

1 **IBD risk loci are enriched in multigenic regulatory modules encompassing**
 2 **causative genes**

3 *Yukihide Momozawa^{1,2*}, Julia Dmitrieva^{1*}, Emilie Théâtre¹, Valérie Deffontaine¹,*
 4 *Souad Rahmouni¹, Benoît Charlotiaux¹, François Crins¹, Elisa Docampo¹, Mahmoud*
 5 *Elansary¹, Ann-Stephan Gori¹, Christelle Lecut³, Rob Mariman¹, Myriam Mni¹,*
 6 *Cécile Oury³, Ilya Altukhov⁴, Dmitry Alexeev⁵, Yuri Aulchenko⁶⁻⁸, Leila Amininejad⁹,*
 7 *G. Bouma¹⁰, F. Hoentjen¹¹, M. Löwenberg¹², B. Oldenburg¹³, M.J. Pierik¹⁴, A.E. van*
 8 *der Meulen-de Jong¹⁵, C.J van der Woude¹⁶, Marijn C Visschedijk¹⁷, The*
 9 *International IBD Genetics Consortium[#], Mark Lathrop¹⁸, Jean-Pierre Hugot¹⁹,*
 10 *Rinse K. Weersma¹⁷, Martine De Vos²⁰, Denis Franchimont⁹, Severine Vermeire²¹,*
 11 *Michiaki Kubo², Edouard Louis²², Michel Georges¹.*

12 ¹Unit of Animal Genomics, WELBIO, GIGA-R & Faculty of Veterinary Medicine, University of Liège
 13 (B34), 1 Avenue de l'Hôpital, 4000-Liège, Belgium. ²Laboratory for Genotyping Development,
 14 RIKEN Center for Integrative Medical Science, 1-7-22, Suehiro-cho, Tsurumi-ku, Yokohama,
 15 Kanagawa 230-0045, Japan. ³Laboratory of Thrombosis and Hemostasis, GIGA-R, University of
 16 Liège (B34), 1 Avenue de l'Hôpital, 4000-Liège, Belgium. ⁴Institute of Physics and Technology,
 17 Institutskiy pereulok 9, Dolgoprudny 141700, Russian Federation. ⁵Novosibirsk State University,
 18 Pirogova ave. 2, 630090 Novosibirsk, Russian Federation, ⁶PolyOmica, Het Vlaggeschip 61, 5237
 19 PA, 's-Hertogenbosch, The Netherlands. ⁷Institute of Cytology and Genetics SD RAS, Lavrentyeva
 20 ave. 10, 630090 Novosibirsk, Russia, ⁸Centre for Global Health Research, Usher Institute of
 21 Population Health Sciences and Informatics, University of Edinburgh, Teviot Place, Edinburgh,
 22 EH8 9AG, UK. ⁹Gastroentérologie Médicale, Faculté de Médecine, Université Libre de Bruxelles,
 23 Route de Lennik 808, 1070 Anderlecht, Belgium. ¹⁰Department of Gastroenterology and
 24 Hepatology, VU University Medical Centre, Amsterdam, The Netherlands. ¹¹ Department of
 25 Gastroenterology and Hepatology, University Medical Centre St. Radboud, Nijmegen, The
 26 Netherlands. ¹²Department of Gastroenterology and Hepatology, Amsterdam Medical Centre,
 27 Amsterdam, The Netherlands. ¹³Department of Gastroenterology and Hepatology, University
 28 Medical Centre Utrecht, Utrecht, The Netherlands. ¹⁴Department of Gastroenterology and
 29 Hepatology, University Medical Centre Maastricht, Maastricht, The Netherlands. ¹⁵Department of
 30 Gastroenterology and Hepatology, Leiden University Medical Centre, Leiden, The Netherlands.
 31 ¹⁶Department of Gastroenterology and Hepatology, Erasmus Medical Centre, the Netherlands,
 32 Rotterdam, The Netherlands. ¹⁷Department of Gastroenterology and Hepatology, University of
 33 Groningen and University Medical Center Groningen, Hanzeplein 1, 9713 GZ Groningen, The
 34 Netherlands. ¹⁸McGill University Centre for Molecular and Computational Genomics, 740 Dr
 35 Penfield Avenue, Montreal, Quebec, Canada. ¹⁹UMR 1149 INSERM/Université Paris-Diderot
 36 Sorbonne Paris-Cité, Assistance Publique Hôpitaux de Paris, 48 Bd Sérurier, 75019 Paris, France.
 37 ²⁰Department of Gastroenterology, University Hospital, De Pintelaan 185, 9000 Gent, Belgium.
 38 ²¹Translational Research in Gastrointestinal Disorders, Department of Clinical and Experimental
 39 Medicine, KU Leuven, UZ Herestraat 49, 3000 Leuven, Belgium. ²²CHU-Liège and Unit of
 40 Gastroenterology, GIGA-R & Faculty of Medicine, University of Liège, 1 Avenue de l'Hôpital, 4000-
 41 Liège, Belgium.

42

43 * Contributed equally.

44 # The complete list of IIBDGC members is provided at the end of the manuscript

45 Correspondence: michel.georges@ulg.ac.be

46 **GWAS have identified >200 risk loci for Inflammatory Bowel Disease (IBD).**
47 **The majority of disease associations are known to be driven by regulatory**
48 **variants. To identify the putative causative genes that are perturbed by**
49 **these variants, we generate a large transcriptome dataset (9 disease-**
50 **relevant cell types) and identify 23,650 cis-eQTL. We show that these are**
51 **determined by ~9,720 regulatory modules, of which ~3,000 operate in**
52 **multiple tissues and ~970 on multiple genes. We identify regulatory**
53 **modules that drive the disease association for 63 of the 200 risk loci, and**
54 **show that these are enriched in multigenic modules. We resequence 45 of**
55 **the corresponding 100 candidate genes in 6,600 Crohn disease (CD) cases**
56 **and 5,500 controls and show that they are significantly enriched in**
57 **causative genes. Our analyses indicate that ≥ 10 -fold larger sample sizes**
58 **will be required to demonstrate the causality of individual genes using**
59 **standard burden tests.**

60

61 **INTRODUCTION**

62 Genome Wide Association Studies (GWAS) scan the entire genome for statistical
63 associations between common variants and disease status in large case-control
64 cohorts. GWAS have identified tens to hundreds of risk loci for nearly all studied
65 common complex diseases of human¹. The study of Inflammatory Bowel Disease
66 (IBD) has been particularly successful, with more than 200 confirmed risk loci
67 reported to date^{2,3}. As a result of the linkage disequilibrium (LD) patterns in the
68 human genome (limiting the mapping resolution of association studies), GWAS-
69 identified risk loci typically span ~ 250 kilobases, encompassing an average of ~
70 5 genes (numbers ranging from zero (“gene deserts”) to more than 50) and
71 hundreds of associated variants. Contrary to widespread misconception, the
72 *causative* variants and genes remain unknown for the vast majority of GWAS-
73 identified risk loci. Yet, this remains a critical goal in order to reap the full
74 benefits of GWAS in identifying new drug targets and developing effective
75 predictive and diagnostic tools. It is the main objective of post-GWAS studies.

76 Distinguishing the few causative variants (i.e. the variants that are directly
77 causing the gene perturbation) from the many neutral variants that are only

78 associated with the disease because they are in LD with the former in the studied
79 population, requires the use of sophisticated fine-mapping methods applied to
80 very large, densely genotyped datasets⁴, ideally followed-up by functional
81 studies⁵. Using such approaches, 18 causative variants for IBD were recently
82 fine-mapped at single base pair resolution, and 51 additional ones at ≤ 10 base
83 pair resolution⁴.

84 A minority of causative variants are coding, i.e. they alter the amino-acid
85 sequence of the encoded protein. In such cases, and particularly if multiple such
86 causative coding variants are found in the same gene (i.e. in case of allelic
87 heterogeneity), the corresponding causative gene is unambiguously identified.
88 In the case of IBD, causative genes have been identified for \sim ten risk loci on the
89 basis of such “independently” (i.e. not merely reflecting LD with other variants)
90 associated coding variants, including *NOD2*, *ATG16L1*, *IL23R*, *CARD9*, *FUT2* and
91 *TYK2*^{4,6-9}.

92 For the majority of risk loci, the GWAS signals are not driven by coding variants.
93 They must therefore be driven by common regulatory variants, i.e. variants that
94 perturb the expression levels of one (or more) target genes in one (or more)
95 disease relevant cell types⁴. Merely reflecting the proportionate sequence space
96 that is devoted to the different layers of gene regulation (transcriptional,
97 posttranscriptional, translational, posttranslational), the majority of regulatory
98 variants are likely to perturb components of “gene switches” (promoters,
99 enhancers, insulators), hence affecting transcriptional output. Indeed, fine-
100 mapped non-coding variants are enriched in known transcription-factor binding
101 sites and epigenetic signatures marking gene switch components⁴. Hence, the
102 majority of common causative variants underlying inherited predisposition to
103 common complex diseases must drive cis-eQTL (expression quantitative trait
104 loci) affecting the causative gene(s) in one or more disease relevant cell types.
105 The corresponding cis-eQTL are expected to operate prior to disease onset, and –
106 driven by common variants – detectable in cohorts of healthy individuals of
107 which most will never develop the disease. The term cis-eQTL refers to the fact
108 that the regulatory variants that drive them only affect the expression of
109 genes/alleles residing on the same DNA molecule, typically no more than one
110 megabase away. Causative variants, whether coding or regulatory, may

111 secondarily perturb the expression of genes/alleles located on different DNA
112 molecules, generating trans-eQTL. Some of these trans-eQTL may participate in
113 the disease process.

114 Cis-eQTL effects are known to be very common, affecting more than 50% of
115 genes¹⁰. Hence, finding that variants associated with a disease are also
116 associated with changes in expression levels of a neighboring gene is not
117 sufficient to incriminate the corresponding genes as causative. Firstly, one has
118 to show that the local association signal for the disease and for the eQTL are
119 driven by the same causative variants. A variety of “colocalisation” methods
120 have been developed to that effect¹¹⁻¹³. Secondly, regulatory variants may affect
121 elements that control the expression of multiple genes¹⁴, which may not all
122 contribute to the development of the disease, i.e. be causative. Thus, additional
123 evidence is needed to obtain formal proof of gene causality. In humans, the only
124 formal test of gene causality that is applicable is the family of “burden” tests, i.e.
125 the search for a differential burden of disruptive mutations in cases and controls,
126 which is expected only for causative genes¹⁵. Burden tests rely on the
127 assumption that – in addition to the common, mostly regulatory variants that
128 drive the GWAS signal – the causative gene will be affected by low frequency and
129 rare causative variants, including coding variants. Thus, the burden test makes
130 the assumption that allelic heterogeneity is common, which is supported by the
131 pervasiveness of allelic heterogeneity of Mendelian diseases in humans¹⁶.
132 Burden tests compare the distribution of rare coding variants between cases and
133 controls¹⁵. The signal-to-noise ratio of the burden test can be increased by
134 restricting the analysis to coding variants that have a higher probability to
135 disrupt protein function¹⁵. In the case of IBD, burden tests have been used to
136 prove the causality of *NOD2*, *IL23R* and *CARD9*^{6,8,9}. A distinct and very elegant
137 genetic test of gene causality is the reciprocal hemizyosity test, and the related
138 quantitative complementation assay^{17,18}. However, with few exceptions^{19,20}, it
139 has only been applied in model organisms in which gene knock-outs can be
140 readily generated²¹.

141 In this paper, we describe the generation of a new and large dataset for eQTL
142 analysis (350 healthy individuals) in nine cell types that are potentially relevant
143 for IBD. We identify and characterize ~24,000 cis-eQTL. By comparing disease

144 and eQTL association patterns using a newly developed statistic, we identify 99
145 strong positional candidate genes in 63 GWAS-identified risk loci. We
146 resequence the 555 exons of 45 of these in 6,600 cases and 5,500 controls in an
147 attempt to prove their causality by means of burden tests. The outcome of this
148 study is relevant to post-GWAS studies of all common complex disease in
149 humans.

150

151 **RESULTS**

152 ***Clustering cis-eQTL into regulatory modules***

153 We generated transcriptome data for six circulating immune cell types (CD4+ T
154 lymphocytes, CD8+ T lymphocytes, CD19+ B lymphocytes, CD14+ monocytes,
155 CD15+ granulocytes, platelets) as well as ileal, colonic and rectal biopsies (IL, TR,
156 RE), collected from 323 healthy Europeans (141 men, 182 women, average age
157 56 years, visiting the clinic as part of a national screening campaign for colon
158 cancer) using Illumina HT12 arrays (CEDAR dataset; Methods). IBD being
159 defined as an inappropriate mucosal immune response to a normal commensal
160 gut flora²², these nine cell types can all be considered to be potentially disease-
161 relevant. Using standard methods based on linear regression and one megabase
162 windows centered on the position of the interrogating probe (Methods), we
163 identified significant cis-eQTL (FDR < 0.05) for 8,804 of 18,580 tested probes
164 (corresponding to 7,216 of 13,615 tested genes) in at least one tissue, amounting
165 to a total of 23,650 cis-eQTL effects (Supplementary Data 1). When a gene shows
166 a cis-eQTL in more than one tissue, the corresponding “eQTL association
167 patterns” (EAP) (i.e. the distribution of association $-\log(p)$ values for all the
168 variants in the region of interest) are expected to be similar if determined by the
169 same regulatory variants, and dissimilar otherwise. Likewise, if several
170 neighboring genes show cis-eQTL in the same or distinct tissues, the
171 corresponding EAP are expected to be similar if determined by the same
172 regulatory variants, and dissimilar otherwise (Fig. 1). We devised the ϑ metric
173 to measure the similarity between association patterns (Methods). ϑ is a
174 correlation measure for paired $-\log(p)$ values (for the two eQTL that are being
175 compared) that ranges between -1 and +1. ϑ shrinks to zero if Pearson’s

176 correlation between paired $-\log(p)$ values does not exceed a chosen threshold (i.e.
177 if the EAP are not similar). ϑ approaches +1 when the two EAP are similar and
178 when variants that increase expression in eQTL 1 consistently increase
179 expression in eQTL 2. ϑ approaches -1 when the two EAP are similar and when
180 variants that increase expression in eQTL 1 consistently decrease expression in
181 eQTL 2. ϑ gives more weight to variants with high $-\log(p)$ for at least one EAP
182 (i.e. it gives more weight to eQTL peaks). Based on the known distribution of ϑ
183 under H_0 (i.e. eQTL determined by distinct variants in the same region) and H_1
184 (i.e. eQTL determined by the same variants), we selected a threshold value
185 $|\vartheta| > 0.60$ to consider that two EAP were determined by the same variant. This
186 corresponds to a false positive rate of 0.05, and a false negative rate of 0.23
187 (Supplementary Fig. 1). We then grouped EAP in “cis-acting regulatory modules”
188 (cRM) using $|\vartheta|$ and a single-link clustering approach (i.e. an EAP needs to have
189 $|\vartheta| > 0.60$ with at least one member of the cluster to be assigned to that cluster).
190 Clusters were visually examined and 29 single edges connecting otherwise
191 unlinked and yet tight clusters manually removed (Supplementary Fig. 2).

192 Using this approach, we clustered the 23,650 effects in 9,720 distinct “cis-
193 regulatory modules” (cRM), encompassing cis-eQTL with similar EAP
194 (Supplementary Data 2). Sixty-eight percent of cRM were gene- and tissue-
195 specific, 22% were gene-specific but operating across multiple tissues (≤ 9
196 tissues, average 3.5), and 10% were multi-genic (≤ 11 genes, average 2.5) and
197 nearly always multi-tissue (Fig. 2&3, Supplementary Fig. 2). In this, cRM are
198 considered gene-specific if the EAPs in the cluster concern only one gene, and
199 tissue-specific if the EAP in the cluster concern only one of the nine cell types.
200 They are, respectively, multigenic and multi-tissue otherwise. cRM operating
201 across multiple tissues tended to affect multiple genes ($r = 0.47$; $p < 10^{-6}$). In
202 such cRM, the direction of the effects tended to be consistent across tissues and
203 genes ($p < 10^{-6}$). Nevertheless, we observed at least 55 probes with effect of
204 opposite sign in distinct cell types ($\vartheta \leq -0.9$), i.e. the corresponding regulatory
205 variants increases transcript levels in one cell type while decreasing them in
206 another (Fig. 4 and Supplementary Data 3). Individual tissues allowed for the
207 detection of 7 to 33% of all cRM, and contributed 3 to 14% unique cRM
208 (Supplementary Fig. 3). Sixty-nine percent of cRM were only detected in one cell

209 type. The rate of cRM sharing between cell types reflects known ontogenic
210 relations. Considering cRM shared by only two cell types (i.e. what jointly
211 differentiates these two cell types from all other), revealed the close proximity of
212 the CD4-CD8, CD14-CD15, ileum-colon, and colon-rectum pairs. Adding
213 information of cRM shared by up to six cell types grouped lymphoid (CD4, CD8,
214 CD19), myeloid (CD14, CD15 but not platelets), and intestinal (ileum, colon and
215 rectum) cells. Adding cRM with up to nine cell types revealed a link between
216 ileum and blood cells, possibly reflecting the presence of blood cells in the ileal
217 biopsies (Fig. 5).

218 ***cRM matching IBD association signals are often multigenic***

219 If regulatory variants affect disease risk by perturbing gene expression, the
220 corresponding “disease association patterns” (DAP) and EAP are expected to be
221 similar, even if obtained in distinct cohorts (yet with same ethnicity) (Fig. 6).
222 We confronted DAP and EAP using the ϑ statistic and threshold ($|\vartheta| > 0.60$)
223 described above for 200 GWAS-identified IBD risk loci. DAP for Crohn’s disease
224 and Ulcerative Colitis were obtained from the International IBD Genetics
225 Consortium (IIBDGC)^{2,3}, EAP from the CEDAR dataset.

226 The probability that two unrelated association signals in a chromosome region of
227 interest are similar (i.e. have high $|\vartheta|$ value) is affected by the degree of LD in the
228 region. If the LD is high it is more likely that two association signals are similar
229 by chance. To account for this, we generated EAP- and locus-specific
230 distributions of $|\vartheta|$ by simulating eQTL explaining the same variance as the
231 studied eQTL, yet driven by 100 variants that were randomly selected in the risk
232 locus (matched for MAF), and computing $|\vartheta|$ with the DAP for all of these. The
233 resulting empirical distribution of $|\vartheta|$ was used to compute the probability to
234 obtain a value of $|\vartheta|$ as high or higher than the observed one, by chance alone
235 (Methods).

236 Strong correlations between DAP and EAP ($|\vartheta| > 0.6$, associated with low
237 empirical p-values) were observed for at least 63 IBD risk loci, involving 99
238 genes (range per locus: 1-6) (Table 1, Fig. 7, Supplementary Data 4). Increased
239 disease risk was associated equally frequently with increased as with decreased
240 expression ($p_{CD} = 0.48$; $p_{UC} = 0.88$). An open-access website has been prepared

241 to visualize correlated DAP-EAP within their genomic context ([http://cedar-](http://cedar-web.giga.ulg.ac.be)
242 [web.giga.ulg.ac.be](http://cedar-web.giga.ulg.ac.be)). Genes with highest $|\vartheta|$ values (≥ 0.9) include known IBD
243 causative genes (f.i. *ATG16L1*, *CARD9*, *FUT2*), known immune regulators (f.i.
244 *IL18R1*, *IL6ST*, *THEMIS*), as well as genes with as of yet poorly defined function in
245 the context of IBD (f.i. *APEH*, *ANKRD55*, *CISD1*, *CPEB4*, *DOCK7*, *ERAP2*, *GNA12*,
246 *GPX1*, *GSDMB*, *ORMDL3*, *SKAP2*, *UBE2L3*, *ZMIZ1*) (Supplementary Note 1).

247 The eQTL link with IBD has not been reported before for at least 47 of the 99
248 reported genes (Table 1). eQTL links with IBD have been previously reported for
249 111 additional genes, not mentioned in Table 1. Our data support these links for
250 19 of them, however, with $|\vartheta| \leq 0.6$ (Supplementary Data 5). We applied SMR¹³
251 as alternative colocalisation method to our data. Using a Bonferroni-corrected
252 threshold of $\leq 2.5 \times 10^{-5}$ for p_{SMR} and ≥ 0.05 for p_{HEIDI} , SMR detected 35 of the 99
253 genes selected with ϑ (Supplementary Data 4). Using the same thresholds, SMR
254 detected nine genes that were not selected by ϑ . Of these, three (*ADAM15*,
255 *AHSA2*, *UBA7*) had previously been reported by others, while six (*FAM189B*,
256 *QRICH1*, *RBM6*, *TAP2*, *ADO*, *LGALS9*) were not. Of these six, three (*RBM6*, *TAP2*,
257 *ADO*) were characterized by $0.45 < |\vartheta| < 0.6$ (Supplementary Data 5).

258 Using an early version of the CEDAR dataset, significant (albeit modest)
259 enrichment of overlapping disease and eQTL signals was reported for CD4, ileum,
260 colon and rectum, focusing on 76 of 97 studied IBD risk loci (MAF of disease
261 variant > 0.05)⁴. By pre-correcting fluorescence intensities with 23 to 53
262 (depending on cell type) principal components to account for unidentified
263 confounders (Methods), we increased the number of significant eQTL from 480
264 to 880 in the corresponding 97 regions (11,964 to 23,650 for the whole genome).
265 We repeated the enrichment analysis focusing on 63 of the same 97 IBD loci (CD
266 risk loci; MAF of disease variant > 0.05), using three colocalisation methods
267 including ϑ (Methods). We observed a systematic excess overlap in all analyzed
268 cell types (2.5-fold on average). The enrichment was very significant with the
269 three methods in CD4 and CD8 (Supplementary Table 1).

270 The 400 analyzed DAP (200 CD and 200 UC) were found to match 76 cRM (in 63
271 risk loci) with $|\vartheta| > 0.6$ (Table 1), of which 25 are multigenic. Knowing that
272 multigenic cRM represent 10% of all cRM (967/9,720), 25/76 (i.e. 33%)

273 corresponds to a highly significant 3-fold enrichment ($p < 10^{-9}$). To ensure that
274 this apparent enrichment was not due to the fact that multigenic cRM have more
275 chance to match DAP (as by definition multiple EAP are tested for multigenic
276 cRM), we repeated the enrichment analysis by randomly sampling only one
277 representative EAP per cRM in the 200 IBD risk loci. The frequency of multigenic
278 cRM amongst DAP-matching cRM averaged 0.22, and was never ≤ 0.10
279 ($p \leq 10^{-5}$) (Supplementary Fig. 4). In loci with high LD, EAP driven by distinct
280 regulatory variants (yet in high LD) may erroneously be merged in the same cRM.
281 To ensure that the observed enrichment in multigenic cRM was not due to higher
282 levels of LD, we compared the LD-based recombination rate of the 63 cRM-
283 matching IBD risk loci with that of the rest of the genome²³. The genome-
284 average recombination rate was 1.23 centimorgan per megabase (cM/Mb), while
285 that of the 63 IBD risk loci was 1.34 cM/Mb, i.e. less LD in the 63 cRM-matching
286 IBD risk loci than in the rest of the genome. We further compared the average
287 recombination rate in the 63 cRM-matching IBD regions with that of sets of 63
288 loci centered on randomly drawn cRM (from the list of 9,720), matched for size
289 and chromosome number (as cM/Mb is affected by chromosome size). The
290 average recombination rate around all cRM was 1.43 cM/Mb, and this didn't
291 differ significantly from the 63 cRM-matching IBD regions ($p=0.46$)
292 (Supplementary Fig. 5). Therefore, the observed enrichment cannot be
293 explained by a higher LD in the 63 studied IBD risk loci. Taken together, EAP
294 that are strongly correlated with DAP ($|\rho| \geq 0.60$), map to regulatory modules
295 that are 2- to 3-fold enriched in multigenic cRM when compared to the genome
296 average and include four of the top 10 (of 9,720) cRM ranked by number of
297 affected genes.

298 ***DAP-matching cRM are enriched in causative genes for IBD***

299 For truly causative genes, the burden of rare disruptive variants is expected to
300 differ between cases and controls²⁴. We therefore performed targeted
301 sequencing for the 555 coding exons (~88 Kb) of 38 genes selected amongst
302 those with strongest DAP-EAP correlations, plus seven genes with suggestive
303 DAP-EAP evidence backed by literature (Table 1), in 6,597 European CD cases
304 and 5,502 matched controls (ref. 25 and Methods). Eighteen of these were part
305 of single-gene cRM and the only gene highlighted in the corresponding locus. The

306 remaining 27 corresponded to multi-gene cRM mapping to 15 risk loci. We
307 added the well-established *NOD2* and *IL23R* causative IBD genes as positive
308 controls. We identified a total of 174 loss-of-function (LoF) variants, 2,567
309 missense variants (of which 991 predicted by SIFT²⁶ to be damaging and
310 Polyphen-2²⁷ to be either possibly or probably damaging), and 1,434
311 synonymous variants (Fig. 8 and Supplementary Data 6). 1,781 of these were
312 also reported in the Genome Aggregation Database²⁸ with nearly identical allelic
313 frequencies (Supplementary Fig. 6). We designed a gene-based burden test to
314 simultaneously evaluate hypothesis (i): all disruptive variants enriched in cases
315 (when $\theta < 0$; risk variants) or all disruptive variants enriched in controls (when θ
316 > 0 ; protective variants), and hypothesis (ii): some disruptive variants enriched
317 in cases and others in controls. Hypothesis (i) was tested with CAST²⁹, and
318 hypothesis (ii) with SKAT³⁰ (Methods). We restricted the analysis to 1,141 LoF
319 and damaging missense variants with minor allele frequency (MAF) ≤ 0.005 to
320 ensure that any new association signal would be independent of the signals from
321 common and low frequency variants having led to the initial identification and
322 fine-mapping of the corresponding loci⁴. For *NOD2* ($p = 6.9 \times 10^{-7}$) and *IL23R* (p
323 $= 1.8 \times 10^{-4}$), LoF and damaging variants were significantly enriched in
324 respectively cases and controls as expected. When considering the 45 newly
325 tested genes as a whole, we observed a significant ($p = 6.9 \times 10^{-4}$) shift towards
326 lower p-values when compared to expectation, while synonymous variants
327 behaved as expected ($p = 0.66$) (Fig. 9 and Supplementary Data 7). This strongly
328 suggests that the sequenced list includes causative genes. *CARD9*, *TYK2* and
329 *FUT2* have recently been shown to be causative genes based on disease-
330 associated low-frequency coding variants (MAF > 0.005)⁴. The shift towards
331 lower p-values remained significant without these ($p = 1.7 \times 10^{-3}$), pointing
332 towards novel causative genes amongst the 42 remaining candidate genes.

333 ***Proving gene causality requires larger case-control cohorts.***

334 Despite the significant shift towards lower p-values when considering the 45
335 genes jointly, none of these were individually significant when accounting for
336 multiple testing ($p \leq \frac{0.05}{2*45} \approx 0.0006$) (Supplemental Data 7). Near identical
337 results were obtained when classifying variants using the Combined Annotation

338 Dependent Depletion (CADD) tool³¹ instead of SIFT/PolyPhen-2 (Supplementary
339 Data 7). We explored three approaches to increase the power of the burden test.
340 The first built on the observation that cRM matching DAP are enriched in
341 multigenic modules. This suggests that part of IBD risk loci harbor multiple co-
342 regulated and hence functionally related genes, of which several (rather than one,
343 as generally assumed) may be causally involved in disease predisposition. To
344 test this hypothesis, we designed a module- rather than gene-based burden test
345 (Methods). However, none of the 30 tested modules reached the experiment-
346 wide significance threshold ($p \leq \frac{0.05}{2*30} \approx 0.0008$). Moreover, the shift towards
347 lower p-values for the 30 modules was not more significant ($p = 2.3 \times 10^{-3}$)
348 than for the gene-based test (Supplementary Fig. 7A and Supplementary Table 7).
349 The second and third approaches derive from the common assumption that the
350 heritability of disease predisposition may be larger in familial and early-onset
351 cases³². We devised orthogonal tests for age-of-onset and familiarity and
352 combined them with our burden tests (Methods). Neither approach would
353 improve the results (Supplementary Fig. 7B&C and Supplementary Data 7).

354 Assuming that *TYK2* and *CARD9* are truly causative and their effect sizes in our
355 data unbiased, we estimated that a case-control cohort ranging from ~ 50,000
356 (*TYK2*) to ~200,000 (*CARD9*) individuals would have been needed to achieve
357 experiment-wide significance (testing 45 candidate genes), and from ~ 78,000
358 (*TYK2*) to >500,000 (*CARD9*) individuals to achieve genome-wide significance
359 (testing 20,000 genes) in the gene-based burden test (Supplementary Fig. 8).

360

361 **DISCUSSION**

362 We herein describe a novel dataset comprising array-based transcriptome data
363 for six circulating immune cell types and intestinal biopsies at three locations
364 collected on ~300 healthy European individuals. We use this CEDAR dataset
365 (“Correlated Expression and Disease Association Research”) to identify 23,650
366 significant cis-eQTL, which fall into 9,720 regulatory modules of which at least
367 ~889 affect more than one gene in more than one tissue. We provide strong
368 evidence that 63 of 200 known IBD GWAS signals reflect the activity of common

369 regulatory variants that preferentially drive multigenic modules. We perform
370 an exon-based burden test for 45 positional candidate CD genes mapping to 33
371 modules, in 5,500 CD cases and 6,500 controls. By demonstrating a significant
372 ($p = 6.9 \times 10^{-4}$) upwards shift of $\log(1/p)$ values for damaging when compared
373 to synonymous variants, we show that the sequenced genes include new
374 causative CD genes.

375 Individually, none of the sequenced genes (other than the positive *NOD2* and
376 *IL23R* controls) exceed the experiment-wide significance threshold, precluding
377 us from definitively pinpointing any novel causative genes. However, we note
378 *IL18R1* amongst the top-ranking genes (see also Supplementary Note 1). *IL18R1*
379 is the only gene in an otherwise relatively gene-poor region (also encompassing
380 *IL1R1* and *IL18RAP*) characterized by robust cis-eQTL in CD4 and CD8 that are
381 strongly correlated with the DAP for CD and UC ($0.68 \leq |\rho| \leq 0.93$). Reduced
382 transcript levels of *IL18R1* in these cell types is associated with increased risk for
383 IBD. Accordingly, rare ($MAF \leq 0.005$) damaging variants were cumulatively
384 enriched in CD cases (CAST $p = 0.05$). The cumulative allelic frequency of rare
385 damaging variants was found to be higher in familial CD cases (0.0027), when
386 compared to non-familial CD cases (0.0016; $p = 0.09$) and controls (0.0010; $p =$
387 0.03). When ignoring carriers of deleterious *NOD2* mutations, average age-of-
388 onset was reduced by ~3 years (25.3 vs 28.2 years) for carriers of rare damaging
389 *IL18R1* variants but this difference was not significant ($p = 0.18$).

390 While the identification of matching cRM for 63/200 DAP points towards a
391 number of strong candidate causative genes, it leaves most risk loci without
392 matching eQTL despite the analysis of nine disease-relevant cell types. This
393 finding is in agreement with previous reports^{4,33}. It suggests that cis-eQTL
394 underlying disease predisposition operate in cell types, cell states (f.i. resting vs
395 activated) or developmental stages that were not explored in this and other
396 studies. It calls for the enlargement and extension of eQTL studies to more
397 diverse and granular cellular panels^{10,34}, possibly by including single-cell
398 sequencing or spatial transcriptomic approaches. By performing eQTL studies in
399 a cohort of healthy individuals, we have made the reasonable assumption that
400 the common regulatory variants that are driving the majority of GWAS signals
401 are acting before disease onset, including in individuals that will never develop

402 the disease. An added advantage of studying a healthy cohort, is that the
403 corresponding dataset is “generic”, usable for the study of perturbation of gene
404 regulation for any common complex disease. However, it is conceivable that
405 some eQTL underlying increased disease risk only manifest themselves once the
406 disease process is initiated, for instance as a result of a modified inflammatory
407 status. Thus, it may be useful to perform eQTL studies with samples collected
408 from affected individuals to see in how far the eQTL landscape is affected by
409 disease status.

410 One of the most striking results of this work is the observation that cRM that
411 match DAP are ≥ 2 -fold enriched in multi-genic modules. We cannot fully
412 exclude that this is due to ascertainment bias. As multi-genic modules tend to
413 also be multi-tissue, multi-genic cRM matching a DAP in a non-explored disease-
414 relevant cell type have a higher probability to be detected in the explored cell
415 types than the equivalent monogenic (and hence more likely cell type specific)
416 cRM. The alternative explanation is that cRM matching DAP are truly enriched in
417 multi-genic cRM. It is tempting to surmise that loci harboring clusters of co-
418 regulated, functionally related causative genes have a higher probability to be
419 detected in GWAS, reflecting a relatively larger target space for causative
420 mutations. We herein tested this hypothesis by applying a module rather than
421 gene-based test. Although this did not appear to increase the power of the
422 burden test in this work, it remains a valuable approach to explore in further
423 studies. Supplementary Data 2 provides a list of >900 multigenic modules
424 detected in this work that could be used in this context.

425 Although we re-sequenced the ORF of 45 carefully selected candidate genes in a
426 total of 5,500 CD cases and 6,600 controls, none of the tested genes exceeded the
427 experiment-wide threshold of significance. This is despite the fact that we used
428 a one-sided, eQTL-informed test to potentially increase power. Established IBD
429 causative genes used as positive control, *NOD2* and *IL23R*, were positive
430 indicating that the experiment was properly conducted. We were not able to
431 improve the signal strength by considering information about regulatory
432 modules, familiarity or age-of-onset. We estimated that ≥ 10 -fold larger sample
433 sizes will be needed to achieve adequate power if using the same approach.

434 Although challenging, these numbers are potentially within reach of
435 international consortia for several common diseases including IBD.

436 It is conceivable that the organ-specificity of nearly all complex diseases (such as
437 the digestive tract for IBD), reflects tissue-specific perturbation of broadly
438 expressed causative genes that may fulfill diverse functions in different organs.
439 If this is true, coding variants may not be the appropriate substrate to perform
440 burden tests, as these will affect the gene across all tissues. In such instances,
441 the disruptive variants of interest may be those perturbing tissue-specific gene
442 switches. Also, it has recently been proposed that the extreme polygenic nature
443 of common complex diseases may reflect the trans-effects of a large proportion
444 of regulatory variants active in a given cell type on a limited number of core
445 genes via perturbation of highly connected gene networks³⁵. Identifying rare
446 regulatory variants is still challenging, however, as tissue-specific gene switches
447 remain poorly catalogued, and the effect of variants on their function difficult to
448 predict. The corresponding sequence space may also be limited in size, hence
449 limiting power. Nevertheless, a reasonable start may be to re-sequence the
450 regions surrounding common regulatory variants that have been fine-mapped at
451 near single base pair resolution⁴.

452 In conclusion, we hereby provide to the scientific community a collection of
453 ~24,000 cis-eQTL in nine cell types that are highly relevant for the study of
454 inflammatory and immune-mediated diseases, particularly of the intestinal tract.
455 The CEDAR dataset advantageously complements existing eQTL datasets
456 including GTEx^{10,34}. We propose a paradigm to rationally organize cis-eQTL
457 effects in co-regulated clusters or regulatory modules. We identify ~100
458 candidate causative genes in 63 out of 200 analyzed risk loci, on the basis of
459 correlated DAP and EAP. We have developed a web-based browser to share the
460 ensuing results with the scientific community (<http://cedar-web.giga.ulg.ac.be>).
461 The CEDAR website will imminently be extended to accommodate additional
462 common complex disease for which GWAS data are publicly available. We show
463 that the corresponding candidate genes are enriched in causative genes, however,
464 that case-control cohorts larger than those used in this study (12,000
465 individuals) are required to formally demonstrate causality by means of
466 presently available burden tests.

467

468 **METHODS**469 **Sample collection in the CEDAR cohort**

470 We collected peripheral blood as well as intestinal biopsies (ileum, transverse
471 colon, rectum) from 323 healthy Europeans visiting the Academic Hospital of the
472 University of Liège as part of a national screening campaign for colon cancer.
473 Participants included 182 women and 141 men, averaging 56 years of age
474 (range: 19-86). Enrolled individuals were not suffering any autoimmune or
475 inflammatory disease and were not taking corticosteroids or non-steroid anti-
476 inflammatory drugs (with the exception of low doses of aspirin to prevent
477 thrombosis). We recorded birth date, weight, height, smoking history, declared
478 ethnicity and hematological parameters (red blood cell count, platelet count,
479 differential white blood cell count) for each individual. The experimental
480 protocol was approved by the ethics committee of the University of Liège
481 Academic Hospital. Informed consent was obtained prior to donation in
482 agreement with the recommendations of the declaration of Helsinki for
483 experiments involving human subjects. We refer to this cohort as CEDAR for
484 Correlated Expression and Disease Association Research.

485 **SNP genotyping and imputation**

486 Total DNA was extracted from EDTA-collected peripheral blood using the
487 MagAttract DNA blood Midi M48 Kit on a QIAcube robot (Qiagen). DNA
488 concentrations were measured using the Quant-iT Picogreen ds DNA Reagents
489 (Invitrogen). Individuals were genotyped for > 700K SNPs using Illumina's
490 Human OmniExpress BeadChips, an iScan system and the Genome Studio
491 software following the guidelines of the manufacturer. We eliminated variants
492 with call rate ≤ 0.95 , deviating from Hardy-Weinberg equilibrium ($p \leq 10^{-4}$), or
493 which were monomorphic. We confirmed European ancestry of all individuals
494 by PCA using the HapMap population as reference. Using the real genotypes of
495 629,570 quality-controlled autosomal SNPs as anchors, we used the Sanger
496 Imputation Services with the UK10K + 1,000 Genomes Phase 3 Haplotype
497 panels⁴³⁻⁴⁶ to impute genotypes at autosomal variants in our population. We

498 eliminated indels, SNPs with $MAF \leq 0.05$, deviating from Hardy-Weinberg
499 equilibrium ($p \leq 10^{-3}$), and with low imputation quality ($INFO \leq 0.4$), leaving
500 6,019,462 high quality SNPs for eQTL analysis.

501 **Transcriptome analysis**

502 Blood samples were kept on ice and treated within one hour after collection as
503 follows. EDTA-collected blood was layered on Ficoll-Paque PLUS (GE
504 Healthcare) to isolate peripheral blood mononuclear cells by density gradient
505 centrifugation. CD4+ T lymphocytes, CD8+ T lymphocytes, CD19+ B lymphocytes,
506 CD14+ monocytes, CD15+ granulocytes were isolated by positive selection using
507 the MACS technology (Miltenyi Biotec). To isolate platelets, blood collected on
508 acid-citrate-dextrose (ACD) anticoagulant was centrifuged at 150g for 10
509 minutes. The platelet rich plasma (PRP) was collected, diluted 2-fold in ACD
510 buffer and centrifuged at 800g for 10 minutes. The platelet pellet was
511 resuspended in MACS buffer (Miltenyi Biotec) and platelets purified by negative
512 selection using CD45 microbeads (Miltenyi Biotec). Intestinal biopsies were
513 flash frozen in liquid nitrogen immediately after collection and kept at -80°C
514 until RNA extraction. Total RNA was extracted from the purified leucocyte
515 populations and intestinal biopsies using the AllPrep Micro Kit and a QIAcube
516 robot (Qiagen). For platelets, total RNA was extracted manually with the RNeasy
517 Mini Kit (Qiagen). Whole genome expression data were generated using HT-12
518 Expression Beadchips following the instructions of the manufacturer (Illumina).
519 Technical outliers were removed using controls recommended by Illumina and
520 the Lumi package⁴⁷. We kept 29,464/47,323 autosomal probes (corresponding
521 to 19,731 genes) mapped by Re-Annotator⁴⁸ to a single gene body with ≤ 2
522 mismatches and not spanning known variants with $MAF > 0.05$. Within cell
523 types, we only considered probes (i.e. “usable” probes) with detection p-value \leq
524 0.05 in $\geq 25\%$ of the samples. Fluorescence intensities were Log_2 transformed
525 and Robust Spline Normalized (RSN) with Lumi⁴⁷. Normalized expression data
526 were corrected for sex, age, smoking status and Satrix Id using ComBat from the
527 SVA R library⁴⁹. We further corrected the ensuing residuals within tissue for the
528 number of Principal Components (PC) that maximized the number of cis-eQTL
529 with $p \leq 10^{-6}$ ⁵⁰. Supplementary Table 2 summarizes the number of usable
530 samples, probes and PC for each tissue type.

531 **Cis-eQTL analysis**

532 Cis-eQTL analyses were conducted with PLINK and using the expression levels
533 precorrected for fixed effects and PC as described above^{51,52}. Analyses were
534 conducted under an additive model, i.e. assuming that the average expression
535 level of heterozygotes is at the midpoint between alternate homozygotes. To
536 identify cis-eQTL we tested all SNPs in a 2Mb window centered around the probe
537 (if “usable”). P-values for individual SNPs were corrected for the multiple
538 testing within the window by permutation (10,000 permutations). For each
539 probe-tissue combination we kept the best (corrected) p-value. Within each
540 individual cell type, the ensuing list of corrected p-values was used to compute
541 the corresponding false discovery rates (FDR or q-value). Supplementary Table
542 3 reports the number of cis-eQTL found in the nine analyzed cell types for
543 different FDR thresholds (see also Supplementary Figure 9).

544 **Comparing EAP with ϑ to identify cis Regulatory Modules**

545 If the transcript levels of a given gene are influenced by the same regulatory
546 variants (one or several) in two tissues, the corresponding EQTL Association
547 Patterns (EAP)(i.e. the $-\log(p)$ values of association for the SNPs surrounding the
548 gene) are expected to be similar. Likewise, if the transcript levels of different
549 genes are influenced by the same regulatory variants in the same or in different
550 tissues, the corresponding EAP are expected to be similar (cfr. main text, Fig. 1).
551 We devised a metric, ϑ , to quantify the similarity between EAP. If two EAP are
552 similar, one can expect the corresponding $-\log(p)$ values to be positively
553 correlated. One particularly wants the EAP peaks, i.e. the highest $-\log(p)$ values,
554 to coincide in order to be convinced that the corresponding cis-eQTL are driven
555 by the same regulatory variants. To quantify the similarity between EAP while
556 emphasizing the peaks we developed a weighted correlation. Imagine two
557 vectors \mathbf{X} and \mathbf{Y} of $-\log(p)$ values for n SNPs surrounding the gene(s) of interest.
558 Using the same nomenclature as in Fig. 1A, \mathbf{X} could correspond to gene A in
559 tissue 1, and \mathbf{Y} to gene A in tissue 2, or \mathbf{X} could correspond to gene A in tissue 1,
560 and \mathbf{Y} to gene B in tissue 2. We only consider for analysis, SNPs within 1Mb of
561 either gene (probe) and for which x_i and/or y_i is superior to 1.3 (i.e. p-value <
562 0.05) hence informative for at least one of the two cis-eQTL. Indeed, the

563 majority of variants with $-\log(p) < 1.3$ ($p > 0.05$) for both EAP are by definition
 564 not associated with either trait. There is therefore no reason to expect that they
 565 could contribute useful information to the correlation metric: their ranking in
 566 terms of $-\log(p)$ values becomes more and more random as the $-\log(p)$
 567 decreases. We define the weight to be given to each SNP in the correlation as:

$$w_i = \left(\text{MAX} \left(\frac{x_i}{x_{\text{MAX}}}, \frac{y_i}{y_{\text{MAX}}} \right) \right)^p$$

568 The larger p , the more weight is given to the top SNPs. In this work, p was set at
 569 one.

570 The weighted correlation between the two EAP, r_w , is then computed as:

$$r_w = \frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i \left(\frac{x_i - \bar{x}_w}{\sigma_x^w} \right) \left(\frac{y_i - \bar{y}_w}{\sigma_y^w} \right)$$

571 in which

$$\bar{x}_w = \frac{\sum_{i=1}^n w_i \times x_i}{\sum_{i=1}^n w_i}$$

$$\bar{y}_w = \frac{\sum_{i=1}^n w_i \times y_i}{\sum_{i=1}^n w_i}$$

$$\sigma_x^w = \sqrt{\frac{\sum_{i=1}^n w_i \times (x_i - \bar{x}_w)^2}{\sum_{i=1}^n w_i}}$$

$$\sigma_y^w = \sqrt{\frac{\sum_{i=1}^n w_i \times (y_i - \bar{y}_w)^2}{\sum_{i=1}^n w_i}}$$

572

573 The larger r_w , the larger the similarity between the EAP, particularly for their
 574 respective peak SNPs.

575 r_w ignores an important source of information. If two EAP are driven by the same
 576 regulatory variant, there should be consistency in the signs of the effects across
 577 SNPs in the region. We will refer to the effect of the “reference” allele of SNP i on
 578 the expression levels for the first and second cis-eQTL as β_i^X and β_i^Y . If the

579 reference allele of the regulatory variant increases expression for both cis-eQTL,
 580 the β_i^X and β_i^Y 's for a SNPs in LD with the regulatory variant are expected to have
 581 the same sign (positive or negative depending on the sign of D for the considered
 582 SNP). If the reference allele of the regulatory variant increases expression for
 583 one cis-eQTL and decreases expression for the other, the β_i^X and β_i^Y 's for a SNPs
 584 in LD with the regulatory variant are expected to have opposite sign. We used
 585 this notion to develop a weighted and signed measure of correlation, r_{ws} . The
 586 approach was the same as for r_w , except that the values of y_i were multiplied by -
 587 1 if the signs of β_i^X and β_i^Y were opposite. r_{ws} is expected to be positive if the
 588 regulatory variant affects the expression of both cis-eQTL in the same direction
 589 and negative otherwise.

590 We finally combined r_w and r_{ws} in a single score referred to as ϑ , as follows:

$$\vartheta = \frac{r_{ws}}{1 + e^{-k(r_w - T)}}$$

591 ϑ penalizes r_{ws} as a function of the value of r_w . The aim is to avoid considering
 592 EAP pairs with strong but negative r_w (which is often the case when the two EAP
 593 are driven by very distinct variants). The link function is a sigmoid-shaped
 594 logistic function with k as steepness parameter and T as sigmoid mid-point. In
 595 this work, we used a value of k of 30, and a value of T of 0.3 (Supplementary
 596 Figure 10).

597 We first evaluated the distribution of ϑ for pairs of EAP driven by the same
 598 regulatory variants by studying 4,693 significant cis-eQTL (FDR < 0.05). For
 599 these, we repeatedly (100 x) split our CEDAR population in two halves,
 600 performed the cis-eQTL analysis separately on both halves and computed ϑ for
 601 the ensuing EAP pairs. Supplementary Figure 1 is showing the obtained results.

602 We then evaluated the distribution of ϑ for pairs of EAP driven by distinct
 603 regulatory variants in the same chromosomal region as follows. We considered
 604 1,207 significant cis-eQTL (mapping to the 200 IBD risk loci described above).
 605 For each one of these, we generated a set of 100 "matching" cis-eQTL effects in
 606 silico, sequentially considering 100 randomly selected SNPs (from the same
 607 locus) as causal. The in silico cis-eQTL were designed such that they would
 608 explain the same fraction of expression variance as the corresponding real cis-

609 eQTL detected with PLINK (cfr. above). When performing cis-eQTL analysis
 610 under an additive model, PLINK estimates β_0 (i.e. the intercept), and β_1 (i.e. the
 611 slope of the regression), including for the top SNP. Assume that the expression
 612 level of the studied gene, Z , for individual i is z_i . Assume that the sample
 613 comprises n_T individuals in total, of which n_{11} are of genotype “11”, n_{12} of
 614 genotype “12”, and n_{22} of genotype “22”, for the top cis-eQTL SNP. The total
 615 expression variance for gene Z equals:

$$\sigma_T^2 = \frac{\sum_{i=1}^{n_T} (z_i - \bar{z}_T)^2}{n_T - 1}$$

616 The variance in expression level due to the cis-eQTL equals:

$$\sigma_{eQTL}^2 = \frac{n_{11}(\beta_0 - \bar{z}_T)^2 + n_{12}(\beta_0 + \beta_1 - \bar{z}_T)^2 + n_{22}(\beta_0 + 2\beta_1 - \bar{z}_T)^2}{n_T}$$

617 The heritability of expression due to the cis-eQTL, i.e. the fraction of the
 618 expression variance that is due to the cis-eQTL is therefore:

$$h_{eQTL}^2 = \frac{\sigma_{eQTL}^2}{\sigma_T^2}$$

619 To simulate cis-eQTL explaining the same h_{eQTL}^2 as the real eQTL in the CEDAR
 620 dataset, we sequentially considered all SNPs in the region. Each one of these
 621 SNPs would be characterized by n_{11} individuals of genotype “11”, n_{12} of genotype
 622 “12”, and n_{22} of genotype “22”, for a total of n_T genotyped individuals. We would
 623 arbitrarily set \bar{z}_{11} , \bar{z}_{12} , and \bar{z}_{22} at -1, 0 and +1. As a consequence, the variance
 624 due to this cis-eQTL equals:

$$\sigma_{eQTL}^2 = \frac{n_{11}(-1 - \bar{z}_T)^2 + n_{12}(0 - \bar{z}_T)^2 + n_{22}(1 - \bar{z}_T)^2}{n_T}$$

625 in which $\bar{z}_T = (n_{22} - n_{11})/n_T$.

626 Knowing σ_{eQTL}^2 and h_{eQTL}^2 , and knowing that

$$h_{eQTL}^2 = \frac{\sigma_{eQTL}^2}{\sigma_{eQTL}^2 + \sigma_{RES}^2}$$

627 the residual variance σ_{RES}^2 can be computed as

$$\sigma_{RES}^2 = \sigma_{eQTL}^2 \left(\frac{1}{h_{eQTL}^2} - 1 \right)$$

628 Individual expression data for the corresponding cis-eQTL (for all individuals of
629 the CEDAR dataset) were hence sampled from the normal distribution

$$z_i \sim N(\bar{z}_{xx}, \sigma_{RES}^2)$$

630 where \bar{z}_{xx} is -1, 0 or +1 depending on the genotype of the individual (11, 12, or
631 22). We then performed cis-eQTL on the corresponding data set using EAP,
632 generating an in silico EAP. Real and in silico EAP were then compared using ϑ .
633 Supplementary Figure 1 shows the corresponding distribution of ϑ values for
634 EAP driven by distinct regulatory variants.

635 The corresponding distributions of ϑ under H_1 and H_0 (Supplementary Figure 1)
636 show that ϑ discriminates very effectively between H_1 and H_0 especially for the
637 most significant cis-eQTL. In the experiment described above, this would yield a
638 false positive rate of 0.05, and a false negative rate of 0.23. We chose a threshold
639 of $|\vartheta| > 0.6$ to cluster EAP in cis-acting regulatory elements or cRM (Fig. 2).
640 Clusters were visually examined as show in Supplementary Figure 2. Twenty-
641 nine edges connecting otherwise unlinked and yet tight clusters were manually
642 removed.

643 **Testing for an excess sharing of cRM between cell types**

644 Assume that cell type 1 is part of n_{1T} cRM, including n_{11} private cRM, n_{12} cRM
645 shared with cell type 2, n_{13} cRM shared with cell type 2, ..., and n_{19} cRM shared
646 with cell type 9. Note that $\sum_{i=1}^9 n_{1i} \geq n_{1T}$, because cRM may include more than
647 two cell types. Assume that $n_{1S} = \sum_{i \neq 1}^9 n_{1i}$ is the sum of pair-wise sharing events
648 for cell type 1. We computed, for each cell type $i \neq 1$, the probability to observe
649 $\geq n_{1i}$ sharing events with cell type 1 assuming that the expected number (under
650 the hypothesis of random assortment) is

$$n_{1S} \times \frac{n_{iT}}{\sum_{j \neq 1}^9 n_{jT}}$$

651 Pair-wise sharing events between tissue 1 and the eight other tissues were
652 generated in silico under this model of random assortment (5,000 simulations).

653 The p-value for n_{1i} was computed as the proportion of simulations that would
654 yield values that would be as large or larger than n_{1i} . The same approach was
655 used for the nine cell types. Thus, two p-values of enrichment are obtained for
656 each pair of cell types i and j , one using i as reference cell type, and the other
657 using j as reference cell type. As can be seen from Fig. 5, the corresponding pairs
658 of p-values were always perfectly consistent.

659 We performed eight distinct analyses. In the first analysis, we only considered
660 cRM involving no more than two tissues (i.e. unique for specific pairs of cell
661 types). In subsequent analyses, we progressively included cRM with no more
662 than three, four, ..., and nine cell types.

663 **Comparing EAP and DAP using ϑ**

664 The approach used to cluster EAP in cRM was also used to assign Disease
665 Association Patterns (DAP) for Inflammatory Bowel Disease (IBD) to EAP-
666 defined cRM. We studied 200 IBD risk loci identified in recent GWAS meta-
667 analyses^{2,3}. The limits of the corresponding risk loci were as defined in the
668 corresponding publications. We measured the similarity between DAP and
669 EAP using the ϑ metric for all cis-eQTL mapping to the corresponding intervals
670 (i.e. for all cis-eQTL for which the top SNP mapped within the interval). To
671 compute the correlations between DAP and EAP we used all SNPs mapping to the
672 disease interval with $-\log(p)$ value ≥ 1.3 either for DAP, EAP or both.

673 In addition to computing ϑ as described in section 5, we computed an empirical
674 p-value for ϑ using the approach (based on in silico generated cis-eQTL)
675 described above to generate the locus-specific distribution of ϑ values for EAP
676 driven by distinct regulatory variants. From this distribution, one can deduce
677 the probability that a randomly generated EAP (explaining as much variance as
678 the real tested EAP) and the DAP would by chance have a $|\vartheta|$ value that is as high
679 or higher than the real EAP. The corresponding empirical p-value accounts for
680 the local LD structure between SNPs.

681 **Evaluating the enrichment of DAP-EAP matching**

682 To evaluate whether DAP matched EAP more often than expected by chance
683 alone, we analyzed 97 IBD risk loci interrogated by the ImmunoChip, (i) in order

684 to allow for convenient comparison with Huang et al.⁴, and (ii) because we
685 needed extensively QC genotypes for the IBDGC data to perform the enrichment
686 analysis with the ϑ -based method (see hereafter). Within these 97 IBD risk loci,
687 we focused on 63 regions affecting CD⁴, encompassing at least one significant
688 eQTL, and for which the lead CD-associated SNP had MAF > 0.05. Indeed, eQTL
689 analyses in the CEDAR dataset were restricted to SNPs with MAF > 0.05 (see
690 above). We used three methods to evaluate whether the observed number of
691 DAP-EAP matches were higher than expected by chance alone: naïve, frequentist
692 and ϑ -based. Analyses were performed separately for the nine cell types.

693 In the “naïve” approach, DAP and EAP were assumed to match if the
694 corresponding lead SNPs were in LD with $r^2 \geq 0.8$. This would yield $n_N \leq 63$
695 risk loci for which the DAP would match at least one EAP. To measure the
696 statistical significance of n_N , we sampled a SNP (MAF > 0.05) at random in each
697 of the 63 risk loci, and counted the number of loci with at least one matching EAP.
698 This “simulation” was repeated 1,000 times. The significance of n_N was
699 measured as the proportion of simulations that would yield $\geq n_N$ matches.

700 The frequentist approach used the method described by Nica et al.⁵³. DAP and
701 EAP were assumed to match if fitting the disease-associated lead SNP in the
702 eQTL analysis caused a larger drop in $-\log(p)$ than 95% of the SNPs with MAF >
703 0.05 in the analyzed risk locus. This would yield $n_F \leq 63$ risk loci for which the
704 DAP would match at least one EAP. To measure the statistical significance of n_F ,
705 we sampled a SNP (MAF > 0.05) at random in each of the 63 risk loci, and
706 counted the number of loci with at least one matching EAP. This “simulation”
707 was repeated 1,000 times. The significance of n_F was measured as the
708 proportion of simulations that would yield $\geq n_F$ matches.

709 Finally, we used our ϑ -based approach in which DAP and EAP were assumed to
710 match if $|\vartheta| > 0.6$. This would yield $n_\vartheta \leq 63$ risk loci for which the DAP would
711 match at least one EAP. To measure the statistical significance of n_ϑ we sampled
712 a SNP (MAF > 0.05) at random in each of the 63 risk loci, and generated a DAP
713 assuming that the corresponding SNPs were causal as follows.

714 Assume a cohort with n_1 cases and n_2 controls (f.i. the IIBDGC cohort). Assume a
 715 SNP with an allelic frequency of p in the cases + controls, an allelic frequency of
 716 $(p + d)$ in cases and $(p + \delta)$ in controls.

717 One can easily show that:

$$718 \quad \delta = -d \frac{n_1}{n_2} \quad (1)$$

719 The odds ratio (OR) for that SNP equals:

$$OR = \frac{(p + d)(1 - p - \delta)}{(p + \delta)(1 - p - d)}$$

720 The ratio between the between-cohort (i.e. cases and controls) variance versus
 721 within-cohort variance (corresponding to an F test) can be shown to equal:

$$F = \frac{d^2 \left(1 + \frac{n_1}{n_2}\right)}{\left(1 + \frac{n_2}{n_1}\right)(p - p^2) - d^2 \left(1 + \frac{n_1}{n_2}\right)}$$

722 If we fix F based on the real top SNP in the IIBDGC data in a given GWAS
 723 identified risk loci, we can determine d (and hence δ using equation 1) for the
 724 randomly selected SNP (that will become an “in silico causative variant”) with
 725 allelic frequency in (cases + controls) of p (different from the real top SNP), by
 726 solving

$$d = \frac{-\beta \pm \sqrt{\beta^2 - 4\alpha\gamma}}{2\alpha}$$

727 where

$$\alpha = \left(1 + \frac{n_1}{n_2}\right)(1 + F)$$

$$\beta = 0$$

$$\gamma = -(p - p^2) \left(1 + \frac{n_2}{n_1}\right) F$$

728 Once we know $(p + d)$ (i.e. the frequency of the SNP in cases), and hence $(p + \delta)$
 729 (i.e. the frequency of the SNP in controls), we can use Hardy–Weinberg to
 730 determine the frequency of the three genotypes in cases ($p_{AA}^{IBD}, p_{AB}^{IBD}, p_{BB}^{IBD}$) and

731 controls ($p_{AA}^{CTR}, p_{AB}^{CTR}, p_{BB}^{CTR}$). We then create an in silico case-control cohort by
732 sampling (with replacement) $n_1 \times p_{AA}^{IBD}$ AA cases, $n_1 \times p_{AB}^{IBD}$ AB cases, ..., and
733 $n_2 \times p_{BB}^{CTR}$ BB controls from the individuals of the IIBDGC (without discriminating
734 real case and control status). Association analysis of the corresponding dataset
735 in the chromosome region of interest generates DAP with $\max - \log(p)$ value
736 similar to the real DAP. This “simulation” was repeated 1,000 times. The
737 significance of n_g was measured as the proportion of simulations that would
738 yield $\geq n_g$ matches.

739 Targeted exon resequencing in CD cases and controls

740 Genes for which EAP match the DAP tightly (high $|\vartheta|$ values) are strong
741 candidate causal genes for the studied disease. In the case of IBD, we identified
742 ~100 such genes (Table 1). Ultimate proof of causality can be obtained by
743 demonstrating a differential burden of rare disruptive variants in cases and
744 controls. Burden tests preferably focus on coding gene segments, in which
745 disruptive variants are most effectively recognized. Analyses are restricted to
746 rare variants to ensure independence from the GWAS signals.

747 To perform burden tests, we collected DNA samples from 7,323 Crohn Disease
748 (CD) cases and 6,342 controls of European descent in France (cases: 1,899 –
749 ctrls: 1,731), the Netherlands (2,002 – 1,923) and Belgium (3,422 – 2,688). The
750 study protocols were approved by the institutional review board at each centre
751 involved with recruitment. Informed consent and permission to share the data
752 were obtained from all subjects, in compliance with the guidelines specified by
753 the recruiting centre’s institutional review board.

754 During the course of this project, we selected 45 genes with high $|\vartheta|$ values for
755 resequencing (Table 1). We designed primers to amplify all corresponding
756 coding exons plus exon-intron boundaries corresponding to all transcripts
757 reported in the CCDS release 15⁵⁴ (Supplemental data 8). Following Momozawa
758 et al.⁵⁵, the primers were merged in five pools to perform a first round of PCR
759 amplification (25 cycles). We then added 8-bp barcodes and common adaptors
760 (for sequencing) to all PCR products by performing a second round of PCR
761 amplification (4 cycles) using primers targeting shared 5’overhangs introduced
762 during the first PCR. The ensuing libraries were purified, quality controlled and

763 sequenced (2 x 150-bp paired-end reads) on a HiSeq 2500 (Illumina) instrument.
764 Sequence reads were sorted by individual using the barcodes, aligned to the
765 human reference sequence (hg19) with the Burrows-Wheeler Aligner (ver.
766 0.7.12)⁵⁶, and further processed using Genome Analysis Toolkit (GATK, ver. 3.2-
767 2)⁵⁷. We only considered individuals for further analyses if $\geq 95\%$ of the target
768 regions was covered ≥ 20 sequence reads. Average sequence depth across
769 individuals and target regions was 1,060. We called variants for each individual
770 separately using the UnifiedGenotyper and HaplotypeCaller of GATK, as well as
771 VCMM (ver. 1.0.2)⁵⁸, and listed all variants detected by either method. Genotypes
772 for all individuals were determined for each variant based on the ratio of
773 reference and alternative alleles amongst sequence reads as determined by
774 Samtools⁵⁹. Individuals were labelled homozygote reference, heterozygote, or
775 homozygote derived when the alternative allele frequency was between 0 and
776 0.15, between 0.25 and 0.75, and between 0.85 and 1, respectively. If the
777 alternative allele frequency was outside these ranges or a variant position was
778 covered with < 20 sequencing reads, the genotype was considered missing. We
779 excluded variants with call rates $< 95\%$ or variants that were not in Hardy-
780 Weinberg equilibrium ($P < 1 \times 10^{-6}$). We excluded 281 individuals with ≥ 2 minor
781 alleles at 23 variants selected to have a MAF ≤ 0.01 in non-Finnish Europeans
782 and ≥ 0.10 in Africans or East-Asians in the Exome Aggregation Consortium⁶⁰.

783 In the end, we used 6,597 cases and 5,502 controls for further analyses, while
784 98.5% of the target regions on average was covered with 20 or more sequence
785 reads.

786 **Gene-based burden test**

787 We first used SIFT⁶¹ and Polyphen-2⁶² to sort the 4,175 variants identified by
788 sequencing in four categories: (i) loss-of-function (LoF) or severe, corresponding
789 to stop gain, stop loss, frameshift and splice-site variants, (ii) damaging,
790 corresponding to missense variants predicted by SIFT to be damaging and
791 Polyphen-2 to be possibly or probably damaging, (iii) benign, corresponding to
792 the other missense variants, and (iv) synonymous. We performed the burden
793 test using the LoF plus damaging variants, and used the synonymous variants as
794 controls. We only considered variants with MAF (computed for the entire

795 dataset, i.e. cases plus controls) ≤ 0.005 . We indeed showed in a previous fine-
796 mapping study that all reported independent effects were driven by variants
797 with $MAF \geq 0.01^4$. By doing so we ensure that the signals of the burden test are
798 independent of previously reported association signals. Thus, 174 LoF, 991
799 damaging, and 1,434 synonymous were ultimately used to perform burden tests.

800 Burden tests come in two main flavors. In the first, one assumes that disruptive
801 variants will be enriched in either cases (i.e. disruptive variants increase risk) or
802 in controls (i.e. disruptive variance decrease risk). In the second, one assumes
803 that - for a given gene - some disruptive variants will be enriched in cases, while
804 other may be enriched in controls (Supplementary Fig. 11). The first was
805 implemented using CAST⁶³. To increase power, we exploited the DAP-EAP
806 information to perform one-sided (rather than two-sided) tests. When $\vartheta < 0$, we
807 tested for an enrichment of disruptive variants in cases; when $\vartheta > 0$, for an
808 enrichment of disruptive variants in controls. P-values were computed by
809 phenotype permutation, i.e. shuffling case-control status. When applying this
810 test on a gene-by-gene basis using synonymous variants ($MAF > 0.005$), the
811 distribution of p-values (QQ-plot) indicated that the CAST test was conservative
812 ($\lambda_{GC} = 0.51$) (Supplementary Fig. 12). The second kind of burden test was
813 implemented with SKAT⁶⁴. It is noteworthy that SKAT ignores information from
814 singletons (Supplementary Fig. 11). Just as for CAST, p-values were computed by
815 phenotype permutation, i.e. shuffling case-control status. When applying this
816 test on a gene-by-gene basis using synonymous variants ($MAF < 0.005$), the
817 distribution of p-values (QQ-plot) indicated that the SKAT test is too permissive
818 ($\lambda_{GC} = 1.73$) (Supplementary Fig. 12). Consequently, gene-based p-values
819 obtained with SKAT were systematically GC corrected using this value of λ_{GC} .
820 We performed the two kinds of analyses for each gene, as one doesn't a priori
821 know what hypothesis will match the reality best for a given gene.

822 We also extracted information from the distribution of p-values (or $-\log(p)$
823 values) across the 45 analyzed genes. Even if individual genes do not yield -
824 $\log(p)$ values that exceed the significance threshold (accounting for the number
825 of analyzed genes and tests performed), the distribution of $-\log(p)$ values may
826 significantly depart from expectations, indicating that the analyzed genes include
827 at least some causative genes. This was done by taking for each gene, the best p-

828 value (whether obtained with CAST or SKAT) and then rank the genes by
829 corresponding $-\log(p)$ value. The same was done for 10^5 phenotype
830 permutations, allowing us to examine the distribution of $-\log(p)$ values for given
831 ranks and compute the corresponding medians and limits of the 95% confidence
832 band, as well as to compute the probability that $-2 \sum_{i=1}^{45} \ln(p_i)$ (Fisher's equation
833 to combine p-values) equals or exceeds the observed. Our results show that
834 there is a significant departure from expectation when analyzing the damaging
835 variants ($p = 6.9 \times 10^{-4}$) but not when analyzing the synonymous variants ($p =$
836 0.66) supporting the presence of genuine causative genes amongst the analyzed
837 list.

838 **cRM-based burden test**

839 The enrichment of multi-genic cRM in IBD risk loci suggests that risk loci may
840 have more than one causative gene belonging to the same cRM. To capitalize on
841 this hypothesis, we developed a cRM-based burden test. Gene-specific p-values
842 were combined within cRM using Fisher's method. For each gene, we considered
843 the best p-value whether obtained with CAST or SKAT. Statistical significance
844 was evaluated by phenotype permutation exactly as described for the gene-
845 based burden test. By doing so we observed a departure from expectation when
846 using the damaging variants ($p = 2.3 \times 10^{-3}$), but not when using the synonymous
847 variants ($p = 0.72$).

848 **Orthogonal tests for age-of-onset and familiarity**

849 It is commonly assumed that the heritability for common complex diseases is
850 higher in familial and early onset cases⁶⁵. To extract the corresponding
851 information from our data in a manner that would be orthogonal to the gene-
852 and module-based tests described above (i.e. the information about age-of-onset
853 and familiarity would be independent of these burden tests), we devised the
854 following approach.

855 For age-of-onset, we summed the age-of-onset of the n_c cases carrying rare
856 disruptive variants for the gene of interest. We then computed the probability
857 that the sum of the age-of-onset of n_c randomly chosen cases was as different
858 from the mean of age-of-onset as the observed one, yielding a gene-specific two-

859 sided p_{SKAT} value. In addition, we used the eQTL information to generate gene-
860 specific one-sided p_{CAST} values, corresponding to the probability that the sum of
861 the age-of-onset of n_C randomly chosen cases was as low or lower than the
862 observed one (for genes for which decrease in expression level as associated
863 with increased risk), or to the probability that the sum of the age-of-onset of n_C
864 randomly chosen cases was as high or higher than the observed one (for genes
865 for which increase in expression level as associated with increased risk). These
866 age-of-onset p-values were then combined with the corresponding p-values from
867 the burden test (CAST with CAST, SKAT with SKAT) using Fisher's method.

868 For familiarity, we determined what fraction of the n_C cases carrying rare
869 disruptive variants for the gene of interest were familial (affected first degree
870 relative). We then computed the probability that the fraction of familial cases
871 amongst n_C randomly chosen cases was as different from the overall proportion
872 of familial cases, yielding a gene-specific two-sided p_{SKAT} value. In addition, we
873 used the eQTL information to generate gene-specific one-sided p_{CAST} values,
874 corresponding to the probability that the fraction of familial cases amongst n_C
875 randomly chosen cases was as high or higher than the observed one (for genes
876 for which decrease in expression level as associated with increased risk), or to
877 the probability that the sum of the age-of-onset of n_C randomly chosen was as
878 low or lower than the observed one (for genes for which increase in expression
879 level as associated with increased risk). These familial p-values were then
880 combined with the corresponding p-values from the burden test (CAST with
881 CAST, SKAT with SKAT) using Fisher's method.

882

883 DATA AVAILABILITY

884 The complete CEDAR eQTL dataset can be downloaded from the Array Express
885 website (<https://www.ebi.ac.uk/arrayexpress/>), accession numbers E-MTAB-
886 6666 (genotypes) and E-MTAB-6667 (expression data). The data, preprocessed
887 as described in Methods, can be downloaded from the CEDAR website
888 (<http://cedar-web.giga.ulg.ac.be>).

889

890

891 **ACKNOWLEDGEMENTS**

892 This work was supported by grants to Michel Georges from WELBIO (CAUSIBD),
893 BELSPO (BeMGI), and Horizon 2020 (SYSCID). Computational resources at ULg
894 have been provided by GIGA and the Consortium des Équipements de Calcul
895 Intensif (CÉCI), funded by the Fonds de la Recherche Scientifique de Belgique
896 (F.R.S.-FNRS) under Grant No. 2.5020.11. This work was conducted as part of
897 the BioBank Japan Project supported by the Japan Agency for Medical Research
898 and Development and by the Ministry of Education, Culture, Sports, Sciences and
899 Technology of the Japanese government. The work of DA and IA was supported
900 by Russian Ministry of Science and Education under 5-100 Excellence
901 Programme. R.K. Weersma is supported by a VIDI grant (016.136.308) from the
902 Netherlands Organisation for Scientific Research (NWO). DNA samples from the
903 Dutch IBD cohort have been collected within the Parelsnoer Institute Project.
904 This nationwide Parelsnoer Institute project is part of and funded by the
905 Netherlands Federation of University Medical Centres and has received initial
906 funding from the Dutch Government (from 2007-2011). The Parelsnoer Institute
907 currently facilitates the uniform nationwide collection of information on and
908 biomaterials of thirteen other diseases. We are grateful to N. Hakozaiki, H. Iijima,
909 N. Maki, and other staff of the Laboratory for Genotyping Development, RIKEN
910 Center for the Integrative Medical Sciences. We thank Wouter Coppieters and the
911 other members of the GIGA genomics platform for their support.

912

913

914 **CONFLICTS OF INTEREST**

915 The authors declare absence of any conflict of interest, whether financial or
916 other.

917

918 **AUTHOR CONTRIBUTIONS**

919 YM, JD and MG conceived experiments, generated data, analyzed data and wrote
920 the manuscript. ET, VD, SR, BC, FC, ED, ME, A-SG,CL,RM,MM,CO generated and
921 analyzed data. IA,DA,YA and MG conceived and generated the CEDAR website. LA,
922 GB, FH, ML, BO, MJP,AEVDMDJ,CJVDW,MVC, ML, JPH, RKW,MDV,DF,SV,MK,EL
923 collected and provided samples.

924

925 REFERENCES

- 926 1. MacArthur, J. *et al.* The new NHGRI-EBI catalog of published genome-wide
927 association studies (GWAS Catalog). *Nucleic Acids Res* **45**, D896-D901 (2017).
- 928 2. Jostins, L. *et al.* Host-microbe interactions have shaped the genetic
929 architecture of inflammatory bowel disease. *Nature* **491**,119-124 (2012).
- 930 3. Liu, J.Z. *et al.* Association analyses identify 38 susceptibility loci for IBD and
931 highlight shared genetic risk across populations. *Nat Genet* **47**, 979-986
932 (2015).
- 933 4. Huang, H. *et al.* Association mapping of IBD loci to single variant resolution.
934 *Nature* **547**, 173-178 (2017).
- 935 5. Claussnitzer, M. *et al.* FTO obesity variant circuitry and adipocyte browning in
936 humans. *N Engl J Med* **373**, 895-907 (2015).
- 937 6. Hugot, J.P. *et al.* Association of NOD2 leucine-rich repeat variants with
938 susceptibility to Crohn's disease. *Nature* **411**, 599-603 (2001).
- 939 7. Hampe, J. *et al.* A genome-wide association scan of nonsynonymous SNPs
940 identifies a susceptibility variant for Crohn disease in ATG16L1. *Nat Genet* **39**,
941 207-211 (2007).
- 942 8. Momozawa, Y. *et al.* Resequencing of positional candidates identifies low
943 frequency IL23R coding variants protecting against inflammatory bowel
944 disease. *Nat Genet* **43**, 43-47 (2011).
- 945 9. Rivas, M.A. *et al.* Deep resequencing of GWAS loci identifies independent rare
946 variants associated with inflammatory bowel disease. *Nat Genet* **43**, 1066-
947 1073 (2011).
- 948 10. The GTEx Consortium. Genetic effects on gene expression across human
949 tissues. *Nature* **550**, 204-213 (2017).

- 950 11. Nica, A.C. *et al.* Candidate causal regulatory effects by integration of
951 expression QTLs with complex trait genetic associations. *PLoS Genet* **6**,
952 e1000895 (2010).
- 953 12. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of
954 genetic association studies using summary statistics. *PLoS Genet* **10**,
955 e1004383 (2014).
- 956 13. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies
957 predicts complex trait gene targets. *Nat Genet* **48**, 481-487 (2016).
- 958 14. The ENCODE Project Consortium. An integrated encyclopedia of DNA
959 elements in the human genome. *Nature* **489**, 57-74 (2012).
- 960 15. Nicolae, D.L. Association tests for rare variants. *Annu Rev Genom Hum Genet*
961 **17**, 117-130 (2016).
- 962 16. Pritchard, J.K. & Cox, N.J. The allelic architecture of human disease genes:
963 common disease-common variant ... or not? *Hum Mol Genet* **11**, 2417-2423
964 (2002).
- 965 17. McGregor, A.P. *et al.* Morphological evolution through multiple cis-regulatory
966 mutations at a single gene. *Nature* **448**, 587-590 (2007).
- 967 18. Mackay, T.F. Quantitative trait loci in *Drosophila*. *Nat Rev Genet* **2**, 11-20
968 (2001).
- 969 19. Yalcin, B. *et al.* Genetic dissection of behavioral QTL shows that *Rgs2*
970 modulates anxiety in mice. *Nat Genet* **36**, 1197-1202 (2004).
- 971 20. Karim, L. *et al.* Variants modulating the expression of a chromosome domain
972 encompassing *PLAG1* influence bovine stature. *Nat Genet* **43**, 405-413 (2011).
- 973 21. Steinmetz, L.M. *et al.* Dissecting the architecture of a QTL in yeast. *Nature* **416**,
974 326-330 (2002).
- 975 22. Khor, B., Gardet, A., Xavier, R. Genetics and pathogenesis of inflammatory
976 bowel disease. *Nature* **474**, 307-317 (2011).
- 977 23. <https://github.com/joepickrell/1000-genomes-genetic-maps>
- 978 24. Fuchsberger, C. *et al.* The genetic architecture of type 2 diabetes. *Nature* **536**,
979 41-47 (2016).
- 980 25. Momozawa, Y. *et al.* Low-frequency coding variants in *CETP* and *CFB* are
981 associated with susceptibility of exudative age-related macular degeneration
982 in the Japanese population. *Hum. Mol. Genet.* **25**, 5027-5034 (2016).

- 983 26. Kumar, P. *et al.* Predicting the effects of coding non-synonymous variants on
984 protein function using the SIFT algorithm. *Nat Protoc* **4**, 1073-1081 (2009).
- 985 27. Adzhubei, I.A. *et al.* A method and server for predicting damaging missense
986 mutations. *Nat Methods* **7**, 248-249 (2010).
- 987 28. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans.
988 *Nature* **536**, 285-291 (2016).
- 989 29. Morgenthaler, S. & Thilly, W.G. A strategy to discover genes that carry multi-
990 allelic or mono-allelic risk for common diseases: a cohort allelic sums test
991 (CAST). *Mutat Res* **615**, 28-56 (2007).
- 992 30. Wu, M.C. *et al.* Rare-variant association testing for sequencing data with the
993 sequence kernel association test. *Am J Hum Genet* **89**, 82-93 (2011).
- 994 31. Richardson, T.G. *et al.* A pathway-centric approach to rare variant association
995 analysis. *Eur J Hum Genet* **25**, 123-129 (2017).
- 996 32. Imielinski, M. *et al.* Common variants at five new loci associated with early-
997 onset inflammatory bowel disease. *Nat Genet* **41**, 1335-1340 (2009).
- 998 33. Chun, S. *et al.* Limited statistical evidence for shared genetic effects of eQTLs
999 and autoimmune-disease-associated loci in three major immune-cell types.
1000 *Nat Genet* **4**, 600-605 (2017).
- 1001 34. The GTEx Consortium. The genotype-tissue expression (GTEx) pilot analysis:
1002 multitissue gene regulation in humans. *Science* **348**, 648-660 (2015).
- 1003 35. Boyle, E.A. *et al.* An expanded view of complex traits: from polygenic to
1004 omnigenic. *Cell* **169**, 1177-1186 (2017).
- 1005 36. Whitehead Pavlides, J.M. *et al.* Predicting targets from integrative analyses of
1006 summary data from GWAS and eQTL studies for 28 human complex traits.
1007 *Genome Medicine* **8**, 84-90 (2016).
- 1008 37. Huler, I. *et al.* Enrichment of inflammatory bowel disease and colorectal
1009 cancer risk variants in colon expression quantitative trait loci. *BMC Genomics*
1010 **16**, 138-153 (2015).
- 1011 38. Gamazon, E.R. *et al.* A gene-based association method for mapping traits
1012 using reference transcriptome data. *Nat Genet* **47**, 1091-1098 (2015).
- 1013 39. De Lange, K.M. *et al.* Genome-wide association study implicates immune
1014 activation of multiple integrin genes in inflammatory bowel disease. *Nat*
1015 *Genet* **49**, 256-261 (2017).

- 1016 40. Libioulle, C. *et al.* Novel Crohn disease locus identified by genome-wide
1017 association maps to a gene desert on 5p13.1 and modulates expression of
1018 PTGER4. *PLoS Genet* **3**, e58 (2007).
- 1019 41. Peltekova, V.D. *et al.* Functional variants of OCTN cation transporter genes are
1020 associated with Crohn disease. *Nat Genet* **39**, 311-318 (2004).
- 1021 42. McCarroll, S.A. *et al.* Deletion polymorphism upstream of IRGM expression
1022 and Crohn's disease. *Nat Genet* **40**, 1107-112 (2008).
- 1023 43. <https://imputation.sanger.ac.uk>
- 1024 44. The 1000 Genomes Project Consortium. A global reference for human genetic
1025 variation. *Nature* **526**, 668-674 (2015).
- 1026 45. Huang, J. *et al.* Improved imputation of low-frequency and rare variants using
1027 the UK10K haplotype reference panel. *Nat Commun* **6**, 8111 (2015).
- 1028 46. McCarthy *et al.* A reference panel of 64,976 haplotypes for genotype
1029 imputation. *Nature Genet.* **48**,1279-1283 (2016).
- 1030 47. Du, P. *et al.* Lumi: a pipeline for processing illumine microarray.
1031 *Bioinformatics* **24**, 1547-1548 (2008).
- 1032 48. Arloth, J. *et al.* Re-Annotator: annotation pipeline for microarray probe
1033 sequences. *PLoS ONE* **10**, e0139516 (2015).
- 1034 49. Johnson, W.E. *et al.* Adjusting batch effects in microarray expression data
1035 using empirical Bayes methods. *Biostatistics* **8**, 118-127 (2007).
- 1036 50. Fairfax, B.P. *et al.* Innate immune activity conditions the effect of regulatory
1037 variants upon monocyte gene expression. *Science* **343**, 1246949 (2014).
- 1038 51. <http://pngu.mgh.harvard.edu/purcell/plink/>
- 1039 52. Purcell, S. *et al.* PLINK: a toolset for whole-genome association and
1040 population-based linkage analysis. *Am. J. Hum. Genet.* **81**, 559-575(2007).
- 1041 53. Nica, A.C. *et al.* Candidate causal regulatory effects by integration of
1042 expression QTLs with complex trait genetic associations. *PLoS Genet* **6**,
1043 e1000895 (2010).
- 1044 54. Farrell, C.M. *et al.* Current status and new features of the Consensus Coding
1045 Sequence database. *Nucleic Acids Res* **42**, D865-72 (2014).
- 1046 55. Momozawa, Y. *et al.* Low-frequency coding variants in CETP and CFB are
1047 associated with susceptibility of exudative age-related macular degeneration
1048 in the Japanese population. *Hum. Mol. Genet.* **25**, 5027-5034 (2016).

- 1049 56. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-
1050 Wheeler transform. *Bioinformatics* **25**, 1754-60 (2009).
- 1051 57. Depristo, M.A. et al. A framework for variation discovery and genotyping
1052 using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491-8 (2011).
- 1053 58. Shigemizu, D. et al. A practical method to detect SNVs and indels from whole
1054 genome and exome sequencing data. *Sci Rep* **3**, 2161 (2013).
- 1055 59. Li, H. et al. The sequence alignment/map (SAM) format and SAMtools.
1056 *Bioinformatics* **25**, 2078-2079 (2009).
- 1057 60. Lek, M. et al. Analysis of protein-coding genetic variation in 60,706 humans.
1058 *Nature* **536**, 285-91 (2016).
- 1059 61. Kumar, P. et al. Predicting the effects of coding non-synonymous variants on
1060 protein function using the SIFT algorithm. *Nat Protoc* **4**, 1073-1081 (2009).
- 1061 62. Adzhubei, I.A. et al. A method and server for predicting damaging missense
1062 mutations. *Nat Methods* **7**, 248-249 (2010).
- 1063 63. Morgenthaler, S. & Thilly, W.G. A strategy to discover genes that carry multi-
1064 allelic or mono-allelic risk for common diseases: a cohort allelic sums test
1065 (CAST). *Mutat Res* **615**, 28-56 (2007).
- 1066 64. Wu, M.C. et al. Rare-variant association testing for sequencing data with the
1067 sequence kernel association test. *Am J Hum Genet* **89**, 82-93 (2011).
- 1068 65. Imielinski, M. et al. Common variants at five new loci associated with early-
1069 onset inflammatory bowel disease. *Nat Genet* **41**, 1335-1340 (2009).
- 1070
- 1071

1072 Figure Legends

1073 **Figure 1: cis Regulatory Module (cRM).** A cis-eQTL affecting gene A in tissue 1
 1074 reveals itself by an “eQTL Association Pattern” ($EAP_{A,1}$), i.e. the pattern of $-\log(p)$
 1075 values for variants in the region. Multiple EAP can be observed in a given
 1076 chromosome region, affecting one or more genes in one or more cell types. EAP
 1077 that are driven by the same underlying variants are expected to be similar, while
 1078 EAP driven by distinct variants (f.i. the green and red regulatory variants in the
 1079 figure) are not. Based on the measure of similarity introduced in this work, ϑ , we
 1080 cluster the EAP in cis-Regulatory Modules (cRM). For EAP in the same module, ϑ
 1081 can be positive or negative, indicating that the variants have the same sign of
 1082 effect (increasing or decreasing expression) for the corresponding EAP pair.

1083 **Figure 2: Single-gene/tissue versus multi-gene/tissue cRM.** Using $|\vartheta| > 0.6$,
 1084 the 23,950 cis-eQTL ($FDR \leq 0.05$) detected in the nine analyzed cell types were
 1085 clustered in 9,691 cis-Regulatory Modules (cRM). 68% of these were single-gene,
 1086 single-tissue cRM (green), 22% were single-gene, multi-tissue cRM (blue), and
 1087 10% were multi-gene, mostly multi-tissue cRM (red). The number of
 1088 observations for single-gene cRM were divided by 10 in the graph for clarity.
 1089 Thus, there are more cases of single-gene, multi-tissue cRM (blue; 2,155) than
 1090 multi-gene cRM (red; 967).

1091 **Figure 3: Example of a multi-gene, multi-tissue cRM.** Gene-tissue
 1092 combinations for which no expression could be detected are marked by “-”, with
 1093 detectable expression but without evidence for cis-eQTL as “→”, with detectable
 1094 expression and evidence for a cis-eQTL as “↑” or “↓” (large arrows: $FDR < 0.05$;
 1095 small arrows: $FDR \geq 0.05$ but high $|\vartheta|$ values). eQTL labelled by the yellow
 1096 arrows constitute the multi-genic and multi-tissular cRM n°57. The
 1097 corresponding regulatory variant(s) increase expression of the *GINM1*, *NUP43*
 1098 and probably *KATNA1* genes (left side of the cRM), while decreasing expression
 1099 of the *PCMT1* and *LRP11* genes (right side of the cRM). The expression of *GINM1*
 1100 in CD15 and *LRP11* in CD4 appears to be regulated in opposite directions by a
 1101 distinct cRM (n°3694, green). The *LATS1* gene, in the same region, is not affected
 1102 by the same regulatory variants in the studied tissues. Inset 1: ϑ values for all

1103 EAP pairs. EAP pairs with $|\vartheta| > 0.6$ are bordered in yellow when corresponding
 1104 to cRM n°57, in green when corresponding to cRM n°3694 (+ green arrow).

1105 **Figure 4: Variant(s) with opposite effects on expression in two cell types.**

1106 Example of a gene (*PNKD*) affected by a cis-eQTL in at least two cell types (CD14
 1107 and platelets) that are characterized by EAP with $\vartheta = -0.97$, indicating that the
 1108 gene's expression level is affected by the same regulatory variant in these two
 1109 cell types, yet with opposite effects, i.e. the variant that is increasing expression
 1110 in platelets is decreasing expression in CD14.

1111 **Figure 5: Significance of the excess sharing of cRM between cell types.** (red:

1112 $p < 0.0002$ (Bonferroni corrected 0.0144), orange: $p < 0.001$ (0.072), rose: $p <$
 1113 0.01 (0.51)). The numbers in the lower-left corner of the squares indicate which
 1114 cRM were used for the analysis: (2) cRM affecting no more than two cell types,
 1115 (3) cRM affecting no more than three cell types, etc. The upper-left square
 1116 indicates the position of the lymphoid cell types (L)(CD4, CD8, CD19), the
 1117 myeloid cell types (M)(CD14,CD15,PLA), and the intestinal cell types (I)(IL, TR,
 1118 RE). For each pair of cell types i and j , we computed two p-values, one using i as
 1119 reference, the other using j as reference (Methods). Pairs of p-values were
 1120 always consistent.

1121 **Figure 6: DAP-matching cRM.** If a regulatory variant (red) affects disease risk

1122 by altering the expression levels of gene B in tissue 2, the $EAP_{B,2}$ is expected to be
 1123 similar (high $|\vartheta|$) to the "disease association pattern" (DAP), both assigned
 1124 therefore to the same cRM. ϑ is positive if increased gene expression is
 1125 associated with increased disease risk, negative otherwise. A cis-eQTL that is
 1126 driven by a regulatory variant (green) that does not directly affect disease risk,
 1127 will be characterized by an EAP (say gene A, tissue 2, $EAP_{A,2}$) that is not similar to
 1128 the DAP (low $|\vartheta|$).

1129 **Figure 7: Screen shots of the CEDAR website,** showing (i) known CD risk loci

1130 on the human karyotype, (ii) a zoom in the HD35 risk locus showing the Refseq
 1131 gene content and summarizing local CEDAR cis-eQTL data (white: no expression
 1132 data, gray: expression data but no evidence for cis-e, black: significant cis-eQTL
 1133 but no correlation with DAP, red: significant cis-eQTL similar to DAP ($\vartheta \leq$
 1134 -0.75), green: significant cis-eQTL similar to DAP ($\vartheta \geq 0.75$)), and (iii) a zoom in

1135 the DAP for Crohn's disease (black) and EAP for *IL18R1* (red), as well as the
1136 signed correlation between DAP and EAP.

1137 **Figure 8: Variants detected by sequencing the coding exons of 45 candidate**
1138 **genes.** Variants are sorted in LoF (Loss-of-Function, i.e. stop gain, frame-shift,
1139 splice site), Damaging MS (missense variants considered as damaging by SIFT⁵
1140 and damaging or possibly damaging by Polyphen-2⁶), Benign MS (other missense
1141 variants), and Synonymous. Blue: variants with $MAF < 0.005$, Red: variants with
1142 $MAF \geq 0.005$.

1143 **Figure 9: QQ-plot for the gene-based burden test.** Ranked $\log(1/p)$ values
1144 obtained when considering LoF and damaging variants (full circles), or
1145 synonymous variants (empty circles). The circles are labeled in blue when the
1146 best p-value for that gene is obtained with CAST, in red when the best p-value is
1147 obtained with SKAT. The black line corresponds to the median $\log(1/p)$ value
1148 obtained (for the corresponding rank) using the same approach on permuted
1149 data (LoF and damaging variants). The grey line marks the upper limit of the
1150 95% confidence band. The name of the genes with nominal p-value ≤ 0.05 are
1151 given. Known causative genes are italicized. The inset p-value corresponds to
1152 the significance of the upwards shift in $\log(1/p)$ values estimated by permutation.

1153

Table 1

Loc	Chr	Beg	End	cRM	Nr	Genes with correlated DAP-EAP	Implicated cell types	Best theta		Best p		Ref
								CD	UC	CD	UC	
HD1	1	2.4	2.8	271	2	<i>TNFRSF14</i>	CD4 CD8 IL TR	-0.74	-0.79	0.02	0.03	4 36
HD2	1	7.7	8.3	2900	1	<i>PARK7</i>	CD15 TR RE	-0.8	-0.82	0.01	0.06	36
N_1_62	1	62.5	63.5	109	3	DOCK7 <i>USP1</i> <i>ATG4C</i>	CD4 CD8 CD19 CD14 CD15	-0.9	0	0.01	1.00	3
N_1_100	1	101.0	102.0	6008	1	<i>SLC30A7</i>	TR	0	-0.71	1.00	0.06	
J_1_119	1	120.2	120.7	9459	1	<i>NOTCH2</i>	CD19	0.68	0	0.13	1.00	
				5	8	<i>GBA</i>	CD4	-0.65	0	0.01	1.00	
HD14	1	155.0	156.1	238	3	<i>THBS3</i> <i>GBA</i> <i>MUC1</i>	CD14 CD15 TR	0	0.81	1.00	0.02	
				4513	1	<i>THBS3</i>	CD4	0	0.66	1.00	0.02	
HD21	1	197.3	198.0	6071	1	<i>DENND1B</i>	CD4	0.7	0.78	0.03	0.02	
HD30	2	62.4	62.7	3716	1	<i>B3GNT2</i>	CD8	-0.63	0	0.01	1.00	
HD35	2	102.8	103.3	1132	1	IL18R1	CD4 CD8	-0.93	-0.87	0.01	0.03	4
				8912	1	(<i>IL18RAP</i>)	CD8	-0.42	0	0.11	0.38	4
J_2_197	2	198.2	199.1	325	2	<i>MARS2</i> <i>PLCL1</i>	CD4 CD14	-0.72	0	0.06	1.00	2 36
J_2_218	2	218.9	219.4	216	3	<i>PNKD</i> <i>GPBAR1</i>	CD14 TR RE	0.72	0.72	0.01	0.06	2 36
HD43	2	234.1	234.6	1177	1	ATG16L1	CD4 CD8 IL TR RE	0.94	0	0.05	1.00	2 39
				2930	1	<i>CCR2</i>	CD19	0.77	0	0.02	1.00	
N_3_45	3	46.0	47.0	1203	1	<i>CCR2</i>	CD4	-0.62	0	0.07	1.00	
				7768	1	<i>CCR9</i>	CD19	0	-0.67	1.00	0.06	
				6798	1	<i>KLHL18</i>	CD14	0	-0.68	1.00	0.03	
				8	7	<i>USP4</i>	CD19	0.64	0.63	0.06	0.07	2
HD50	3	48.4	51.4	217	3	GPX1 <i>APEH</i> <i>IP6K1</i>	CD19 CD14 TR RE	0.91	0.97	0.01	0.01	2 39
				122	3	<i>FAM212A</i>	CD19	0	0.61	1.00	0.05	
J_3_52	3	52.8	53.3	3190	1	<i>SFMBT1</i>	TR RE	0	-0.88	1.00	0.01	37
J_4_73	4	74.6	75.1	1271	1	<i>CXCL5</i>	CD4 CD8 CD19 CD14 PLA	0	-0.84	1.00	0.01	2
HD60	5	40.0	40.7			(<i>PTGER4</i>)	CD15	0	0	0.28	0.15	40
HD61	5	55.4	55.5	360	2	ANKRD55 <i>IL6ST</i>	CD4 CD8	0.9	0	0.02	1.00	4
HD62	5	72.4	72.6	6625	1	<i>FOXD1</i>	IL	-0.74	0	0.03	1.00	4
HD63	5	95.9	96.5	365	2	ERAP2 <i>LNPEP</i>	CD4 CD8 CD19 CD14 CD15	0.94	0.71	0.01	0.02	2 4 37
							PLA IL TR RE					
HD65	5	130.4	132.0	55	4	(<i>SLC22A4</i>) (<i>SLC22A5</i>)	CD4 CD15	-0.55	0	0.06	0.07	4 41
HD66	5	141.4	141.7	2389	1	<i>NDFIP1</i>	CD8 PLA	0.87	0.88	0.04	0.01	2
HD67	5	149.0	151.0	-	-	(<i>JRGM</i>)	-	-	-	-	-	42
HD71	5	173.2	173.6	1349	1	<i>CPEB4</i>	CD4 CD8 CD19 CD14 CD15	-0.92	0	0.01	1.00	2 4
							PLA TR					
J_66_32	6	32.3	32.9	7853	1	<i>HLA-DQA2</i>	IL	0	-0.62	1.00	0.02	
HD76	6	90.8	91.1	1404	1	<i>BACH2</i>	CD4	0.67	0	0.14	1.00	
HD78	6	111.3	112.0	9603	1	<i>SLC16A10</i>	IL	0	-0.71	1.00	0.11	
HD80	6	127.9	128.4	707	2	THEMIS <i>PTPRK</i>	CD8	-0.92	0	0.01	1.00	
HD83	6	167.3	167.6	1425	1	<i>RNASET2</i>	CD4 CD8 CD15 PLA	-0.87	0	0.02	1.00	4
J_7_1	7	2.5	3.0	2729	1	GNA12	CD19 CD14 TR	0	-0.94	1.00	0.02	2
HD84	7	26.6	27.3	1441	1	SKAP2	CD4 CD8 CD19	0.97	0	0.01	1.00	4
HD85	7	28.1	28.3	6438	1	<i>JAZF1</i>	CD4	0.78	0	0.01	1.00	2
HD92	7	128.5	128.8	401	2	<i>IRF5</i> <i>TNPO3</i>	CD15 IL	0	-0.64	1.00	0.02	2 36
				7046	1	<i>TSPAN33</i>	CD19	-0.64	0	0.01	1.00	
				5869	1	<i>PTK2B</i>	CD14	-0.69	0	0.01	1.00	
N_8_26	8	26.7	27.7	5841	1	<i>TRIM35</i>	CD4	0	0.66	1.00	0.01	
HD106	9	139.1	139.5	64	4	CARD9 <i>INPP5E</i> <i>SEC16A</i> <i>SDCCAG3</i>	CD4 CD8 CD19 CD14 CD15	0.95	0.86	0.01	0.02	2 4 37
							IL TR RE					
HD109	10	30.6	30.9	1603	1	<i>MTPAP</i>	TR	-0.62	0	0.11	1.00	
HD112	10	59.8	60.2	1609	1	CISD1	CD4 CD8 CD19 CD14 CD15	0.94	0.83	0.04	0.01	2 4 36
							TR RE					
J_10_74	10	75.4	75.9	436	2	<i>VCL</i>	CD4 CD8 CD19 CD14 RE	0	-0.79	1.00	0.04	
				4279	1	<i>CAM2KG</i>	CD4	-0.67	0	0.04	1.00	
HD114	10	81.0	81.2	5476	1	ZMIZ1	CD8	-0.91	-0.86	0.03	0.01	
J_10_80	10	82.0	82.5	712	2	<i>TSPAN14</i>	TR	-0.71	0	0.01	1.00	
				2216	1	<i>TSPAN14</i>	CD4 CD14	0.76	0	0.01	1.00	2
HD116	10	101.2	101.4	5439	1	<i>SLC25A28</i>	CD14	-0.61	0	0.22	1.00	
J_11_57	11	58.1	58.6	7164	1	<i>ZFP91</i>	PLA	-0.64	-0.75	0.02	0.07	
J_11_59	11	61.3	61.8	1670	1	<i>TMEM258</i>	CD4 CD8 CD19	0.83	0	0.04	1.00	
J_11_65	11	65.4	65.9	451	2	<i>CTSW</i> <i>FIBP</i>	CD4 CD8	-0.73	0	0.01	1.00	2
HD122	11	114.2	114.6	268	3	<i>REXO2</i> <i>NXPE1</i> <i>NXPE4</i>	TR RE	0	-0.89	1.00	0.02	4 37
HD123	11	118.3	118.8	8200	1	<i>TREH</i>	IL	0	0.7	1.00	0.05	
HD142	14	88.2	88.7	8940	1	<i>GPR65</i>	CD14	0.8	0.79	0.01	0.01	
				6353	1	(<i>GALC</i>)	CD14	-0.52	-0.23	0.06	0.06	4
J_15_40	15	41.3	41.8	9109	1	<i>CHP1</i>	IL	0.62	0	0.01	1.00	
J_16_22	16	23.6	24.1	2672	1	<i>PRKCB</i>	CD14	0	0.64	1.00	0.05	2
HD150	16	28.2	29.1	6	8	<i>TUFM</i> <i>SBK1</i> <i>APOBR</i> <i>SGF29</i> <i>CLN3</i>	CD4 CD8 CD19 CD14 CD15 IL	0.81	0.86	0.05	0.03	4
						<i>SPNS1</i>	TR RE					
HD151	16	30.4	31.4	2673	1	<i>RNF40</i>	CD15	-0.63	0	0.02	1.00	
				1886	1	<i>ITGAL</i>	CD4 CD8 CD19	0	0.74	1.00	0.01	39
HD153	16	68.4	68.9	1894	1	<i>ZFP90</i>	CD4 CD8 CD19 CD14 TR	0	0.83	1.00	0.07	2 36
HD156	16	85.9	86.1	3328	1	<i>IRF8</i>	TR RE	0	0.72	1.00	0.01	
HD159	17	37.3	38.3	37	5	GSDMB <i>ORMDL3</i> <i>PGAP3</i> (<i>GSDMA</i>)	CD4 CD8 CD19 CD14 IL TR	-0.98	-0.92	0.02	0.01	2 4
							RE					
HD161	17	40.3	41.0	836	2	<i>STAT3</i>	PLA	0.67	0	0.10	1.00	
HD164	18	67.4	67.6	1988	1	<i>CD226</i>	CD4 CD8 PLA	0	-0.86	1.00	0.01	2
N_18_76	18	76.7	77.7	7292	1	<i>PQLC1</i>	PLA	-0.68	0	0.01	1.00	
HD166	19	10.3	10.7	9232	1	(<i>TYK2</i>)	CD14	-0.44	-0.09	0.10	0.10	
HD168	19	47.1	47.4	581	2	<i>GNG8</i>	CD4	0	-0.63	1.00	0.06	
HD169	19	49.0	49.3	3128	1	FUT2	IL TR RE	-0.95	0	0.01	1.00	4
J_20_31	20	31.1	31.6	593	2	<i>COMMD7</i>	CD14	0	0.61	1.00	0.01	
				7	8	<i>UQCC1</i>	CD19	-0.69	0	0.02	1.00	2
J_20_32	20	33.6	34.1	3369	1	<i>MMP24-AS1</i>	RE	-0.63	-0.71	0.03	0.03	
HD175	20	62.2	62.5	2322	1	<i>LIME1</i>	CD4 CD19	-0.86	0	0.01	1.00	2
HD176	21	16.6	16.9	9578	1	<i>NRIP1</i>	CD4	0	-0.69	1.00	0.02	
HD180	22	21.7	22.1	2130	1	UBE2L3	CD4 CD8 CD19 CD14 CD15	0.97	0.92	0.01	0.07	2 4
							IL TR RE					
N_22_41	22	41.4	42.4	2149	1	<i>EP300</i>	CD8 CD19 CD15	0	0.71	1.00	0.02	

1154

1155 **Table 1:** IBD risk loci for which at least one cis-eQTL association pattern (EAP)
1156 was found to match the disease association pattern (DAP). Given are (i) the
1157 name and chromosomal coordinates of the corresponding loci (Locus, Chr, Beg,
1158 End)(GRCh37/hg19 in Mb), (ii) the identifier and total number of genes in the
1159 matching cis-acting regulatory module (cRM, Nr), (iii) the genes and tissues
1160 involved in matching DAP-EAP ($|\vartheta| > 0.6$) (bold when $|\vartheta| \geq 0.9$), (iv) the best ϑ -
1161 values and corresponding empirical p-values obtained for CD and UC,
1162 respectively, and (vi) references reporting a link between one or more of the
1163 same genes and IBD on the basis of eQTL information. Genes that were
1164 resequenced are shown in italics. Genes that were resequenced despite $|\vartheta| \leq$
1165 0.6 are bracketed, and the supporting references provided in “Ref”. The higher
1166 number of matching DAP-EAP in this study when compared to Huang et al. ⁴ are
1167 primarily due to the fact that (i) we herein study 200 IBD risk loci (vs 97), and
1168 (ii) we increase the number of detected cis-eQTL approximately two-fold by
1169 correcting for hidden confounders using PCs.

1170

1171

1172

1173

1174 *The International IBD Genetics Consortium*

1175

1176 Clara Abraham²³, Jean-Paul Achkar^{24,25}, Tariq Ahmad²⁶, Ashwin N Ananthakrishnan^{27,28}, Vibeke
 1177 Andersen^{29,30,31}, Carl A Anderson³², Jane M Andrews³³, Vito Annese^{34,35}, Guy Aumais^{36,37}, Leonard
 1178 Baidoo³⁸, Robert N Baldassano³⁹, Peter A Bampton⁴⁰, Murray Barclay⁴¹, Jeffrey C Barrett³²,
 1179 Theodore M Bayless⁴², Johannes Bethge⁴³, Alain Bitton⁴⁴, Gabrielle Boucher⁴⁵, Stephan Brand⁴⁶,
 1180 Berenice Brandt⁴³, Steven R Brant⁴², Carsten Büning⁴⁷, Angela Chew^{48,49}, Judy H Cho⁵⁰, Isabelle
 1181 Cleynen²¹, Ariella Cohain⁵¹, Anthony Croft⁵², Mark J Daly^{53,54}, Mauro D'Amato^{55,56,57}, Silvio
 1182 Danese⁵⁸, Dirk De Jong¹¹, Goda Denapiene⁵⁹, Lee A Denson⁶⁰, Kathy L Devaney²⁷, Olivier Dewit⁶¹,
 1183 Renata D'Inca⁶², Marla Dubinsky⁶³, Richard H Duerr^{38,64}, Cathryn Edwards⁶⁵, David Ellinghaus⁶⁶,
 1184 Jonah Essers^{67,68}, Lynnette R Ferguson⁶⁹, Eleonora A Festen¹⁷, Philip Fleshner⁷⁰, Tim Florin⁷¹,
 1185 Andre Franke⁶⁶, Karin Fransen⁷², Richard Gearry^{41,73}, Christian Gieger⁷⁴, Jürgen Glas^{46,75}, Philippe
 1186 Goyette⁴⁵, Todd Green^{54,67}, Anne M Griffiths⁷⁶, Stephen L Guthery⁷⁷, Hakon Hakonarson⁷⁸, Jonas
 1187 Halfvarson⁷⁸, Katherine Hanigan⁵², Talin Haritunians⁷⁰, Ailsa Hart⁷⁹, Chris Hawkey⁸⁰, Nicholas K
 1188 Hayward⁸¹, Matija Hedl²³, Paul Henderson^{82,83}, Xinli Hu⁸⁴, Hailiang Huang^{53,54}, Ken Y Hui⁵⁰, Marcin
 1189 Imielinski³⁹, Andrew Ippoliti⁷⁰, Laimas Jonaitis⁸⁵, Luke Jostins^{86,87}, Tom H Karlsen^{88,89,90}, Nicholas
 1190 A Kennedy⁹¹, Mohammed Azam Khan^{92,93}, Gediminas Kiudelis⁸⁵, Krupa Krishnaprasad⁹⁴, Subra
 1191 Kugathasan⁹⁵, Limas Kupcinskas⁹⁶, Anna Latiano³⁴, Debby Laukens²⁰, Ian C Lawrance^{48,97}, James C
 1192 Lee⁹⁸, Charlie W Lees⁹¹, Marcis Leja⁹⁹, Johan Van Limbergen⁷⁶, Paolo Lionetti¹⁰⁰, Jimmy Z Liu³²,
 1193 Gillian Mahy¹⁰¹, John Mansfield¹⁰², Dunecan Massey⁹⁸, Christopher G Mathew^{103,104}, Dermot PB
 1194 McGovern⁷⁰, Raquel Milgrom¹⁰⁵, Mitja Mitrovic^{72,106}, Grant W Montgomery⁸¹, Craig Mowat¹⁰⁷,
 1195 William Newman^{92,93}, Aylwin Ng^{27,108}, Siew C Ng¹⁰⁹, Sok Meng Evelyn Ng²³, Susanna Nikolaus⁴³,
 1196 Kaida Ning²³, Markus Nöthen¹¹⁰, Ioannis Oikonomou²³, Orazio Palmieri³⁴, Miles Parkes⁹⁸, Anne
 1197 Phillips¹⁰⁷, Cyriel Y Ponsioen¹², Urös Potocnik^{106,111}, Natalie J Prescott¹⁰³, Deborah D Proctor²³,
 1198 Graham Radford-Smith^{52,112}, Jean-Francois Rahier¹¹³, Soumya Raychaudhuri⁸⁴, Miguel Regueiro³⁸,
 1199 Florian Rieder²⁴, John D Rioux^{36,45}, Stephan Ripke^{53,54}, Rebecca Roberts⁴¹, Richard K Russell⁸²,
 1200 Jeremy D Sanderson¹¹⁴, Miquel Sans¹¹⁵, Jack Satsangi⁹¹, Eric E Schadt⁵¹, Stefan Schreiber^{43,66},
 1201 Dominik Schulte⁴³, L Philip Schumm¹¹⁶, Regan Scott³⁸, Mark Seielstad^{117,118}, Yashoda Sharma²³,
 1202 Mark S Silverberg¹⁰⁵, Lisa A Simms⁵², Jurgita Skieceviciene⁸⁵, Sarah L Spain^{32,119}, A. Hillary
 1203 Steinhart¹⁰⁵, Joanne M Stempak¹⁰⁵, Laura Stronati¹²⁰, Jurgita Sventoraityte⁹⁴, Stephan R Targan⁷⁰,
 1204 Kirstin M Taylor¹¹⁴, Anje ter Velde¹², Leif Torkvist¹²¹, Mark Tremelling¹²², Suzanne van
 1205 Sommeren¹⁷, Eric Vasiliauskas⁷⁰, Hein W Verspaget¹⁵, Thomas Walters^{76,123}, Kai Wang³⁹, Ming-
 1206 Hsi Wang^{24,42}, Zhi Wei¹²⁴, David Whiteman⁸¹, Cisca Wijmenga⁷², David C Wilson^{82,83}, Juliane
 1207 Winkelmann^{125,1266}, Ramnik J Xavier^{27,54}, Bin Zhang⁵¹, Clarence K Zhang¹²⁷, Hu Zhang^{128,129}, Wei
 1208 Zhang²³, Hongyu Zhao¹²⁷, Zhen Z Zhao⁸¹

1209

1210 ²³Section of Digestive Diseases, Department of Internal Medicine, Yale School of Medicine, New
 1211 Haven, Connecticut, USA. ²⁴Department of Gastroenterology and Hepatology, Digestive Disease
 1212 Institute, Cleveland Clinic, Cleveland, Ohio, USA. ²⁵Department of Pathobiology, Lerner Research
 1213 Institute, Cleveland Clinic, Cleveland, Ohio, USA. ²⁶Peninsula College of Medicine and Dentistry,

1214 Exeter, UK. ²⁷Gastroenterology Unit, Massachusetts General Hospital, Harvard Medical School,
1215 Boston, Massachusetts 02114, USA. ²⁸Division of Medical Sciences, Harvard Medical School,
1216 Boston, Massachusetts, USA. ²⁹Focused research unit for Molecular Diagnostic and Clinical
1217 Research (MOK), IRS-Center Sonderjylland, Hospital of Southern Jutland, 6200 Åbenrå, Denmark.
1218 ³⁰Institute of Molecular Medicine, University of Southern Denmark, 5000 Odense, Denmark.
1219 ³¹Institute of Regional Health Research , University of Southern Denmark, Odense, Denmark.
1220 ³²Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10
1221 1SA, UK. ³³Inflammatory Bowel Disease Service, Department of Gastroenterology and Hepatology,
1222 Royal Adelaide Hospital, Adelaide, Australia. ³⁴Unit of Gastroenterology, Istituto di Ricovero e
1223 Cura a Carattere Scientifico-Casa Sollievo della Sofferenza (IRCCS-CSS) Hospital, San Giovanni
1224 Rotondo, Italy. ³⁵Strutture Organizzative Dipartimentali (SOD) Gastroenterologia 2, Azienda
1225 Ospedaliero Universitaria (AOU) Careggi, Florence, Italy. ³⁶Facult de Médecine, Universit de
1226 Montréal, Montréal, Québec H3C 3J7, Canada. ³⁷Department of Gastroenterology, Hôpital
1227 Maisonneuve-Rosemont, Montréal, Québec, Canada. ³⁸Division of Gastroenterology, Hepatology
1228 and Nutrition, Department of Medicine, University of Pittsburgh School of Medicine, Pittsburgh,
1229 Pennsylvania 15213, USA. ³⁹Center for Applied Genomics, The Children's Hospital of Philadelphia,
1230 Philadelphia, Pennsylvania, USA. ⁴⁰Department of Gastroenterology and Hepatology, Flinders
1231 Medical Centre and School of Medicine, Flinders University, Adelaide, Australia. ⁴¹Department of
1232 Medicine, University of Otago, Christchurch, New Zealand. ⁴²Meyerhoff Inflammatory Bowel
1233 Disease Center, Department of medicine, Johns Hopkins University School of Medicine, Baltimore,
1234 Maryland, USA. ⁴³Department for General Internal Medicine, Christian-Albrechts-University, Kiel,
1235 Germany. ⁴⁴Division of Gastroenterology, Royal Victoria Hospital, Montréal, Québec, Canada.
1236 ⁴⁵Research Center, Montreal Heart Institute, Montréal, Québec H1T 1C8, Canada. ⁴⁶Department of
1237 Medicine II, Ludwig-Maximilians-University Hospital Munich-Grosshadern, Munich, Germany.
1238 ⁴⁷Department of Gastroenterology, Campus Charité Mitte, Universitätsmedizin Berlin, Berlin,
1239 Germany. ⁴⁸Harry Perkins Institute for Medical Research, School of Medicine and Pharmacology,
1240 University of Western Australia, Murdoch, Western Australia 6150, Australia. ⁴⁹IBD unit ,
1241 Fremantle Hospital, Fremantle, Australia. ⁵⁰Department of Genetics, Yale School of Medicine, New
1242 Haven, Connecticut 06510, USA. ⁵¹Department of Genetics and Genomic Sciences, Mount Sinai
1243 School of Medicine, New York, New York, USA. ⁵²Inflammatory Bowel Diseases, Genetics and
1244 Computational Biology, Queensland Institute of Medical Research, Brisbane, Australia. ⁵³Analytic
1245 and Translational Genetics Unit, Massachusetts General Hospital, Harvard Medical School,
1246 Boston, Massachusetts 02114, USA. ⁵⁴Broad Institute of MIT and Harvard, Cambridge,
1247 Massachusetts 02141, USA. ⁵⁵Clinical Epidemiology Unit, Department of Medicine Solna,
1248 Karolinska Institutet, 17176 Stockholm, Sweden. ⁵⁶Department of Gastrointestinal and Liver
1249 Diseases, BioDonostia Health Research Institute, 20014 San Sebastián, Spain. ⁵⁷IKERBASQUE,
1250 Basque Foundation for Science, 48013 Bilbao, Spain. ⁵⁸IBD Center, Department of
1251 Gastroenterology, Istituto Clinico Humanitas, Milan, Italy. ⁵⁹Center of hepatology,
1252 Gastroenterology and Dietetics, Vilnius University, Vilnius, Lithuania. ⁶⁰Pediatric

1253 Gastroenterology, Cincinnati Children's Hospital Medical Center, Cincinnati, Ohio, USA.
1254 ⁶¹Department of Gastroenterology, Université Catholique de Louvain (UCL) Cliniques
1255 Universitaires Saint-Luc, Brussels, Belgium. ⁶²Division of Gastroenterology, University Hospital
1256 Padua, Padua, Italy. ⁶³Department of Pediatrics, Cedars Sinai Medical Center, Los Angeles,
1257 California, USA. ⁶⁴Department of Human Genetics, University of Pittsburgh Graduate School of
1258 Public Health, Pittsburgh, Pennsylvania 15261, USA. ⁶⁵Department of Gastroenterology, Torbay
1259 Hospital, Torbay, Devon, UK. ⁶⁶Institute of Clinical Molecular Biology, Christian-Albrechts-
1260 University of Kiel, 24118 Kiel, Germany. ⁶⁷Center for Human Genetic Research, Massachusetts
1261 General Hospital, Harvard Medical School, Boston, Massachusetts, USA. ⁶⁸Pediatrics, Harvard
1262 Medical School, Boston, Massachusetts, USA. ⁶⁹Faculty of Medical & Health Sciences, School of
1263 Medical Sciences, The University of Auckland, Auckland, New Zealand. ⁷⁰F. Widjaja Foundation
1264 Inflammatory Bowel and Immunobiology Research Institute, Cedars-Sinai Medical Center, Los
1265 Angeles, California 90048, USA. ⁷¹Department of Gastroenterology, Mater Health Services,
1266 Brisbane, Australia. ⁷²Department of Genetics, University Medical Center Groningen, Groningen,
1267 The Netherlands. ⁷³Department of Gastroenterology, Christchurch Hospital, Christchurch, New
1268 Zealand. ⁷⁴Institute of Genetic Epidemiology, Helmholtz Zentrum München - German Research
1269 Center for Environmental Health, Neuherberg, Germany. ⁷⁵Department of Preventive Dentistry
1270 and Periodontology, Ludwig-Maximilians-University Hospital Munich-Grosshadern, Munich,
1271 Germany. ⁷⁶Division of Pediatric Gastroenterology, Hepatology and Nutrition, Hospital for Sick
1272 Children, Toronto, Ontario, Canada. ⁷⁷Department of Pediatrics, University of Utah School of
1273 Medicine, Salt Lake City, Utah, USA. ⁷⁸Department of Gastroenterology, Faculty of Medicine and
1274 Health, Örebro University, SE-70182 Örebro, Sweden. ⁷⁹Department of Medicine, St Mark's
1275 Hospital, Harrow, Middlesex, UK. ⁸⁰Nottingham Digestive Diseases Centre, Queens Medical Centre,
1276 Nottingham, UK. ⁸¹Molecular Epidemiology, Genetics and Computational Biology, Queensland
1277 Institute of Medical Research, Brisbane, Australia. ⁸²Paediatric Gastroenterology and Nutrition,
1278 Royal Hospital for Sick Children, Edinburgh, UK. ⁸³Child Life and Health, University of Edinburgh,
1279 Edinburgh, Scotland, UK. ⁸⁴Division of Rheumatology Immunology and Allergy, Brigham and
1280 Women's Hospital, Boston, Massachusetts, USA. ⁸⁵Academy of Medicine, Lithuanian University of
1281 Health Sciences, Kaunas, Lithuania. ⁸⁶Wellcome Trust Centre for Human Genetics, University of
1282 Oxford, Headington OX3 7BN, UK. ⁸⁷Christ Church, University of Oxford, St Aldates OX1 1DP, UK.
1283 ⁸⁸Research Institute of Internal Medicine, Department of Transplantation Medicine, Division of
1284 Cancer, Surgery and Transplantation, Oslo University Hospital Rikshospitalet, Oslo, Norway.
1285 ⁸⁹Norwegian PSC Research Center, Department of Transplantation Medicine, Division of Cancer,
1286 Surgery and Transplantation, Oslo University Hospital Rikshospitalet, Oslo, Norway. ⁹⁰K.G. Jebsen
1287 Inflammation Research Centre, Institute of Clinical Medicine, University of Oslo, Oslo,
1288 Norway. ⁹¹Gastrointestinal Unit, Western General Hospital University of Edinburgh, Edinburgh,
1289 UK.
1290 ⁹²Genetic Medicine, Manchester Academic Health Science Centre, Manchester, UK. ⁹³The
1291 Manchester Centre for Genomic Medicine, University of Manchester, Manchester, UK. ⁹⁴QIMR
1292 Berghofer Medical Research Institute, Royal Brisbane Hospital, Brisbane, Australia. ⁹⁵Department

1293 of Pediatrics, Emory University School of Medicine, Atlanta, Georgia, USA. ⁹⁶Department of
1294 Gastroenterology, Kaunas University of Medicine, Kaunas, Lithuania. ⁹⁷Centre for Inflammatory
1295 Bowel Diseases, Saint John of God Hospital, Subiaco, Western Australia 6008, Australia.
1296 ⁹⁸Inflammatory Bowel Disease Research Group, Addenbrooke's Hospital, Cambridge CB2 0QQ,
1297 UK. ⁹⁹Faculty of medicine, University of Latvia, Riga, Latvia. ¹⁰⁰Dipartimento di Neuroscienze,
1298 Psicologia, Area del Farmaco e Salute del Bambino, Università di Firenze Strutture Organizzative
1299 Dipartimentali (SOD) Gastroenterologia e Nutrizione Ospedale pediatrico Meyer, Firenze, Italy.
1300 ¹⁰¹Department of Gastroenterology, The Townsville Hospital, Townsville, Australia.
1301 ¹⁰²Institute of Human Genetics, Newcastle University, Newcastle upon Tyne, UK. ¹⁰³Department of
1302 Medical and Molecular Genetics, King's College London, London SE1 9RT, UK. ¹⁰⁴Sydney Brenner
1303 Institute for Molecular Bioscience, University of the Witwatersrand, Johannesburg 2193, South
1304 Africa. ¹⁰⁵Inflammatory Bowel Disease Centre, Mount Sinai Hospital, Toronto, Ontario, Canada.
1305 ¹⁰⁶Center for Human Molecular Genetics and Pharmacogenomics, Faculty of Medicine, University
1306 of Maribor, Maribor, Slovenia. ¹⁰⁷Department of Medicine, Ninewells Hospital and Medical School,
1307 Dundee, UK. ¹⁰⁸Center for Computational and Integrative Biology, Massachusetts General
1308 Hospital, Harvard Medical School, Boston, Massachusetts, USA. ¹⁰⁹Department of Medicine and
1309 Therapeutics, Institute of Digestive Disease, Chinese University of Hong Kong, Hong Kong.
1310 ¹¹⁰Department of Genomics Life & Brain Center, University Hospital Bonn, Bonn, Germany.
1311 ¹¹¹Faculty for Chemistry and Chemical Engineering, University of Maribor, Maribor, Slovenia.
1312 ¹¹²Department of Gastroenterology, Royal Brisbane and Womens Hospital, Brisbane, Australia.
1313 ¹¹³Department of Gastroenterology, Université Catholique de Louvain (UCL) Centre Hospitalier
1314 Universitaire (CHU) Mont-Godinne, Mont-Godinne, Belgium. ¹¹⁴Department of Gastroenterology,
1315 Guy's & St Thomas' NHS Foundation Trust, St-Thomas Hospital, London, UK. ¹¹⁵Department of
1316 Digestive Diseases, Hospital Quiron Teknon, Barcelona, Spain. ¹¹⁶Department of Public Health
1317 Sciences, University of Chicago, Chicago, Illinois, USA. ¹¹⁷Human Genetics, Genome Institute of
1318 Singapore, Singapore. ¹¹⁸Institute for Human Genetics, University of California, San Francisco, San
1319 Francisco, California, USA. ¹¹⁹Open Targets, Wellcome Trust Genome Campus, Hinxton,
1320 Cambridgeshire CB10 1SD, UK. ¹²⁰Department of Biology of Radiations and Human Health,
1321 Agenzia nazionale per le nuove tecnologie l'energia e lo sviluppo economico sostenibile (ENEA),
1322 Rome, Italy. ¹²¹Department of Clinical Science Intervention and Technology, Karolinska Institutet,
1323 Stockholm, Sweden. ¹²²Gastroenterology & General Medicine, Norfolk and Norwich University
1324 Hospital, Norwich, UK. ¹²³Faculty of medicine, University of Toronto, Toronto, Ontario, Canada.
1325 ¹²⁴Department of Computer Science, New Jersey Institute of Technology, Newark, New Jersey,
1326 USA. ¹²⁵Institute of Human Genetics, Technische Universität München, Munich,
1327 Germany. ¹²⁶Department of Neurology, Technische Universität München, Munich, Germany.
1328 ¹²⁷Department of Biostatistics, School of Public Health, Yale University, New Haven, Connecticut,
1329 USA. ¹²⁸Department of Gastroenterology, West China Hospital, Chengdu, Sichuan, China. ¹²⁹State
1330 Key Laboratory of Biotherapy, Sichuan University West China University of Medical Sciences
1331 (WCUMS), Chengdu, Sichuan, China.
1332

1

Supplementary Information

2

3 **IBD risk loci are enriched in multigenic regulatory modules encompassing**
4 **causative genes.**

5 Momozawa et al.

6 **Supplementary note 1: Genes with strong DAP-EAP correlation**

7 ***IL18R1*** encodes the IL-18r1, the receptor of IL-18, a potent proinflammatory
8 cytokine governing host-microorganism homeostasis and is postulated to play a
9 role in IBD^{1,2}. However, IL-18/IL-18r1 precise contribution to the disease remains
10 controversial. Indeed, compared to wild-type mice, *Il18*^{-/-} and *Il18r1*^{-/-} full KO mice
11 are more susceptible to AOM/DSS-induced colitis and polyp formation³. However,
12 targeted deletion of *Il18*^{-/-} and *Il18r1*^{-/-} in intestinal epithelial cells confers
13 protection from colitis and mucosal damage in mice⁴. In human, several studies
14 have associated circulating or local IL-18 with IBD severity, suggesting that IL-18
15 could be an effector cytokine in IBD⁵.

16 ***IL6ST*** encodes the interleukin 6 signal transducer protein (IL6ST), also called IL6
17 beta, GP130 or CD130. IL6ST is a common transmembrane receptor for all family
18 members of IL6 that include IL-6, IL-11, ciliary neurotrophic factor (CNTF),
19 cardiotrophin-1 (CT-1), cardiotrophin like cytokine (CLC), leukaemia inhibitory
20 factor (LIF), oncostatin M (OSM), neuropoitin (NPN) and interleukin-27 (IL-27)⁶.
21 IL6 family members / IL6ST signaling pathways involve the activation of JAK
22 (Janus kinase) family members, leading to the activation of STAT (signal
23 transducers and activators of transcription) family, as well as the activation of
24 MAPK (mitogen-activated protein kinase) pathway. These pathways are involved
25 in cell survival, apoptosis, differentiation and proliferation⁶. The involvement of
26 IL6/IL6ST/STAT3 in the pathophysiology of IBD is well documented⁷. Indeed,
27 high circulating levels of IL6 is associated with increased severity of the disease⁷.
28 T cells from IBD patient show increased STAT3 activation with increased
29 expression of IL6ST and enhanced resistance to apoptosis⁸. A pilot clinical trial
30 (phase I) targeting of IL6/IL6ST pathway in patients with CD has shown that
31 blocking this pathway has effects similar to the inhibition of TNF^{9,10}.

32 ***THEMIS*** encodes the thymocyte-expressed molecule involved in selection
33 (THEMIS), the expression of which is limited to lymphoid tissues. In mice, THEMIS
34 is highly expressed in pre-TcR thymocytes and plays an important role in T-cell
35 development and TCR activation signaling^{11,12}. Its expression is reduced in
36 differentiated T lymphocytes¹². THEMIS deficiency in mice is associated with the

37 presence of higher percent of T_{reg} cells, with reduced TCR-mediated T cell
38 response, increased proportion of memory CD4 and CD8 T cells and reduced
39 proportions of naïve-phenotype populations¹². Interestingly, all these T cells
40 associated features are implicated in the pathogenesis of IBD. Indeed, lamina
41 propria T cells in IBD are hypo-responsive to TCR stimulation and high number of
42 effector T cells are present in the inflamed bowel¹³. As for T_{reg}, only moderate
43 expansion was seen in intestinal lesions of Crohn's patients suggesting that their
44 suppressive activity is probably not sufficient against the overwhelming effector
45 T cells activity¹³.

46 **APEH** encodes the acylpeptide hydrolase (APEH) enzyme that contributes to
47 protein degradation processes in concert with the proteasome. It catalyzes the
48 removal of *N*-acylated amino acids from acetylated peptides¹⁴. Its physiological
49 role is not well understood. SNPs in APEH gene have been associated with both CD
50 and UC¹⁵. Like other ubiquitin proteasome systems (UPS) such as USP40 or CYLD,
51 APEH may also regulate the NF-κB pathway. Under this scenario, an alteration of
52 NF-κB signaling may lead to aberrant immune response and inflammation.

53 **ANKRD55** encodes an Ankyrin repeat domain-containing protein 55 with
54 unknown function. Ankyrin repeats are composed of 33-34 aa and are the most
55 abundant motifs in nature with highly diverse cellular functions¹⁶. SNPs at the
56 ANKRD55 locus have also been associated with multiple sclerosis¹⁷ and RA¹⁸.

57 **CISD1** gene encodes a highly conserved iron-sulfur domain-containing protein A,
58 known as mitoNEET. This iron-containing protein is a dynamic redox-sensitive
59 molecule that serves an important role in mitochondrial functions. It participates
60 in critical process such as electron shuttling through the electron transport chain,
61 regulation of enzymatic activity, and synthesis of heme and iron-sulfur
62 clusters^{19,20}. Deregulation of iron metabolism and associated anemia has been
63 associated with IBD²¹. The role that mitoNEET plays in the etiology of IBD remains
64 to be determined.

65 **CPEB4** gene encodes the cytoplasmic polyadenylation element-binding protein 4
66 (CEBP4), which belongs to a family of proteins that bind mRNAs and contain a
67 cytoplasmic polyadenylation element (CPE) in their 3'-UTR. Binding results in 3'-

68 poly(A) tail extension and translational upregulation of target mRNAs. *Cpeb4*
69 mRNA is rhythmically regulated in mouse liver, conferring temporal translational
70 regulation. In the absence of CPEB4, a large number of mRNAs are transcribed,
71 but remain untranslated until needed²². A recent study, using knockout mice
72 models, showed that CPEB4 was required for translation of numerous proteins
73 involved in ER homeostasis and CPEB4 loss resulted in mitochondrial dysfunction
74 and defective lipid metabolism, two hallmarks of ER stress. *Cpeb4* KO livers were
75 highly susceptible to ER stress-induced apoptosis and to development of NAFLD²³.
76 In CD, reduced CPEB4 may also lead to ER stress and mitochondrial dysfunction.

77 **DOCK7** encodes dedicator of cytokinesis 7 protein (Dock7), a member of Dock
78 proteins family and an activator of Rac GTPases. DOCK7 plays an important role
79 in axon outgrowth, Schwann cell migration, and axon myelination²⁴. Mutation in
80 this gene in mice leads to hypopigmentation suggesting a non-redundant role in
81 the distribution and function of dermal and follicular melanocytes. However,
82 mutant mice show normal neuronal function despite the high expression of DOCK7
83 in the developing brain, suggesting redundancy with other Docks²⁵. The role of
84 DOCK7 in IBD and immune cells function is totally unknown.

85 **ERAP2** gene encodes an endoplasmic reticulum aminopeptidase (ERAP2), an
86 enzyme involved in trimming of peptides for MHC-I loading. Aberrant ERAP2
87 function could influence peptide-HLA-B27 stability, formation of MHC-I free heavy
88 chains and ER stress^{26,27,28}. SNPs in *ERAP2* gene have been associated with CD²⁹.
89 Although the underlying mechanisms are not known, it is possible that ERAP2
90 modification contributes to the reported reduction of MHCI on CD4 T cells from
91 CD patients³⁰. ERAP2 modification may also contribute to the epithelial ER stress
92 associated with CD and UC.

93 **GNA12** encodes Guanine nucleotide-binding protein subunit alpha-12 or $G\alpha_{12}$,
94 which belongs to the heterotrimeric G proteins. $G\alpha_{12}$ is found in tight junctions
95 (TJ) where it interacts with ZO-1³¹ and plays important roles in para-cellular
96 permeability^{32,33}. $G\alpha_{12}$ is ubiquitously expressed and interacts, upon receptor-
97 mediated activation, with certain Rho guanine nucleotide exchange factors
98 (RhoGEFs) which in turn mediate activation of the small GTPase RhoA³⁴. Intestinal

99 permeability and barrier dysfunction is a hallmark of CD and UC. Several studies
100 reported changes in the expression of several TJ proteins in both diseases³⁵. It is
101 conceivable that modifications in the $G\alpha_{12}$ pool leads to alteration of intestinal
102 permeability. Tissue-specific $G\alpha_{12}$ -deficient mice revealed important functions of
103 this protein in modulating T cell trafficking and proliferation, as well as in the
104 response to foreign and self antigens³⁶, important processes that may affect
105 susceptibility for T cell-mediated diseases.

106 **GPX1** encodes the glutathione peroxidase 1 (GPX1), a highly abundant and
107 ubiquitously expressed cytosolic enzyme. Like all glutathione peroxidases family
108 members, GPX1 catalyzes the reduction of H₂O₂ by glutathione and consequently,
109 protects cells from oxidative damage. In IBD, it is believed that intestinal and
110 colonic injuries and dysfunction is at least partially due to elevation of reactive
111 metabolites of oxygen and nitrogen³⁷. Although the role of GPX1 is not known in
112 IBD, deficiency of both GPX1 and GPX2 in mice lead to spontaneous ileo-colitis and
113 intestinal cancer³⁸. A protective role of GPX1 and GPX2 against oxidative stress has
114 also been suggested by studies reporting elevated *Gpx1/2* gene expression in
115 gastric mucosa after *H. pylori* infection³⁹. Association of the elevated expression of
116 *Gpx1/2* gene with tumorigenesis could be due to its anti-apoptotic activity⁴⁰.

117 **GSDMB** encodes Gasdermin-B protein (GSDMB) the function of which is largely
118 unknown. The expression of *GSDMB* has been associated with differentiated
119 epithelial cells and with regions containing proliferating cells or stem cells,
120 respectively, of the esophagus and the gastric mucosa^{41,42}.

121 **JAZF1**, also known as *TIP27*, encodes a transcriptional repressor of *NR2C2*, also
122 known as *TAK1* or *TR4*⁷⁶. Mice deficient in *NR2C2* show low IGF1 serum
123 concentrations and perinatal and early postnatal hypoglycemia, as well as growth
124 retardation⁷⁷. *JAZF1* also affects variation in human height⁷⁸. SNPs in *JAZ1F* have
125 been associated with type II diabetes⁷⁹, prostate⁸⁰ and endometrial cancer⁸¹ and
126 with systemic lupus erythematosus⁸². However, the role of *JAZF1* in immune
127 response and autoimmunity remains to be elucidated.

128 **LSP1** encodes a leukocyte-specific protein 1 (LSP1), a Ca²⁺-activated, intracellular

129 filamentous actin-binding protein that interacts with the cytoskeleton and is
130 expressed in hematopoietic lineage and in endothelial cells⁷⁰. Evidence from mice
131 model studies suggest that LSP1 plays a negative regulatory role on neutrophil
132 and T cell migration^{71,72}. A recent study identified a novel *LSP1* deletion variant
133 for RA susceptibility through CNV GWAS⁷³. The copy number of *LSP1* was found
134 to be significantly lower in RA patients and was associated with increased T cell
135 migration⁷³. We found a positive correlation of LSP1 expression (in CD14⁺ cells)
136 with UC, but not with CD. UC, as well as CD, is characterized by an increased
137 infiltration of immune cells in inflamed tissues. Our finding is therefore surprising
138 if we consider the concept of an association between increased cell migration
139 with *LSP1* CNVs and LSP1 insufficiency. It is possible that LSP1 plays an additional,
140 yet unknown role in monocytes. On the other hand, if LSP1 participates actively in
141 the cross-talk between leukocytes and endothelial cells during leukocyte
142 transmigration, the physiological differences in microvasculature and the
143 integrins involved may dictate organ-specific roles for LSP1 in leukocyte
144 recruitment into the inflammatory sites.

145 ***NXPE1***: Encodes Neuroexophilin and PC-esterase domain family member 1
146 (*NXPE1*). A human gastrointestinal tract (GIT) specific transcriptome and
147 proteome study validate the expression pattern of this gene and protein in the
148 intestine⁷⁴. *NXPE1* was recently identified as a novel target gene for IBD-
149 associated variants⁷⁵. Its function remains largely unknown.

150 ***ORMDL3*** encodes ORM1-like protein 3, a negative regulator of sphingolipid
151 synthesis and a regulator of endoplasmic reticulum-mediated calcium signaling⁴⁵.
152 *ORMDL3* is involved in the regulation of eosinophil and T cell functions^{46,47}. It also
153 facilitate B cells survival and regulates autophagy through the ATF6 signaling
154 pathway⁴⁸. Genetic variants regulating *ORMDL3* expression have been associated
155 with susceptibility to asthma⁴⁹, T1D⁵⁰, atherosclerosis⁵¹, ankylosing spondylitis⁵²
156 and IBD⁵³. *ORMDL3* might be associated with IBDs and other autoimmune and
157 inflammatory diseases by activating ERS, inducing autophagy and/or promoting
158 immune cells activation.

159 **REXO2** encodes an oligoribonuclease protein. Its depletion, using RNAi, causes a
160 significant decrease of mtDNA and mtRNA and impaired *de novo* mitochondrial
161 protein synthesis⁸³. REXO2's function remains unknown but it may be involved in
162 the well documented mitochondrial defects associated with IBD⁸⁴.

163 **RNASET2** is the only RNase T2 family member in humans and is potentially
164 involved in the inhibition of tumorigenesis, metastasis and angiogenesis^{85,86}. Loss-
165 of-function of RNASET2 protects fibroblasts from oxidative stress⁸⁹ while its
166 overexpression in melanocytes and keratinocytes sensitizes these cells to
167 oxidative-stress-induced apoptosis⁹⁰. Interestingly, CD is characterized by an
168 impaired immune cells apoptosis associated with elevated H₂O₂ in PBMC during
169 the active phase of the disease⁹¹. Although speculative, it is possible that reduced
170 **RNASET2** contributes to the altered oxidative stress in CD.

171 **SKAP2** encodes the Src kinase-associated phosphoprotein 2 (Skap2), a cytosolic
172 adaptor protein expressed in a variety of cell types including hematopoietic
173 cells^{54,55,56}. Skap2 has been implicated in cell adhesion through association to
174 integrins and cytoplasmic actin⁵⁵, and is required for global actin reorganization.
175 It interacts with different molecules implicated in integrin signaling events^{54,56,57}.
176 Loss of Skap2 in mice results in reduced inflammation in experimental
177 autoimmune encephalomyelitis as well as defects in macrophage migration into
178 tumor metastasis, suggesting a physiologically important role of Skap2 for
179 leukocyte recruitment *in vivo*^{55,58}.

180 **UBE2L3** gene encodes an atypical Ubiquitin E2 Conjugase (UBE2L3) the role of
181 which has been recently uncovered. It is an indirect human and mouse Caspase-1
182 target and plays an important role in the maturation of IL-1 β . UBE2L3 depletion
183 in mice increases pro-IL-1 β levels and mature-IL-1 β secretion by
184 inflammasomes⁶¹. Several GWAS identified polymorphisms in the genomic locus
185 of **UBE2L3** that are associated with multiple autoimmune diseases⁶² including
186 CD²⁹. Decreased secretion of the inflammasome cytokine IL-1 β was noted in
187 monocytes of Crohn's disease patients⁶³. It is therefore tempting to speculate that
188 **UBE2L3** contributes to disease at least partially by modulating IL-1 β secretion.

189 **ZMIZ1** encodes Zmiz1, a member of the protein inhibitor of activated STAT (PIAS)-
190 like family of coregulators⁶⁴. Zmiz1 is widely and variably expressed⁶⁵. In GWAS,
191 a SNP within ZMIZ1 gene was associated with early-onset Crohn's disease and
192 IBD⁶⁶. *ZMIZ1* is co-expressed with activated *NOTCH1* across a broad range of T-
193 ALL oncogenic subgroups. Its inhibition slows human T-ALL cell proliferation
194 and/or sensitizes them to γ -Secretase inhibitors (GSI)⁶⁷. Evidence from Zmiz1-
195 deficient mice demonstrated that Zmiz1 is a direct Notch1 cofactor that
196 heterogeneously regulates Notch1 target genes and plays an important role in T
197 cells development⁶⁸. Altered expression of *ZMIZ1* has been reported to affect
198 Smad3-mediated transcription⁶⁹. Interestingly, our analysis shows that increased
199 UC disease risk was associated with decrease of both *SMAD3* and *ZMIZ1*
200 expression while no association was observed with *NOTCH1*. This association was
201 observed in different tissues/cell types suggesting a possible trans effect of *ZMIZ1*
202 on *SMAD3* expression.

203

204

205 **Supplementary Table 1**

206

Cell type	Naive (r^2 based)			Frequentist (Nica et al., 2010)			Theta-based		
	Overlaps observed	Overlaps expected	P value	Overlaps observed	Overlaps expected	P value	Overlaps observed	Overlaps expected	P value
CD4	12	3.3	< 0.01	14	4.9	< 0.01	17	8.4	< 0.01
CD8	12	3.5	< 0.01	18	4.3	< 0.01	16	6.9	< 0.01
CD14	8	3.3	0.061	9	4.7	0.211	10	7.1	0.720
CD15	4	1.9	0.646	4	2	0.720	7	5.1	0.909
CD19	7	2	0.010	7	3.6	0.410	12	5.8	0.044
PLA	4	0.9	0.010	3	0.9	0.475	5	1.8	0.119
IL	4	1.6	0.432	7	2.1	0.027	8	4.1	0.281
TR	6	2.6	0.211	5	3.5	0.928	11	6	0.086
RE	5	1.5	0.103	6	2.4	0.204	9	5.5	0.509

207

208 Enrichment of DAP-EAP matching in 63 of 97 CD risk loci covered by the ImmunoChip. For each cell type, we provide the number of
209 matches (or overlaps) observed with the top disease-associated SNPs ($MAF > 0.05$), as well as the number of matches expected with the
210 same number of SNPs ($MAF > 0.05$) sampled at random in the same 63 risk loci. The analyses were conducted using three "colocalisation"
211 methods (Naive, Frequentist and Theta-based). The p-values were determined by simulation (1,000 sets of 63 randomly sampled SNPs)
212 and Bonferroni corrected for the analysis of 9 cell types. < 0.01 means that the number of matches observed with the real disease-
213 associated SNPs was never observed with any set (out of 1,000) of randomly sampled SNPs.

214 **Supplementary Table 2**

215

Tissue	Nr of samples	Nr of probes	Nr of PCs
CD4	303	13,466	38
CD8	294	13,317	35
CD19	282	12,648	40
CD14	286	13,170	36
CD15	289	11,069	27
PLA	251	6,565	23
IL	200	15,401	59
TR	271	15,082	50
RE	267	14,844	53

216

217 Number of usable samples, probes and PC for each tissue type.

218

219 **Supplementary Table 3**

220

Tissue	Nr of probes	FDR≤0.25	FDR≤0.10	FDR≤0.05	FDR≤0.01
CD4	13,466	7,417	4,957	4,176	3,247
CD8	13,317	6,760	4,309	3,599	2,779
CD19	12,648	4,984	3,138	2,549	1,953
CD14	13,170	7,118	4,728	3,961	3,106
CD15	11,069	3,611	2,396	1,983	1,512
PLA	6,565	1,404	996	854	653
IL	15,401	2,769	1,728	1,426	1,031
TR	15,082	5,183	3,391	2,807	2,160
RE	14,844	4,180	2,726	2,295	1,731

221

222 Number of cis-eQTL found in the nine analyzed cell types for different FDR
 223 thresholds (see also Suppl. Figure 7).

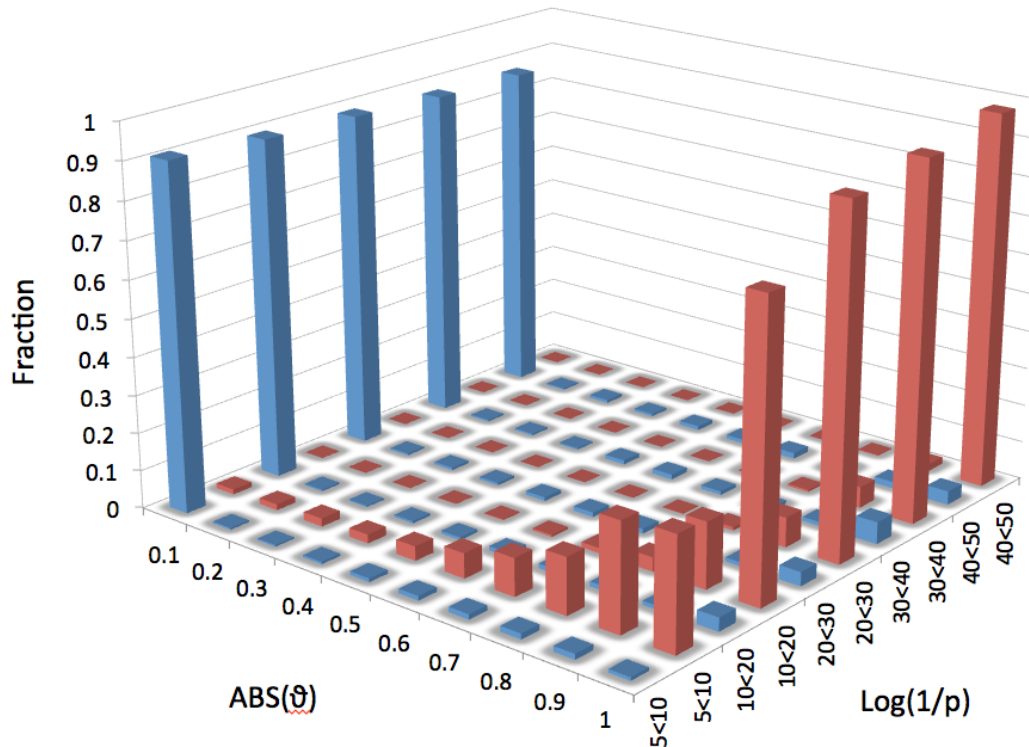
224

Supplementary Figures

225

226

1. Supplementary Figure 1



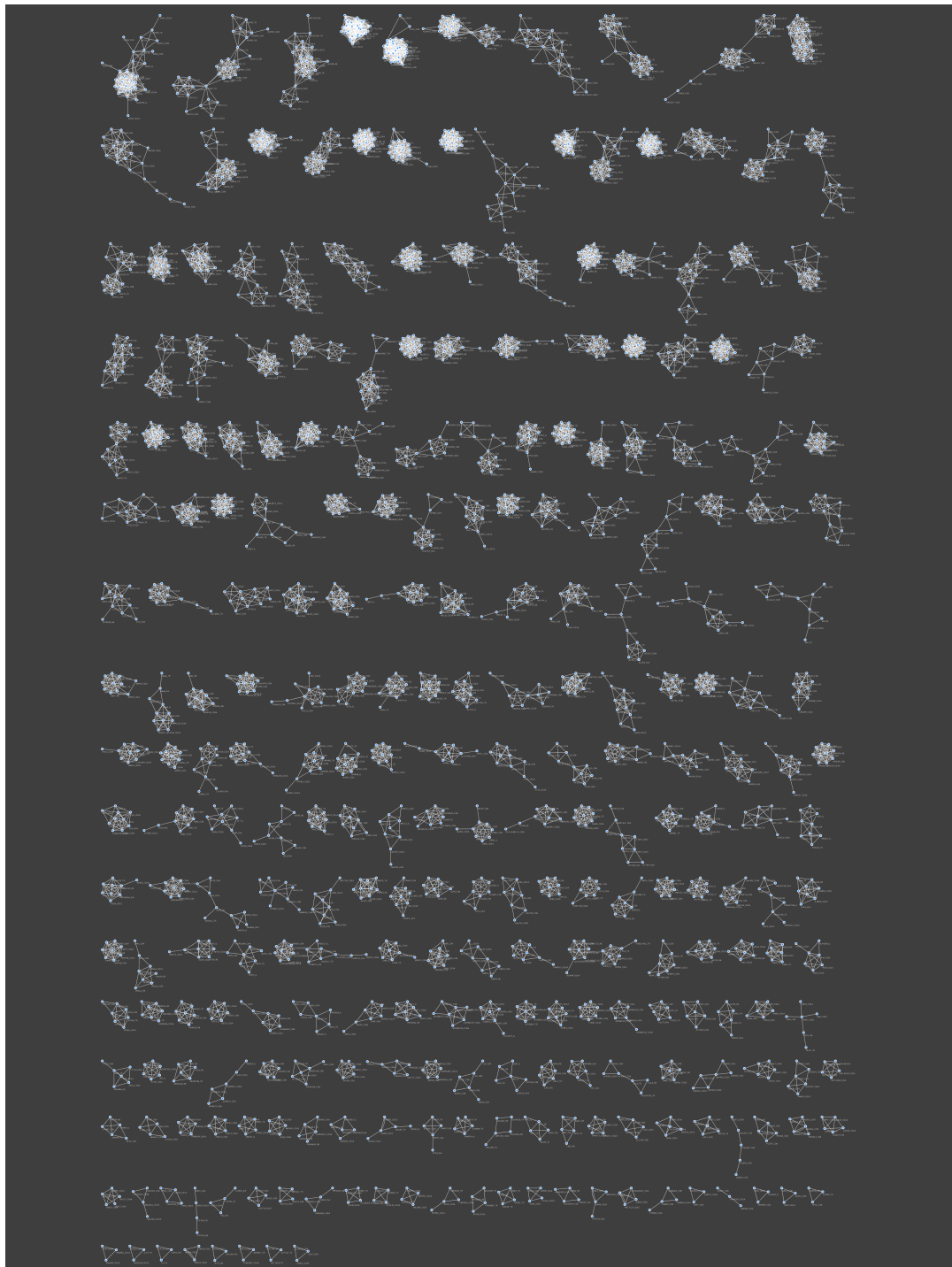
227

228 Absolute values of ϑ for pairs of eQTL driven by distinct regulatory variants (blue),
 229 and for pairs of eQTL driven by the same regulatory variants (red). The first (blue)
 230 were obtained by confronting real cis-eQTL with in silico simulated eQTL
 231 explaining the same variance as the real eQTL but driven by a randomly chosen
 232 SNPs in a 2Mb window centered around the probe. The second (red) were
 233 obtained by confronting eQTL obtained by reanalyzing two mutually exclusive
 234 halves of the CEDAR population separately in a region harboring a real cis-eQTL.
 235 It can be seen that ϑ very effectively discriminates between pairs of eQTL driven
 236 by distinct (blue) vs the same (red) regulatory variants. By choosing 0.6 as
 237 threshold value for ϑ , one captures most red pairs (~88%) with minimum
 238 contamination of blue pairs (~5%). $\text{Log}(1/p)$: eQTL are sorted by the smallest
 239 $\text{log}(1/p)$ value of the two eQTL being compared.

240

241 **2. Supplementary Figure 2**

242



243

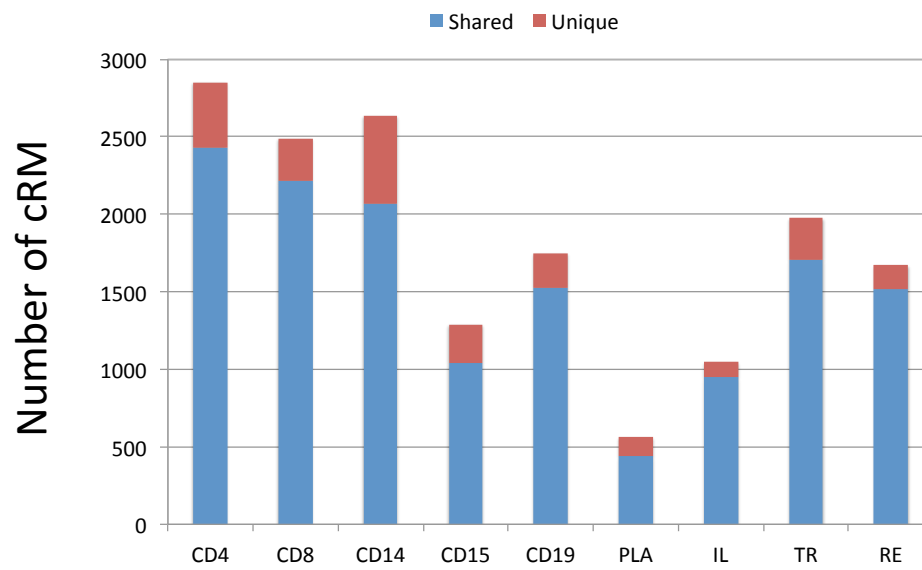
244 Graphical representation (using Cytoscape¹) of 269 cis acting regulatory modules
245 (cRM) including at least three genes (see Suppl. Table 2). Every node corresponds
246 to a cis-eQTL involving a specific gene-tissue combination. Edges connect pairs of
247 cis-eQTL for which $|\vartheta| \geq 0.6$.

248 1. Shannon, P. et al. Genome Res. 13, 2498-2504 (2003).

249

250 **3. Supplementary Figure 3**

251



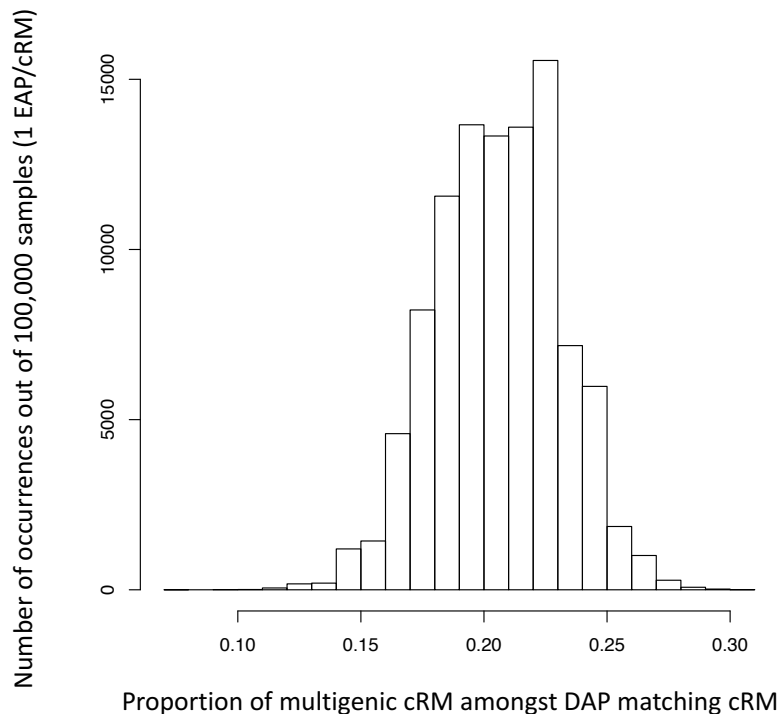
252

253 Number of cRM detected in each cell type. Blue: shared cRM (i.e. also detected
254 in at least one other cell type). Red: Unique (i.e. detected only in that cell type).

255

256 4. Supplementary Figure 4

257

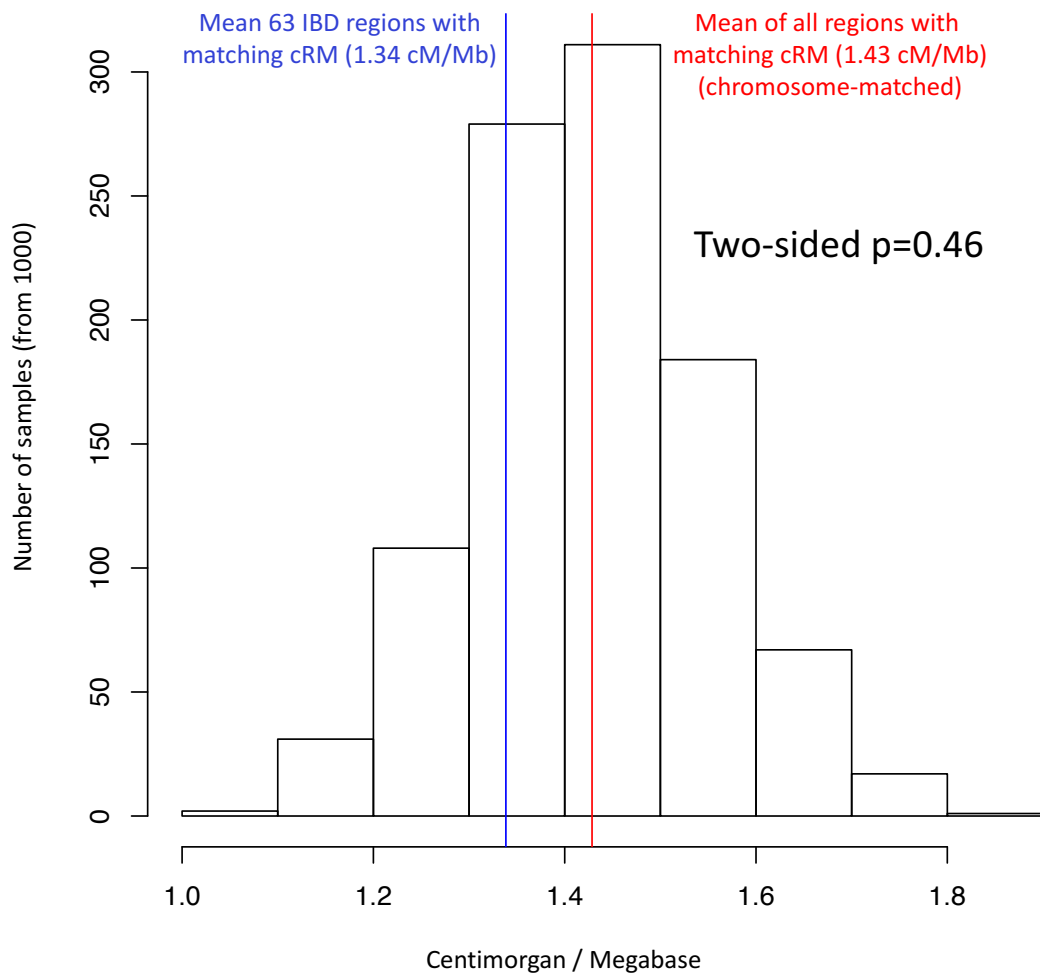


258

259 Across the entire genome, the proportion of multigenic cRM was shown to be
260 0.10 (see also main text, figure 1B). Amongst DAP matching cRM (mapping to
261 63 of 200 studied IBD risk loci; main text Table 1) this proportion was shown
262 to be 0.33, hence a highly significant enrichment. To ensure that this
263 enrichment was not only due to the fact that matching between DAP and EAP
264 was de facto tested multiple times for multigenic cRM and only once for other
265 cRM, we only tested one randomly sampled EAP per cRM (whether monogenic
266 or multigenic). This was repeated 100,000 times and yielded the distribution
267 of the proportion of multigenic cRM amongst DAP matching cRM shown above.
268 The average was 0.22, and we never observed values ≤ 0.11 , i.e. the genome-
269 wide average.

270

271 5. Supplementary Figure 5

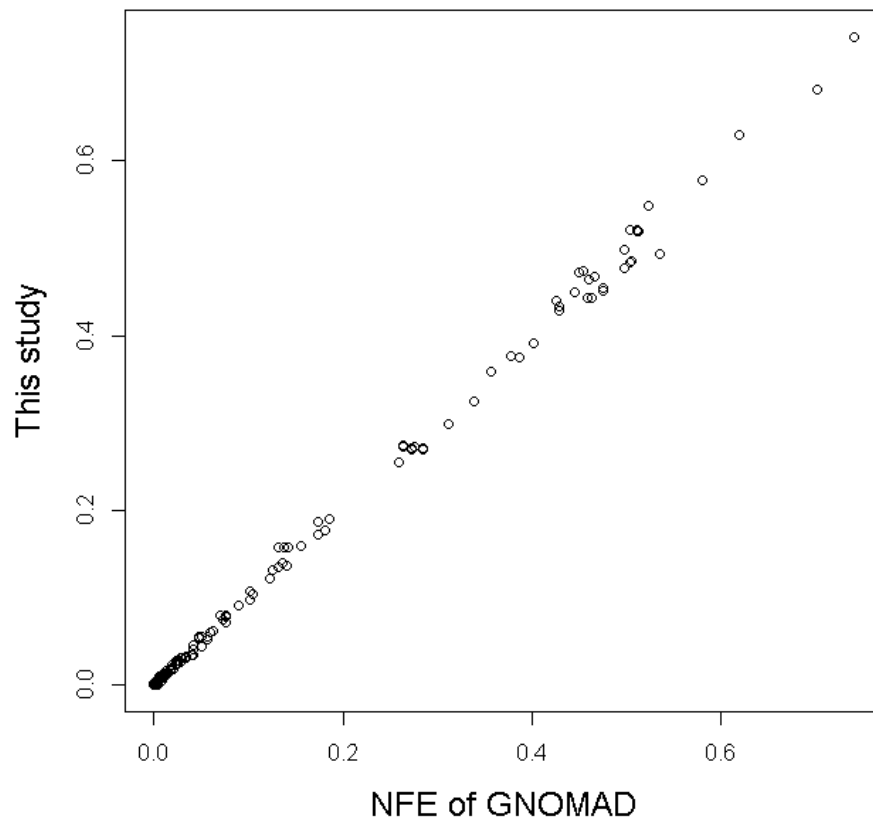


272

273

274 The 63 IBD risk loci with matching cRM are 2- to 3-fold enriched in multigenic
 275 cRM ($p \leq 10^{-5}$). This could be due to the fact that the LD is higher in IBD
 276 regions than in the rest of the genome. To test this, we downloaded LD-based
 277 recombination maps of the human genome from
 278 <https://github.com/joepickrell/1000-genomes-genetic-maps>. The average
 279 recombination rate across the human genome was 1.23 centimorgan per
 280 megabase (cM/Mb). The average recombination rate for the 63 IBD risk loci
 281 with matching cRM was 1.34 cM/Mb, i.e. less LD than in the rest of the genome.
 282 Regions encompassing eQTL (and hence cRM) may differ from the rest of the
 283 genome with regards to LD. Thus, we further sampled 1,000 sets of 63 loci
 284 centered on cRM (from our list of 9,720) that were matched for size and
 285 chromosomal location with the 63 cRM-matching IBD risk loci. The mean
 286 recombination rate for the cRM-centered genome was 1.43 cM/Mb. The figure
 287 shows the frequency distribution of the corresponding mean cRM/Mb per set
 288 (black), the mean of means of the 1,000 sets of 63 randomly drawn loci (red),
 289 and the mean of the 3 IBD risk loci (blue). The mean of the 63 IBD risk loci did
 290 not differ significantly from the rest of the cRM centered portion of the genome
 291 (two-tailed p-value: 0.46).

292 **6. Supplemental Figure 6**
293

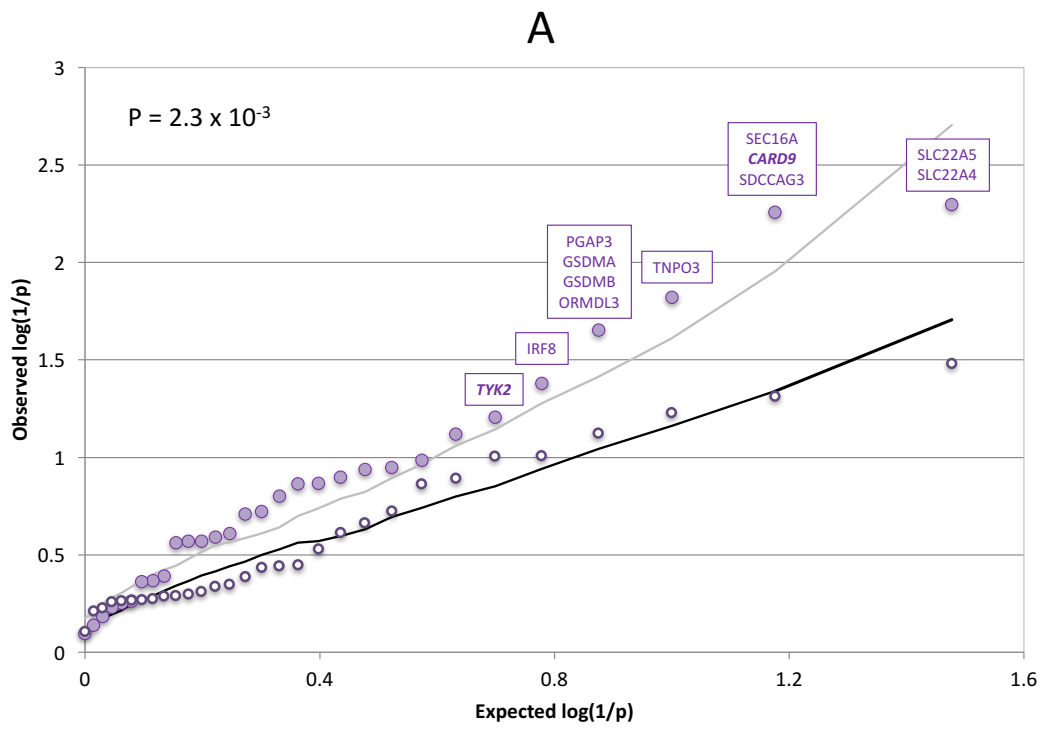


294

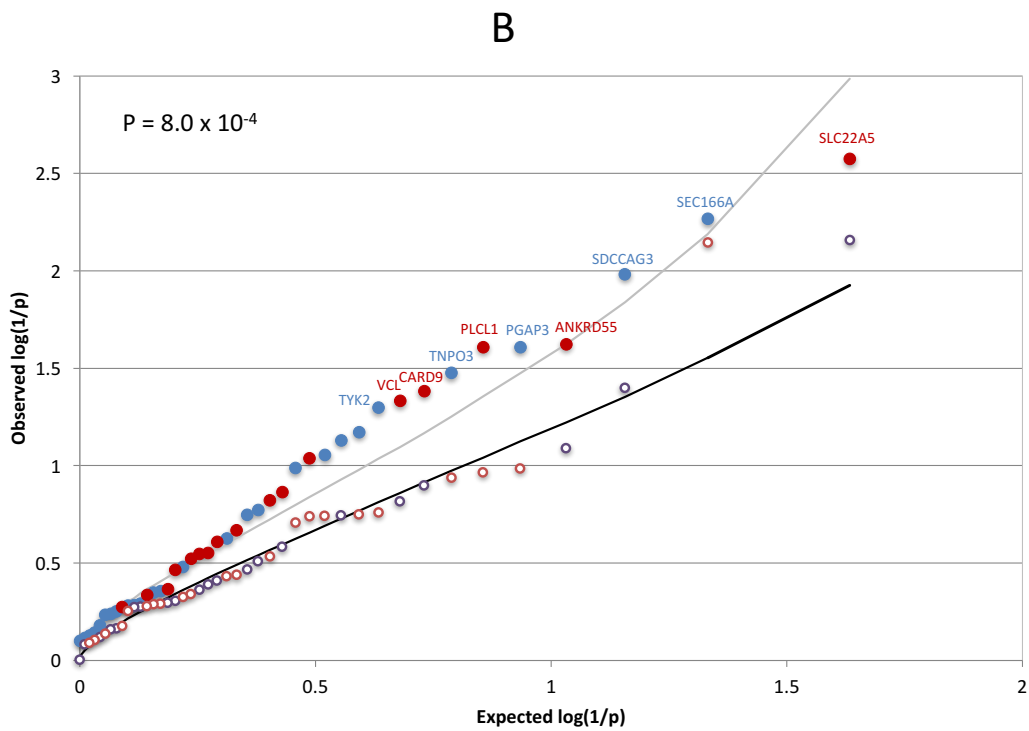
295 Comparison of the alternative allele frequency for 1,781 variants observed in this
296 study and in 55,860 non-Finnish European samples from the GNOMAD study.

297

298 7. Supplementary Figure 7
 299

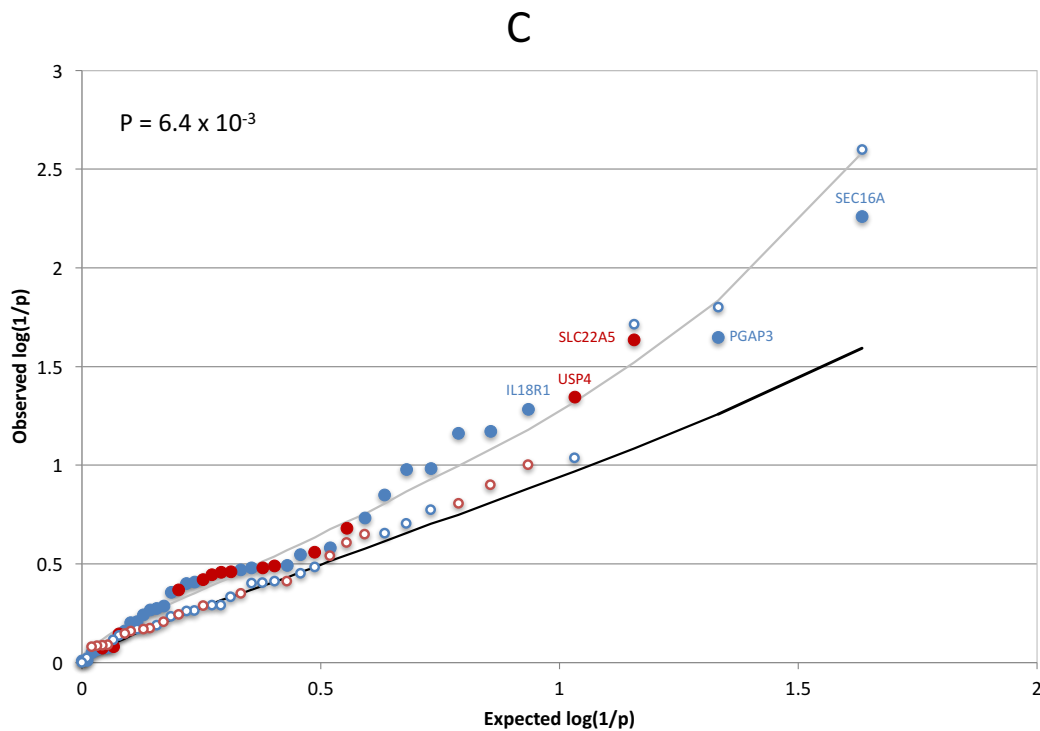


300



301

302

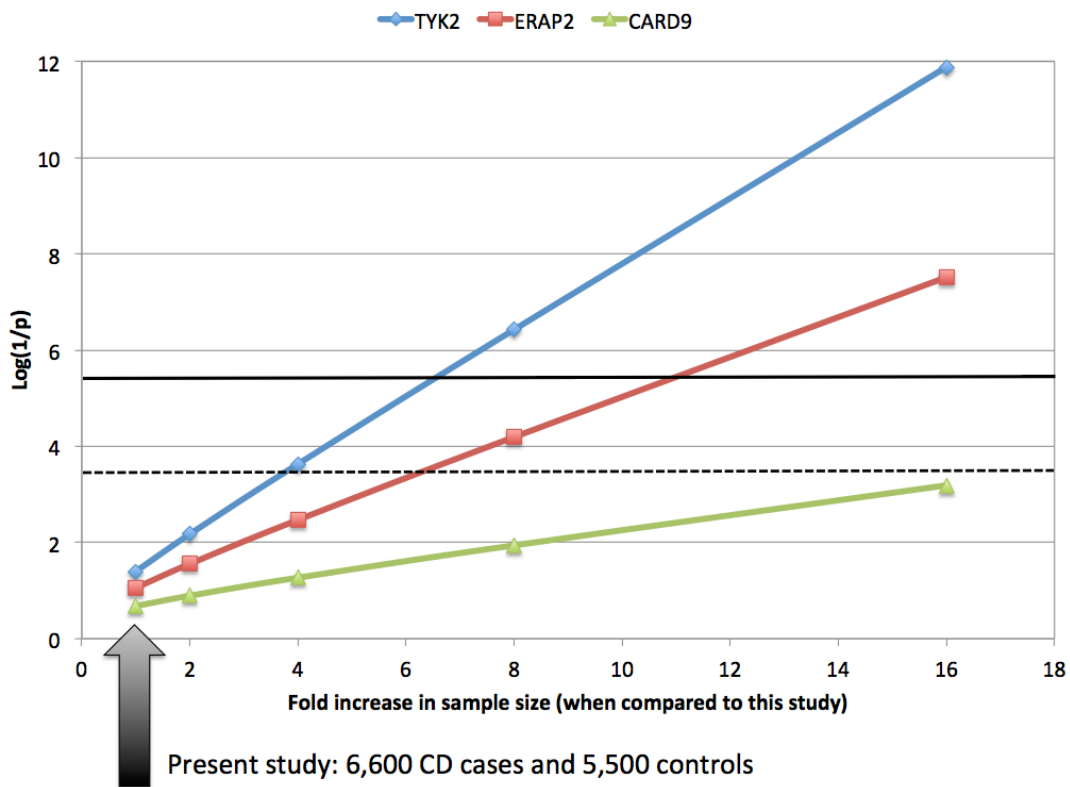


303

304 QQ-plot for the module-based burden test (A), disease plus age-of-onset-based
 305 burden test (B), and disease plus familiarity-based burden test (C). Ranked
 306 $\log(1/p)$ values obtained when considering LoF and damaging variants (full
 307 circles), or synonymous variants (empty circles). The circles are labeled in blue
 308 when the best p-value for that gene is obtained with CAST, in red when the best
 309 p-value is obtained with SKAT, or in purple for the module-based test (as some
 310 genes in the module may have their best p-value with CAST and other with
 311 SKAT). The black line corresponds to the median $\log(1/p)$ value obtained (for
 312 the corresponding rank) using the same approach on permuted data (LoF and
 313 damaging variants). The grey line marks the upper limit of the 95% confidence
 314 band. The name of the genes/modules exceeding the nominal p-value of 0.05
 315 are given. The inset p-values correspond to the significance of the upwards
 316 shift in $\log(1/p)$ values estimated by permutation.

317

318 **8. Supplementary Figure 8**
 319



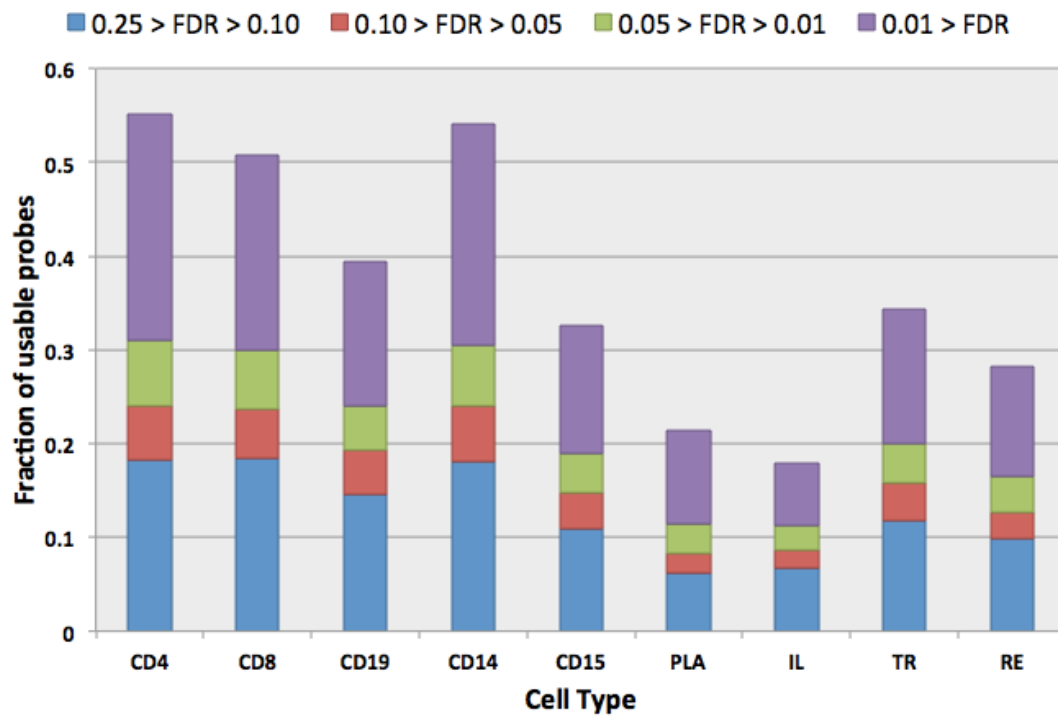
320

321 Effect of increasing sample size on the $\log(1/p)$ values of a one-sided burden
 322 test assuming that the effects observed for *TYK2* (blue), *ERAP2* (red) and
 323 *CARD9* (green) observed in this study are real unbiased. The dotted horizontal
 324 black line corresponds to an hypothetical experiment-wide significance
 325 threshold assuming the realization of 200 independent tests (targeting for
 326 instance 100-200 genes selected on the basis of coincident DAP-EAP patterns).
 327 The plain horizontal black line corresponds to an hypothetical genome-wide
 328 significance threshold assuming the realization of 20,000 independent tests
 329 (targeting all genes). It can be seen that an at least 4-fold increase in sample
 330 size is needed to achieve significance in the first scenario and at least 7-fold
 331 increase in the second scenario.

332

333 9. Supplementary Figure 9

334



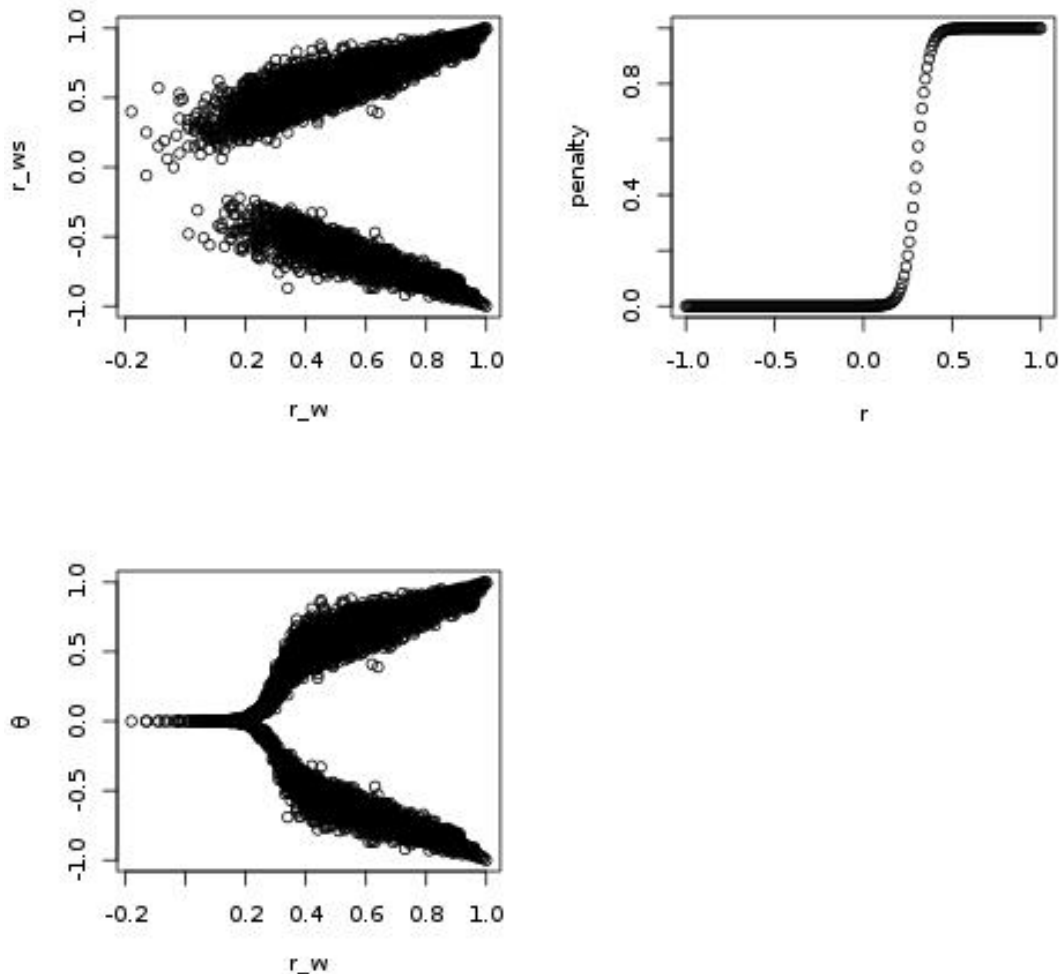
335

336 Proportion of usable probes with cis-eQTL at various levels of FDR in the nine
 337 analyzed cell types.

338

339 10. Supplementary Figure 10

340



341

342 Graphical illustration of the relationship between r_w , r_{ws} and ϑ . The penalty
 343 function applied to r_{ws} to generate ϑ , corresponds to $\frac{1}{1+e^{-k(r_w-T)}}$. The graph is
 344 shown for $k = 30$ and $T=0.3$, the values used in this study.

345 The point here is that if two association patterns are “similar” (driven by the same
 346 variants), the correlation (r_w in Suppl. Methods) between $-\log(1/p)$ values is
 347 expected to be positive. If two association patterns are different (driven by
 348 distinct variants) they may generate strong negative correlations (r_w). The first
 349 part of the method aims at weeding out such instances (negative r_w). One way to
 350 do this is to choose a simple threshold value for r_w . We herein propose an
 351 approach that offers more flexibility: it generates a penalty that increases when
 352 the correlation decreases with an adaptable rate. As shown in Suppl. Fig. 8, the
 353 values of $k=30$ and $T=0.3$ essentially correspond to a threshold value of 0.3. As
 354 can also be seen from Suppl. Fig. 8, there is (as expected) a strong linear
 355 relationship with slope 1 between r_w and $|r_{ws}|$ (and hence between r_w and $|\vartheta|$ for
 356 pairs with $r_w > 0.3$). Because we subsequently use a threshold value $|\vartheta| \geq 0.6$, the
 357 choice T has very little impact on the outcome unless one approaches 0.6.

358

359

360

11. Supplementary Figure 11

361
362

363

364

365

366

367

368

369

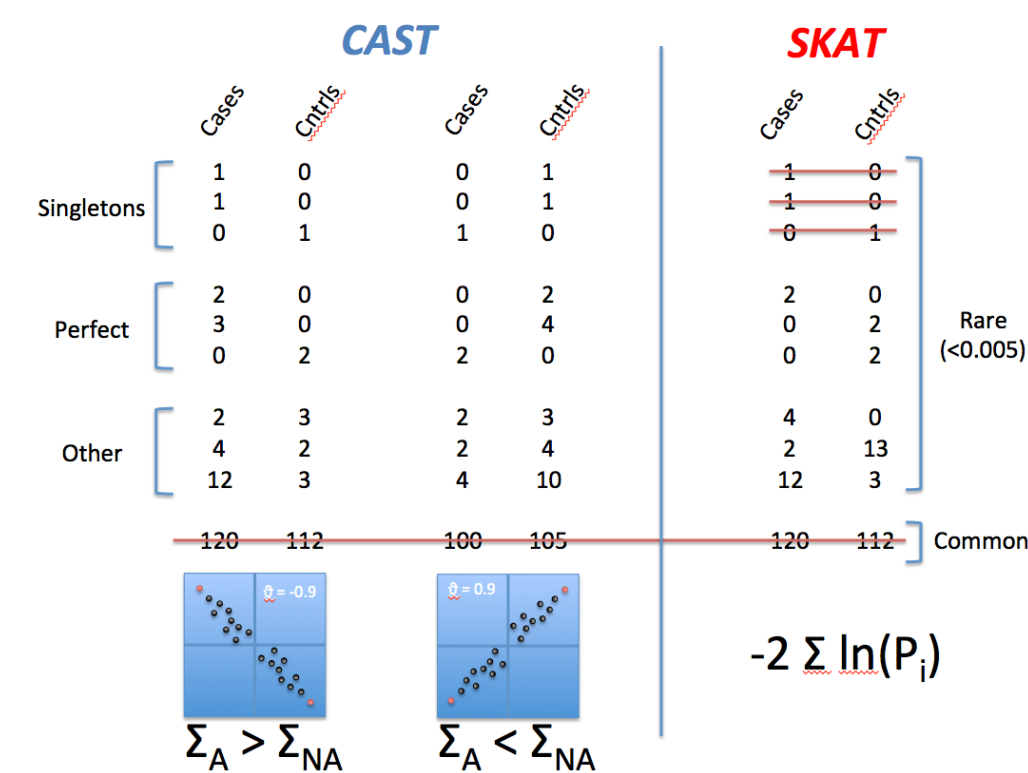
370

371

372

373

374



375

Schematic representation of the key features of the implemented “burden test”.

376

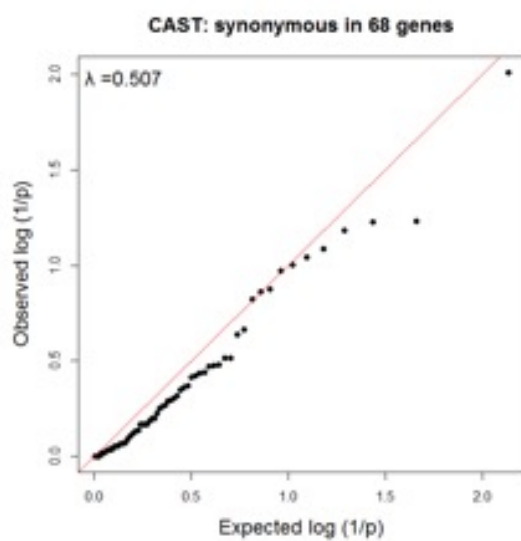
The analysis is restricted to rare variants with MAF < 0.005 to ensure that the new signal is independent of the one that lead to the identification of the corresponding risk loci by GWAS (based on common and low frequency variants). Variants can be sorted in (i) singletons (i.e. observed only ones in the analyzed samples), (ii) perfect (i.e. observed more than ones in the sample but perfectly associated with disease status), and (iii) other (i.e. observed more than ones in the sample in both cases and controls).

383

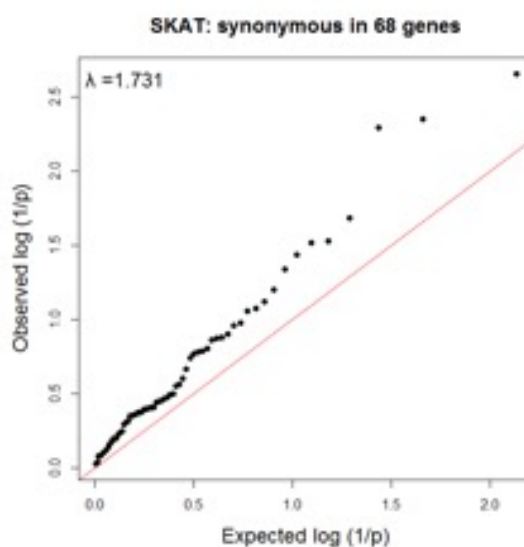
We test two hypotheses. The first assumes that disruptive variants are either enriched in cases or in controls as a function of the sign of the correlation between DAP and EAP (if decreased expression is associated with increased risk, disruptive “risk” variants are expected to be enriched in cases; if increased expression is associated with increased risk, disruptive “protective” variants are expected to be enriched in controls). The test is implemented with CAST and in essence performs a one-sided test of independence (what is the probability to observe the excess of disruptive variants in cases (respectively controls) by chance alone?). The second hypothesis tests whether the distribution of the variants in cases and controls is characterized by too many variants that tend to be overrepresented either in cases or in controls. Thus, this hypothesis allows some disruptive variants to increase risk and others to be protective. This hypothesis does not use information from singletons. Testing this hypothesis is implemented with SKAT. It can be seen in simplified form as combining the p-values (from a test of independence) across variants (without considering the sign of the effect) using for instance Fisher’s method.

399

400 12. Supplementary Figure 12



401



402

403

404 Distribution of permutation-based $-\log(p)$ values obtained for 68 analyzed genes
405 with synonymous variants using CAST (A), and SKAT (B), indicating that CAST is
406 conservative ($\lambda_{GC} = 0.507$), while SKAT is too permissive ($\lambda_{GC} = 1.73$). The 68 genes
407 correspond to the 47 genes reported in this study, plus 21 genes sequenced in the
408 same cohort as part of another study.

409

410

- 411 1. Dinarello, C. A., Novick, D., Kim, S. & Kaplanski, G. Interleukin-18 and IL-18
412 binding protein. *Frontiers in Immunology* **4**, (2013).
- 413 2. Elinav, E. *et al.* NLRP6 inflammasome regulates colonic microbial ecology and
414 risk for colitis. *Cell* **145**, 745–757 (2011).
- 415 3. Salcedo, R. *et al.* MyD88-mediated signaling prevents development of
416 adenocarcinomas of the colon: role of interleukin 18. *J. Exp. Med.* **207**, 1625–36
417 (2010).
- 418 4. Nowarski, R. *et al.* Epithelial IL-18 Equilibrium Controls Barrier Function in
419 Colitis. *Cell* **163**, 1444–1456 (2015).
- 420 5. Dinarello, C. A. Interleukin-18 and the Pathogenesis of Inflammatory Diseases.
421 *Semin. Nephrol.* **27**, 98–114 (2007).
- 422 6. Heinrich, P. C. *et al.* Principles of interleukin (IL)-6-type cytokine signalling and
423 its regulation. *Biochem. J.* **374**, 1–20 (2003).
- 424 7. Mitsuyama, K. *et al.* Therapeutic strategies for targeting the IL-6/STAT3
425 cytokine signaling pathway in inflammatory bowel disease. *Anticancer Res.* **27**,
426 3749–56
- 427 8. Atreya, R. *et al.* Blockade of interleukin 6 trans signaling suppresses T-cell
428 resistance against apoptosis in chronic intestinal inflammation: evidence in
429 Crohn disease and experimental colitis in vivo. *Nat. Med.* **6**, 583–588 (2000).
- 430 9. Ito, H. *et al.* A Pilot Randomized Trial of a Human Anti-Interleukin-6 Receptor
431 Monoclonal Antibody in Active Crohn's Disease. *Gastroenterology* **126**, 989–
432 996 (2004).
- 433 10. Jones, S. A., Scheller, J. & Rose-John, S. Therapeutic strategies for the clinical
434 blockade of IL-6/gp130 signaling. *J. Clin. Invest.* **121**, 3375–3383 (2011).
- 435 11. Lesourne, R. *et al.* Themis, a T cell-specific protein important for late thymocyte
436 development. *Nat. Immunol.* **10**, 840–7 (2009).
- 437 12. Fu, G. *et al.* Themis controls thymocyte selection through regulation of T cell
438 antigen receptor-mediated signaling. *Nat. Immunol.* **10**, 848–56 (2009).
- 439 13. Neurath, M. F. Cytokines in inflammatory bowel disease. *Nat. Rev. Immunol.* **14**,

- 440 329–342 (2014).
- 441 14. Palumbo, R. *et al.* APEH Inhibition Affects Osteosarcoma Cell Viability via
442 Downregulation of the Proteasome. *Int. J. Mol. Sci.* **17**, 1614 (2016).
- 443 15. Cleynen, I. *et al.* Genetic and microbial factors modulating the ubiquitin
444 proteasome system in inflammatory bowel disease. *Gut* **63**, 1265–74 (2014).
- 445 16. Li, J., Mahajan, A. & Tsai, M.-D. Ankyrin Repeat: A Unique Motif Mediating
446 Protein–Protein Interactions †. *Biochemistry* **45**, 15168–15178 (2006).
- 447 17. International Multiple Sclerosis Genetics Consortium (IMSGC) *et al.* Analysis
448 of immune-related loci identifies 48 new susceptibility variants for multiple
449 sclerosis. *Nat. Genet.* **45**, 1353–60 (2013).
- 450 18. Viatte, S. *et al.* Genetic markers of rheumatoid arthritis susceptibility in anti-
451 citrullinated peptide antibody negative patients. *Ann. Rheum. Dis.* **71**, 1984–90
452 (2012).
- 453 19. Wiley, S. E., Murphy, A. N., Ross, S. A., van der Geer, P. & Dixon, J. E.
454 MitoNEET is an iron-containing outer mitochondrial membrane protein that
455 regulates oxidative capacity. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 5318–5323
456 (2007).
- 457 20. Geldenhuys, W. J., Leeper, T. C. & Carroll, R. T. mitoNEET as a novel drug
458 target for mitochondrial dysfunction. *Drug Discovery Today* **19**, 1601–1606
459 (2014).
- 460 21. Goldberg, N. D. Iron deficiency anemia in patients with inflammatory bowel
461 disease. *Clinical and Experimental Gastroenterology* **6**, 61–70 (2013).
- 462 22. Kojima, S., Sher-Chen, E. L. & Green, C. B. Circadian control of mRNA
463 polyadenylation dynamics regulates rhythmic protein expression. *Genes Dev.* **26**,
464 2724–2736 (2012).
- 465 23. Maillo, C. *et al.* Circadian- and UPR-dependent control of CPEB4 mediates a
466 translational response to counteract hepatic steatosis under ER stress. *Nat. Cell*
467 *Biol.* **19**, 94–105 (2017).
- 468 24. Watabe-Uchida, M., John, K. A., Janas, J. A., Newey, S. E. & Van Aelst, L. The
469 Rac Activator DOCK7 Regulates Neuronal Polarity through Local
470 Phosphorylation of Stathmin/Op18. *Neuron* **51**, 727–739 (2006).

- 471 25. Blasius, A. L. *et al.* Mice with mutations of Dock7 have generalized
472 hypopigmentation and white-spotting but show normal neurological function.
473 *Proc. Natl. Acad. Sci.* **106**, 2706–2711 (2009).
- 474 26. Zhang, Z. *et al.* Functional interaction of ERAP2 and HLA-B27 activates the
475 unfolded protein response. *Arthritis Rheumatol.* (2016). doi:10.1002/art.40033
- 476 27. Andrés, A. M. *et al.* Balancing selection maintains a form of ERAP2 that
477 undergoes nonsense-mediated decay and affects antigen presentation. *PLoS*
478 *Genet.* **6**, 1–13 (2010).
- 479 28. Saveanu, L. *et al.* Concerted peptide trimming by human ERAP1 and ERAP2
480 aminopeptidase complexes in the endoplasmic reticulum. *Nat. Immunol.* **6**, 689–
481 97 (2005).
- 482 29. Franke, A. *et al.* Genome-wide meta-analysis increases to 71 the number of
483 confirmed Crohn’s disease susceptibility loci. *Nat. Genet.* **42**, 1118–25 (2010).
- 484 30. Damjanovich, L., Volkó, J., Forgács, A., Hohenberger, W. & Bene, L. Crohn’s
485 disease alters MHC-rafts in CD4 + T-cells. *Cytom. Part A* **81 A**, 149–164 (2012).
- 486 31. Dodane, V. & Kachar, B. Identification of isoforms of G proteins and PKC that
487 colocalize with tight junctions. *J. Membr. Biol.* **149**, 199–209 (1996).
- 488 32. Meyer, T. N., Schwesinger, C. & Denker, B. M. Zonula occludens-1 is a
489 scaffolding protein for signaling molecules: G α 12 directly binds to the Src
490 homology 3 domain and regulates paracellular permeability in epithelial cells. *J.*
491 *Biol. Chem.* **277**, 24855–24858 (2002).
- 492 33. Sabath, E. *et al.* G α 12 regulates protein interactions within the MDCK cell tight
493 junction and inhibits tight-junction assembly. *J. Cell Sci.* **121**, 814–824 (2008).
- 494 34. Fukuhara, S., Chikumi, H. & Gutkind, J. S. RGS-containing RhoGEFs: the
495 missing link between transforming G proteins and Rho? *Oncogene* **20**, 1661–
496 1668 (2001).
- 497 35. Landy, J. *et al.* Tight junctions in inflammatory bowel diseases and inflammatory
498 bowel disease associated colorectal cancer. *World Journal of Gastroenterology*
499 **22**, 3117–3126 (2016).
- 500 36. Herroeder, S. *et al.* Guanine Nucleotide-Binding Proteins of the G12 Family
501 Shape Immune Functions by Controlling CD4+ T Cell Adhesiveness and

- 502 Motility. *Immunity* **30**, 708–720 (2009).
- 503 37. Pavlick, K. P. *et al.* Role of reactive metabolites of oxygen and nitrogen in
504 inflammatory bowel disease 1,2 1This article is part of a series of reviews on
505 ‘Reactive Oxygen and Nitrogen in Inflammation.’ The full list of papers may be
506 found on the homepage of the journal. 2Gue. *Free Radic. Biol. Med.* **33**, 311–
507 322 (2002).
- 508 38. Chu, F. F. *et al.* Bacteria-Induced Intestinal Cancer in Mice with Disrupted Gpx1
509 and Gpx2 Genes. *Cancer Res.* **64**, 962–968 (2004).
- 510 39. Maeda, S. *et al.* cDNA microarray analysis of Helicobacter pylori-mediated
511 alteration of gene expression in gastric cancer cells. *Biochem. Biophys. Res.*
512 *Commun.* **284**, 443–9 (2001).
- 513 40. Hockenbery, D. M., Oltvai, Z. N., Yin, X. M., Milliman, C. L. & Korsmeyer, S.
514 J. Bcl-2 functions in an antioxidant pathway to prevent apoptosis. *Cell* **75**, 241–
515 51 (1993).
- 516 41. Saeki, N. *et al.* GASDERMIN, suppressed frequently in gastric cancer, is a target
517 of LMO1 in TGF-beta-dependent apoptotic signalling. *Oncogene* **26**, 6488–98
518 (2007).
- 519 42. Saeki, N. *et al.* Distinctive expression and function of four GSDM family genes
520 (GSDMA-D) in normal and malignant upper gastrointestinal epithelium. *Genes,*
521 *Chromosom. Cancer* **48**, 261–271 (2009).
- 522 43. Pal, L. R. & Moulton, J. Genetic basis of common human disease: Insight into the
523 role of missense SNPs from genome-wide association studies. *J. Mol. Biol.* **427**,
524 2271–2289 (2015).
- 525 44. Jostins, L. *et al.* Host-microbe interactions have shaped the genetic architecture
526 of inflammatory bowel disease. *Nature* **491**, 119–24 (2012).
- 527 45. Breslow, D. K. *et al.* Orm family proteins mediate sphingolipid homeostasis.
528 *Nature* **463**, 1048–1053 (2010).
- 529 46. Ha, S. G. *et al.* ORMDL3 promotes eosinophil trafficking and activation via
530 regulation of integrins and CD48. *Nat. Commun.* **4**, 2479 (2013).
- 531 47. Schmiedel, B. J. *et al.* 17q21 asthma-risk variants switch CTCF binding and
532 regulate IL-2 production by T cells. *Nat. Commun.* **7**, 13426 (2016).

- 533 48. Dang, J. *et al.* ORMDL3 Facilitates the Survival of Splenic B Cells via an
534 ATF6 α -Endoplasmic Reticulum Stress-Beclin1 Autophagy Regulatory Pathway.
535 *J. Immunol.* (2017). doi:10.4049/jimmunol.1602124
- 536 49. Zhai, W. H., Song, C. Y., Huang, Z. G. & Sha, H. Correlation between the
537 genetic polymorphism of ORMDL3 gene and asthma risk: a meta-analysis.
538 *Genet. Mol. Res.* **14**, 7101–12 (2015).
- 539 50. Saleh, N. M. *et al.* Genetic association analyses of atopic illness and
540 proinflammatory cytokine genes with type 1 diabetes. *Diabetes. Metab. Res. Rev.*
541 **27**, 838–43 (2011).
- 542 51. Ma, X. *et al.* ORMDL3 contributes to the risk of atherosclerosis in Chinese Han
543 population and mediates oxidized low-density lipoprotein-induced autophagy in
544 endothelial cells. *Sci. Rep.* **5**, 17194 (2015).
- 545 52. Laukens, D. *et al.* Evidence for significant overlap between common risk
546 variants for Crohn's disease and ankylosing spondylitis. *PLoS One* **5**, e13795
547 (2010).
- 548 53. McGovern, D. P. B. *et al.* Genome-wide association identifies multiple
549 ulcerative colitis susceptibility loci. *Nat. Genet.* **42**, 332–7 (2010).
- 550 54. Asazuma, N. *et al.* Interaction of linker for activation of T cells with multiple
551 adapter proteins in platelets activated by the glycoprotein VI-selective ligand,
552 convulxin. *J. Biol. Chem.* **275**, 33427–33434 (2000).
- 553 55. Togni, M. *et al.* Regulation of In Vitro and In Vivo Immune Functions by the
554 Cytosolic Adaptor Protein SKAP-HOM. *Mol. Cell. Biol.* **25**, 8052–8063 (2005).
- 555 56. Alenghat, F. J. *et al.* Macrophages require Skap2 and Sirp α for integrin-
556 stimulated cytoskeletal rearrangement. *J. Cell Sci.* **125**, 5535–5545 (2012).
- 557 57. Königsberger, S. *et al.* HPK1 associates with SKAP-HOM to negatively regulate
558 Rap1-mediated B-lymphocyte adhesion. *PLoS One* **5**, 1–9 (2010).
- 559 58. Tanaka, M. *et al.* SKAP2 Promotes Podosome Formation to Facilitate Tumor-
560 Associated Macrophage Infiltration and Metastatic Progression. *Cancer Res.* **76**,
561 358–369 (2016).
- 562 59. Sartor, R. B. Mechanisms of disease: pathogenesis of Crohn's disease and
563 ulcerative colitis. *Nat. Clin. Pract. Gastroenterol. Hepatol.* **3**, 390–407 (2006).

- 564 60. Ghosh, S. *et al.* Natalizumab for active Crohn's disease. *N Engl J Med* **348**, 24–
565 32. (2003).
- 566 61. Eldridge MJ, Sanchez-Garrido J, Hoben GF, Goddard PJ, S. A. The Atypical
567 Ubiquitin E2 Conjugase UBE2L3 Is an Indirect Caspase-1 Target and Controls
568 IL-1 β Secretion by Inflammasomes. *Cell Rep.* **18**, 1285–1297 (2017).
- 569 62. Alpi, A. F., Chaugule, V. & Walden, H. Mechanism and disease association of
570 E2-conjugating enzymes: lessons from UBE2T and UBE2L3. *Biochem. J.* **473**,
571 3401–3419 (2016).
- 572 63. Lamkanfi, M., Walle, L. Vande & Kanneganti, T. D. Deregulated inflammasome
573 signaling in disease. *Immunological Reviews* **243**, 163–173 (2011).
- 574 64. Shuai, K. & Liu, B. Regulation of gene-activation pathways by PIAS proteins in
575 the immune system. *Nat. Rev. Immunol.* **5**, 593–605 (2005).
- 576 65. Sharma, M. *et al.* hZimp10 is an androgen receptor co-activator and forms a
577 complex with SUMO-1 at replication foci. *EMBO J.* **22**, 6101–6114 (2003).
- 578 66. Imielinski, M. *et al.* Common variants at five new loci associated with early-
579 onset inflammatory bowel disease. *Nat. Genet.* **41**, 1335–40 (2009).
- 580 67. Rakowski, L. A. *et al.* Convergence of the ZMIZ1 and NOTCH1 Pathways at C-
581 MYC in Acute T Lymphoblastic Leukemias. *Cancer Res.* **73**, 930–941 (2013).
- 582 68. Pinnell, N. *et al.* The PIAS-like Coactivator Zmiz1 Is a Direct and Selective
583 Cofactor of Notch1 in T Cell Development and Leukemia. *Immunity* **43**, 870–
584 883 (2015).
- 585 69. Li, X., Thyssen, G., Beliakoff, J. & Sun, Z. The Novel PIAS-like Protein
586 hZimp10 Enhances Smad Transcriptional Activity. *J. Biol. Chem.* **281**, 23748–
587 23756 (2006).
- 588 70. Jongstra-Bilen, J. & Jongstra, J. Leukocyte-specific protein 1 (LSP1): a regulator
589 of leukocyte emigration in inflammation. *Immunol. Res.* **35**, 65–74 (2006).
- 590 71. Wang, C. *et al.* Modulation of Mac-1 (CD11b/CD18)-Mediated Adhesion by the
591 Leukocyte-Specific Protein 1 Is Key to Its Role in Neutrophil Polarization and
592 Chemotaxis. *J. Immunol.* **169**, 415–423 (2002).
- 593 72. Wang, J. *et al.* Accelerated wound healing in leukocyte-specific, protein 1-
594 deficient mouse is associated with increased infiltration of leukocytes and

- 595 fibrocytes. *J. Leukoc. Biol.* **82**, 1554–63 (2007).
- 596 73. Hwang, S.-H. *et al.* Leukocyte-specific protein 1 regulates T-cell migration in
597 rheumatoid arthritis. *Proc. Natl. Acad. Sci.* **112**, E6535–E6543 (2015).
- 598 74. Gremel, G. *et al.* The human gastrointestinal tract-specific transcriptome and
599 proteome as defined by RNA sequencing and antibody-based profiling. *J.*
600 *Gastroenterol.* **50**, 46–57 (2015).
- 601 75. Hular, I. *et al.* Enrichment of inflammatory bowel disease and colorectal cancer
602 risk variants in colon expression quantitative trait loci. *BMC Genomics* **16**, 138
603 (2015).
- 604 76. Nakajima, T. TIP27: a novel repressor of the nuclear orphan receptor TAK1/TR4.
605 *Nucleic Acids Res.* **32**, 4194–4204 (2004).
- 606 77. Collins, L. L. *et al.* Growth retardation and abnormal maternal behavior in mice
607 lacking testicular orphan nuclear receptor 4. *Proc. Natl. Acad. Sci. U. S. A.* **101**,
608 15058–63 (2004).
- 609 78. Johansson, ??sa *et al.* Common variants in the JAZF1 gene associated with
610 height identified by linkage and genome-wide association analysis. *Hum. Mol.*
611 *Genet.* **18**, 373–380 (2009).
- 612 79. Cooper, J. D. *et al.* Meta-analysis of genome-wide association study data
613 identifies additional type 1 diabetes risk loci. *Nat. Genet.* **40**, 1399–401 (2008).
- 614 80. Thomas, G. *et al.* Multiple loci identified in a genome-wide association study of
615 prostate cancer. *Nat. Genet.* **40**, 310–5 (2008).
- 616 81. Koontz, J. I. *et al.* Frequent fusion of the JAZF1 and JJAZ1 genes in endometrial
617 stromal tumors. *Proc. Natl. Acad. Sci.* **98**, 6348–6353 (2001).
- 618 82. Martin, J. E. *et al.* A systemic sclerosis and systemic lupus erythematosus pan-
619 meta-GWAS reveals new shared susceptibility loci. *Hum. Mol. Genet.* **22**, 4021–
620 4029 (2013).
- 621 83. Bruni, F., Gramegna, P., Oliveira, J. M. A., Lightowers, R. N. & Chrzanowska-
622 Lightowers, Z. M. A. REXO2 is an oligoribonuclease active in human
623 mitochondria. *PLoS One* **8**, e64670 (2013).
- 624 84. Matondo, A. & Kim, S. S. Targeted-mitochondria antioxidants therapeutic
625 implications in inflammatory bowel disease. *J. Drug Target.* 1–8 (2017).

- 626 doi:10.1080/1061186X.2017.1339196
- 627 85. Acquati, F. *et al.* Microenvironmental control of malignancy exerted by
628 RNASET2, a widely conserved extracellular RNase. *Proc. Natl. Acad. Sci. U. S.*
629 *A.* **108**, 1104–9 (2011).
- 630 86. Acquati, F. *et al.* Loss of function of Ribonuclease T2, an ancient and
631 phylogenetically conserved RNase, plays a crucial role in ovarian tumorigenesis.
632 *Proc. Natl. Acad. Sci.* 1222079110- (2013). doi:10.1073/pnas.1222079110
- 633 87. Gabrielsen, I. S. M. *et al.* Genetic risk variants for autoimmune diseases that
634 influence gene expression in thymus. *Hum. Mol. Genet.* ddw152 (2016).
635 doi:10.1093/hmg/ddw152
- 636 88. Chu, X. *et al.* A genome-wide association study identifies two new risk loci for
637 Graves' disease. *Nat. Genet.* **43**, 897–901 (2011).
- 638 89. Caputa, G. *et al.* RNASET2 is required for ROS propagation during oxidative
639 stress-mediated cell death. *Cell Death Differ.* **23**, 347–57 (2016).
- 640 90. Wang, Q. *et al.* Stress-induced RNASET2 overexpression mediates melanocyte
641 apoptosis via the TRAF2 pathway in vitro. *Cell Death Dis.* **5**, e1022 (2014).
- 642 91. Moret-Tatay, I. *et al.* Possible Biomarkers in Blood for Crohn's Disease:
643 Oxidative Stress and MicroRNAs—Current Evidences and Further Aspects to
644 Unravel. *Oxid. Med. Cell. Longev.* **2016**, 1–9 (2016).
- 645