regulates the expression of thousands of genes after binding to nuclear  
77 receptors including the retinoic-acid receptor and retinoid-X receptor subgroups<sup>5,6</sup>.

78  
79 The majority of dietary retinol is delivered to the liver, which is the primary organ responsible  
80 for its storage and metabolism. Retinol binding protein 4 (RBP4) is the major systemic  
81 transporter of retinol after hepatic secretion, facilitating delivery of retinol throughout the  
82 body<sup>3,7,8</sup>. RBP4 in turn complexes with the tetramer protein transthyretin (TTR), which  
83 stabilises circulating RBP4 and reduces renal filtration<sup>9</sup>. Notably, retinol can also be delivered  
84 directly to target tissues through other mechanisms, such as its postprandial packaging into  
85 lipid chylomicrons, as reviewed elsewhere<sup>3</sup>.

86  
87 The role of retinoid-related interventions in human disease for individuals who are not retinol  
88 deficient has been of long-standing interest. Synthetic retinoids that are structurally similar to  
89 retinol/retinoic acid are approved for dermatological indications (e.g., adapalene) and some  
90 cancers (e.g., bexarotene)<sup>10,11</sup>, with continued interest in repurposing these compounds across  
91 a range of other indications, including neuropsychiatry<sup>2,12</sup>. Currently, retinol supplementation  
92 is not specifically indicated unless an individual is deficient, which is rare in high-income  
93 countries, though much more common in low-income countries. Numerous observational or  
94 randomised controlled studies have explored the effects of supplementation, a high vitamin A  
95 diet, and/or measured circulating retinol in a variety of disease contexts. However, the data  
96 from these efforts have often either been null or conflicting between studies<sup>12-16</sup>. Despite this,  
97 recent observational evidence suggests a relationship between a greater serum retinol  
98 abundance and lower mortality in a large, prospective 30-year follow up study<sup>17</sup>.

99

100 Genetics provides a powerful tool to better characterise factors that influence the abundance of  
101 circulating retinol in serum. Moreover, estimated genetic effects on retinol can be utilised to  
102 understand potential causal relationships with human health and disease, which may be  
103 informative for supplementation, dietary intervention, or drug repurposing<sup>18,19</sup>. Family studies  
104 have suggested that circulating retinol is significantly heritable<sup>20</sup>, albeit estimated in small  
105 sample sizes. Similarly, dedicated genome-wide association studies (GWAS) of circulating  
106 retinol have also been limited to very modest sample sizes<sup>21,22</sup>. In 2011, Mondul *et al.* published  
107 findings of two genome-wide significant loci associated with retinol (N = 5006). These loci  
108 were plausibly mapped to the genes *RBP4* and *TTR*, respectively, which form the primary  
109 retinol transport complex<sup>22</sup>. There has been comparatively little progress in further  
110 characterising the genetic architecture of retinol since that time relative to other micronutrients  
111 like vitamin D, for which large sample size GWAS (N > 400,000) have been released<sup>23,24</sup>. The  
112 recent adoption of untargeted high-throughput metabolomics platforms with coverage for  
113 retinol in some existing genotyped cohorts presents a new opportunity to boost statistical  
114 power. As a result, in the current study, we aim to perform the largest GWAS of circulating  
115 retinol to-date to identify novel loci and leverage these data to study how retinol relates to  
116 health and disease.

117

## 118 **RESULTS**

119

### 120 **The common and rare variant genetic architecture of circulating retinol**

121 We integrated common and rare variant data from up to 22,274 individuals of European  
122 ancestry in our discovery meta-analyses to estimate genetic effects on circulating retinol  
123 (Figure 1A-B, Online Methods). Firstly, variants from the INTERVAL and METSIM studies  
124 were meta-analysed (N<sub>Meta</sub> = 17,268 – termed *METSIM+INTERVAL*, Figure 1, Online  
125 Methods). After harmonisation, there were 8,173,975 common and 5,091,050 rare [minor allele  
126 frequency (MAF) < 1%] overlapping variants between INTERVAL and METSIM,  
127 respectively. As retinol effect sizes were estimated in the same units in both studies (plasma  
128 SD units, quantified by the same instrument), we conducted both a fixed-effects inverse-  
129 variance weighted (IVW) meta-analysis and a sample-size weighted meta-analysis of Z scores  
130 (Stouffer's method). Secondly, we conducted an additional meta-analysis (N<sub>Meta</sub> = 22,274)  
131 which also included data from two other studies (ATBC and PLCO), termed the  
132 *METSIM+INTERVAL+ATBC+PLCO* meta-analysis. However, there were markedly fewer  
133 variants available in this meta-analysis after imputing the ATBC and PLCO summary statistics

134 and harmonising with METSIM+INTERVAL ( $N_{\text{var}} = 3,896,351$ , Online Methods). As a result,  
135 we focused on the *METSIM+INTERVAL* meta-analysis as the primary discovery dataset.

136

137 Considering common variants from the HapMap3 panel ( $\text{MAF} > 0.05$ , outside of the major  
138 histocompatibility complex (MHC) region), we observed a relatively subtle inflation of retinol  
139 signals across the genome as indexed by the mean  $\chi^2$  statistic, with mean  $\chi^2$  values around  
140 1.03-1.04, regardless of the meta-analysis considered (Supplementary Table 1). Common  
141 variant SNP heritability ( $h^2_{\text{SNP}}$ ) estimated using the linkage disequilibrium (LD) score  
142 regression (LDSR) approach and the 1000 genomes European reference panel was between  
143 6%-7% and nominally statistically significant, although with somewhat large standard errors  
144 (2.6%-3.2%) (Figure 1C, Supplementary Table 2). SNP heritability estimates of circulating  
145 retinol increased to between 10%-13%, but were still noisy, using two alternate models to  
146 estimate  $h^2_{\text{SNP}}$  and the UK Biobank (UKBB) as the LD reference (Figure 1C, Online Methods,  
147 Supplementary Table 2). Partitioned  $h^2_{\text{SNP}}$  across tissues and cell-types demonstrated nominal  
148 enrichment in biologically logical contexts such as liver, adipose, pancreas, and blood  
149 (Supplementary Figure 1). The somewhat large  $h^2_{\text{SNP}}$  standard errors are likely a product of  
150 sample size; however, we then conducted further analysis to explore the extent of the polygenic  
151 signal associated with retinol across the genome using an Empirical Bayes' method (Online  
152 Methods). This method was utilised to model the number of non-null effects on retinol genome-  
153 wide stratified by bins of LD scores that index the extent of LD a variant exhibits with other  
154 variants (Figure 1D). Across all LD score bins, we estimated that the mean fraction of common  
155 variants across the genome with non-null effects on retinol was between 1.4% to 2.4%,  
156 depending on the modelling parameters used. In line with expectation, the proportion of non-  
157 null retinol effects was very high ( $> 50\%$ ) when considering the variants that display the most  
158 extensive LD (highest LD score bins). The application of these analyses to two GWAS of  
159 another vitamin (25-hydroxyvitamin D<sub>3</sub>) with either comparable or much larger sample  
160 sizes<sup>23,25</sup>, suggested that the less-polygenic architecture of retinol observed in this study may  
161 become more diffuse across the genome with greater sample sizes, in line with many other  
162 quantitative traits (Supplementary Figure 2). However, we caution that further investigation  
163 will be required as these data become available to confirm this.

164

165 Next, we processed the common variant results of both meta-analyses to identify genome-wide  
166 significant loci associated with circulating retinol ( $P_{\text{GWAS}} < 5 \times 10^{-8}$ ). In the primary discovery  
167 meta-analysis (*METSIM+INTERVAL*, Stouffer's method), we uncovered eight genome-wide

168 significant loci, six of which were not reported in the previous Mondul *et al.* retinol GWAS  
169 (Table 1). The absolute effect sizes of these lead SNPs were between 0.066 and 0.172 SD in  
170 circulating retinol per effect allele, as derived from the IVW meta-analysis. We observed  
171 minimal heterogeneity between the two cohorts for these lead SNPs, although heterogeneity  
172 was slightly more marked for rs6601299. Replication was then attempted in the TwinsUK  
173 cohort (N up to 1621, Online methods). Considering the mean association across all timepoints  
174 retinol was measured, as well as the twin pairs separately, 7 out of the 8 lead SNPs in the loci  
175 had effect sizes in the same direction, which was greater than expected by chance alone  
176 (Binomial  $P = 0.035$ , Supplementary Table 3).

177

178 We also investigated the effect of retinol associated lead SNPs on factors associated with  
179 dietary intake of retinoids. Using a GWAS of retinol intake derived from a self-reported 24-  
180 hour dietary recall in the UK Biobank (UKBB, N = 62,991)<sup>26</sup>, we found no evidence to suggest  
181 that any of the effect of these genetic signals on circulating retinol is mediated through  
182 influencing dietary intake behaviours (Supplementary Table 4, Supplementary Figure 3).

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

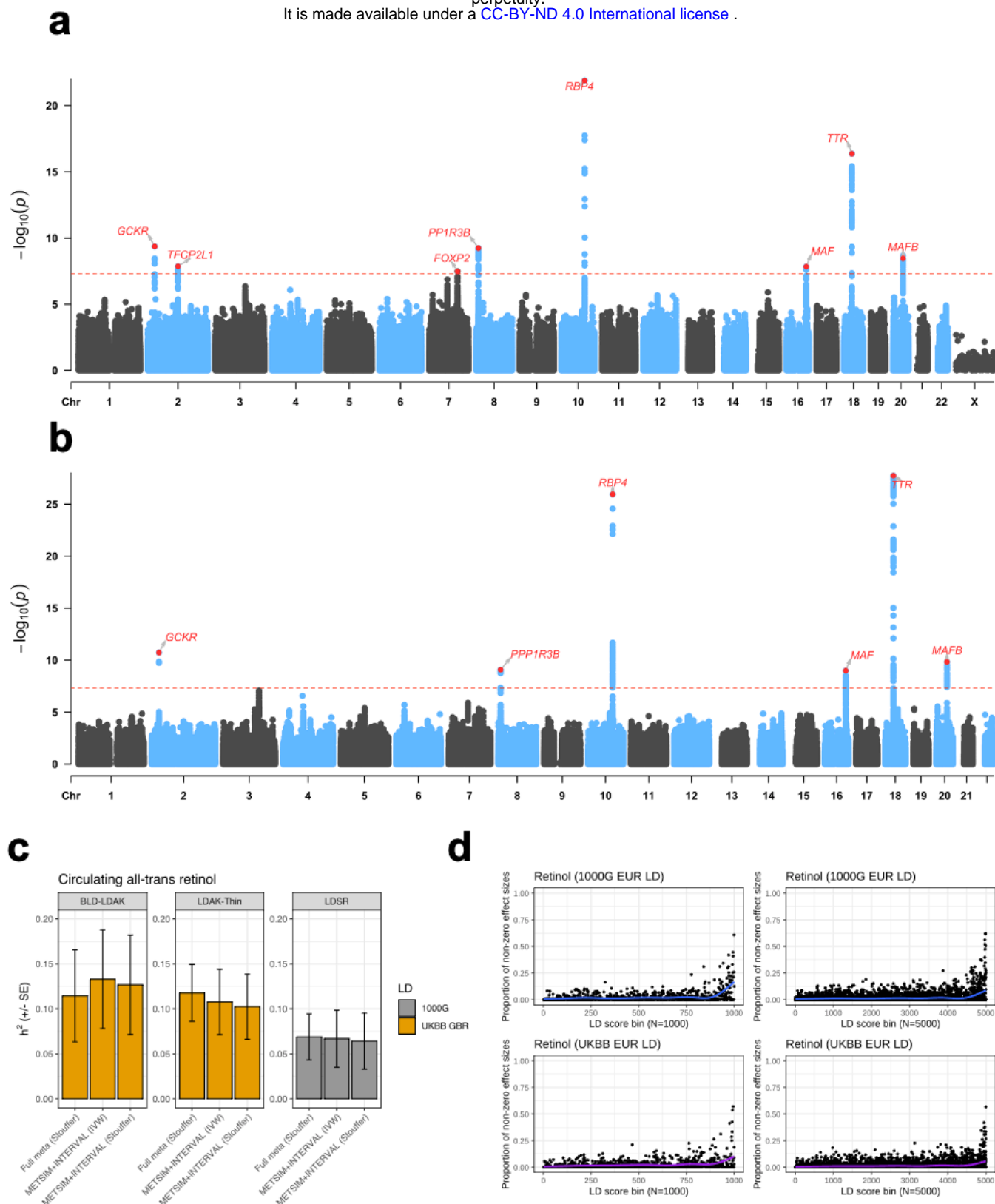
200

201 **Table 1. Genome-wide significant common loci associated with retinol**  
202 **(METSIM+INTERVAL)**

Lead SNP (Stouffer)	Locus	EA/ NEA	EAF (NFE)	EAF (FIN)	Beta (IVW)	Z-score (Stouffer)	$P_{GWAS}$	$P_{Het}$
rs1260326	2:27598097- 27752871	T/C	0.409	0.358	0.071	6.242	$4.32e^{-10}$	0.023
rs34898035	2:122078406- 122084285	A/G	0.042	0.027	-0.172	-5.677	$1.37e^{-8}$	0.190
rs11762406	7:114014488- 114286611	A/C	0.092	0.085	-0.107	-5.526	$3.28e^{-8}$	0.104
rs6601299	8:9167797- 9224907	T/C	0.17	0.10	-0.105	-6.197	$5.76e^{-10}$	0.004
<b>rs10882283</b>	<b>10:95295876- 95360964</b>	A/C	<b>0.622</b>	<b>0.664</b>	<b>0.109</b>	<b>9.789</b>	<b><math>1.26e^{-22}</math></b>	<b>0.192</b>
rs12149203	16:79696939- 79756197	C/G	0.708	0.726	0.069	5.668	$1.45e^{-8}$	0.841
<b>rs1667226</b>	<b>18:29134171- 29190174</b>	A/T	<b>0.481</b>	<b>0.507</b>	<b>-0.092</b>	<b>-8.405</b>	<b><math>4.27e^{-17}</math></b>	<b>0.06</b>
rs6029188	20:39142516- 39234223	A/G	0.637	0.662	-0.066	-5.908	$3.47e^{-9}$	0.147

203 Lead SNP based on statistical significance from sample size weighted (Stouffer) meta-analysis.  
204 Loci boundaries as defined by FUMA (hg19 coordinates). EA = effect allele, that is, allele to  
205 which the effect size relates, NEA = non-effect allele. The effect allele frequency (EAF) is given  
206 from gnomAD v2.1.1 for non-Finnish Europeans (NFE) and Finnish Europeans (FIN). The effect  
207 size in standard deviation units (beta) is denoted for each lead SNP from the IVW meta-analysis,  
208 as well as the Z-score from the Stouffer meta-analysis. The  $P_{GWAS}$  (statistical significance of  
209 association) and  $P_{Het}$  (Heterogeneity between input effect sizes  $P$  value derived from Cochran's  
210  $Q$ ) are also from the Stouffer meta-analysis. Bolded loci are known retinol signals.

211  
212  
213  
214  
215  
216



217 **Figure 1. Common variant influences on circulating serum retinol.** (a) Manhattan plot of the  
 218 meta-analysis of common variants shared between the INTERVAL and METSIM cohorts  
 219 (Stouffer's sample size weighted meta-analysis). Variant-wise  $-\log_{10} P$ -values for association  
 220 are plotted, with the dotted red line denoting genome-wide significance. The closest genic  
 221 transcription start site is labelled for each lead SNP. (b) Manhattan plot, as above, for the  
 222 larger sample size meta-analysis that includes the ATBC and PLCO cohorts, but with fewer



223 *variants available for meta-analysis (c) Estimates of SNP heritability of retinol ( $h^2$ ), with the*  
224 *error bars denoting the standard errors of the estimates. The first two panels denote estimates*  
225 *using the BLD-LDAK model and the LDAK-thin model, respectively, both using LD tagging*  
226 *files from the Great British ancestry participants in the UK Biobank. The last panel estimates*  
227 *heritability using the LDSR model with LD from the 1000 genomes European participants.*  
228 *Estimates were for the METSIM+INTERVAL meta-analyses (Stouffer's and IVW), as well as*  
229 *the larger meta-analysis including ATBC/PLCO. (d) Empirical Bayes' estimation of non-null*  
230 *effects on retinol genome-wide, stratified by bins of ascendingly sorted LD score by magnitude.*  
231 *The LD score bins were different for each panel – 1000 bins, 1000 genomes European LD*  
232 *scores (top left); 5000 bins, 1000 genomes European LD scores (top right); 1000 bins, UKBB*  
233 *white British LD scores (bottom left); 5000 bins, white British LD scores (bottom right). Each*  
234 *point represents the proportion of non-null effect sizes for that bin, with the trendline estimated*  
235 *using a generalised additive model for the relationship between the LD score bin and the*  
236 *proportion of non-null effects.*

237

238 In the larger sample-size meta-analysis with fewer variants available across all input datasets  
239 (METSIM+INTERVAL+ATBC+PLCO), six of the eight genome-wide significant loci from the  
240 smaller meta-analysis were available to test for association. Five of the six loci available  
241 became more statistically significant in this larger meta-analysis, whilst the chromosome 8  
242 locus obtained a very similar level of statistical significance in both meta-analyses (Table 2).  
243 It should be noted that there was relatively large heterogeneity for the lead SNP at the *TTR*  
244 locus on chromosome 18 locus between cohorts. This was due to a more significant association  
245 in the ATBC+PLCO GWAS, although this region was still very strongly associated in  
246 METSIM and INTERVAL in the same direction.

247

248

249

250

251

252

253

254

255

256 **Table 2. Genome-wide significant common loci associated with retinol in the larger**  
 257 **meta-analysis (METSIM+INTERVAL+ATBC+PLCO)**

Lead SNP (Stouffer)	Locus	EA/NEA	Z-score (METSIM+INTERVAL)	Z-score (METSIM+INTERVAL+PLCO+ATBC)	$P_{GWAS}$	$P_{Het}$
rs1260326	2:27598097-27752871	T/C	6.242	6.714	$1.90e^{-11}$	0.06
rs6601299	8:9167797-9224907	T/C	-6.197	-6.137	$8.41e^{-10}$	0.004
<b>rs11187547</b>	<b>10:95279771-95360964</b>	<b>A/G</b>	<b>8.769</b>	<b>10.691</b>	<b><math>1.12e^{-26}</math></b>	<b>0.11</b>
rs11865979	16:79696939-79756197	T/C	5.575	6.103	$1.04e^{-9}$	0.887
<b>rs4799581</b>	<b>18:29068068-29230411</b>	<b>T/C</b>	<b>-7.998</b>	<b>-11.065</b>	<b><math>1.85e^{-28}</math></b>	<b><math>9.72e^{-6}</math></b>
rs6029188	20:39152458-39234223	A/G	-5.908	-6.407	$1.48e^{-10}$	0.299

258 Lead SNP based on statistical significance from sample size weighted (Stouffer) meta-analysis for  
 259 SNPs available in this extended analysis. Loci boundaries as defined by FUMA (hg19 coordinates).  
 260 EA = effect allele, that is, allele to which the effect size relates, NEA = non-effect allele. The Z  
 261 score is given for the full extended meta-analysis as well as the smaller METSIM+INTERVAL  
 262 analysis in terms of sample size. The  $P_{GWAS}$  (statistical significance of association) and  $P_{Het}$   
 263 (Heterogeneity between input effect sizes  $P$  value derived from Cochran's  $Q$ ) are also from the  
 264 Stouffer meta-analysis from the extended meta-analysis. Bolded loci are known retinol signals.

265  
 266 We also estimated rare variant (frequency < 1%) effects on circulating retinol using variants  
 267 available in the METSIM+INTERVAL meta-analysis. Despite relatively low power to detect  
 268 rare variant association, we identified a genome-wide significant rare variant signal on  
 269 chromosome five (chr5:86765041:T:C, dbSNP ID: rs138675130) associated with reduced  
 270 circulating plasma retinol per C allele (-0.441 SD,  $SE = 0.0709$ ,  $P = 6.37 \times 10^{-9}$ ) with no  
 271 significant heterogeneity between the contributing studies. The frequency of this C allele in  
 272 Europeans (gnomAD v.3.1.2) ranges from 0.5% in non-Finnish Europeans to 0.8% in Finns,

273 and whilst it is marginally rarer in Africans and South Asians, it is entirely absent in the East  
274 Asian and Middle Eastern populations in that database. The variant is intergenic and in  
275 FinnGen release 8, the C allele was associated at phenome-wide significance ( $P < 1 \times 10^{-5}$ )  
276 with increased odds of benign neoplasm of the eye and adnexa, as well as whooping cough.  
277 The closest canonical transcription start site to this variant is that of *COX7C*, which encodes a  
278 subunit of a terminal component of the mitochondrial respiratory chain. We also uncovered  
279 several suggestively significant rare variant signals ( $P < 1 \times 10^{-5}$ ), including three non-  
280 synonymous variants in the genes *FREM2*, *NAXD*, and *CHDI1* (Supplementary Table 5). Of  
281 these, the rare *NAXD* non-synonymous allele suggestively associated with lower retinol  
282 (rs3742192) had some *in silico* evidence to suggest deleteriousness (Supplementary Table 5),  
283 although it is classed as benign in ClinVar. Finally, to boost power, we also statistically  
284 aggregated rare variants to genes (Online Methods). There were no significant retinol genes  
285 after Bonferroni correction of these gene level association results ( $P < 3.02 \times 10^{-7}$ ); however,  
286 there were two novel suggestively associated genes ( $P < 3.02 \times 10^{-5}$ ), *GALM* and *ZDHHC18*  
287 (Supplementary Table 6).

288

### 289 **Prioritisation of retinol-associated genes reveals novel mechanisms influencing** 290 **circulating retinol**

291 In common variant loci, prioritising causal genes can be difficult due to confounding factors  
292 like linkage disequilibrium. For circulating retinol, we employed a multi-faceted approach to  
293 prioritise genes that are confidently associated. Firstly, we sought to interrogate the eight  
294 genome-wide significant loci uncovered in the main discovery meta-analysis  
295 (*METSIM+INTERVAL*) to uncover plausible causal genes. This was achieved by adapting an  
296 integrative pipeline developed in previous work that considers annotation, probabilistic  
297 finemapping, integrative scoring, and *in silico* prediction (Online Methods, Supplementary  
298 Table 7). We describe these results further for each locus in the supplementary text but  
299 summarise the prioritisation evidence forthwith. There was quite consistent evidence in four  
300 loci for a likely causal gene (*GCKR*, *FOXP2*, *RBP4*, and *TTR*). *RBP4* (chromosome 10 locus)  
301 and *TTR* (chromosome 18 locus) were previously associated with retinol in the only other  
302 dedicated GWAS of this trait and form the complex that transports retinol in serum<sup>22</sup>, thereby,  
303 having a direct biological link to retinol abundance. *GCKR*, a gene encoding a protein that  
304 binds to and regulates the key metabolic enzyme glucokinase, was confidently the causal gene  
305 for the locus on chromosome 2 with the rs1260326 lead SNP. This gene is known to have a large

306 and varied metabolic association profile due to the role of glucokinase in glycaemic and lipid  
307 related processes, amongst others<sup>27-29</sup>. Interestingly, the lead SNP, and most likely causal  
308 variant derived from probabilistic finemapping, was a common missense allele (rs1260326),  
309 whereby the retinol increasing *T* allele corresponds to a substitution of leucine for proline in  
310 the GCKR protein. Previous experimental work suggests that this variant impacts GCKR  
311 affinity for glucokinase<sup>27,30</sup>.

312

313 The transcription factor gene *FOXP2* was also strongly supported by multiple lines of evidence  
314 as a causal gene for the locus on chromosome seven. The role of this gene in the brain and in  
315 relation to neurological phenotypes like language has been extensively studied<sup>31</sup>, but less so in  
316 the periphery despite relatively high expression across many different systemic tissues.  
317 Therefore, we analysed RNA-sequencing data of *FOXP2* overexpression in a cell-line not  
318 derived from the central nervous system (human osteosarcoma epithelial cell line) and revealed  
319 the transcriptional correlates of *FOXP2* overexpression were enriched for a broad range of  
320 pathways related to factors including extracellular matrix biology, glycosylation, and  
321 interleukin signalling, amongst many others (Supplementary Tables 8-11, Online Methods).

322

323 The remaining four loci exhibited less clear evidence of which gene to prioritise, although all  
324 point to potentially interesting functional mechanisms. On chromosome eight, there is some  
325 evidence to support *PPP1R3B* as a gene that influences retinol, which encodes a catalytic  
326 subunit of the phosphatase PP1 that is implicated in relevant metabolic processes like glycogen  
327 synthesis<sup>32</sup>. However, other lines of evidence point to a role of long-noncoding RNA in this  
328 locus. The loci on chromosomes 16 and 20 are noteworthy as the closest transcription start sites  
329 to the respective lead SNPs are two transcription factors (TF) from the *Maf* family (*MAF* and  
330 *MAFB*). Interestingly, *MAFB* has been shown to regulate both *TTR* and *RBP4* expression in  
331 various tissue contexts from human or murine studies<sup>33,34</sup>. As only some lines of evidence  
332 support these two TF, further functional characterisation of these two loci is warranted.  
333 Interestingly, in the locus on chromosome 16, two of the other genes with some evidence for a  
334 retinol-related function (*MAFTRR* and *LINC01229*) have been shown experimentally to  
335 regulate *MAF* expression and are also associated with other biochemical traits like urate, further  
336 supporting the role of the *Maf* family on retinol abundance in serum<sup>35</sup>. Finally, the remaining  
337 locus on chromosome 2 (2:122078406-122084285) had the least interpretable functional  
338 prioritisation results. The closest transcription start site to the lead SNP was another TF  
339 (*TFCP2LI*) that has broad physiological roles including in the kidney<sup>36</sup>.

340 We then sought to expand our scope for gene discovery beyond genome-wide significant  
341 retinol-associated loci through further integration of genetics with transcriptomics and  
342 proteomics (Online Methods, Supplementary Tables 12-13). Firstly, we leveraged multivariate  
343 models of genetically regulated expression (GR<sub>E</sub>X) to perform a transcriptome-wide (liver,  
344 whole blood, adipose, small intestine, pancreas, and breast mammary tissue) and proteome-  
345 wide (plasma) association study (TWAS/PWAS) of circulating retinol. Tissues for the TWAS  
346 were selected based on prior knowledge of retinol biology and the results of the partitioned  
347 heritability analyses (Online Methods), whilst plasma was the only relevant tissue available for  
348 PWAS. After applying multiple-testing correction to the TWAS and PWAS individually (FDR  
349 < 0.05) and testing whether there was a shared causal variant via colocalisation [Posterior  
350 probability (*PP*) of a shared causal variant (H<sub>4</sub>),  $PP_{H_4} > 0.8$ ], we identified strong evidence of  
351 four additional retinol-associated genes outside of genome-wide significant loci (at least +/- 1  
352 megabase away). These were as follows: *MLXIPL*, which binds to carbohydrate response  
353 elements to regulate triglycerides<sup>37-39</sup>; *GSK3B*, a gene that encodes a member of the glycogen  
354 synthase kinase family involved in metabolism and glycaemic homeostasis<sup>40</sup>; the tankyrase  
355 gene (*TNKS*) implicated in processes like Wnt signalling<sup>41</sup>; and *INHBC*, part of the inhibin  
356 family of proteins with important endocrine functionality<sup>42</sup>. Genetically predicted mRNA  
357 expression of *MLXIPL* in adipose, pancreas, and breast mammary tissue was inversely  
358 associated with circulating retinol levels. Conversely, TWAS analyses revealed that genetically  
359 predicted expression of *GSK3B* and *TNKS* was positively associated with circulating retinol  
360 levels. Finally, genetically predicted plasma protein expression of *INHBC* showed a positive  
361 correlation with retinol levels ( $Z_{PWAS} = 4.72$ ). We then used a more conservative approach  
362 whereby finemapped variants that influence protein expression (pQTLs) were used as  
363 instrumental variables (IV) to estimate the causal effect of plasma proteins on retinol using  
364 Mendelian randomisation (MR, Online Methods, Supplementary Table 14). We applied the  
365 same filters to the results (FDR < 0.05 and  $PP_{H_4} > 0.8$ ). These analyses further supported that  
366 upregulated *INHBC* likely increases serum retinol, with each SD increase in plasma protein  
367 expression associated with a small but highly statistically significant [0.05 per SD in  
368 expression, 95% CI: 0.03, 0.07] impact on circulating retinol. In line with expectation, pQTL-  
369 MR, and subsequent colocalisation, further genetically validates that elevated *RBP4* protein  
370 abundance correspondingly increases serum retinol with somewhat large effect (0.6 [95% CI:  
371 0.48, 0.72] SD in retinol per SD in plasma *RBP4* expression). There was also evidence that  
372 *RBP4* protein expression and retinol colocalise under the hypothesis of a single causal variant  
373 ( $PP_{H_4} = 1$ ). Considering the eight genes prioritised in this and the previous section (*RBP4*,

374 *GCKR, FOXP2, TTR, MLXIPL, GSK3B, TNKS, INHBC*), we found that these genes exhibited  
375 upregulated expression in the liver ( $P_{\text{Adjusted}} < 0.05$ ) upon analysing data from 54 GTEx tissues.  
376 This further consolidates the salience of hepatic processing to genetic influences on circulating  
377 retinol abundance. Pathway analyses of these genes demonstrated that they were enriched  
378 amongst factors involved in carbohydrate metabolism (Supplementary Table 15).

379

### 380 **Wide-ranging evidence of causal effects of retinol across the human clinical phenome**

381 The role of retinol in human health and disease has been of long-standing interest. However,  
382 most evidence has been observational, limiting the ability for causal inference. Moreover,  
383 randomised controlled trials (RCT) of interventions like retinol supplementation and synthetic  
384 retinoids have only been performed for a small fraction of the traits implicated through  
385 observational studies. We sought to increase our understanding of causal effects of retinol on  
386 human health by leveraging genetic variants associated with retinol uncovered in this study as  
387 IVs. Given certain assumptions are met, these genetic proxies of retinol can be utilised to  
388 estimate causal effects of circulating retinol at scale using Mendelian randomisation (MR)  
389 (Online Methods). Firstly, we utilised a single IV in *RBP4* (rs10882283), as this gene has a  
390 clear and unambiguous association with circulating retinol levels, and therefore, is less likely  
391 to be prone to horizontal pleiotropy than other retinol-associated loci. We do caution, however,  
392 that *RBP4* does exert some other functionality that may not be directly related to retinol  
393 transport<sup>43</sup>, although this gene is still likely the best available single IV associated with  
394 circulating retinol at genome-wide significance. We utilised this *RBP4* IV  
395 (*METSIM+INTERVAL* meta-analysis effect size) to estimate the effect of retinol on over  
396 19,500 outcomes in the IEUGWAS database (IEUGWASdb, Online Methods, Supplementary  
397 Table 16). After multiple-testing correction ( $FDR < 0.01$ ), retinol was found to putatively exert  
398 causal effects on outcomes including several lipid traits, leukocyte counts (total leukocytes and  
399 neutrophils), reticulocytes, optic disc area, and resting-state connectivity of a functional MRI  
400 (fMRI) derived network edge. For instance, each SD in circulating retinol was associated with  
401 a -0.1 SD [95% CI: -0.14, -0.06] decrease in leukocyte count, whilst this unit retinol increase  
402 was estimated to increase optic disc area by 0.22 SD [95% CI: 0.12, 0.32]. We further  
403 interrogated a representative subset of these associations using colocalisation to test if the  
404 signals are driven by a shared causal variant in *RBP4* (Supplementary Table 17). Colocalisation  
405 strongly supported that the effect of circulating retinol on leukocyte count, optic disc area and  
406 the edge in the fMRI network arises from a shared causal variant in *RBP4*. In contrast, the  
407 effect of retinol on lipids through *RBP4* was shown to likely arise from linkage due to the

408 proximal gene *FFAR4* that encodes a free-fatty acid receptor. Therefore, the lipid findings most  
409 likely do not represent a causal impact of circulating retinol, at least through *RBP4*, but rather  
410 the influence of *FFAR4*.

411

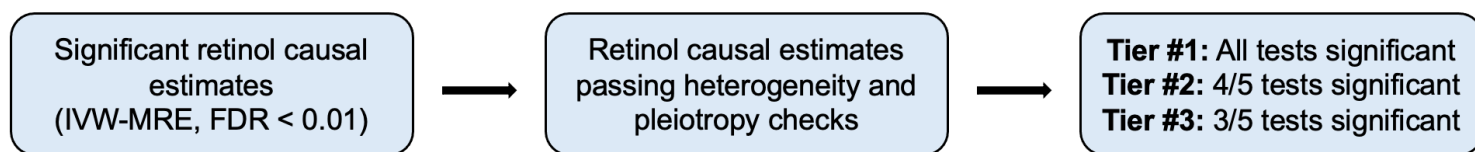
412 To boost power for discovery of causal retinol effects, we then utilised all independent (LD  $r^2$   
413  $< 0.001$ ) genome-wide significant SNPs as IVs (Online Methods). Firstly, we used the inverse-  
414 variance weighted estimator with multiplicative random effects (IVW-MRE) to estimate the  
415 effect of retinol on IEUGWASdb outcomes for which at least six of the IVs were available ( $>$   
416 17,000 outcome phenotypes). There was moderate positive correlation between the IVW-MRE  
417 and single *RBP4* estimates across all traits tested ( $r = 0.43$ , Supplementary Figure 4). Whilst  
418 well powered, the IVW-MRE assumes all IVs are valid, which is unlikely in practice.  
419 Therefore, we developed a pipeline to prioritise the most confident causal relationships that  
420 survive multiple-testing correction considering the IVW-MRE estimates (FDR  $< 0.01$ , Online  
421 Methods, Figure 2A, Supplementary Tables 18-20). This was comprised of three tiers, with  
422 Tier #1 being the highest level of evidence. Retinol/outcome pairings that were assigned a tier  
423 had to exhibit no significant heterogeneity between IV-outcome effects, a non-significant MR-  
424 Egger intercept (which screens for unbalanced pleiotropy), and not be driven by a single IV.  
425 Four additional MR methods with differing assumptions were then applied in this study (Online  
426 Methods). From the trait pairings that passed the above heterogeneity and pleiotropy filters,  
427 Tier #1 traits were those for which all five methods were statistically significant, whilst Tier  
428 #2 traits had 4/5 significant methods, and Tier #3 3/5 methods significant. There were no Tier  
429 #1 retinol/outcome trait pairings, but several Tier #2 and Tier #3 trait pairings (Figure 2B-C),  
430 with all of them directionally consistent with the estimates from the single *RBP4* IV, supporting  
431 their validity. Broadly, we found that retinol increased body fat related measures, resting-state  
432 fMRI connectivity of several network edges, as well as food consumption phenotypes related  
433 to carbohydrates. Retinol also exhibited evidence of a relationship with the cortical thickness  
434 and surface area of several brain regions, as well as microbiome composition and keratometry  
435 measurements. These results were dominated by continuous traits, for which we are better  
436 powered, however, there were some binary traits assigned Tier #2 or Tier #3 evidence.  
437 Specifically, retinol was associated with decreased odds of tuberculosis sequelae and  
438 coxarthrosis (arthrosis of the hip), whilst it was associated with increased odds of non-specific  
439 skin eruptions, adverse asthma/COPD medication effects, and dental problems. We caution  
440 that all these inferred associations require further investigation, and should be treated with  
441 requisite caution, as reviewed elsewhere<sup>19,44,45</sup>. A consideration of this approach that leverages

442 multiplicative random effects with relatively few IVs ( $< 10$ ), is the potential influence of  
443 residual standard errors below 1 on the estimation of the IVW standard errors. We see some  
444 evidence of this impact on the IVW-MRE estimates for Tier #2 and Tier #3 traits – as these  
445 traits exhibit no significant heterogeneity between IV estimates. Specifically, whilst all fixed  
446 effect IVW results are still highly statistically significant for these traits, they have larger  
447 standard errors than the IVW-MRE, indicative of residual standard errors  $< 1$ . This is a function  
448 of the MRE not scaling the standard error of the IVW by the model's residual standard error  
449 like in the fixed effects model. We discuss this further in the supplementary text and in  
450 supplementary figure 5. However, these issues only impact the  $P$ -value of the IVW-MRE  
451 relative to that of the fixed effects, and Tier #2/Tier #3 traits still exhibit non-zero evidence  
452 across multiple methods and no indication of a single IV driving the association. Moreover,  
453 using instead a fixed effects IVW estimator as the primary test for trait pairings with no  
454 heterogeneity (Cochran's  $Q$   $P < 0.05$ ) yields similar outcomes being prioritised as most  
455 statistically significant after FDR correction (Supplementary Text). It also important to  
456 consider when interpreting these estimates that MR approaches are only interpretable under the  
457 assumptions they make, and as a result, any potential causal relationship reported here requires  
458 further validation, ideally using a randomised control trial design.

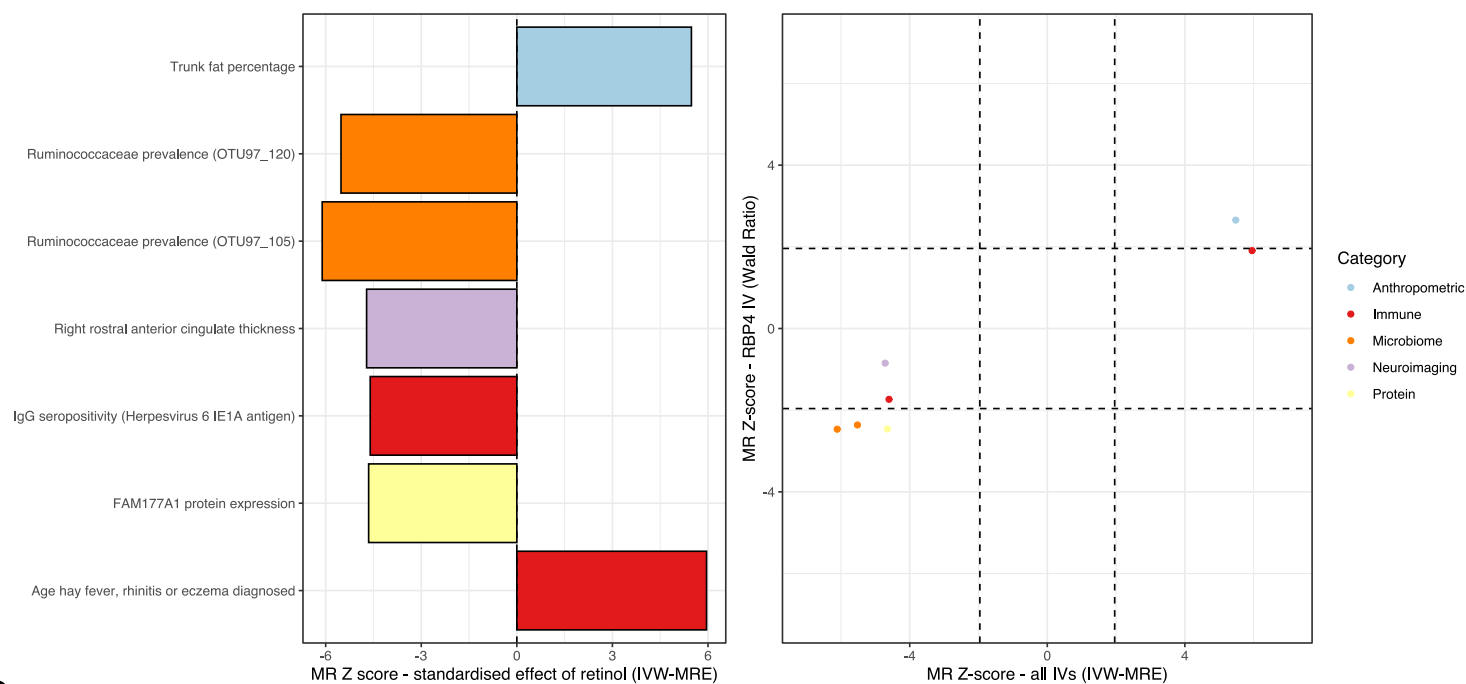
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475



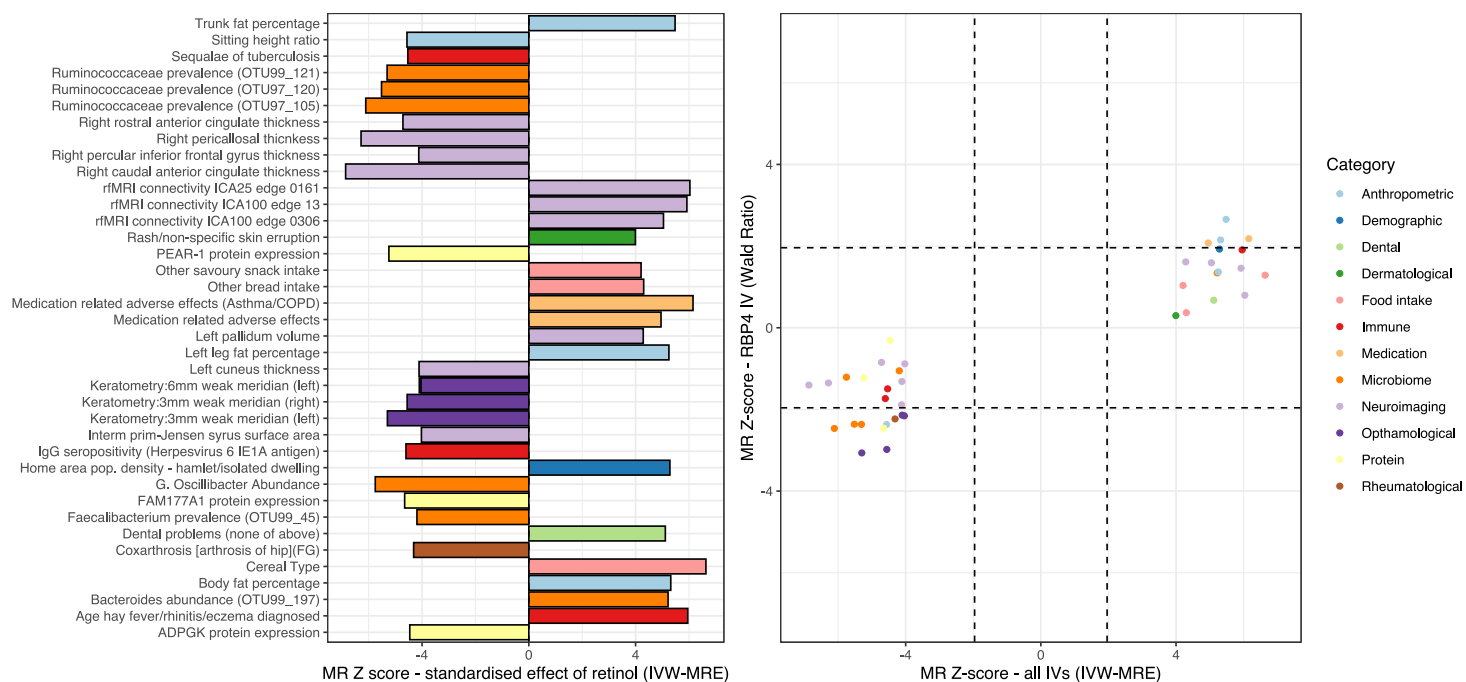
**a** 476



**b** 481



**c**



508 **Figure 2. Estimated causal effects of circulating retinol across the human clinical phenome.**

509 **(a) Prioritisation pipeline overview for retinol causal estimates [inverse-variance weighted**

510 *estimator with multiplicative random effects (IVW-MRE)] that survive multiple-testing*  
511 *correction ( $FDR < 0.01$ ). These estimates are subjected to tests for heterogeneity and*  
512 *pleiotropy (Online Methods), with a tier then assigned based on how many of the five*  
513 *Mendelian randomisation (MR) methods applied are at least nominally statistically significant.*  
514 *In panels (b) and (c), the left-hand plot denotes the Z-score ( $\beta/SE$ ) from the MR IVW-MRE*  
515 *estimates. Positive Z scores denote a positive IVW-MRE estimate of the effect of circulating*  
516 *retinol on that trait, and vice versa. The traits are coloured by their broad phenotypic category.*  
517 *The right-hand plot visualises the Z score using the IVW-MRE model verses that of the MR*  
518 *estimate using the RBP4 IV alone (Wald Ratio). The dotted lines approximately represent*  
519 *nominal statistical significance ( $P < 0.05$ ). In panel (b), just tier #2 traits are plotted (Online*  
520 *Methods), whilst panel (c) plots both tier #2 and tier #3 traits.*

521

522 We hypothesised that the putative effect of retinol on body fat could be one explanation for its  
523 relationship in this study to brain phenotypes beyond a direct effect of retinoid signalling,  
524 particularly given that obesity and adiposity have been linked with MRI related indices<sup>46</sup>.  
525 However, using IVs for body fat percentage, we did not find any strong evidence that it is  
526 causally related to any of the retinol-associated brain regions (Supplementary Table 21),  
527 suggesting a direct effect of retinol on these regions/networks or a relationship induced through  
528 some other unobserved confounder. Reverse causality for these Tier #2 and Tier #3 exposures  
529 was then also considered, although for binary traits this should be treated purely as a test of the  
530 null hypothesis given the difficulties in using binary traits as IVs<sup>47</sup>. There was no strong  
531 evidence for reverse causality of any of these traits (Supplementary Table 22). One exception  
532 to this was in relation to expression of the protein PEAR1, for which there was very nominal  
533 evidence of bidirectional effects.

534

535 A limitation of the above phenome-wide analyses is that the multiple-testing burden that arises  
536 from the inclusion of over 17,000 traits may obscure retinol effects on binary disease  
537 phenotypes, as these are usually less powered than GWAS of continuous traits. To overcome  
538 this, we also applied the above pipeline using all retinol IVs to 1141 binary endpoints with at  
539 least 1000 cases in FinnGen release 8 (not featured in IEUGWASdb), allowing a phenome-  
540 wide analysis of electronic health record derived binary outcomes. There were eight disease  
541 phenotypes that retinol was associated with after multiple testing correction (Figure 3A,  $FDR$   
542  $< 0.01$ ), which increased to 19 with an exploratory  $FDR < 0.1$  threshold (Supplementary Table  
543 23). After applying the above pipeline to these results that considers heterogeneity, pleiotropy,

544 and consistency across MR methods, there were four disease endpoints with Tier #3 evidence  
545 (Figure 3B, Supplementary Table 24). Specifically, retinol was estimated to increase the odds  
546 of congenital malformations of the heart and great arteries, whilst it was protective for type 2  
547 diabetes with coma and inflammatory liver disease. As these traits had tier #3 evidence, there  
548 was some inconsistency in the strength of the results across different MR methods, and  
549 therefore, these relationships should be interpreted cautiously. One of the most active areas of  
550 research in retinol epidemiology is the relationship between retinol and cancer risk<sup>48</sup>. The  
551 estimated effect of circulating retinol on the odds of any malignant neoplasm was not  
552 significantly different than one - OR = 0.97 [95% CI: 0.91, 1.04],  $P = 0.423$  (Supplementary  
553 Figure 6). However, there was some indication of a protective effect of retinol on squamous  
554 non-small cell lung cancer, which approached the threshold for statistical significance after  
555 FDR correction - OR = 0.64 [95% CI: 0.51, 0.80],  $P = 8.19 \times 10^{-5}$ ,  $q = 0.01$ . There was also  
556 some data to support retinol having effects on other respiratory neoplasm endpoints  
557 (Supplementary Figure 6). Given previous observational evidence that retinol is protective for  
558 lung cancer<sup>48</sup>, as well as some data supporting the use of synthetic retinoids like bexarotene in  
559 lung neoplasms<sup>49</sup>, this relationship warrants further exploration.

560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

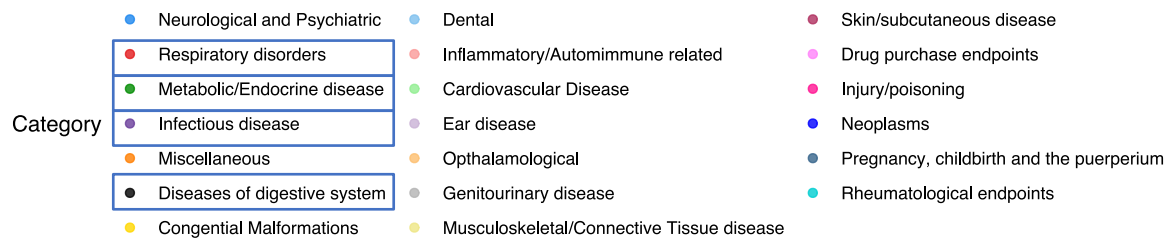
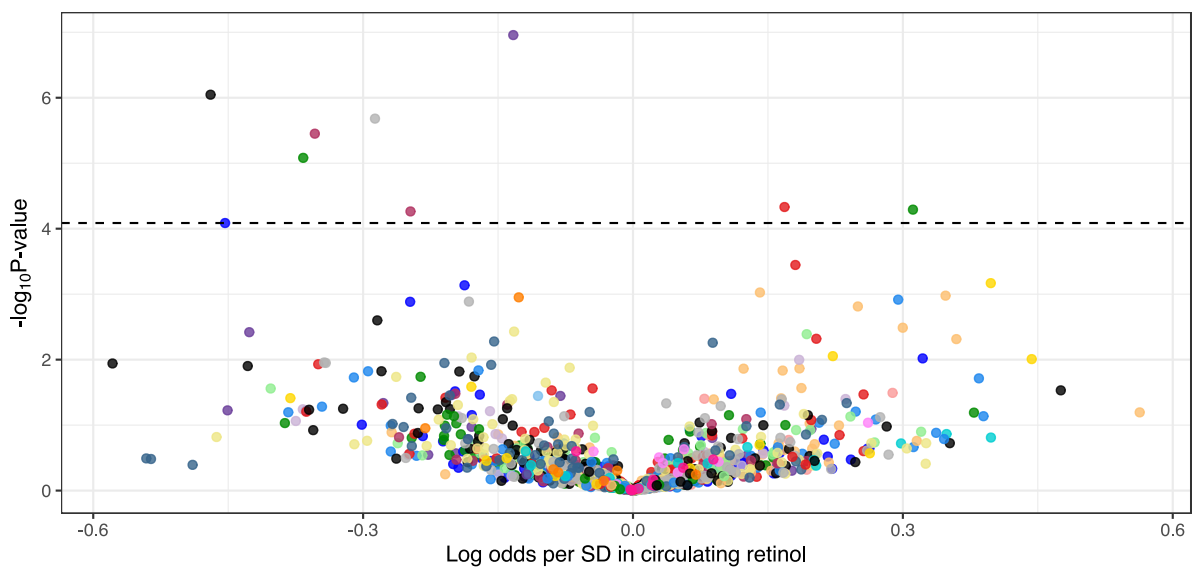
607

608

609

610

611



594

595

596

597

598

599

600

601

602

603

604

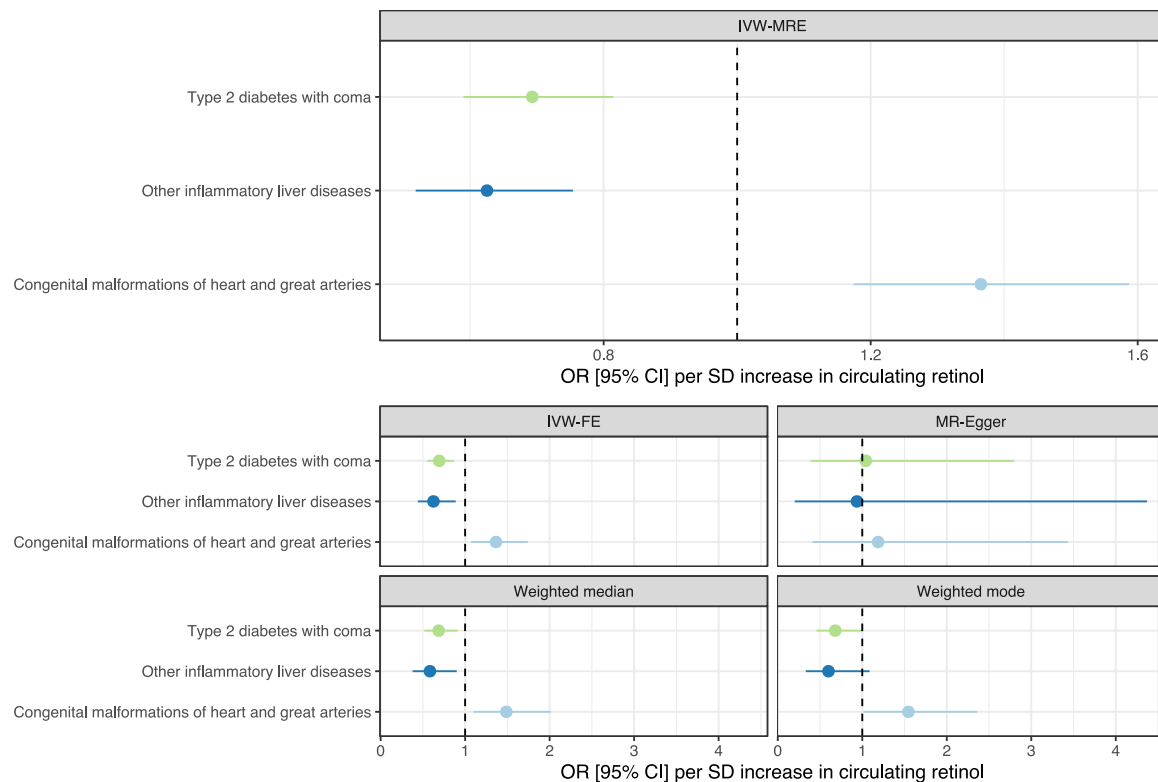
605

606

607

608

609



**Figure 3. Genetic estimates of the relationship between circulating retinol and binary disease endpoints in the Finnish population. (a) Volcano plot denoting the phenome-wide MR**

612 *estimates of retinol (SD units) and the log odds of each binary endpoint in FinnGen release 8.*  
613 *The x-axis denotes the log odds (IVW-MRE), whilst the y-axis denotes the  $-\log_{10}$  P-value of the*  
614 *MR estimate. The dotted horizontal line is approximately equivalent to the magnitude of P-*  
615 *value that is estimated at an FDR of 1%. Categories highlighted by a blue box on the legend*  
616 *indicate that an endpoint in that category survived multiple-testing correction. (b) Retinol*  
617 *effects on binary FinnGen endpoints after multiple-testing correction ( $FDR < 0.01$ ) with the*  
618 *strongest evidence from all phenotypes tested based on heterogeneity, pleiotropy, and*  
619 *consistency across MR estimators (Tier #3, Online Methods). Each panel denotes the effect*  
620 *sizes (odds ratio with 95% confidence intervals per SD in circulating retinol) for the IVW-*  
621 *MRE, the IVW with fixed effects (IVW-FE), the MR-Egger, Weighted median, and Weighted*  
622 *Mode, respectively.*

623

624 We also investigated evidence for bidirectional effects involving these Tier #3 traits that retinol  
625 is genetically predicted to influence. As described above, using binary traits as exposures in  
626 MR should be treated with suitable caution and are often underpowered. The use of either  
627 genome-wide significant or suggestively significant SNPs ( $P < 1 \times 10^{-5}$ ) as IVs did not indicate  
628 evidence of reverse causality of these diseases to retinol (Supplementary Table 25). However,  
629 given some of the statistical limitations of these analyses, such effects warrant further  
630 consideration.

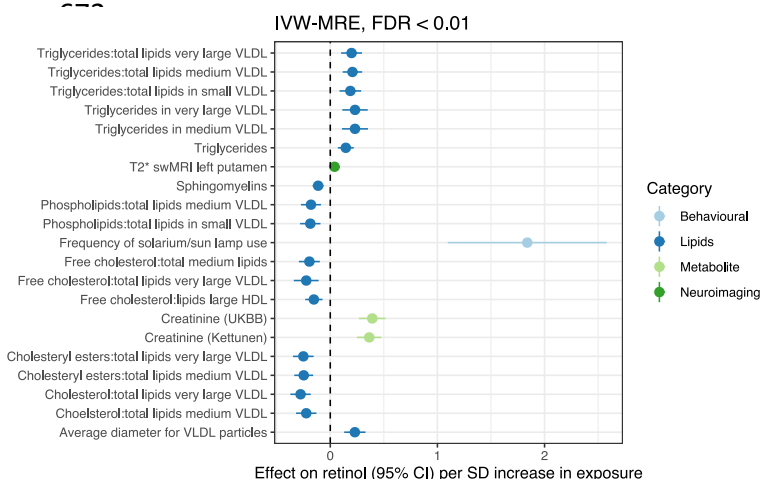
631

### 632 **Genetic evidence that lipids and kidney function influence circulating retinol**

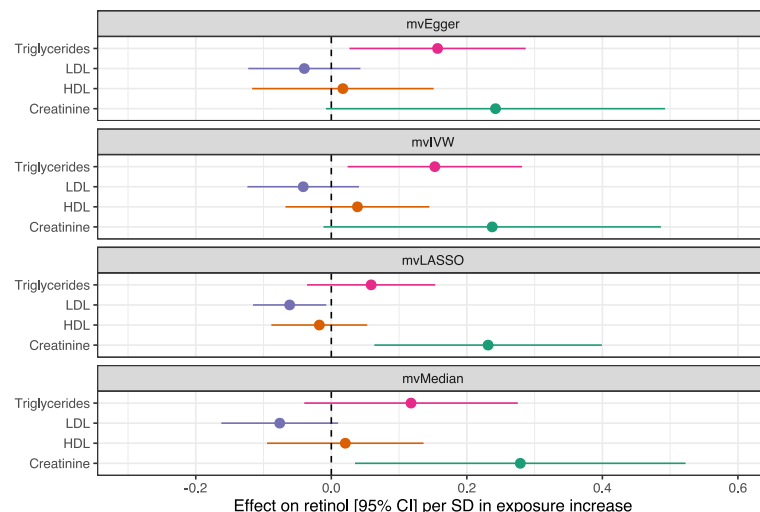
633 It is also clinically valuable to understand exposures and diseases that impact circulating retinol  
634 abundance. To explore this in greater detail, we leveraged retinol as an outcome trait in MR  
635 analyses. We utilised a diverse range of thousands of continuous and ordinal phenotypes from  
636 IEUGWASdb as exposures in a similar pipeline described above (Online Methods). Several  
637 lipid species were demonstrated to putatively influence retinol abundance after multiple-testing  
638 correction ( $FDR < 0.01$ , Figure 4A, Supplementary Table 26); for example, triglycerides were  
639 implicated to increase circulating retinol whilst cholesteryl ester related traits decreased  
640 circulating retinol. We also saw a positive effect of the frequency of solarium and sun lamp use  
641 on retinol, which may arise from behavioural related mechanisms. Furthermore, our findings  
642 suggest that increased retinol levels are associated with a susceptibility-weighted MRI measure  
643 in the left putamen, called T2\*. T2\*, reflecting magnetic susceptibility relative to tissue water,  
644 can be influenced by factors such as iron and calcium content. This association may also be  
645 influenced by behavioural or other pathways that warrant further investigation.

646 After applying the same tiering system used for retinol as an exposure, we observed Tier #1  
 647 evidence strongly supporting a causal effect of creatinine on circulating retinol. This  
 648 relationship is biologically plausible and likely represents an association with kidney  
 649 function<sup>50,51</sup>. Due to the biological complexity of lipid traits, we observed significant  
 650 heterogeneity between IV effects, and as a result, they were not assigned a tier in our analysis.  
 651 However, the effect of triglycerides on increasing circulating retinol levels is consistent with  
 652 established knowledge of retinol biology, as well as this study implicating two genes that are  
 653 mechanistically confirmed to impact triglycerides (*GCKR* and *MLXIPL*). We then leveraged  
 654 the CAUSE model to distinguish causal effects of creatinine on circulating retinol from  
 655 correlated pleiotropy that may arise between these two traits due to the extensive polygenicity  
 656 of creatinine<sup>52</sup>. We found that a model that includes a causal effect of creatinine on retinol was  
 657 more parsimonious than a model of pleiotropy alone ('sharing model') through comparison of  
 658 these models using the Bayesian expected log pointwise posterior density (ELPD) method  
 659 (Supplementary Figure 7,  $\Delta\text{ELPD}_{\text{Sharing vs Causal}} = -4.34$ ,  $P = 8.9 \times 10^{-3}$ ). Given that IVs for  
 660 creatinine could plausibly act through lipid species like triglycerides to influence circulating  
 661 retinol, we then constructed multivariable MR (MVMR) models that estimated the creatinine  
 662 to retinol relationship conditioned on high density lipoprotein (HDL), low density lipoprotein  
 663 (LDL), and triglycerides (Online Methods). While there was some evidence that the effect of  
 664 creatinine on retinol could arise due to triglycerides, there was also evidence to suggest an  
 665 independent effect of both triglycerides and creatinine on increasing circulating retinol,  
 666 depending on the modelling parameters used (Figure 4B, Supplementary Text). In addition to  
 667 investigating the causal effects of exposures on retinol, we also examined the possibility of  
 668 bidirectional effects. We found very weak evidence that retinol has a negative effect on  
 669 creatinine ( $P = 0.023$ ), but there was no significant evidence to suggest bidirectional  
 670 relationships between retinol and the other implicated exposures.

671  
 672  
**a**



**b**



680 **Figure 4. Exploring the causal effects of continuous exposures on circulating retinol. (a)**  
681 *Exposure traits that demonstrated a significant causal estimate (IVW-MRE) on circulating*  
682 *retinol after multiple-testing correction (FDR < 0.01). Traits are coloured relative to their*  
683 *broad phenotypic category. (b) Multivariable MR (MVMR) models investigating the effect of*  
684 *creatinine and major lipid species on circulating retinol. Each panel represents the results*  
685 *from a different MVMR model (each with different underlying assumptions (Online Methods)).*  
686 *The exposure – retinol relationship plotted is conditional on the three other traits in the model.*  
687

688 Finally, we explored pharmacological agents and molecular perturbagens that may influence  
689 circulating retinol (Online Methods). Considering the novel genes prioritised in this study with  
690 an assigned direction of expression (TWAS/PWAS, pQTL MR), it was found that *GSK3B* is a  
691 drug-target known to be inhibited by lithium and related compounds. This may be of clinical  
692 interest as it suggests that lithium, utilised as a therapy in mood disorders, may decrease  
693 circulating retinol via its inhibition of *GSK3B* given that genetically predicted expression of  
694 this gene was positively associated with retinol. We then employed computational signature  
695 mapping to further characterise pharmacological agents related to retinol (Online Methods).  
696 However, these analyses did not yield any compounds for which the *in vitro* transcriptomic  
697 signature significantly matched or opposed genetically predicted expression associated with  
698 retinol after multiple-testing correction (Supplementary Table 27). We then considered  
699 perturbagen signatures aggregated to biological pathways or overall mechanisms of action  
700 (MOA) groups of compounds (Online Methods). After multiple-testing correction (FDR <  
701 0.05), there were 13 gene-set based perturbagen signatures that were significantly similar to  
702 the directionality of genetically predicted expression associated with retinol (Supplementary  
703 Table 28). For example, the expression signature of compounds in the HDAC inhibitor MOA  
704 opposed expression genetically predicted to increase serum retinol. This can be interpreted as  
705 while no single HDAC inhibitor was significantly associated with retinol, there was at least  
706 some evidence for the overall relationship with this MOA. This accords with the suggested  
707 effect of HDAC inhibitors like valproic acid on downregulating expression of *RBP4*<sup>53–55</sup>, a  
708 gene not included in the signature mapping analyses given the large effect of its encoded  
709 protein on retinol.

710  
711  
712

713 **Genetically proxied retinol can identify individuals outside of the normative range of**  
714 **circulating retinol for a given age**

715 We were also interested in evaluating the performance of a genetically proxied index of  
716 circulating retinol, that is, a circulating retinol polygenic score (PGS). The independent  
717 TwinsUK cohort was utilised to tune and evaluate retinol PGS (Online Methods). We used  
718 several methods to evaluate the performance of a retinol PGS in this cohort. Firstly, we test  
719 different retinol PGS configurations in a random selected training subset of the model (70% of  
720 cohort), using a linear mixed model to account for relatedness between the twin pairs. The best  
721 performing retinol PGS configuration in the training subset explained approximately 2.12% of  
722 the phenotypic variance of retinol when applied to the remaining 30% of the cohort (mean  
723 variance explained across three retinol measurement timepoints). We also found similar  
724 performance when applied to each subset of twin pairs (Supplementary Table 29). A limitation  
725 of this approach for using the same cohort for tuning and testing the retinol PGS is that the  
726 estimated effect sizes may not be representative. As a result, we employed an approach to tune  
727 weights for the PGS using the summary statistics alone. This was achieved through leveraging  
728 the principles of probabilistic finemapping to update variant weights by their posterior  
729 probability of association (Online Methods). Due to the modest polygenicity of retinol, this  
730 method upweights a small number of variants. Despite this, these scores were still significantly  
731 associated with circulating retinol in TwinsUK (Supplementary Table 29).

732

733 Like most micronutrients, circulating retinol has been shown previously to have a complex  
734 relationship with age<sup>56</sup>. However, population-level approaches investigating these effects do  
735 not account for inter-individual variability. We hypothesised that normative modelling could  
736 be used to characterise individual patterns of circulating retinol, and to evaluate the  
737 contribution of genetics to these individualised effects. Normative modelling, derived from the  
738 application of growth charts in paediatric medicine, aims to estimate normative reference  
739 ranges of variation in the population (e.g., of circulating retinol) based on age and/or other  
740 relevant variables. Here, we established reference ranges for circulating retinol as a function of  
741 age using generalised-additive models for location, scale, and shape (GAMLSS) frameworks  
742 in the TwinsUK dataset (Online Methods, Supplementary Text). These models were  
743 constructed in the full cohort, as well as in one twin subset only for comparison. Briefly, this  
744 involved identifying the optimal distribution for the GAMLSS model using all samples,  
745 followed by splitting the data into two partitions. One partition was utilised to estimate centiles  
746 (5<sup>th</sup>, 25<sup>th</sup>, 50<sup>th</sup>, 75<sup>th</sup> and 95<sup>th</sup>), and the left-out subjects were benchmarked against this reference



747 chart to determine the position of their individual retinol measurement (Online Methods). This  
748 process was then repeated using the opposite subject allocation. An individual's retinol level  
749 was classified as infra-normal if their measured retinol fell below the 5<sup>th</sup> percentile for their  
750 age, and as supra-normal if their retinol value exceeded the 95<sup>th</sup> percentile for their age.

751

752 We then investigated the extent to which retinol PGS was associated with these individual  
753 profiles of circulating retinol with respect to age (Supplementary Table 30). By way of  
754 example, we report results forthwith from the more conservative modelling approach which  
755 only used one half of the twins. Retinol PGS was at least nominally significantly associated  
756 with supra and infra-normal deviations except for supra-normal deviations at the first (youngest  
757 visit) for which there was only a trend observed (Supplementary Figures 8-9). For example, at  
758 the second visit each SD in retinol PGS was associated with an approximately 70% [95% CI:  
759 23%, 136%] increase in the odds of exhibiting supra-normal retinol levels for a given age  
760 relative to all remaining participants, whilst conversely reducing the odds of displaying infra-  
761 normal levels by approximately 37% [95% CI: 12%, 55%]. These data suggest that genetics is  
762 a non-zero contributor to circulating retinol levels that fall outside the normative range for a  
763 given age. In future, a normative modelling approach could be utilized to examine additional  
764 factors, including dietary intake, as well as the interplay between genetics and other influences,  
765 that might contribute to individualized deviations in retinol levels relative to the population  
766 benchmarks.

767

## 768 **DISCUSSION**

769 We conducted the largest GWAS of circulating retinol to-date, revealing important novel  
770 insights into genetic influences on this trait. The sample sizes in this study facilitated the first  
771 published estimate of SNP heritability for circulating retinol, plausibly between 5-10%.  
772 However, the large standard errors accompanying these estimates reinforces that greater  
773 sample sizes are still needed. Moreover, we were able to uncover confident genetic signals  
774 associated with retinol at genome-wide significance outside of the RBP4:TTR transport  
775 complex. The gene prioritisation pipeline applied both within and beyond genome-wide  
776 significant loci prioritised eight genes with high-confidence for a role in retinol biology. These  
777 genes were highly expressed in the liver and overrepresented amongst biological pathways  
778 including carbohydrate metabolism. The liver is known to be the key organ responsible for  
779 retinol storage and processing<sup>1</sup>, which is represented strongly by the genetic data in this study.  
780 Further, the prioritised genes assigned as overrepresented in the *regulation of carbohydrate*

781 *metabolic process* pathway (*GCKR*, *GSK3B*, and *MLXIPL1*) are all broadly known to be related  
782 to hepatic energy metabolism. As lipids are directly mechanistically linked to retinol  
783 absorption, storage, and delivery<sup>1,3</sup>, it is plausible that the varied metabolic roles of these genes  
784 converge on changes in the abundance of different lipid species. The role of lipids in circulating  
785 retinol abundance is also highlighted by our Mendelian randomisation analyses. However, it is  
786 still likely that glycaemic homeostasis may impact circulating retinol via mechanisms not  
787 directly linked to lipid biology; for example, expression of the insulin-controlled glucose  
788 transporter GLUT4 is postulated to be related to RBP4 protein levels<sup>57</sup>. In summary, our results  
789 suggest that the most identifiable common variant influences on circulating retinol are either  
790 mediated through direct effects on transport or metabolic factors, particularly related to lipids.  
791 We also prioritised genes like *FOXP2* for which a mechanistic relationship to retinol is less  
792 inherently clear. Our analyses of transcriptomic correlates of *FOXP2* supported the immense  
793 biological pleiotropy associated with this transcription factor, reinforcing its significance  
794 outside of its traditionally conceived association in the literature with neurological phenotypes  
795 like language. Work is now needed to disentangle the mechanisms which specifically underlie  
796 this relationship between *FOXP2* and circulating retinol that were infer from these genetic  
797 findings.

798  
799 Our study also represents a significant advancement as it is the first to perform a high-  
800 throughput, hypothesis free, analysis investigating the potential causal effects of retinol across  
801 a wide range of human clinical phenotypes using Mendelian randomisation. This work  
802 recapitulated known influences of retinol on ophthalmological measures<sup>58</sup>, the innate and  
803 adaptive immune response<sup>59</sup>, and congenital heart malformations<sup>60</sup>. However, we also  
804 uncovered some more novel relationships that may be of direct clinical relevance. We highlight  
805 forthwith the example of circulating retinol being genetically predicted to impact the thickness  
806 and surface area of several brain regions, as well as indices of brain connectivity. Retinoic acid,  
807 a downstream metabolite of retinol, is considered one of the most intrinsic central nervous  
808 system signalling molecules, particularly as it exerts control over processes like neuronal  
809 differentiation *in utero* and adult neurogenesis, as reviewed elsewhere<sup>2</sup>. It is, therefore, logical  
810 that retinol would plausibly influence brain structure and connectivity throughout the lifespan.  
811 However, the regions implicated in this study require further examination with respect to their  
812 clinical significance. By way of example, we associated increased circulating retinol with a  
813 reduction in thickness in the right rostral anterior cingulate cortex. This cortical region has been  
814 identified by a large international mega-analysis from the ENIGMA consortium to exhibit

815 increased thickness in individuals with the neuropsychiatric disorder schizophrenia compared  
816 to controls<sup>61</sup>, suggesting a potential protective effect of retinol in this region with respect to  
817 schizophrenia. This accords with previous evidence linking retinoids to schizophrenia<sup>2,62</sup>. It is  
818 known clinically that both retinol deficiency and toxicity can have harmful neurological  
819 effects, highlighting the complexity of the relationship of retinol to the brain throughout the  
820 human lifespan. This complexity is also seen with synthetic retinoids. For instance, isotretinoin  
821 (13-*cis* retinoic acid), indicated for conditions like acne, has been shown to have opposing  
822 effects on adult neurogenesis relative to all-*trans* retinoic acid and putatively increases the risk  
823 of suicide<sup>63,64</sup>, although evidence for this association is mixed<sup>65</sup>. Conversely, another synthetic  
824 retinoid, bexarotene, with different receptor affinities to isotretinoin, has demonstrated some  
825 promise as a potential adjuvant to antipsychotics in schizophrenia<sup>66</sup>. Future work should  
826 attempt to understand these relationships with greater fidelity by investigating genetic  
827 influences on other retinoids beyond retinol, as well as how tissue-specific abundance can  
828 differ from what circulates in serum<sup>67</sup>. Emerging methods for non-linear Mendelian  
829 randomisation would also be useful in this context given that retinol often exerts dose  
830 dependent effects<sup>68</sup>. The causal estimates generated in this study also need to be treated with  
831 appropriate caution due to the limitations of MR and require further validation in study designs  
832 that can enable causal inference, such as randomised control trials. As reviewed  
833 previously<sup>18,19,69</sup>, MR tests are only unbiased when their assumptions are plausibly satisfied,  
834 which is why we implement a suite of different methods and sensitivity analyses with quite  
835 distinct underlying assumptions in this study. Genetic estimates on circulating retinol used as  
836 IVs could also be confounded by factors including uncontrolled population stratification,  
837 selection bias, and measurement error. In summary, we provide a large resource to the literature  
838 of putative effects of circulating retinol across the human clinical phenome that will be  
839 informative for future investigation of this trait.

840

841 Finally, we make some recommendations for future GWAS of retinoid molecules. An  
842 important limitation of our analyses is that we only investigated genetic effects on circulating  
843 retinol, rather than retinol availability within target tissues. Given the complexities of retinol  
844 homeostasis, it is plausible that genetic effects on factors like retinol entry into the cell (e.g.,  
845 STRA6 receptor) and esterification for storage (e.g., lecithin retinol acyltransferase) are  
846 obscured when considering only retinol present in serum or plasma. Therefore, future studies  
847 should attempt to measure retinol abundance in different tissues from genotyped samples,  
848 although this will pose a challenge in terms of obtaining sufficient sample sizes. Moreover, it

849 would be of interest to characterise the genetic overlap between effects on retinol versus other  
850 retinoids like retinaldehyde and all-*trans* retinoic acid. Despite these limitations, and the need  
851 for concerted efforts to collect more data, we believe this study demonstrates the value in  
852 conducting retinol GWAS to both better characterise retinoid associated biology and its clinical  
853 significance.

854

## 855 **ONLINE METHODS**

### 856 **Study cohorts**

857 The proceeding section outlines the datasets included in the genome-wide meta-analysis of  
858 circulating retinol, as well as the replication cohort.

859

### 860 *INTERVAL*

861 The largest constituent cohort of the meta-analysis was drawn from the INTERVAL study,  
862 comprised of recruited blood donors from the United Kingdom<sup>70</sup>. Retinol abundance was  
863 measured from plasma using the high-throughput metabolomics platform DiscoveryHD4<sup>®</sup>  
864 (Metabolon, Inc., Durham, USA), as outlined in the supplementary text. Briefly, after  
865 adjustment for various technical/biological confounders and outlier effects, residualised  
866 plasma retinol was inverse-rank normal transformed before association testing. Whole-genome  
867 sequencing of this cohort was performed as described elsewhere<sup>71</sup>. A GWAS of the normalised  
868 residuals was performed in HAIL via multiple-linear regression adjusted for INTERVAL  
869 metabolon batch and 10 genetic PCs<sup>72</sup>. The final GWAS sample size was 11,132 European  
870 ancestry participants.

871

### 872 *METSIM*

873 Plasma retinol was also measured using the DiscoveryHD4<sup>®</sup> high-throughput platform in a  
874 recent metabolome-wide GWAS of the METSIM (Metabolic Syndrome in Men) study<sup>73</sup>. The  
875 METSIM cohort consists of middle-aged men recruited from Northern Finland between 2005-  
876 2010<sup>74</sup>. As described by Yin *et al.*<sup>73</sup>, METSIM participants were genotyped using the Human  
877 OmniExpress-12v1\_C BeadChip and imputed using a custom METSIM panel of whole  
878 genome-sequenced participants in the study. A linear mixed model implemented in EPACTS  
879 v.3.2.6 was then leveraged to perform GWAS on residualised retinol, subjected to inverse-rank  
880 transformation after adjustment for technical and biological confounders. The final GWAS had  
881 a sample size of 6136 METSIM participants (European ancestry).

882

883 *ATBC+PLCO*

884 The largest previous dedicated GWAS of circulating retinol from 2011 was also included in  
885 this study<sup>22</sup>. This GWAS comprised data from two studies that measured serum retinol: the  
886 Alpha-Tocopherol, Beta-Carotene Cancer Prevention (ATBC) Study, a Finnish randomised  
887 control trial of beta-carotene/alpha-tocopherol supplementation for cancer prevention<sup>75</sup>, and  
888 the Prostate, Lung, Colorectal, and Ovarian (PLCO) Cancer Screening Trial, a United States  
889 trial of cancer screening effectiveness<sup>76</sup>. The inclusion criteria, measurement of serum retinol,  
890 and genotyping have been outlined by Mondul *et al.*<sup>22</sup>. Briefly, ATBC/PLCO samples were  
891 genotyped using the Illumina HumanHap550/610 arrays and imputed to the HapMap Central  
892 European reference panel, whilst serum retinol concentrations were estimated using reversed-  
893 phase liquid chromatography. The retinol GWAS (N=5006) was performed in R (version  
894 2.10.1) using multiple linear regression adjusted for age at sample collection, SNP derived PCs,  
895 cancer status, serum cholesterol, and body mass index. The retinol units for the GWAS effect  
896 sizes were in natural log transformed  $\mu\text{g/L}$ . A limitation of using summary statistics from this  
897 study is that only  $< 600,000$  variants were available for inclusion in the GWAS, as was common  
898 at the time before more recent advances in imputation pipelines that result in larger post-  
899 imputation yield. Therefore, to increase the number of variants available for meta-analysis, we  
900 applied a summary statistics-based imputation procedure to boost the number of variants  
901 available for meta-analysis. After harmonisation with the 1000 genomes phase 3 reference  
902 panel, we applied Gaussian summary statistics imputation (ImpG) as implemented by the FIZI  
903 v0.7.2 python package (<https://github.com/bogdanlab/fizi>) with the default window size of 250  
904 kb<sup>77</sup>. The ImpG model leverages the assumed Gaussian distribution of GWAS  $Z$  scores with  
905 mean zero and variance that arises due the LD-induced correlation between variants. As  
906 outlined elsewhere<sup>77</sup>,  $Z$  scores of unobserved (imputed) variants can be estimated given the LD  
907 correlation matrix derived from the reference panel, along with a metric of imputation accuracy  
908 ( $R^2$ ) using the conditional variance. We retained only confidently imputed variants ( $R^2 > 0.8$ ).

909

910 *TwinsUK*

911 We performed a serum retinol GWAS in the TwinsUK cohort to serve as a replication dataset.  
912 High-throughput metabolomics profiling in this cohort has been described extensively  
913 elsewhere<sup>78</sup>. TwinsUK is a prospective population-based study of mostly female twin pairs  
914 which has been profiled using a variety of multiomic technologies<sup>79</sup>. As outlined by Shin *et al.*,  
915 genotyping was performed with a combination of Illumina arrays (HumanHap300,  
916 HumanHap610Q, 1M-Duo and 1.2MDuo 1M)<sup>80</sup>. This was followed by imputation to the 1000

917 genomes phase reference panel after quality control and retaining individuals of predominantly  
918 European ancestry after PCA. We retained physically genotyped variants and those with  
919 moderate imputation accuracy for GWAS ( $R^2 > 0.3$ ) in 5654 samples. Our strategy for GWAS  
920 in this cohort given the limited sample size with measured retinol available was to split the  
921 twin pairs into two separate cohorts and ensure individuals in each sub-cohort were unrelated.  
922 Relatedness testing in both sub-cohorts containing one of the two possible twin pairs was  
923 performed separately using KING as implemented by plink2 (PLINK v2.00a3LM AVX2  
924 Intel), with one participant from third-degree relative or greater pairs randomly removed<sup>81,82</sup>.  
925 Kinship estimation via KING was performed for autosomal variants physically genotyped on  
926 the array, with a MAF  $> 0.05$ , outside of regions of long-range LD like the MHC<sup>83</sup>, and in  
927 relative linkage equilibrium ( $r^2 < 0.05$ ). PCA was then applied in each sub-cohort using plink2  
928 to calculate eigenvectors for use as downstream covariates. Retinol was measured from  
929 samples at three timepoints. The mean age of participants at each timepoint was 51.5 (SD =  
930 8.41), 58.6 (SD = 8.38), and 64.7 (SD = 8.41), respectively. There were only a very small  
931 number of males in this cohort (~ 3%), so only females were retained for further analysis (N =  
932 1696) due to this imbalance in sex composition. After merging with the genotyped split twin  
933 cohorts, as described above, there were up to 916 and 717 genotyped participants with  
934 measured retinol at three timepoints in each subset, respectively. Six GWAS were performed:  
935 in each sub-cohort (one or two), multiple linear regression was utilised to test the additive effect  
936 of each variant on measured retinol at one of the three measured time-points covaried for age,  
937 five SNP derived PCs, and metabolomics batch. These GWAS was performed using the --glm  
938 flag in plink2, resulting a 9,051,192 by three matrix of estimated retinol effect sizes for both  
939 sub-cohorts.

940

### 941 **Genome-wide meta-analyses**

942 We conducted genome-wide meta-analysis of common variants (MAF  $\geq 0.01$ ) using METAL  
943 (version March 2011), followed by a rare variant (MAF  $< 0.01$ ) meta-analysis also with  
944 METAL. The METSIM and INTERVAL cohorts were integrated for the primary meta-analysis  
945 as they both had expansive genome-wide coverage of common and rare variants. A sample size  
946 weighted meta-analysis of Z scores (Stouffer's method) was utilised for this purpose. We also  
947 conducted the METSIM+INTERVAL meta-analysis via an inverse-variance weighted  
948 estimator with fixed effects to estimate the effect sizes of effect alleles in SD units of plasma  
949 retinol given this was the unit of both the METSIM and INTERVAL GWAS, as well as both  
950 studies using the same Metabolon Inc. platform for metabolite quantification. Heterogeneity

951 between the studies was assessed using Cochran's  $Q$  test. We then conducted a meta-analysis  
952 with fewer available variants that also included ATBC+PLCO using Stouffer's method (as the  
953 unit for this GWAS differed from METSIM and INTERVAL). This larger sample-size meta-  
954 analysis was also restricted to common variants as there was very limited rare variant coverage  
955 in ATBC+PLCO.

956

957 The FUMA v1.4.1 (Functional Mapping and Annotation of Genome-Wide Association  
958 Studies) platform was utilised to annotate variants, define lead SNPs, and infer loci boundaries  
959 for genome-wide significant signals ( $P < 5 \times 10^{-8}$ )<sup>84</sup>. We utilised the default settings for  
960 defining independent significant SNPs ( $r^2 \leq 0.6$ ) and lead SNPs ( $r^2 \leq 0.1$ ). LD estimation was  
961 achieved using the 1000 genomes phase 3 European reference panel, with LD blocks within  
962 250 kb of each other merged into a single locus. We then attempted to replicate the lead SNPs  
963 from the eight genome-wide significant loci (METSIM+INTERVAL) in TwinsUK, as  
964 described in the previous section. Given the small sample size of TwinsUK, we sought to  
965 ascertain if the lead SNPs were directionally consistent with the meta-analysis. This was  
966 achieved by taking the mean SNP  $Z$  score across the three timepoints for the two sub-cohorts  
967 of unrelated participants, with a binomial test utilised to infer whether the number of lead SNPs  
968 (mean  $Z$ ) that were directionally consistent was greater than chance alone (Binomial  $P < 0.05$ ).

969

970 In the rare-variant meta-analysis, we annotated variants using the Functional Annotation of  
971 Variants (FAVOR) online resource<sup>85</sup>. Phenome-wide association profiles of selected variants  
972 were also investigated using the pheweb browser collated from FinnGen release 8  
973 (<https://r8.finnngen.fi/>)<sup>86</sup>. Rare variants were then aggregated to genes through leveraging the  
974 characteristics of the Cauchy distribution<sup>87,88</sup>. In this approach, gene-wise  $P$  values are summed  
975 and then transformed to approximate a Cauchy distribution, which due to its heavy tail is  
976 insensitive to correlations amongst the  $P$  values. This behaviour of the Cauchy distribution is  
977 important as covariance amongst rare variants is difficult to estimate, and therefore, this  
978 approach guards against inflated type I error due to potential unknown covariance/LD between  
979 rare variants. Code for implementing the Cauchy aggregation was adapted from  
980 <https://github.com/yaowuliu/ACAT>.

981

## 982 **SNP heritability estimation**

983 Summary statistics for the meta-analyses were 'munged' using the munge\_sumstats.py script  
984 from the *ldsc* repository of scripts (<https://github.com/bulik/ldsc>) and only common (MAF >

985 0.05) HapMap3 variants outside of the MHC retained. SNP heritability was then estimated  
986 using the LDSR model and the 1000 genomes phase 3 reference panel<sup>89</sup>. By way of  
987 comparison, we then estimated SNP heritability using the LDAK model via the SumHer  
988 implemented in LDAK v.5.2 (<https://dougsspeed.com/>)<sup>90</sup>. Pre-computed tagging files derived  
989 from 2000 white British individuals in the UKBB for HapMap3 SNPs were utilised to calculate  
990 SNP heritability using the LDAK-thin and BLD-LDAK models, as described elsewhere<sup>90,91</sup>.  
991 We also estimated partitioned SNP heritability via LDSR using a multi-tissue and cell-type  
992 panel<sup>92</sup>.

993

### 994 **Empirical Bayes' modelling of the genetic architecture of retinol**

995 We investigated the polygenicity of the genetic architecture of circulating retinol using an  
996 Empirical Bayes' adaptive shrinkage method termed ashR<sup>93</sup>. Functions to perform this method  
997 were implemented via the ashR R package v2.2-54 (<https://github.com/stephens999/ashr>).  
998 Briefly, this approach models effect sizes, along with their standard error, as a mixture of zero  
999 and non-zero effects. Empirical Bayes' inference is performed under the assumption that the  
1000 distribution of these variant effects is unimodal, which is a realistic assumption for genetic  
1001 effects on complex traits. The METSIM+INTERVAL IVW meta-analysis was utilised for this  
1002 as it provides an interpretable effect size and standard error for each variant (plasma SD units).  
1003 In line with previous work<sup>94</sup>, we annotated each HapMap3 variant with its corresponding LD  
1004 score from the 1000 genomes phase 3 European reference panel, as well LD scores from the  
1005 UKBB White Great British samples for comparison, and sorted these into bins of similar LD  
1006 scores ( $N_{\text{Bins}}=1000$  and  $N_{\text{Bins}}=5000$ ). The ashR Empirical Bayes' inference of the proportion  
1007 of non-zero effects was undertaken in each LD score bin, followed by calculating the mean  
1008 across all bins. We utilised a generalised additive model to plot a smoothed trend line of the  
1009 relationship between increasing LD score bin (higher LD score) and the proportion of non-zero  
1010 effects.

1011

### 1012 **Gene prioritisation**

1013 Gene prioritisation was performed within genome-wide significant loci, as well as outside of  
1014 loci that obtained genome-wide significance. Due to the better coverage of common variants,  
1015 we focused on the METSIM+INTERVAL meta-analysis for this analysis. The pipeline for  
1016 prioritising putative causal genes within the eight genome-wide significant loci was adapted  
1017 from a previous GWAS performed by our group<sup>95</sup>. The following criteria were utilised in this  
1018 study: 1) closest transcription start site (TSS) to the lead SNP, 2) closest gene (any) to the lead



1019 SNP, 3) gene encoding retinoid transporter or enzyme in locus, 4) non-synonymous variant in  
1020 locus, 5) the most statistically significant GTEx eGene [expression quantitative trait loci  
1021 (eQTL) signal] in locus, 6) most significant GTEx eGene in finemapped credible set for eQTL  
1022 signal [posterior inclusion probability ( $PIP$ ) > 0.1, DAP-G method]<sup>96,97</sup>, 7) strongest plasma  
1023 pGene [protein quantitative trait loci (pQTL) signal] drawn from finemapped pQTLs ( $PIP$  >  
1024 0.5)<sup>98</sup>, 8) highest CADD score<sup>99</sup>, 9) lowest RegulomeDB score<sup>100</sup>, 10) the OpenTargets  $V2G$   
1025 predicted gene for the lead SNP<sup>101</sup>, and 11) genes physically mapped to SNPs in the 95%  
1026 credible set derived from probabilistic finemapping (assuming a single causal variant such that  
1027 LD did not have to be modelled)<sup>102</sup>. We used a prior variance of 0.15 to approximate Bayes'  
1028 factors from variant-wise effect sizes (METSIM+INTERVAL IVW meta-analysis), in line  
1029 with previous work finemapping association signals with quantitative traits<sup>103</sup>. We  
1030 characterised which genes satisfied the greatest number of the above criteria on a per locus  
1031 basis.

1032

1033 *FOXP2* was one of our confidently prioritised genes but its biological significance outside of  
1034 the brain is less well understood. To investigate this, we analysed RNA sequencing data from  
1035 an *in vitro* experiment that overexpressed *FOXP2* in a human osteosarcoma epithelial cell line  
1036 (U2OS) via transfection of wild-type *FOXP2* expressing plasmids. Raw read counts from five  
1037 control cell line replicates versus five plasmid transfected *FOXP2* overexpression replicates  
1038 were downloaded from the Gene Expression Omnibus (GEO) resource (GEO Accession:  
1039 GSE138938). Data normalisation, filtration, and differential expression analyses were  
1040 performed using the edgeR package version 3.34.0<sup>104</sup>. Specifically, raw counts were firstly  
1041 normalised to library size and lowly expressed genes with fewer than 10 raw counts in the  
1042 smallest library removed via a counts-per-million thresholding approach. Data were inspected  
1043 before and after the filtration step via coefficient of variation (BCV) and multidimensional  
1044 scaling (MDS) plots (Supplementary Figure 10). Differential expression for each gene that  
1045 survived quality control was then performed using exact tests for differences in the means  
1046 between two groups of negative-binomially distributed counts. We defined a differentially  
1047 expressed gene as those which survived multiple-testing correction using the Bonferroni  
1048 method ( $P_{\text{Corrected}} < 0.05$ ), with three different absolute  $\log_2$  fold change (FC) cut-offs  
1049 considered:  $|\log_2\text{FC}| > 1.5$ ,  $|\log_2\text{FC}| > 2$ , and  $|\log_2\text{FC}| > 5$ . The use of Bonferroni correction and  
1050 large absolute FC thresholds is very conservative; however, given the volume of differentially  
1051 expressed genes, and a relatively large number of replicates for a cell line experiment boosting  
1052 power, we believe these strict parameters are warranted to prioritise the most salient *FOXP2*

1053 associated signals. The overrepresentation of each set of candidate genes amongst biological  
1054 pathways and other ontology sets was tested using g:Profiler<sup>105</sup>.

1055

1056 To identify potential causal genes that have not reached genome-wide significance at our  
1057 current sample size, we integrated the circulating retinol GWAS with genetic effects on mRNA  
1058 and protein expression. Firstly, we conducted a transcriptome and proteome-wide association  
1059 study (TWAS/PWAS) of circulating retinol using the FUSION approach<sup>106</sup>. As outlined  
1060 previously<sup>107–110</sup>, FUSION leverages models of *cis*-acting genetically regulated expression  
1061 (GReX) that exhibit statistically significant non-zero heritability. Variant weights from GReX  
1062 models are integrated with the effect of those same variants on retinol to estimate the direction  
1063 of genetically regulated expression associated with increasing circulating retinol. TWAS  
1064 GReX (mRNA) were estimated previously using GTEx v8  
1065 (<http://gusevlab.org/projects/fusion/>). We selected the following biologically informative  
1066 tissues to perform TWAS based on known retinol biology or tissues that exhibited at least  
1067 nominally significant ( $P < 0.01$ ) enrichment of SNP heritability in the partitioned-LDSR  
1068 model. The selected tissues were: small intestine terminal ileum, pancreas, liver, adipose  
1069 (visceral omentum), adipose (subcutaneous), breast (mammary tissue), and whole blood. It has  
1070 been suggested previously that TWAS signals from tissues that are less directly trait relevant  
1071 can induce spurious associations, which is why we limited our hypothesis space to these  
1072 tissues<sup>111</sup>. Protein GReX were derived from plasma (ARIC study), as outlined elsewhere<sup>98</sup>. We  
1073 applied Benjamini-Hochberg false discovery rate (FDR) correction across all TWAS Z,  
1074 followed by all PWAS Z. Colocalisation between GReX models and retinol was performed for  
1075 all TWAS/PWAS signals that were at least nominally significant via the coloc package as  
1076 implemented by FUSION (single shared variant hypothesis)<sup>103</sup>. We then considered a more  
1077 conservative approach to prioritise proteins for whom expression could be causally linked to  
1078 circulating retinol through leveraging finemapped plasma pQTLs ( $PIP > 0.5$ , ARIC) as  
1079 instrumental variables (IV) for Mendelian randomisation<sup>19,44,98</sup>. The Wald ratio method was  
1080 implemented for proteins with single IVs, whilst an inverse-variance weighted estimator with  
1081 fixed effects was utilised for proteins with more than 1 IV. Fixed effects were used for the IVW  
1082 rather than multiplicative-random effects in this instance as no protein had  $> 4$  IVs. The use of  
1083 IV based approaches for identifying trait-associated genes versus GReX has been discussed  
1084 extensively elsewhere<sup>109</sup>. Colocalisation was also performed for proteins that survived FDR  
1085 correction. Mendelian randomisation was performed using the TwoSampleMR package v0.5.6,

1086 with IVs clumped through leveraging LD from the 1000 genomes phase 3 European reference  
1087 panel to retain only independent pQTLs ( $r^2 < 0.001$ ).

1088

1089 Finally, we investigated the tissue specificity of prioritised genes from GWAS loci and the  
1090 TWAS/PWAS/MR approach using FUMA. To do this, we compared the expression of these  
1091 prioritised genes using a *t*-test (one-sided and two-sided) against all other available genes in  
1092 GTEx v8 on a per tissue basis (54 tissues), followed by applying Bonferroni correction to these  
1093 *P*-values<sup>84</sup>. We also conducted pathway analyses of these genes using g:Profiler with default  
1094 parameters<sup>105</sup>.

1095

### 1096 **Causal inference**

1097 We developed and implemented a comprehensive pipeline to leverage this retinol GWAS to  
1098 identify putative causal effects of retinol on traits across the human clinical phenome, as well  
1099 as traits that causally influence retinol in the reverse direction. This was achieved using  
1100 Mendelian randomisation (MR), which has been reviewed extensively elsewhere<sup>19,45,69</sup>. Firstly,  
1101 we considered circulating retinol as the MR exposure. The lead SNP in *RBP4* was first chosen  
1102 as a single IV to proxy circulating retinol as out of all the genes implicated in genome-wide  
1103 significant loci, *RBP4* has exhibits the most specificity in terms of its relationship with serum  
1104 retinol. We utilised the Wald ratio method to estimate the effect of the circulating retinol  
1105 increasing rs10882283-A allele (METSIM+INTERVAL IVW effect size) on over 19,000  
1106 outcomes in the IEUGWASdb v6.9.2 resource via the ieugwasr package version 0.1.5<sup>26</sup>. To  
1107 correct for multiple testing, we applied the false discovery rate (FDR) method, and we retained  
1108 only those retinol/outcome pairs with estimates that were significant below the 1% FDR  
1109 threshold ( $q < 0.01$ ). Subsequently, we investigated whether significant trait pairs colocalised  
1110 using *coloc* v5.1.0 with default priors.

1111

1112 Although the single IV approach is more conservative, power to detect causal effects can be  
1113 boosted by using all available independent ( $r^2 < 0.001$ , 1000 genomes phase 3 European  
1114 reference panel) genome-wide significant SNPs as IVs. The inverse-variance weighted  
1115 estimator with multiplicative random effects (IVW-MRE) was utilised to estimate the effect of  
1116 retinol on outcomes from IEUGWASdb for which at least 6 IVs were available in that  
1117 GWAS<sup>112</sup>. The IVW approach has a zero percent breakdown level as it assumes all IVs are  
1118 valid for use in Mendelian randomisation, which is often unrealistic in practice. We developed  
1119 pipeline to prioritise the most reliable causal estimates from the IVW-MRE that survived

1120 multiple-testing correction ( $q < 0.01$ ). This involved identifying retinol MR estimates on traits  
1121 for which the following applied: i) no significant heterogeneity ( $P < 0.05$ ) between IV  
1122 exposure-outcome effects tested with Cochran's  $Q^{113}$ , ii) a non-significant intercept of an MR-  
1123 Egger model that does not constrain the intercept to pass through the origin<sup>114</sup>, and iii) no  
1124 evidence that leaving out any single IV ablates the statistical significance (at least  $P < 0.05$ ) of  
1125 the estimate<sup>115</sup>. Traits satisfying these criteria were then assigned a tier (Tier #1, Tier #2, Tier  
1126 #3) based on the statistical significance of the retinol causal estimate using other MR methods  
1127 besides the IVW-MRE that have different assumptions regarding IV validity. These were: the  
1128 IVW with fixed effects (does not model heterogeneity like with MRE), the weighted median  
1129 method<sup>116</sup>, the weighted mode method<sup>117</sup>, and the MR-Egger method<sup>114</sup>. The underlying  
1130 assumptions and methodological considerations of using these methods have discussed  
1131 extensively elsewhere<sup>118,119</sup>. Traits for which the effect of retinol was at least nominally  
1132 statistically significant using all five methods were assigned as Tier #1, whilst four tests being  
1133 statistically significant was Tier #2, and three tests being statistically significant was Tier #3.  
1134 The  $F$ -statistic and  $I^2$  of the IVs was also assessed to ensure they were well powered ( $F > 10$ )  
1135 and suited for MR-Egger ( $I^2 > 0.9$ ), respectively<sup>120</sup>. To investigate the effect of body fat  
1136 percentage on the retinol associated MRI indices, we used IVs from a non-overlapping GWAS  
1137 of that trait ( $N=65,831$ )<sup>121</sup>. We then repeated the above process of binary outcomes with at least  
1138 1000 cases from FinnGen release 8, which is not included in the current version of  
1139 IEUGWASdb at time of analysis, to increase power to detect effects on disease endpoints.

1140

1141 We also systematically investigated continuous outcomes that may causally impact circulating  
1142 retinol. Continuous outcomes were our focus as causal estimates from binary exposures are  
1143 difficult to interpret and are often less powered in the context of MR<sup>47</sup>. Exposures were filtered  
1144 from phenotypes available in IEUGWASdb to retain continuous traits, with further filtering to  
1145 identify traits with  $\geq 5$  genome-wide significant, independent ( $r^2 < 0.001$ , 1000 genomes phase  
1146 3 European reference panel) IVs available in the retinol GWAS. The same pipeline as above  
1147 was then applied (IVW-MRE FDR  $< 0.01$ , followed by sensitivity analyses and tier  
1148 assignment). The causal estimate of creatinine on retinol was a Tier #1 trait, and due to the  
1149 pervasive polygenicity of creatinine afforded by its large sample size, we followed up this  
1150 relationship using the MR model "Causal Analysis Using Summary Effect Estimates"  
1151 (CAUSE), as described elsewhere using the CAUSE R package v1.2.0<sup>52</sup>. Briefly, this method  
1152 is a more polygenic approach and seeks to distinguish casual effects from correlated pleiotropy  
1153 by fitting competing models that account for these terms and comparing them using ELPD.

1154 After using 1 million random variants to estimate nuisance parameters, LD clumping and  
1155 thresholding was applied to the serum creatinine summary statistics (IEUGWASdb trait ID:  
1156 met-d-creatinine) in line with the original CAUSE publication ( $P < 0.001$ ,  $r^2 < 0.01$ , 1000  
1157 genomes phase 3 European reference panel). The competing CAUSE models were then fit  
1158 (null, sharing, and causal), ensuring that all Pareto  $k$  estimates were  $< 0.5$  during the model  
1159 comparison using ELPD.

1160

1161 We then performed a multivariable MR (MVMR) analysis to estimate causal effects of  
1162 creatinine on retinol conditioned on three major lipid species using GWAS from the global  
1163 lipids genetics consortium (LDL, HDL, and triglycerides)<sup>122</sup>. MVMR was undertaken in  
1164 accordance with previous work using the R packages MVMR v0.3 and  
1165 MendelianRandomization v0.6.0<sup>118</sup>. Briefly, this entailed identifying variants associated at  
1166 genome-wide significance with at least one of the four exposures that are independent ( $r^2 <$   
1167  $0.001$ ), calculating a conditional  $F$ -statistic for multivariable instruments<sup>123</sup>, and applying four  
1168 MVMR models (IVW, Egger regression, Weighted Median, and a LASSO based penalised  
1169 regression approach for selecting the optimal IV configuration)<sup>124</sup>.

1170

### 1171 **Drugs and perturbagens associated with circulating retinol**

1172 We also considered drugs that may influence circulating retinol. We searched genes prioritised  
1173 from our pipeline with an assigned direction of retinol-associated expression using DGIdb  
1174 v4.2.0 and DrugBank v5 to identify retinol associated genes targeted by drugs<sup>125,126</sup>, outside of  
1175 known drugs that target RBP4 and TTR. Retained drug-gene interactions were restricted to  
1176 those with known mechanism of action and  $> 2$  lines of supporting evidence. We also utilised  
1177 computational signature mapping to identify pharmacological agents that may enhance or  
1178 inhibit the expression of genes associated with circulating retinol. To boost power for signature  
1179 mapping, we considered all gene that were nominally associated with retinol in the  
1180 TWAS/PWAS ( $P < 0.05$ ) that exhibited moderate colocalisation of a shared causal variant  
1181 ( $PP_{H4} > 0.4$ ). These genes were uploaded to the Connectivity Map Query online tool to quantify  
1182 the similarity, termed connectivity, with drug perturbagen associated expression profiles, as  
1183 outlined in detail elsewhere<sup>127</sup>.

1184

### 1185 **Polygenic scoring**

1186 We used the independent TwinsUK replication cohort, described in a preceding section, to  
1187 investigate retinol polygenic scores (PGS). PGS were applied to genotyped variants and high

1188 confidence imputed variants ( $R^2 > 0.8$ ) in TwinsUK. There were two different methodologies  
1189 implemented to construct PGS: LD clumping and thresholding (LD C+T) and a probabilistic  
1190 finemapping based method that scales variant effect sizes based on their posterior probability  
1191 of causality, thereby upweighting signals more likely to be causal (RápidoPGS)<sup>128,129</sup>. In the  
1192 LD C+T approach, variants were clumped using the within sample LD of TwinsUK at the  
1193 following  $P$ -value thresholds:  $5 \times 10^{-8}$ ,  $1 \times 10^{-5}$ ,  $1 \times 10^{-3}$ , 0.01, 0.05, 0.1, 0.5, and 1. Additive PGS  
1194 were then profiled using PRSice2 v2.3.5<sup>130</sup>. The RápidoPGS applies probabilistic finemapping  
1195 (with Bayes' factors approximated using Wakefield's method) to independent LD blocks  
1196 genome-wide such that variant-wise posterior probabilities of causality can be estimated. This  
1197 in essence is a 'shrinkage' approach to PGS to account for double-counting effects that arise  
1198 due to correlated effect sizes induced by LD; however, does not inherently require an  
1199 independent genotyped sample from the GWAS for tuning. Approximate Bayes' factors in  
1200 each LD block were derived assuming a prior variance of 0.15, conventionally used for  
1201 quantitative traits. This parameter choice was compared to a data driven approach to estimating  
1202 the prior variance based on SNP heritability, as outlined elsewhere<sup>129</sup>. Variant-wise effect sizes  
1203 are multiplied by their posterior probabilities before PGS calculation, as above.

1204

1205 As we are using a twin cohort, there were two different approaches we utilised to tune the  
1206 optimal PGS configuration from the LD C+T approach. Firstly, we randomly split the cohort  
1207 into a training (70% of participants) and test (30% of participants) partition. A linear mixed  
1208 model was then fit in the training partition with fixed effects of PGS, age, and metabolomics  
1209 batch and a random effect of family ID to account for twin relatedness. This was applied for  
1210 retinol measured at each of the three visits for all the  $P$ -value thresholds. The marginal  $R^2$  from  
1211 a null model with no PGS was subtracted from the full model to infer the best performing PGS  
1212 in the training partition (mean marginal  $R^2$  across three visits). The variance explained of the  
1213 best performing  $P$ -value threshold was then estimated using the same approach in the test  
1214 partition. By way of comparison, we also split the twins into separate unrelated cohorts and  
1215 used one set as training and one as testing. Fixed effects instead of mixed effects linear  
1216 regression was then implemented in a similar fashion to above, additionally covaried for five  
1217 SNP derived principal components. We then repeated all the above for the two probabilistic  
1218 finemapping weighted PGS (prior variance = 0.15 and data driven prior variance).

1219

1220

1221

## 1222 Normative modelling

1223 We built a normative model of retinol as a function of age per study visit in TwinsUK. This  
1224 was achieved using a generalised additive model for location ( $\mu$ ), scale ( $\sigma$ ), and shape  
1225 (GAMLSS)<sup>131,132</sup>, implemented in R v4.4.1. The GAMLSS approach is useful in this  
1226 application as it is semi-parametric and able to account for factors such as heteroskedasticity  
1227 and non-Gaussian distributions. Our modelling approach can be summarised as follows: we  
1228 firstly fit a GAMLSS model in the full sample for a variety of GAMLSS distribution families  
1229 implemented by the *gamlss* R package v5.4.12. The model on retinol at visit  $i, i \in \{1,2,3\}$  for  
1230 each of the GAMLSS families set the  $\mu$  term as the first order fractional polynomial of age,  
1231 along with metabolomics batch as an additional covariate, with the term  $\sigma$  just the first order  
1232 fractional polynomial of age to model the scale of the distribution. The model fit of each of the  
1233 tested families was then assessed using the Bayesian information criterion (BIC) and the  
1234 Akaike information criterion (AIC). We repeated the above also including measured body mass  
1235 index (BMI) at that timepoint in the  $\mu$  term by way of comparison. We chose the GAMLSS  
1236 family (Box-Cox  $t$  distribution) through considering the model performance (minimum AIC  
1237 and BIC) over all three visits (Supplementary Text, Supplementary Figures 11-12). We then  
1238 split the cohort in half, separating the twins, and fit normative centile curves using the selected  
1239 GAMLSS family to one half of each batch of retinol measurement, and computed deviations  
1240 on the other independent sample half on a per batch basis. We repeated this process with the  
1241 modelling and deviation subsets reversed to compute deviations for all samples. Individuals  
1242 whose measured retinol was above the model derived 95th percentile were classified as having  
1243 supra-normal retinol for their age, while those below the 5th percentile were classified as  
1244 having infra-normal retinol. To guard against overfitting due to the relatedness of twins  
1245 between the two subsets, we then performed all of the above just using one half of the twins as  
1246 the full cohort for model fitting, followed by splitting this subset as described above. We tested  
1247 the relationship between scaled retinol PGS (mean = 0, SD = 1) and infra-normal retinol at  
1248 each visit was tested in half of the cohort (unrelated) using binomial logistic regression  
1249 additionally covaried for five SNP derived PCs. The same models were also constructed for  
1250 supra-normal individuals.

1251

## 1252 Software and operating systems

1253 The primary analyses in this manuscript were performed either on a MacBook Pro (OS X:  
1254 Ventura 13.3), an in-house linux cluster (Ubuntu 18.04.5 LTS), or the High-Performance

1255 Computing Research Compute Grid of the University of Newcastle [Red Hat Enterprise Linux  
1256 release 8.1 (Ootpa)]. The primary R version utilised was version 4.1.1 (2021-08-10), with some  
1257 additional analyses using R version 4.0.3 (2020-10-10) (linux cluster). The Python version  
1258 utilised was either Python 2.7.17 or Python 3.6.9, depending on the requirements of the  
1259 analyses.

1260

#### 1261 **DATA AVAILABILITY**

1262 Genome-wide summary statistics will be uploaded to GWAS catalog upon final publication.  
1263 In the interim, summary statistics can also be found at the following: 10.5281/zenodo.7905523.  
1264 TwinsUK data can be accessed by approve researchers upon application  
1265 (<https://twinsuk.ac.uk/resources-for-researchers/our-data/>).

1266

#### 1267 **CODE AVAILABILITY**

1268 Code used in this study is freely available at the following GitHub repository -  
1269 [https://github.com/Williamreay/Retinol\\_GWAS\\_code](https://github.com/Williamreay/Retinol_GWAS_code).

1270

#### 1271 **ACKNOWLEDGEMENTS**

1272 TwinsUK is funded by the Wellcome Trust, Medical Research Council, Versus Arthritis,  
1273 European Union Horizon 2020, Chronic Disease Research Foundation (CDRF), Zoe Ltd, the  
1274 National Institute for Health and Care Research (NIHR) Clinical Research Network (CRN) and  
1275 Biomedical Research Centre based at Guy's and St Thomas' NHS Foundation Trust in  
1276 partnership with King's College London. P.S. was supported by a Rutherford Fund Fellowship  
1277 from the Medical Research Council (grant no. MR/S003746/1). Participants in the INTERVAL  
1278 trial were recruited with the active collaboration of NHS Blood and Transplant England  
1279 ([www.nhsbt.nhs.uk](http://www.nhsbt.nhs.uk)), which has supported field work and other elements of the trial. DNA  
1280 extraction and genotyping were co-funded by the National Institute for Health and Care  
1281 Research (NIHR), the NIHR BioResource (<http://bioresource.nihr.ac.uk>) and the NIHR  
1282 Cambridge Biomedical Research Centre (BRC-1215-20014). The academic coordinating  
1283 centre for INTERVAL was supported by core funding from the: NIHR Blood and Transplant  
1284 Research Unit (BTRU) in Donor Health and Genomics (NIHR BTRU-2014-10024), NIHR  
1285 BTRU in Donor Health and Behaviour (NIHR203337), UK Medical Research Council  
1286 (MR/L003120/1), British Heart Foundation (SP/09/002; RG/13/13/30194; RG/18/13/33946)  
1287 and NIHR Cambridge BRC (BRC-1215-20014; NIHR203312) [\*].



1288 Metabolon metabolomics assays in INTERVAL were funded by the: NIHR BioResource,  
1289 NIHR Cambridge Biomedical Research Centre (BRC-1215-20014) [\*], Wellcome Trust grant  
1290 number 206194 and BioMarin Pharmaceutical, Inc. The academic coordinating centre thank  
1291 blood donor centre staff and blood donors for participating in the INTERVAL trial. This work  
1292 was supported by Health Data Research UK, which is funded by the UK Medical Research  
1293 Council, Engineering and Physical Sciences Research Council, Economic and Social Research  
1294 Council, Department of Health and Social Care (England), Chief Scientist Office of the  
1295 Scottish Government Health and Social Care Directorates, Health and Social Care Research  
1296 and Development Division (Welsh Government), Public Health Agency (Northern Ireland),  
1297 British Heart Foundation and Wellcome. For the purpose of open access, the author(s) has  
1298 applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript  
1299 version arising from this submission.\*The views expressed are those of the author(s) and not  
1300 necessarily those of the NIHR, NHSBT or the Department of Health and Social Care. Access  
1301 to the TwinsUK cohort for this study was funded by a Hunter Medical Research Institute  
1302 (Newcastle, Australia) Precision Medicine Research Program Pilot Grant (W.R.R). C.E.C. is  
1303 supported by an investigator grant from the National Health and Medical Research Council  
1304 (NHMRC). M.A.D is supported by an investigator grant from the National Health and Medical  
1305 Research Council (NHMRC). P.S. was supported by a Rutherford Fund Fellowship from the  
1306 Medical Research Council (grant no. MR/S003746/1). We also acknowledge and thank all  
1307 participants in the genetic studies included in this paper.

1308

### 1309 **ETHICS DECLARATION**

1310 P.S. is now a full-time employee of GlaxoSmithKline. The remaining authors declare no  
1311 competing financial interests.

1312

### 1313 **AUTHOR CONTRIBUTIONS**

1314 W.R.R designed the study and was the primary analyst. D.J.K. assisted with developing the  
1315 phenome-wide Mendelian randomisation pipeline. M.A.D developed the normative modelling  
1316 pipeline. Z.F.G. performed the signature mapping analyses. P.S performed quality control on  
1317 the INTERVAL data prior to GWAS. K.K performed the INTERVAL GWAS. E.D.C and  
1318 C.E.C provided input into the TwinsUK retinol quality control, and the interpretation of the  
1319 clinical associations. L.A.G assisted with data curation and preparation of the manuscript.  
1320 A.M.M and D.A provided the ATBC/PLCO data. M.J.C contributed to drafting manuscript and

1321 provided funding. All authors contributed to interpretation of the results and the final  
1322 manuscript.

1323

## 1324 REFERENCES

1325 1. Blomhoff, R. & Blomhoff, H. K. Overview of retinoid metabolism and function. *J*

1326 *Neurobiol* **66**, 606–630 (2006).

1327 2. Reay, W. R. & Cairns, M. J. The role of the retinoids in schizophrenia: genomic and

1328 clinical perspectives. *Mol Psychiatry* **25**, 706–718 (2020).

1329 3. D’Ambrosio, D. N., Clugston, R. D. & Blaner, W. S. Vitamin A metabolism: an update.

1330 *Nutrients* **3**, 63–103 (2011).

1331 4. Britton, G. Carotenoid research: History and new perspectives for chemistry in biological

1332 systems. *Biochimica et Biophysica Acta (BBA) - Molecular and Cell Biology of Lipids*

1333 **1865**, 158699 (2020).

1334 5. Duester, G., Mic, F. A. & Molotkov, A. Cytosolic retinoid dehydrogenases govern

1335 ubiquitous metabolism of retinol to retinaldehyde followed by tissue-specific metabolism

1336 to retinoic acid. *Chem Biol Interact* **143–144**, 201–210 (2003).

1337 6. Cunningham, T. J. & Duester, G. Mechanisms of retinoic acid signalling and its roles in

1338 organ and limb development. *Nat Rev Mol Cell Biol* **16**, 110–123 (2015).

1339 7. Li, Y., Wongsiriroj, N. & Blaner, W. S. The multifaceted nature of retinoid transport and

1340 metabolism. *Hepatobiliary Surg Nutr* **3**, 126–139 (2014).

1341 8. Kanai, M., Raz, A. & Goodman, D. S. Retinol-binding protein: the transport protein for

1342 vitamin A in human plasma. *J Clin Invest* **47**, 2025–2044 (1968).

1343 9. van Bennekum AM, null *et al.* Biochemical basis for depressed serum retinol levels in

1344 transthyretin-deficient mice. *J Biol Chem* **276**, 1107–1113 (2001).

1345 10. Gudas, L. J. Synthetic Retinoids Beyond Cancer Therapy. *Annu Rev Pharmacol*

1346 *Toxicol* **62**, 155–175 (2022).

- 1347 11. Chisholm, D. R. & Whiting, A. Design of synthetic retinoids. *Methods Enzymol* **637**,  
1348 453–491 (2020).
- 1349 12. Behl, T. *et al.* Therapeutic insights elaborating the potential of retinoids in  
1350 Alzheimer’s disease. *Front. Pharmacol.* **13**, 976799 (2022).
- 1351 13. Kim, J. A., Jang, J.-H. & Lee, S.-Y. An Updated Comprehensive Review on Vitamin  
1352 A and Carotenoids in Breast Cancer: Mechanisms, Genetics, Assessment, Current  
1353 Evidence, and Future Clinical Implications. *Nutrients* **13**, 3162 (2021).
- 1354 14. Harirchian, M. H., Mohammadpour, Z., Fatehi, F., Firoozeh, N. & Bitarafan, S. A  
1355 systematic review and meta-analysis of randomized controlled trials to evaluating the trend  
1356 of cytokines to vitamin A supplementation in autoimmune diseases. *Clinical Nutrition* **38**,  
1357 2038–2044 (2019).
- 1358 15. Leelakanok, N., D’Cunha, R. R., Sutamtewagul, G. & Schweizer, M. L. A systematic  
1359 review and meta-analysis of the association between vitamin A intake, serum vitamin A,  
1360 and risk of liver cancer. *Nutr Health* **24**, 121–131 (2018).
- 1361 16. O’Connor, E. A. *et al.* Vitamin and Mineral Supplements for the Primary Prevention  
1362 of Cardiovascular Disease and Cancer: Updated Evidence Report and Systematic Review  
1363 for the US Preventive Services Task Force. *JAMA* **327**, 2334 (2022).
- 1364 17. Huang, J., Weinstein, S. J., Yu, K., Männistö, S. & Albanes, D. Association between  
1365 serum retinol and overall and cause-specific mortality in a 30-year prospective cohort  
1366 study. *Nat Commun* **12**, 6418 (2021).
- 1367 18. Bennett, D. A. & Du, H. An Overview of Methods and Exemplars of the Use of  
1368 Mendelian Randomisation in Nutritional Research. *Nutrients* **14**, 3408 (2022).
- 1369 19. Reay, W. R. & Cairns, M. J. Advancing the use of genome-wide association studies  
1370 for drug repurposing. *Nat Rev Genet* **22**, 658–671 (2021).

- 1371 20. Gueguen, S. *et al.* Genetic and environmental contributions to serum retinol and  
1372 alpha-tocopherol concentrations: the Stanislas Family Study. *Am J Clin Nutr* **81**, 1034–  
1373 1044 (2005).
- 1374 21. Ferrucci, L. *et al.* Common variation in the beta-carotene 15,15'-monooxygenase 1  
1375 gene affects circulating levels of carotenoids: a genome-wide association study. *Am J Hum*  
1376 *Genet* **84**, 123–133 (2009).
- 1377 22. Mondul, A. M. *et al.* Genome-wide association study of circulating retinol levels.  
1378 *Human Molecular Genetics* **20**, 4724–4731 (2011).
- 1379 23. Revez, J. A. *et al.* Genome-wide association study identifies 143 loci associated with  
1380 25 hydroxyvitamin D concentration. *Nat Commun* **11**, 1647 (2020).
- 1381 24. Manousaki, D. *et al.* Genome-wide Association Study for Vitamin D Levels Reveals  
1382 69 Independent Loci. *The American Journal of Human Genetics* **106**, 327–337 (2020).
- 1383 25. Manousaki, D. *et al.* Low-Frequency Synonymous Coding Variation in CYP2R1 Has  
1384 Large Effects on Vitamin D Levels and Risk of Multiple Sclerosis. *Am J Hum Genet* **101**,  
1385 227–238 (2017).
- 1386 26. Elsworth, B. *et al.* *The MRC IEU OpenGWAS data infrastructure*.  
1387 <http://biorxiv.org/lookup/doi/10.1101/2020.08.10.244293> (2020)  
1388 doi:10.1101/2020.08.10.244293.
- 1389 27. Perez-Martinez, P. *et al.* Association between glucokinase regulatory protein (GCKR)  
1390 and apolipoprotein A5 (APOA5) gene polymorphisms and triacylglycerol concentrations  
1391 in fasting, postprandial, and fenofibrate-treated states. *The American Journal of Clinical*  
1392 *Nutrition* **89**, 391–399 (2009).
- 1393 28. Fernandes Silva, L., Vangipurapu, J., Kuulasmaa, T. & Laakso, M. An intronic  
1394 variant in the GCKR gene is associated with multiple lipids. *Sci Rep* **9**, 10240 (2019).

- 1395 29. Yeh, K.-H. *et al.* Pleiotropic Effects of Common and Rare GCKR Exonic Mutations  
1396 on Cardiometabolic Traits. *Genes* **13**, 491 (2022).
- 1397 30. Zelent, B. *et al.* Analysis of the co-operative interaction between the allosterically  
1398 regulated proteins GK and GKRP using tryptophan fluorescence. *Biochemical Journal*  
1399 **459**, 551–564 (2014).
- 1400 31. Lai, C. S. L. FOXP2 expression during brain development coincides with adult sites  
1401 of pathology in a severe speech and language disorder. *Brain* **126**, 2455–2462 (2003).
- 1402 32. Mehta, M. B. *et al.* Hepatic protein phosphatase 1 regulatory subunit 3B (Ppp1r3b)  
1403 promotes hepatic glycogen synthesis and thereby regulates fasting energy homeostasis. *J*  
1404 *Biol Chem* **292**, 10444–10454 (2017).
- 1405 33. Russell, R. *et al.* Loss of the transcription factor MAFB limits  $\beta$ -cell derivation from  
1406 human PSCs. *Nat Commun* **11**, 2742 (2020).
- 1407 34. Artner, I. *et al.* MafA and MafB regulate genes critical to beta-cells in a unique  
1408 temporal manner. *Diabetes* **59**, 2530–2539 (2010).
- 1409 35. Leask, M. *et al.* Functional Urate-Associated Genetic Variants Influence Expression  
1410 of lincRNAs LINC01229 and MAFTRR. *Front Genet* **9**, 733 (2018).
- 1411 36. Werth, M. *et al.* Transcription factor TFCEP2L1 patterns cells in the mouse kidney  
1412 collecting ducts. *eLife* **6**, e24265 (2017).
- 1413 37. Ortega-Azorín, C. *et al.* Amino Acid Change in the Carbohydrate Response Element  
1414 Binding Protein Is Associated With Lower Triglycerides and Myocardial Infarction  
1415 Incidence Depending on Level of Adherence to the Mediterranean Diet in the PREDIMED  
1416 Trial. *Circ Cardiovasc Genet* **7**, 49–58 (2014).
- 1417 38. Dentin, R. *et al.* Liver-Specific Inhibition of ChREBP Improves Hepatic Steatosis and  
1418 Insulin Resistance in *ob/ob* Mice. *Diabetes* **55**, 2159–2170 (2006).

- 1419 39. Kathiresan, S. *et al.* Six new loci associated with blood low-density lipoprotein  
1420 cholesterol, high-density lipoprotein cholesterol or triglycerides in humans. *Nat Genet* **40**,  
1421 189–197 (2008).
- 1422 40. Beurel, E., Grieco, S. F. & Jope, R. S. Glycogen synthase kinase-3 (GSK3):  
1423 regulation, actions, and diseases. *Pharmacol Ther* **148**, 114–131 (2015).
- 1424 41. Mariotti, L., Pollock, K. & Guettler, S. Regulation of Wnt/ $\beta$ -catenin signalling by  
1425 tankyrase-dependent poly(ADP-ribosyl)ation and scaffolding. *Br J Pharmacol* **174**, 4611–  
1426 4636 (2017).
- 1427 42. Namwanje, M. & Brown, C. W. Activins and Inhibins: Roles in Development,  
1428 Physiology, and Disease. *Cold Spring Harb Perspect Biol* **8**, a021881 (2016).
- 1429 43. Nono Nankam, P. A. & Blüher, M. Retinol-binding protein 4 in obesity and metabolic  
1430 dysfunctions. *Molecular and Cellular Endocrinology* **531**, 111312 (2021).
- 1431 44. Davies, N. M., Holmes, M. V. & Davey Smith, G. Reading Mendelian randomisation  
1432 studies: a guide, glossary, and checklist for clinicians. *BMJ* k601 (2018)  
1433 doi:10.1136/bmj.k601.
- 1434 45. Sanderson, E. *et al.* Mendelian randomization. *Nat Rev Methods Primers* **2**, 6 (2022).
- 1435 46. McWhinney, S. R. *et al.* Obesity and brain structure in schizophrenia – ENIGMA  
1436 study in 3021 individuals. *Mol Psychiatry* **27**, 3731–3737 (2022).
- 1437 47. Burgess, S. & Labrecque, J. A. Mendelian randomization with a binary exposure  
1438 variable: interpretation and presentation of causal estimates. *Eur. J. Epidemiol.* **33**, 947–  
1439 952 (2018).
- 1440 48. Hada, M., Mondul, A. M., Weinstein, S. J. & Albanes, D. Serum Retinol and Risk of  
1441 Overall and Site-Specific Cancer in the ATBC Study. *American Journal of Epidemiology*  
1442 **189**, 532–542 (2020).

- 1443 49. Dragnev, K. H. *et al.* A phase I/II study of bexarotene with carboplatin and weekly  
1444 paclitaxel for the treatment of patients with advanced non-small cell lung cancer. *J Thorac*  
1445 *Dis* **10**, 5531–5537 (2018).
- 1446 50. Olsen, T. *et al.* Creatinine, total cysteine and uric acid are associated with serum  
1447 retinol in patients with cardiovascular disease. *Eur J Nutr* **59**, 2383–2393 (2020).
- 1448 51. Kučerová, K. *et al.* Determination of urinary retinol and creatinine as an early  
1449 sensitive marker of renal dysfunction. *J Chromatogr A* **1607**, 460390 (2019).
- 1450 52. Morrison, J., Knoblauch, N., Marcus, J. H., Stephens, M. & He, X. Mendelian  
1451 randomization accounting for correlated and uncorrelated pleiotropic effects using  
1452 genome-wide summary statistics. *Nat Genet* **52**, 740–747 (2020).
- 1453 53. Shinde, V. *et al.* Definition of transcriptome-based indices for quantitative  
1454 characterization of chemically disturbed stem cell development: introduction of the STOP-  
1455 Toxukn and STOP-Toxukk tests. *Arch Toxicol* **91**, 839–864 (2017).
- 1456 54. Krug, A. K. *et al.* Human embryonic stem cell-derived test systems for developmental  
1457 neurotoxicity: a transcriptomics approach. *Arch Toxicol* **87**, 123–143 (2013).
- 1458 55. Meganathan, K. *et al.* Neuronal developmental gene and miRNA signatures induced  
1459 by histone deacetylase inhibitors in human embryonic stem cells. *Cell Death Dis* **6**, e1756  
1460 (2015).
- 1461 56. Stephensen, C. B. & Gildengorin, G. Serum retinol, the acute phase response, and the  
1462 apparent misclassification of vitamin A status in the third National Health and Nutrition  
1463 Examination Survey,,. *The American Journal of Clinical Nutrition* **72**, 1170–1178 (2000).
- 1464 57. Inoue, E. *et al.* Identification of Glucose Transporter 4 Knockdown-dependent  
1465 Transcriptional Activation Element on the Retinol Binding Protein 4 Gene Promoter and  
1466 Requirement of the 20 S Proteasome Subunit for Transcriptional Activity. *Journal of*  
1467 *Biological Chemistry* **285**, 25545–25553 (2010).

- 1468 58. Kiser, P. D., Golczak, M. & Palczewski, K. Chemistry of the retinoid (visual) cycle.  
1469 *Chem Rev* **114**, 194–232 (2014).
- 1470 59. Hall, J. A., Grainger, J. R., Spencer, S. P. & Belkaid, Y. The Role of Retinoic Acid in  
1471 Tolerance and Immunity. *Immunity* **35**, 13–22 (2011).
- 1472 60. Stefanovic, S. & Zaffran, S. Mechanisms of retinoic acid signaling during  
1473 cardiogenesis. *Mechanisms of Development* **143**, 9–19 (2017).
- 1474 61. van Erp, T. G. M. *et al.* Cortical Brain Abnormalities in 4474 Individuals With  
1475 Schizophrenia and 5098 Control Subjects via the Enhancing Neuro Imaging Genetics  
1476 Through Meta Analysis (ENIGMA) Consortium. *Biol Psychiatry* **84**, 644–654 (2018).
- 1477 62. Reay, W. R. *et al.* Polygenic disruption of retinoid signalling in schizophrenia and a  
1478 severe cognitive deficit subtype. *Mol Psychiatry* **25**, 719–731 (2020).
- 1479 63. Crandall, J. *et al.* 13-cis-retinoic acid suppresses hippocampal cell division and  
1480 hippocampal-dependent learning in mice. *Proc Natl Acad Sci U S A* **101**, 5111–5116  
1481 (2004).
- 1482 64. Bremner, J. D., Shearer, K. D. & McCaffery, P. J. Retinoic acid and affective  
1483 disorders: the evidence for an association. *J Clin Psychiatry* **73**, 37–50 (2012).
- 1484 65. Li, C. *et al.* Use of isotretinoin and risk of depression in patients with acne: a  
1485 systematic review and meta-analysis. *BMJ Open* **9**, e021549 (2019).
- 1486 66. Lerner, V. *et al.* The retinoid X receptor agonist bexarotene relieves positive  
1487 symptoms of schizophrenia: a 6-week, randomized, double-blind, placebo-controlled  
1488 multicenter trial. *J Clin Psychiatry* **74**, 1224–1232 (2013).
- 1489 67. Tanprasertsuk, J. *et al.* Serum Carotenoids, Tocopherols, Total n-3 Polyunsaturated  
1490 Fatty Acids, and n-6/n-3 Polyunsaturated Fatty Acid Ratio Reflect Brain Concentrations in  
1491 a Cohort of Centenarians. *The Journals of Gerontology: Series A* **74**, 306–314 (2019).



- 1492 68. Staley, J. R. & Burgess, S. Semiparametric methods for estimation of a nonlinear  
1493 exposure-outcome relationship using instrumental variables with application to Mendelian  
1494 randomization. *Genet. Epidemiol.* **41**, 341–352 (2017).
- 1495 69. Burgess, S. *et al.* Guidelines for performing Mendelian randomization investigations.  
1496 *Wellcome Open Res* **4**, 186 (2019).
- 1497 70. Moore, C. *et al.* The INTERVAL trial to determine whether intervals between blood  
1498 donations can be safely and acceptably decreased to optimise blood supply: study protocol  
1499 for a randomised controlled trial. *Trials* **15**, 363 (2014).
- 1500 71. Bomba, L. *et al.* Whole-exome sequencing identifies rare genetic variants associated  
1501 with human plasma metabolites. *The American Journal of Human Genetics* **109**, 1038–  
1502 1054 (2022).
- 1503 72. Hail Team. Hail 0.2.
- 1504 73. Yin, X. *et al.* Genome-wide association studies of metabolites in Finnish men identify  
1505 disease-relevant loci. *Nat Commun* **13**, 1644 (2022).
- 1506 74. Laakso, M. *et al.* The Metabolic Syndrome in Men study: a resource for studies of  
1507 metabolic and cardiovascular diseases. *J Lipid Res* **58**, 481–493 (2017).
- 1508 75. The alpha-tocopherol, beta-carotene lung cancer prevention study: design, methods,  
1509 participant characteristics, and compliance. The ATBC Cancer Prevention Study Group.  
1510 *Ann Epidemiol* **4**, 1–10 (1994).
- 1511 76. Hayes, R. B. *et al.* Methods for etiologic and early marker investigations in the PLCO  
1512 trial. *Mutat Res* **592**, 147–154 (2005).
- 1513 77. Pasaniuc, B. *et al.* Fast and accurate imputation of summary statistics enhances  
1514 evidence of functional enrichment. *Bioinformatics* **30**, 2906–2914 (2014).
- 1515 78. Long, T. *et al.* Whole-genome sequencing identifies common-to-rare variants  
1516 associated with human blood metabolites. *Nat Genet* **49**, 568–578 (2017).

- 1517 79. Moayyeri, A., Hammond, C. J., Hart, D. J. & Spector, T. D. The UK Adult Twin  
1518 Registry (TwinsUK Resource). *Twin Res Hum Genet* **16**, 144–149 (2013).
- 1519 80. Shin, S.-Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat*  
1520 *Genet* **46**, 543–550 (2014).
- 1521 81. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and  
1522 richer datasets. *GigaSci* **4**, 7 (2015).
- 1523 82. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association  
1524 studies. *Bioinformatics* **26**, 2867–2873 (2010).
- 1525 83. Price, A. L. *et al.* Long-range LD can confound genome scans in admixed  
1526 populations. *Am. J. Hum. Genet.* **83**, 132–135; author reply 135-139 (2008).
- 1527 84. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping  
1528 and annotation of genetic associations with FUMA. *Nat Commun* **8**, 1826 (2017).
- 1529 85. Zhou, H. *et al.* FAVOR: functional annotation of variants online resource and  
1530 annotator for variation across the human genome. *Nucleic Acids Res* **51**, D1300–D1311  
1531 (2023).
- 1532 86. Kurki, M. I. *et al.* FinnGen provides genetic insights from a well-phenotyped isolated  
1533 population. *Nature* **613**, 508–518 (2023).
- 1534 87. Liu, Y. & Xie, J. Cauchy Combination Test: A Powerful Test With Analytic  $p$  -Value  
1535 Calculation Under Arbitrary Dependency Structures. *Journal of the American Statistical*  
1536 *Association* **115**, 393–402 (2020).
- 1537 88. Liu, Y. *et al.* ACAT: A Fast and Powerful  $p$  Value Combination Method for Rare-  
1538 Variant Analysis in Sequencing Studies. *Am. J. Hum. Genet.* **104**, 410–421 (2019).
- 1539 89. Schizophrenia Working Group of the Psychiatric Genomics Consortium *et al.* LD  
1540 Score regression distinguishes confounding from polygenicity in genome-wide association  
1541 studies. *Nat Genet* **47**, 291–295 (2015).

- 1542 90. Speed, D. & Balding, D. J. SumHer better estimates the SNP heritability of complex  
1543 traits from summary statistics. *Nat Genet* **51**, 277–284 (2019).
- 1544 91. Speed, D., Holmes, J. & Balding, D. J. Evaluating and improving heritability models  
1545 using summary statistics. *Nat Genet* **52**, 458–462 (2020).
- 1546 92. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-  
1547 wide association summary statistics. *Nat Genet* **47**, 1228–1235 (2015).
- 1548 93. Stephens, M. False discovery rates: a new deal. *Biostatistics* **18**, 275–294 (2017).
- 1549 94. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An Expanded View of Complex Traits: From  
1550 Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).
- 1551 95. Reay, W. R. *et al.* The genetic architecture of pneumonia susceptibility implicates  
1552 mucin biology and a relationship with psychiatric illness. *Nat Commun* **13**, 3756 (2022).
- 1553 96. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across  
1554 human tissues. *Science* **369**, 1318–1330 (2020).
- 1555 97. Wen, X., Pique-Regi, R. & Luca, F. Integrating molecular QTL data into genome-  
1556 wide genetic association analysis: Probabilistic assessment of enrichment and  
1557 colocalization. *PLoS Genet* **13**, e1006646 (2017).
- 1558 98. Zhang, J. *et al.* Plasma proteome analyses in individuals of European and African  
1559 ancestry identify cis-pQTLs and models for proteome-wide association studies. *Nat Genet*  
1560 **54**, 593–602 (2022).
- 1561 99. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD:  
1562 predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids*  
1563 *Research* **47**, D886–D894 (2019).
- 1564 100. Boyle, A. P. *et al.* Annotation of functional variation in personal genomes using  
1565 RegulomeDB. *Genome Res* **22**, 1790–1797 (2012).

- 1566 101. Ghoussaini, M. *et al.* Open Targets Genetics: systematic identification of trait-  
1567 associated genes using large-scale genetics and functional genomics. *Nucleic Acids*  
1568 *Research* **49**, D1311–D1320 (2021).
- 1569 102. Wakefield, J. Bayes factors for genome-wide association studies: comparison with P-  
1570 values. *Genet Epidemiol* **33**, 79–86 (2009).
- 1571 103. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic  
1572 association studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
- 1573 104. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package  
1574 for differential expression analysis of digital gene expression data. *Bioinformatics* **26**,  
1575 139–140 (2010).
- 1576 105. Raudvere, U. *et al.* g:Profiler: a web server for functional enrichment analysis and  
1577 conversions of gene lists (2019 update). *Nucleic Acids Research* **47**, W191–W198 (2019).
- 1578 106. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association  
1579 studies. *Nat. Genet.* **48**, 245–252 (2016).
- 1580 107. Reay, W. R. & Cairns, M. J. Pairwise common variant meta-analyses of  
1581 schizophrenia with other psychiatric disorders reveals shared and distinct gene and gene-  
1582 set associations. *Transl Psychiatry* **10**, 134 (2020).
- 1583 108. Adams, D. M., Reay, W. R. & Cairns, M. J. Multiomic prioritisation of risk genes for  
1584 anorexia nervosa. *Psychol Med* 1–9 (2023) doi:10.1017/S0033291723000235.
- 1585 109. Reay, W. R. *et al.* Genetics-informed precision treatment formulation in  
1586 schizophrenia and bipolar disorder. *The American Journal of Human Genetics* **109**, 1620–  
1587 1637 (2022).
- 1588 110. Reay, W. R. *et al.* Genetic association and causal inference converge on  
1589 hyperglycaemia as a modifiable factor to improve lung function. *eLife* **10**, e63115 (2021).

- 1590 111. Wainberg, M. *et al.* Opportunities and challenges for transcriptome-wide association  
1591 studies. *Nat Genet* **51**, 592–599 (2019).
- 1592 112. Burgess, S., Butterworth, A. & Thompson, S. G. Mendelian randomization analysis  
1593 with multiple genetic variants using summarized data. *Genet. Epidemiol.* **37**, 658–665  
1594 (2013).
- 1595 113. Bowden, J., Hemani, G. & Davey Smith, G. Invited Commentary: Detecting  
1596 Individual and Global Horizontal Pleiotropy in Mendelian Randomization-A Job for the  
1597 Humble Heterogeneity Statistic? *Am J Epidemiol* **187**, 2681–2685 (2018).
- 1598 114. Bowden, J., Davey Smith, G. & Burgess, S. Mendelian randomization with invalid  
1599 instruments: effect estimation and bias detection through Egger regression. *Int J Epidemiol*  
1600 **44**, 512–525 (2015).
- 1601 115. Burgess, S., Bowden, J., Fall, T., Ingelsson, E. & Thompson, S. G. Sensitivity  
1602 Analyses for Robust Causal Inference from Mendelian Randomization Analyses with  
1603 Multiple Genetic Variants. *Epidemiology* **28**, 30–42 (2017).
- 1604 116. Bowden, J., Davey Smith, G., Haycock, P. C. & Burgess, S. Consistent Estimation in  
1605 Mendelian Randomization with Some Invalid Instruments Using a Weighted Median  
1606 Estimator. *Genet. Epidemiol.* **40**, 304–314 (2016).
- 1607 117. Burgess, S., Zuber, V., Gkatzionis, A. & Foley, C. N. Modal-based estimation via  
1608 heterogeneity-penalized weighting: model averaging for consistent and efficient estimation  
1609 in Mendelian randomization when a plurality of candidate instruments are valid. *Int J*  
1610 *Epidemiol* **47**, 1242–1254 (2018).
- 1611 118. Reay, W. R. *et al.* Genetic estimates of correlation and causality between blood-based  
1612 biomarkers and psychiatric disorders. *Sci. Adv.* **8**, eabj8969 (2022).
- 1613 119. Slob, E. A. W. & Burgess, S. A comparison of robust Mendelian randomization  
1614 methods using summary data. *Genetic Epidemiology* **44**, 313–329 (2020).

- 1615 120. Bowden, J. *et al.* Assessing the suitability of summary data for two-sample Mendelian  
1616 randomization analyses using MR-Egger regression: the role of the I<sup>2</sup> statistic. *Int J*  
1617 *Epidemiol* **45**, 1961–1974 (2016).
- 1618 121. Lu, Y. *et al.* New loci for body fat percentage reveal link between adiposity and  
1619 cardiometabolic disease risk. *Nat Commun* **7**, 10495 (2016).
- 1620 122. Willer, C. J. *et al.* Discovery and refinement of loci associated with lipid levels. *Nat*  
1621 *Genet* **45**, 1274–1283 (2013).
- 1622 123. Sanderson, E., Spiller, W. & Bowden, J. *Testing and Correcting for Weak and*  
1623 *Pleiotropic Instruments in Two-Sample Multivariable Mendelian Randomisation.*  
1624 <http://biorxiv.org/lookup/doi/10.1101/2020.04.02.021980> (2020)  
1625 doi:10.1101/2020.04.02.021980.
- 1626 124. Grant, A. J. & Burgess, S. Pleiotropy robust methods for multivariable Mendelian  
1627 randomization. *Statistics in Medicine* **40**, 5813–5830 (2021).
- 1628 125. Freshour, S. *et al.* *Integration of the Drug-Gene Interaction Database (DGIdb) with*  
1629 *open crowdsourcing efforts.* <http://biorxiv.org/lookup/doi/10.1101/2020.09.18.301721>  
1630 (2020) doi:10.1101/2020.09.18.301721.
- 1631 126. Wishart, D. S. *et al.* DrugBank 5.0: a major update to the DrugBank database for  
1632 2018. *Nucleic Acids Res* **46**, D1074–D1082 (2018).
- 1633 127. Subramanian, A. *et al.* A Next Generation Connectivity Map: L1000 Platform and the  
1634 First 1,000,000 Profiles. *Cell* **171**, 1437–1452.e17 (2017).
- 1635 128. International Schizophrenia Consortium *et al.* Common polygenic variation  
1636 contributes to risk of schizophrenia and bipolar disorder. *Nature* **460**, 748–752 (2009).
- 1637 129. Reales, G., Vigorito, E., Kelemen, M. & Wallace, C. RápidoPGS: a rapid polygenic  
1638 score calculator for summary GWAS data without a test dataset. *Bioinformatics* **37**, 4444–  
1639 4450 (2021).

- 1640 130. Choi, S. W. & O'Reilly, P. F. PRSice-2: Polygenic Risk Score software for biobank-  
1641 scale data. *GigaScience* **8**, giz082 (2019).
- 1642 131. Stasinopoulos, M. D., Rigby, R. A., Heller, G. Z., Voudouris, V. & Bastiani, F. D.  
1643 *Flexible Regression and Smoothing: Using GAMLSS in R*. (Chapman and Hall/CRC,  
1644 2017). doi:10.1201/b21973.
- 1645 132. Rigby, R. A. & Stasinopoulos, D. M. Generalized Additive Models for Location,  
1646 Scale and Shape. *Journal of the Royal Statistical Society Series C: Applied Statistics* **54**,  
1647 507–554 (2005).
- 1648