

ARTICLE

DOI: 10.1038/s41467-018-04365-8

OPEN

# IBD risk loci are enriched in multigenic regulatory modules encompassing putative causative genes

Yukihide Momozawa, Julia Dmitrieva et al.<sup>#</sup>

GWAS have identified >200 risk loci for Inflammatory Bowel Disease (IBD). The majority of disease associations are known to be driven by regulatory variants. To identify the putative causative genes that are perturbed by these variants, we generate a large transcriptome data set (nine disease-relevant cell types) and identify 23,650 *cis*-eQTL. We show that these are determined by ~9720 regulatory modules, of which ~3000 operate in multiple tissues and ~970 on multiple genes. We identify regulatory modules that drive the disease association for 63 of the 200 risk loci, and show that these are enriched in multigenic modules. Based on these analyses, we resequence 45 of the corresponding 100 candidate genes in 6600 Crohn disease (CD) cases and 5500 controls, and show with burden tests that they include likely causative genes. Our analyses indicate that  $\geq 10$ -fold larger sample sizes will be required to demonstrate the causality of individual genes using this approach.

---

Correspondence and requests for materials should be addressed to M.G. (email: [michel.georges@ulg.ac.be](mailto:michel.georges@ulg.ac.be))

<sup>#</sup>A full list of authors and their affiliations appears at the end of the paper.

Genome Wide Association Studies (GWAS) scan the entire genome for statistical associations between common variants and disease status in large case–control cohorts. GWAS have identified tens to hundreds of risk loci for nearly all studied common complex diseases of human<sup>1</sup>. The study of Inflammatory Bowel Disease (IBD) has been particularly successful, with more than 200 confirmed risk loci reported to date<sup>2,3</sup>. As a result of the linkage disequilibrium (LD) patterns in the human genome (limiting the mapping resolution of association studies), GWAS-identified risk loci typically span ~250 kb, encompassing an average of ~5 genes (numbers ranging from zero (“gene deserts”) to more than 50) and hundreds of associated variants. Contrary to widespread misconception, the causative variants and genes remain unknown for the vast majority of GWAS-identified risk loci. Yet, this remains a critical goal in order to reap the full benefits of GWAS in identifying new drug targets and developing effective predictive and diagnostic tools. It is the main objective of post-GWAS studies.

Distinguishing the few causative variants (i.e., the variants that are directly causing the gene perturbation) from the many neutral variants that are only associated with the disease because they are in LD with the former in the studied population, requires the use of sophisticated fine-mapping methods applied to very large, densely genotyped data sets<sup>4</sup>, ideally followed-up by functional studies<sup>5</sup>. Using such approaches, 18 causative variants for IBD were recently fine-mapped at single base pair resolution, and 51 additional ones at ≤10 base pair resolution<sup>4</sup>.

A minority of causative variants are coding, i.e., they alter the amino-acid sequence of the encoded protein. In such cases, and particularly if multiple such causative coding variants are found in the same gene (i.e., in case of allelic heterogeneity), the corresponding causative gene is unambiguously identified. In the case of IBD, causative genes have been identified for approximately ten risk loci on the basis of such “independently” (i.e., not merely reflecting LD with other variants) associated coding variants, including *NOD2*, *ATG16L1*, *IL23R*, *CARD9*, *FUT2*, and *TYK2*<sup>4,6–9</sup>.

For the majority of risk loci, the GWAS signals are not driven by coding variants. They must therefore be driven by common regulatory variants, i.e., variants that perturb the expression levels of one (or more) target genes in one (or more) disease relevant cell types<sup>4</sup>. Merely reflecting the proportionate sequence space that is devoted to the different layers of gene regulation (transcriptional, posttranscriptional, translational, posttranslational), the majority of regulatory variants are likely to perturb components of “gene switches” (promoters, enhancers, insulators), hence affecting transcriptional output. Indeed, fine-mapped non-coding variants are enriched in known transcription-factor binding sites and epigenetic signatures marking gene switch components<sup>4</sup>. Hence, the majority of common causative variants underlying inherited predisposition to common complex diseases must drive *cis*-eQTL (expression quantitative trait loci) affecting the causative gene(s) in one or more disease relevant cell types. The corresponding *cis*-eQTL are expected to operate prior to disease onset, and—driven by common variants—detectable in cohorts of healthy individuals of which most will never develop the disease. The term *cis*-eQTL refers to the fact that the regulatory variants that drive them only affect the expression of genes/alleles residing on the same DNA molecule, typically no more than one megabase away. Causative variants, whether coding or regulatory, may secondarily perturb the expression of genes/alleles located on different DNA molecules, generating *trans*-eQTL. Some of these *trans*-eQTL may participate in the disease process.

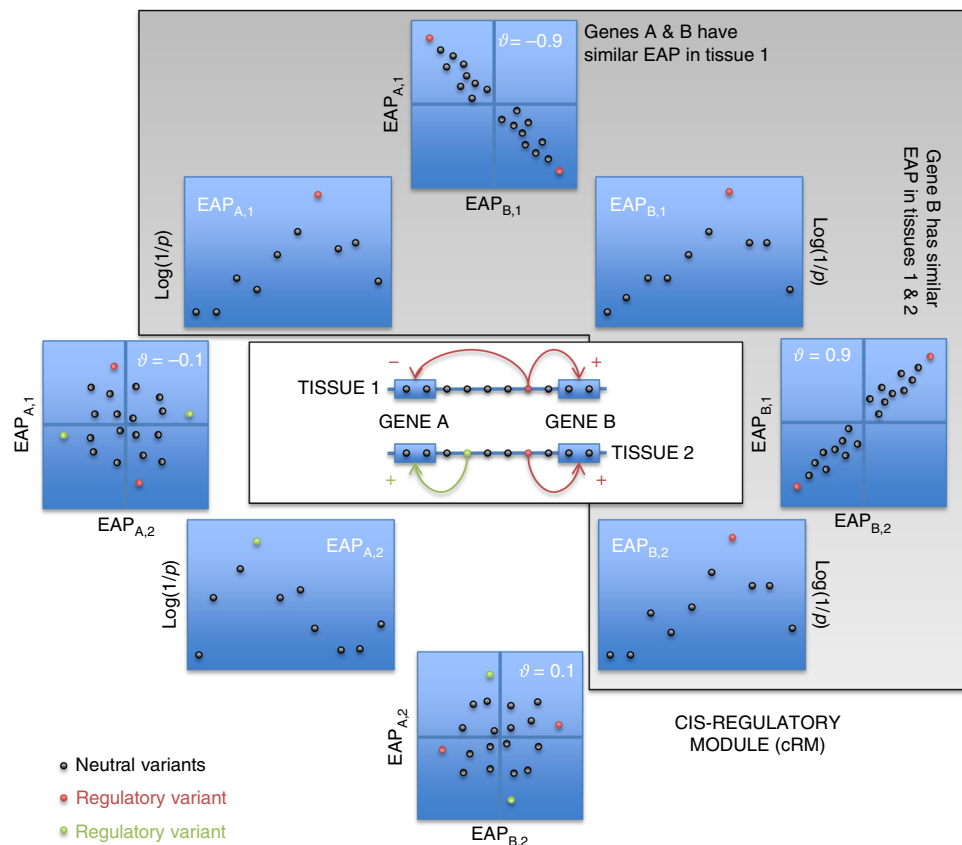
*cis*-eQTL effects are known to be very common, affecting more than 50% of genes<sup>10</sup>. Hence, finding that variants associated with

a disease are also associated with changes in expression levels of a neighboring gene is not sufficient to incriminate the corresponding genes as causative. Firstly, one has to show that the local association signal for the disease and for the eQTL are driven by the same causative variants. A variety of “colocalisation” methods have been developed to that effect<sup>11–13</sup>. Secondly, regulatory variants may affect elements that control the expression of multiple genes<sup>14</sup>, which may not all contribute to the development of the disease, i.e., be causative. Thus, additional evidence is needed to obtain formal proof of gene causality. In humans, the only formal test of gene causality that is applicable is the family of “burden” tests, i.e., the search for a differential burden of disruptive mutations in cases and controls, which is expected only for causative genes<sup>15</sup>. Burden tests rely on the assumption that—in addition to the common, mostly regulatory variants that drive the GWAS signal—the causative gene will be affected by low frequency and rare causative variants, including coding variants. Thus, the burden test makes the assumption that allelic heterogeneity is common, which is supported by the pervasiveness of allelic heterogeneity of Mendelian diseases in humans<sup>16</sup>. Burden tests compare the distribution of rare coding variants between cases and controls<sup>15</sup>. The signal-to-noise ratio of the burden test can be increased by restricting the analysis to coding variants that have a higher probability to disrupt protein function<sup>15</sup>. In the case of IBD, burden tests have been used to prove the causality of *NOD2*, *IL23R*, and *CARD9*<sup>6,8,9</sup>. A distinct and very elegant genetic test of gene causality is the reciprocal hemizygosity test, and the related quantitative complementation assay<sup>17,18</sup>. However, with few exceptions<sup>19,20</sup>, it has only been applied in model organisms in which gene knock-outs can be readily generated<sup>21</sup>.

In this paper, we describe the generation of a new and large data set for eQTL analysis (350 healthy individuals) in nine cell types that are potentially relevant for IBD. We identify and characterize ~24,000 *cis*-eQTL. By comparing disease and eQTL association patterns (EAP) using a newly developed statistic, we identify 99 strong positional candidate genes in 63 GWAS-identified risk loci. We resequence the 555 exons of 45 of these in 6600 cases and 5500 controls in an attempt to prove their causality by means of burden tests. The outcome of this study is relevant to post-GWAS studies of all common complex disease in humans.

## Results

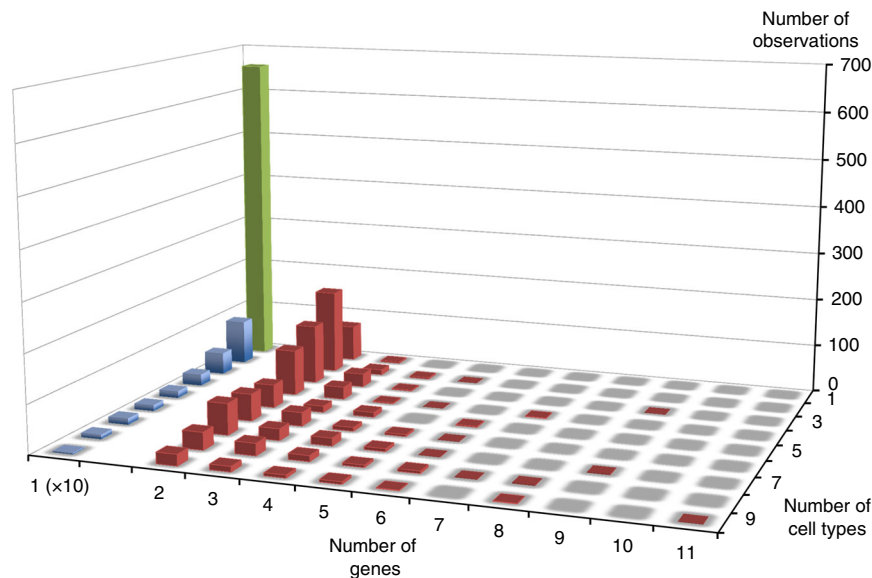
**Clustering *cis*-eQTL into regulatory modules.** We generated transcriptome data for six circulating immune cell types (CD4+ T lymphocytes, CD8+ T lymphocytes, CD19+ B lymphocytes, CD14+ monocytes, CD15+ granulocytes, platelets) as well as ileal, colonic, and rectal biopsies (IL, TR, RE), collected from 323 healthy Europeans (141 men, 182 women, average age 56 years, visiting the clinic as part of a national screening campaign for colon cancer) using Illumina HT12 arrays (CEDAR data set; Methods). IBD being defined as an inappropriate mucosal immune response to a normal commensal gut flora<sup>22</sup>, these nine cell types can all be considered to be potentially disease relevant. Using standard methods based on linear regression and two megabase windows centered on the position of the interrogating probe (Methods), we identified significant *cis*-eQTL (FDR < 0.05) for 8804 of 18,580 tested probes (corresponding to 7216 of 13,615 tested genes) in at least one tissue, amounting to a total of 23,650 *cis*-eQTL effects (Supplementary Data 1). When a gene shows a *cis*-eQTL in more than one tissue, the corresponding “eQTL association patterns” (EAP) (i.e., the distribution of association  $-\log(p)$  values for all the variants in the region of interest) are expected to be similar if determined by the same regulatory variants, and dissimilar otherwise. Likewise, if several



**Fig. 1** *cis*-Regulatory Module (cRM). A *cis*-eQTL affecting gene A in tissue 1 reveals itself by an “eQTL Association Pattern” (EAP<sub>A,1</sub>), i.e., the pattern of  $-\log(p)$  values for variants in the region. Multiple EAP can be observed in a given chromosome region, affecting one or more genes in one or more cell types. EAP that are driven by the same underlying variants are expected to be similar, while EAP driven by distinct variants (for instance, the green and red regulatory variants in the figure) are not. Based on the measure of similarity introduced in this work,  $\vartheta$ , we cluster the EAP in cRM. For EAP in the same module,  $\vartheta$  can be positive or negative, indicating that the variants have the same sign of effect (increasing or decreasing expression) for the corresponding EAP pair

neighboring genes show *cis*-eQTL in the same or distinct tissues, the corresponding EAP are expected to be similar if determined by the same regulatory variants, and dissimilar otherwise (Fig. 1). We devised the  $\vartheta$  metric to measure the similarity between association patterns (Methods).  $\vartheta$  is a correlation measure for paired  $-\log(p)$  values (for the two eQTL that are being compared) that ranges between  $-1$  and  $+1$ .  $\vartheta$  shrinks to zero if Pearson’s correlation between paired  $-\log(p)$  values does not exceed a chosen threshold (i.e., if the EAP are not similar).  $\vartheta$  approaches  $+1$  when the two EAP are similar and when variants that increase expression in eQTL 1 consistently increase expression in eQTL 2.  $\vartheta$  approaches  $-1$  when the two EAP are similar and when variants that increase expression in eQTL 1 consistently decrease expression in eQTL 2.  $\vartheta$  gives more weight to variants with high  $-\log(p)$  for at least one EAP (i.e., it gives more weight to eQTL peaks). Based on the known distribution of  $\vartheta$  under  $H_0$  (i.e., eQTL determined by distinct variants in the same region) and  $H_1$  (i.e., eQTL determined by the same variants), we selected a threshold value  $|\vartheta| > 0.60$  to consider that two EAP were determined by the same variant. This corresponds to a false positive rate of 0.05, and a false negative rate of 0.23 (Supplementary Fig. 1). We then grouped EAP in “*cis*-acting regulatory modules” (cRM) using  $|\vartheta|$  and a single-link clustering approach (i.e., an EAP needs to have  $|\vartheta| > 0.60$  with at least one member of the cluster to be assigned to that cluster). Clusters were visually examined and 29 single edges connecting otherwise unlinked and yet tight clusters manually removed (Supplementary Fig. 2).

Using this approach, we clustered the 23,650 effects in 9720 distinct “*cis*-regulatory modules” (cRM), encompassing *cis*-eQTL with similar EAP (Supplementary Data 2). Sixty-eight percent of cRM were gene- and tissue-specific, 22% were gene-specific but operating across multiple tissues ( $\leq 9$  tissues, average 3.5), and 10% were multi-genic ( $\leq 11$  genes, average 2.5) and nearly always multi-tissue (Figs. 2 and 3, Supplementary Fig. 2). In this, cRM are considered gene specific if the EAPs in the cluster concern only one gene, and tissue specific if the EAP in the cluster concern only one of the nine cell types. They are, respectively, multigenic and multi-tissue otherwise. cRM operating across multiple tissues tended to affect multiple genes ( $r = 0.47$ ;  $p < 10^{-6}$ ). In such cRM, the direction of the effects tended to be consistent across tissues and genes ( $p < 10^{-6}$ ). Nevertheless, we observed at least 55 probes with effect of opposite sign in distinct cell types ( $\vartheta \leq -0.9$ ), i.e., the corresponding regulatory variants increases transcript levels in one cell type while decreasing them in another (Fig. 4 and Supplementary Data 3). Individual tissues allowed for the detection of 7–33% of all cRM, and contributed 3–14% unique cRM (Supplementary Fig. 3). Sixty-nine percent of cRM were only detected in one cell type. The rate of cRM sharing between cell types reflects known ontogenic relations. Considering cRM shared by only two cell types (i.e., what jointly differentiates these two cell types from all other), revealed the close proximity of the CD4–CD8, CD14–CD15, ileum–colon, and colon–rectum pairs. Adding information of cRM shared by up to six cell types grouped lymphoid (CD4, CD8, CD19), myeloid (CD14, CD15 but not platelets), and intestinal (ileum, colon and rectum) cells.



**Fig. 2** Single-gene/tissue versus multi-gene/tissue cRM. Using  $|\rho| > 0.6$ , the 23,950 *cis*-eQTL ( $FDR \leq 0.05$ ) detected in the nine analyzed cell types were clustered in 9720 *cis*-Regulatory Modules (cRM). 68% of these were single-gene, single-tissue cRM (green), 22% were single-gene, multi-tissue cRM (blue), and 10% were multi-gene, mostly multi-tissue cRM (red). The number of observations for single-gene cRM were divided by 10 in the graph for clarity. Thus, there are more cases of single-gene, multi-tissue cRM (blue; 2155) than multi-gene cRM (red; 967)

Adding cRM with up to nine cell types revealed a link between ileum and blood cells, possibly reflecting the presence of blood cells in the ileal biopsies (Fig. 5).

**cRM matching IBD association signals are often multigenic.** If regulatory variants affect disease risk by perturbing gene expression, the corresponding “disease association patterns” (DAP) and EAP are expected to be similar, even if obtained in distinct cohorts (yet with same ethnicity) (Fig. 6). We confronted DAP and EAP using the  $\vartheta$  statistic and threshold ( $|\vartheta| > 0.60$ ) described above for 200 GWAS-identified IBD risk loci. DAP for Crohn’s disease and ulcerative colitis were obtained from the International IBD Genetics Consortium (IIBDGC)<sup>2,3</sup>, EAP from the CEDAR data set.

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

Strong correlations between DAP and EAP ( $|\vartheta| > 0.6$ , associated with low empirical  $p$  values) were observed for at least 63 IBD risk loci, involving 99 genes (range per locus: 1–6) (Table 1, Fig. 7, Supplementary Data 4). Increased disease risk was associated equally frequently with increased as with decreased expression ( $p_{CD} = 0.48$ ;  $p_{UC} = 0.88$ ). An open-access website has been prepared to visualize correlated DAP–EAP within their genomic context (<http://cedar-web.giga.ulg.ac.be>). Genes with highest  $|\vartheta|$  values ( $\geq 0.9$ ) include known IBD causative genes (for instance, *ATG16L1*, *CARD9*, and *FUT2*), known immune regulators (for instance, *IL18R1*, *IL6ST*, and *THEMIS*), as well as genes with as of yet poorly defined function in the context of IBD (for instance, *APEH*, *ANKRD55*, *CISD1*, *CPEB4*, *DOCK7*,

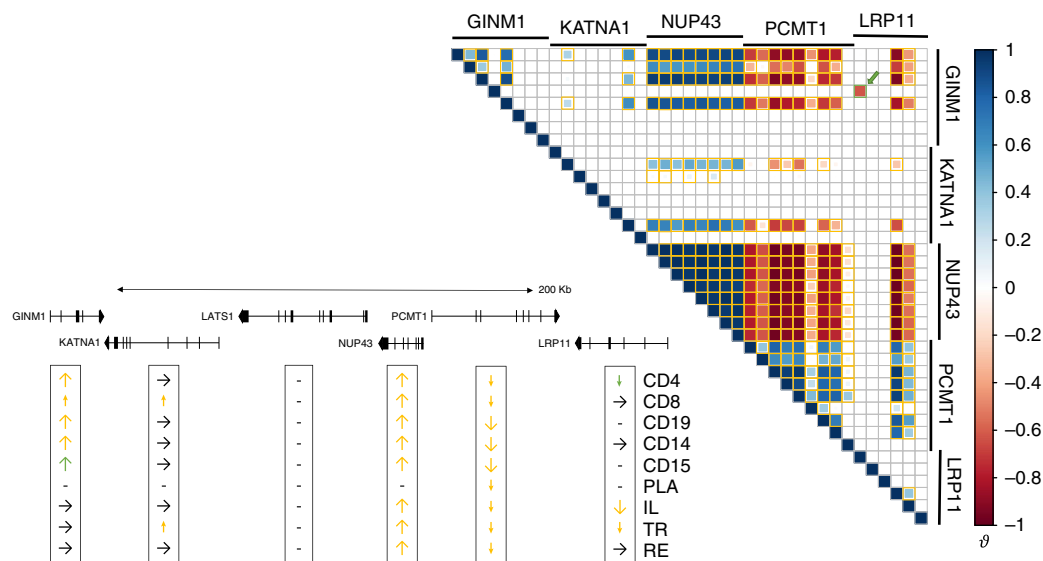
*ERAP2*, *GNA12*, *GPX1*, *GSDMB*, *ORMDL3*, *SKAP2*, *UBE2L3*, and *ZMIZ1*) (Supplementary Note 1).

The eQTL link with IBD has not been reported before for at least 47 of the 99 reported genes (Table 1). eQTL links with IBD have been previously reported for 111 additional genes, not mentioned in Table 1. Our data support these links for 19 of them, however, with  $|\vartheta| \leq 0.6$  (Supplementary Data 5). We applied SMR<sup>13</sup> as alternative colocalisation method to our data. Using a Bonferroni-corrected threshold of  $\leq 2.5 \times 10^{-5}$  for  $p_{SMR}$  and  $\geq 0.05$  for  $p_{HEIDI}$ , SMR detected 35 of the 99 genes selected with  $\vartheta$  (Supplementary Data 4). Using the same thresholds, SMR detected nine genes that were not selected by  $\vartheta$ . Of these, three (*ADAM15*, *AHSA2*, *UBA7*) had previously been reported by others, while six (*FAM189B*, *QRICH1*, *RBM6*, *TAP2*, *ADO*, *LGALS9*) were not. Of these six, three (*RBM6*, *TAP2*, *ADO*) were characterized by  $0.45 < |\vartheta| < 0.6$  (Supplementary Data 5).

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

The 400 analyzed DAP (200 CD and 200 UC) were found to match 76 cRM (in 63 risk loci) with  $|\vartheta| > 0.6$  (Table 1), of which 25 are multigenic. Knowing that multigenic cRM represent 10% of all cRM (967/9720), 25/76 (i.e., 33%) corresponds to a highly significant three-fold enrichment ( $p < 10^{-9}$ ). To ensure that this apparent enrichment was not due to the fact that multigenic cRM have more chance to match DAP (as by definition multiple EAP are tested for multigenic cRM), we repeated the enrichment





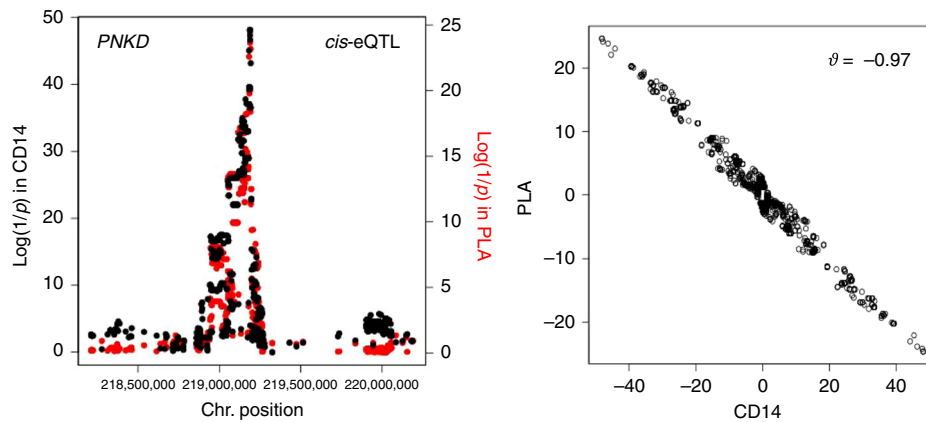
**Fig. 3** Example of a multi-gene, multi-tissue cRM. Gene-tissue combinations for which no expression could be detected are marked by “-”, with detectable expression but without evidence for *cis*-eQTL as “→”, with detectable expression and evidence for a *cis*-eQTL as “↑” or “↓” (large arrows: FDR < 0.05; small arrows: FDR ≥ 0.05 but high  $|\vartheta|$  values). eQTL labeled by the yellow arrows constitute the multi-genic and multi-tissular cRM no. 57. The corresponding regulatory variant(s) increase expression of the *GINM1*, *NUP43* and probably *KATNA1* genes (left side of the cRM), while decreasing expression of the *PCMT1* and *LRP11* genes (right side of the cRM). The expression of *GINM1* in CD15 and *LRP11* in CD4 appears to be regulated in opposite directions by a distinct cRM (no. 3694, green). The *LATS1* gene, in the same region, is not affected by the same regulatory variants in the studied tissues. Inset 1:  $\vartheta$  values for all EAP pairs. EAP pairs with  $|\vartheta| > 0.6$  are bordered in yellow when corresponding to cRM no. 57, in green when corresponding to cRM no. 3694 (+green arrow)

analysis by randomly sampling only one representative EAP per cRM in the 200 IBD risk loci. The frequency of multigenic cRM amongst DAP-matching cRM averaged 0.22, and was never  $\leq 0.10$  ( $p \leq 10^{-5}$ ) (Supplementary Fig. 4). In loci with high LD, EAP driven by distinct regulatory variants (yet in high LD) may erroneously be merged in the same cRM. To ensure that the observed enrichment in multigenic cRM was not due to higher levels of LD, we compared the LD-based recombination rate of the 63 cRM-matching IBD risk loci with that of the rest of the genome (<https://github.com/joepickrell/1000-genomes-genetic-maps>). The genome-average recombination rate was 1.23 centimorgan per megabase (cM/Mb), while that of the 63 IBD risk loci was 1.34 cM/Mb, i.e., less LD in the 63 cRM-matching IBD risk loci than in the rest of the genome. We further compared the average recombination rate in the 63 cRM-matching IBD regions with that of sets of 63 loci centered on randomly drawn cRM (from the list of 9720), matched for size and chromosome number (as cM/Mb is affected by chromosome size). The average recombination rate around all cRM was 1.43 cM/Mb, and this didn't differ significantly from the 63 cRM-matching IBD regions ( $p = 0.46$ ) (Supplementary Fig. 5). Therefore, the observed enrichment cannot be explained by a higher LD in the 63 studied IBD risk loci. Taken together, EAP that are strongly correlated with DAP ( $|\vartheta| \geq 0.60$ ), map to regulatory modules that are 2- to 3-fold enriched in multigenic cRM when compared to the genome average and include four of the top 10 (of 9720) cRM ranked by number of affected genes.

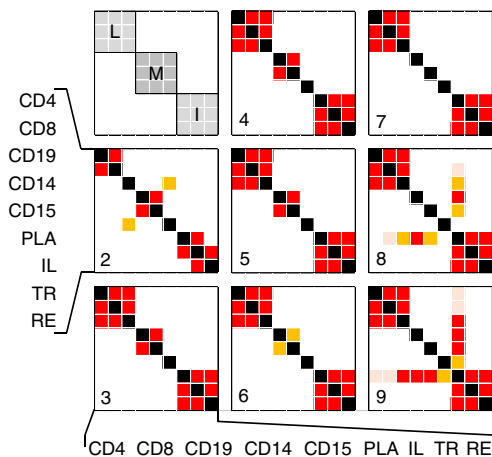
#### DAP-matching cRM are enriched in causative genes for IBD.

For truly causative genes, the burden of rare disruptive variants is expected to differ between cases and controls<sup>23</sup>. We therefore performed targeted sequencing for the 555 coding exons (~88 kb) of 38 genes selected amongst those with strongest DAP-EAP

correlations, plus seven genes with suggestive DAP-EAP evidence backed by literature (Table 1), in 6597 European CD cases and 5502 matched controls (ref.<sup>24</sup> and Methods). Eighteen of these were part of single-gene cRM and the only gene highlighted in the corresponding locus. The remaining 27 corresponded to multi-gene cRM mapping to 15 risk loci. We added the well-established *NOD2* and *IL23R* causative IBD genes as positive controls. We identified a total of 174 loss-of-function (LoF) variants, 2567 missense variants (of which 991 predicted by SIFT<sup>25</sup> to be damaging and Polyphen-2<sup>26</sup> to be either possibly or probably damaging), and 1434 synonymous variants (Fig. 8 and Supplementary Data 6). 1781 of these were also reported in the Genome Aggregation Database<sup>27</sup> with nearly identical allelic frequencies (Supplementary Fig. 6). We designed a gene-based burden test to simultaneously evaluate hypothesis (i): all disruptive variants enriched in cases (when  $\vartheta < 0$ ; risk variants) or all disruptive variants enriched in controls (when  $\vartheta > 0$ ; protective variants), and hypothesis (ii): some disruptive variants enriched in cases and others in controls. Hypothesis (i) was tested with CAST<sup>28</sup>, and hypothesis (ii) with SKAT<sup>29</sup> (Methods). We restricted the analysis to 1141 LoF and damaging missense variants with minor allele frequency (MAF)  $\leq 0.005$  to ensure that any new association signal would be independent of the signals from common and low frequency variants having led to the initial identification and fine-mapping of the corresponding loci<sup>4</sup>. For *NOD2* ( $p = 6.9 \times 10^{-7}$ ) and *IL23R* ( $p = 1.8 \times 10^{-4}$ ), LoF and damaging variants were significantly enriched in respectively cases and controls as expected. When considering the 45 newly tested genes as a whole, we observed a significant ( $p = 6.9 \times 10^{-4}$ ) shift towards lower  $p$  values when compared to expectation, while synonymous variants behaved as expected ( $p = 0.66$ ) (Fig. 9 and Supplementary Data 7). This strongly suggests that the sequenced list includes causative genes. *CARD9*, *TYK2*, and *FUT2* have recently been shown to be causative genes based on disease-associated low-frequency coding variants (MAF > 0.005)<sup>4</sup>. The



**Fig. 4** Variant(s) with opposite effects on expression in two cell types. Example of a gene (*PNKD*) affected by a *cis*-eQTL in at least two cell types (CD14 and platelets) that are characterized by EAP with  $\theta = -0.97$ , indicating that the gene's expression level is affected by the same regulatory variant in these two cell types, yet with opposite effects, i.e., the variant that is increasing expression in platelets is decreasing expression in CD14



**Fig. 5** Significance of the excess sharing of cRM between cell types. (red:  $p < 0.0002$  (Bonferroni corrected 0.0144), orange:  $p < 0.001$  (0.072), rose:  $p < 0.01$  (0.51)). The numbers in the lower-left corner of the squares indicate which cRM were used for the analysis: (2) cRM affecting no more than two cell types, (3) cRM affecting no more than three cell types, etc. The upper-left square indicates the position of the lymphoid cell types (L) (CD4, CD8, CD19), the myeloid cell types (M) (CD14, CD15, PLA), and the intestinal cell types (I) (IL, TR, RE). For each pair of cell types  $i$  and  $j$ , we computed two  $p$  values, one using  $i$  as reference, the other using  $j$  as reference (Methods). Pairs of  $p$  values were always consistent

shift towards lower  $p$  values remained significant without these ( $p = 1.7 \times 10^{-3}$ ), pointing towards novel causative genes amongst the 42 remaining candidate genes.

### Proving gene causality requires larger case-control cohorts.

Despite the significant shift towards lower  $p$  values when considering the 45 genes jointly, none of these were individually significant when accounting for multiple testing ( $p \leq \frac{0.05}{2 \times 45} \approx 0.0006$ ) (Supplementary Data 7). Near identical results were obtained when classifying variants using the Combined Annotation Dependent Depletion (CADD) tool<sup>30</sup> instead of SIFT/PolyPhen-2 (Supplementary Data 7). We explored three approaches to increase the power of the burden test. The first built on the observation that cRM matching DAP are enriched in multigenic modules. This suggests that part of IBD risk loci harbor multiple co-regulated and hence functionally related

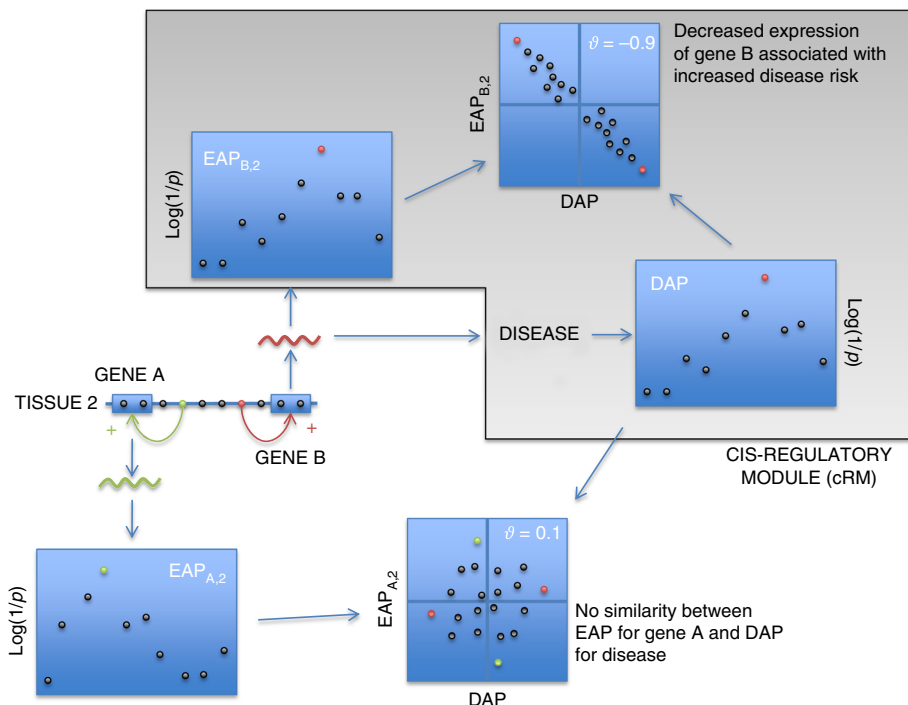
genes, of which several (rather than one, as generally assumed) may be causally involved in disease predisposition. To test this hypothesis, we designed a module- rather than gene-based burden test (Methods). However, none of the 30 tested modules reached the experiment-wide significance threshold ( $p \leq \frac{0.05}{2 \times 30} \approx 0.0008$ ). Moreover, the shift towards lower  $p$  values for the 30 modules was not more significant ( $p = 2.3 \times 10^{-3}$ ) than for the gene-based test (Supplementary Fig. 7a and Supplementary Table 7). The second and third approaches derive from the common assumption that the heritability of disease predisposition may be larger in familial and early-onset cases<sup>31</sup>. We devised orthogonal tests for age-of-onset and familiarity and combined them with our burden tests (Methods). Neither approach would improve the results (Supplementary Fig. 7b, c and Supplementary Data 7).

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

### Discussion

We herein describe a novel dataset comprising array-based transcriptome data for six circulating immune cell types and intestinal biopsies at three locations collected on ~300 healthy European individuals. We use this ‘‘Correlated Expression and Disease Association Research’’ data set (CEDAR) to identify 23,650 significant *cis*-eQTL, which fall into 9720 regulatory modules of which at least ~889 affect more than one gene in more than one tissue. We provide strong evidence that 63 of 200 known IBD GWAS signals reflect the activity of common regulatory variants that preferentially drive multigenic modules. We perform an exon-based burden test for 45 positional candidate CD genes mapping to 33 modules, in 5500 CD cases and 6500 controls. By demonstrating a significant ( $p = 6.9 \times 10^{-4}$ ) upwards shift of  $\log(1/p)$  values for damaging when compared to synonymous variants, we show that the sequenced genes include new causative CD genes.

Individually, none of the sequenced genes (other than the positive *NOD2* and *IL23R* controls) exceed the experiment-wide significance threshold, precluding us from definitively pinpointing any novel causative genes. However, we note *IL18R1* amongst



**Fig. 6** DAP-matching cRM. If a regulatory variant (red) affects disease risk by altering the expression levels of gene B in tissue 2, the  $EAP_{B,2}$  is expected to be similar (high  $|\theta|$ ) to the “disease association pattern” (DAP), both assigned therefore to the same cRM.  $\theta$  is positive if increased gene expression is associated with increased disease risk, negative otherwise. A *cis*-eQTL that is driven by a regulatory variant (green) that does not directly affect disease risk, will be characterized by an EAP (say gene A, tissue 2,  $EAP_{A,2}$ ) that is not similar to the DAP (low  $|\theta|$ )

the top-ranking genes (see also Supplementary Note 1). *IL18R1* is the only gene in an otherwise relatively gene-poor region (also encompassing *IL1R1* and *IL18RAP*) characterized by robust *cis*-eQTL in CD4 and CD8 that are strongly correlated with the DAP for CD and UC ( $0.68 \leq |\theta| \leq 0.93$ ). Reduced transcript levels of *IL18R1* in these cell types is associated with increased risk for IBD. Accordingly, rare ( $MAF \leq 0.005$ ) damaging variants were cumulatively enriched in CD cases (CAST  $p = 0.05$ ). The cumulative allelic frequency of rare damaging variants was found to be higher in familial CD cases (0.0027), when compared to non-familial CD cases (0.0016;  $p = 0.09$ ) and controls (0.0010;  $p = 0.03$ ). When ignoring carriers of deleterious *NOD2* mutations, average age-of-onset was reduced by  $\sim 3$  years (25.3 vs 28.2 years) for carriers of rare damaging *IL18R1* variants but this difference was not significant ( $p = 0.18$ ).

While the identification of matching cRM for 63/200 DAP points towards a number of strong candidate causative genes, it leaves most risk loci without matching eQTL despite the analysis of nine disease-relevant cell types. This finding is in agreement with previous reports<sup>4,32</sup>. It suggests that *cis*-eQTL underlying disease predisposition operate in cell types, cell states (for instance, resting vs activated) or developmental stages that were not explored in this and other studies. It calls for the enlargement and extension of eQTL studies to more diverse and granular cellular panels<sup>10,33</sup>, possibly by including single-cell sequencing or spatial transcriptomic approaches. By performing eQTL studies in a cohort of healthy individuals, we have made the reasonable assumption that the common regulatory variants that are driving the majority of GWAS signals are acting before disease onset, including in individuals that will never develop the disease. An added advantage of studying a healthy cohort, is that the corresponding dataset is “generic”, usable for the study of perturbation of gene regulation for any common complex disease. However, it is conceivable that some eQTL underlying increased disease risk only manifest themselves once

the disease process is initiated, for instance as a result of a modified inflammatory status. Thus, it may be useful to perform eQTL studies with samples collected from affected individuals to see in how far the eQTL landscape is affected by disease status.

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

Although we re-sequenced the ORF of 45 carefully selected candidate genes in a total of 5500 CD cases and 6600 controls, none of the tested genes exceeded the experiment-wide threshold of significance. This is despite the fact that we used a one-sided, eQTL-informed test to potentially increase power. Established IBD causative genes used as positive control, *NOD2* and *IL23R*, were positive indicating that the experiment was properly conducted. We were not able to improve the signal strength by considering information about regulatory modules, familiarity or age-of-onset. We estimated that  $\geq 10$ -fold larger sample sizes will be needed to achieve adequate power if using the same approach.

**Table 1 IBD risk loci for which at least one *cis*-eQTL association pattern (EAP) was found to match the disease association pattern (DAP)**

Loci	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	4,48
HD2	1	7.7	8.3	2900	1	<i>PARK7</i>	CD15 TR RE	-0.8	-0.82	0.01	0.06	48
N_1_62	1	62.5	63.5	109	3	<b>DOCK7</b> <i>USP1 ATG4C</i>	CD4 CD8 <b>CD19</b> 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	
HD14	1	155.0	156.1	5	8	<i>GBA</i>	CD4	-0.65	0	0.01	1.00	
				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	<b>IL18R1</b>	<b>CD4 CD8</b>	-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	2,48
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	2,48
HD43	2	234.1	234.6	1177	1	<b>ATG16L1</b>	CD4 CD8 IL TR RE	0.94	0	0.05	1.00	2,49
N_3_45	3	46.0	47.0	2930	1	<i>CCR2</i>	CD19	0.77	0	0.02	1.00	
				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	
HD50	3	48.4	51.4	8	7	<i>USP4</i>	CD19	0.64	0.63	0.06	0.07	2
				217	3	<b>GPX1</b> <i>APEH</i> <i>IP6K1</i>	<b>CD19</b> CD14 <b>TR RE</b>	0.91	0.97	0.01	0.01	2,49
				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	50
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	51
HD61	5	55.4	55.5	360	2	<b>ANKRD55</b> <b>IL6ST</b>	<b>CD4 CD8</b>	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	<b>ERAP2</b> <i>LNPEP</i>	<b>CD4 CD8 CD19 CD14 CD15</b> PLA	0.94	0.71	0.01	0.02	2,4,50
							<b>IL TR RE</b>					
HD65	5	130.4	132.0	55	4	( <i>SLC22A4</i> ) ( <i>SLC22A5</i> )	CD4 CD15	-0.55	0	0.06	0.07	4,52
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>IRGM</i> )	-	-	-	-	-	53
HD71	5	173.2	173.6	1349	1	<i>CPEB4</i>	<b>CD4</b> CD8 CD19 <b>CD14</b> CD15 PLA	-0.92	0	0.01	1.00	2,4
							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	<b>THEMIS</b> <i>PTPRK</i>	<b>CD8</b>	-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	<b>GNA12</b>	<b>CD19</b> CD14 TR	0	-0.94	1.00	0.02	2
HD84	7	26.6	27.3	1441	1	<b>SKAP2</b>	<b>CD4 CD8 CD19</b>	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	2,48
				7046	1	<i>TSPAN33</i>	CD19	-0.64	0	0.01	1.00	
N_8_26	8	26.7	27.7	5869	1	<i>PTK2B</i>	CD14	-0.69	0	0.01	1.00	
				5841	1	<i>TRIM35</i>	CD4	0	0.66	1.00	0.01	
HD106	9	139.1	139.5	64	4	<b>CARD9</b> <i>INPP5E</i> <i>SEC16A</i>	CD4 CD8 CD19 CD14 <b>CD15</b> IL TR	0.95	0.86	0.01	0.02	2,4,50
						<i>SDCCAG3</i>	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	<b>CISD1</b>	<b>CD4 CD8 CD19 CD14 CD15</b> TR	0.94	0.83	0.04	0?	2,4,48
							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	<b>ZMIZ1</b>	<b>CD8</b>	-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	4,50
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	



**Table 1** (continued)

Loc	Chr	Beg	End	cRM	Nr	Genes with correlated DAP-EAP	Implicated cell types	Best $\theta$		Best $p$		Ref
								CD	UC	CD	UC	
J_16_22	16	23.6	24.1	2672	1	PRKCB	CD14	0	0.64	1.00	0.05	<sup>2</sup>
HD150	16	28.2	29.1	6	8	<i>TUFM</i> SBK1 APOBR <i>SGF29</i> <i>CLN3</i> <i>SPNS1</i>	CD4 CD8 CD19 CD14 CD15 IL TR RE	0.81	0.86	0.05	0.03	<sup>4</sup>
HD151	16	30.4	31.4	2673	1	RNF40	CD15	-0.63	0	0.02	1.00	
				1886	1	<i>ITGAL</i>	CD4 CD8 CD19	0	0.74	1.00	0.01	<sup>54</sup>
HD153	16	68.4	68.9	1894	1	ZFP90	CD4 CD8 CD19 CD14 TR	0	0.83	1.00	0.07	<sup>2,48</sup>
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	<b>GSDMB ORMDL3</b> <i>PGAP3</i> ( <i>GSDMA</i> )	<b>CD4 CD8 CD19</b> CD14 <b>IL TR RE</b>	-0.98	-0.92	0.02	0.01	<sup>2,4</sup>
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	CD226	CD4 CD8 PLA	0	-0.86	1.00	0.01	<sup>2</sup>
N_18_76	18	76.7	77.7	7292	1	PQLC1	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	GNG8	CD4	0	-0.63	1.00	0.06	
HD169	19	49.0	49.3	3128	1	<b>FUT2</b>	<b>IL TR RE</b>	-0.95	0	0.01	1.00	<sup>4</sup>
J_20_31	20	31.1	31.6	593	2	COMMD7	CD14	0	0.61	1.00	0.01	
J_20_32	20	33.6	34.1	7	8	UQCC1	CD19	-0.69	0	0.02	1.00	<sup>2</sup>
				3369	1	MMP24-AS1	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	<sup>2</sup>
HD176	21	16.6	16.9	9578	1	NRIP1	CD4	0	-0.69	1.00	0.02	
HD180	22	21.7	22.1	2130	1	<b>UBE2L3</b>	<b>CD4 CD8 CD19 CD14 CD15 IL TR RE</b>	0.97	0.92	0.01	0.07	<sup>2,4</sup>
N_22_41	22	41.4	42.4	2149	1	EP300	CD8 CD19 CD15	0	0.71	1.00	0.02	

Given are (i) the name and chromosomal coordinates of the corresponding loci (Locus, Chr, Beg, End) (GRCh37/hg19 in Mb), (ii) the identifier and total number of genes in the matching cis-acting regulatory module (cRM, Nr), (iii) the genes and tissues involved in matching DAP-EAP ( $|\theta| > 0.6$ ) (bold when  $|\theta| \geq 0.9$ ), (iv) the best  $\theta$ -values and corresponding empirical  $p$  values obtained for CD and UC, respectively, and (v) references reporting a link between one or more of the same genes and IBD on the basis of eQTL information. Genes that were resequenced are shown in italics. Genes that were resequenced despite  $|\theta| \leq 0.6$  are bracketed, and the supporting references provided in "Ref". The higher number of matching DAP-EAP in this study when compared to Huang et al.<sup>4</sup> are primarily due to the fact that (i) we herein study 200 IBD risk loci (vs 97), and (ii) we increase the number of detected cis-eQTL approximately two-fold by correcting for hidden confounders using PCs

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

It is conceivable that the organ-specificity of nearly all complex diseases (such as the digestive tract for IBD), reflects tissue-specific perturbation of broadly expressed causative genes that may fulfill diverse functions in different organs. If this is true, coding variants may not be the appropriate substrate to perform burden tests, as these will affect the gene across all tissues. In such instances, the disruptive variants of interest may be those perturbing tissue-specific gene switches. Also, it has recently been proposed that the extreme polygenic nature of common complex diseases may reflect the trans-effects of a large proportion of regulatory variants active in a given cell type on a limited number of core genes via perturbation of highly connected gene networks<sup>34</sup>. Identifying rare regulatory variants is still challenging, however, as tissue-specific gene switches remain poorly cataloged, and the effect of variants on their function difficult to predict. The corresponding sequence space may also be limited in size, hence limiting power. Nevertheless, a reasonable start may be to re-sequence the regions surrounding common regulatory variants that have been fine-mapped at near single base pair resolution<sup>4</sup>.

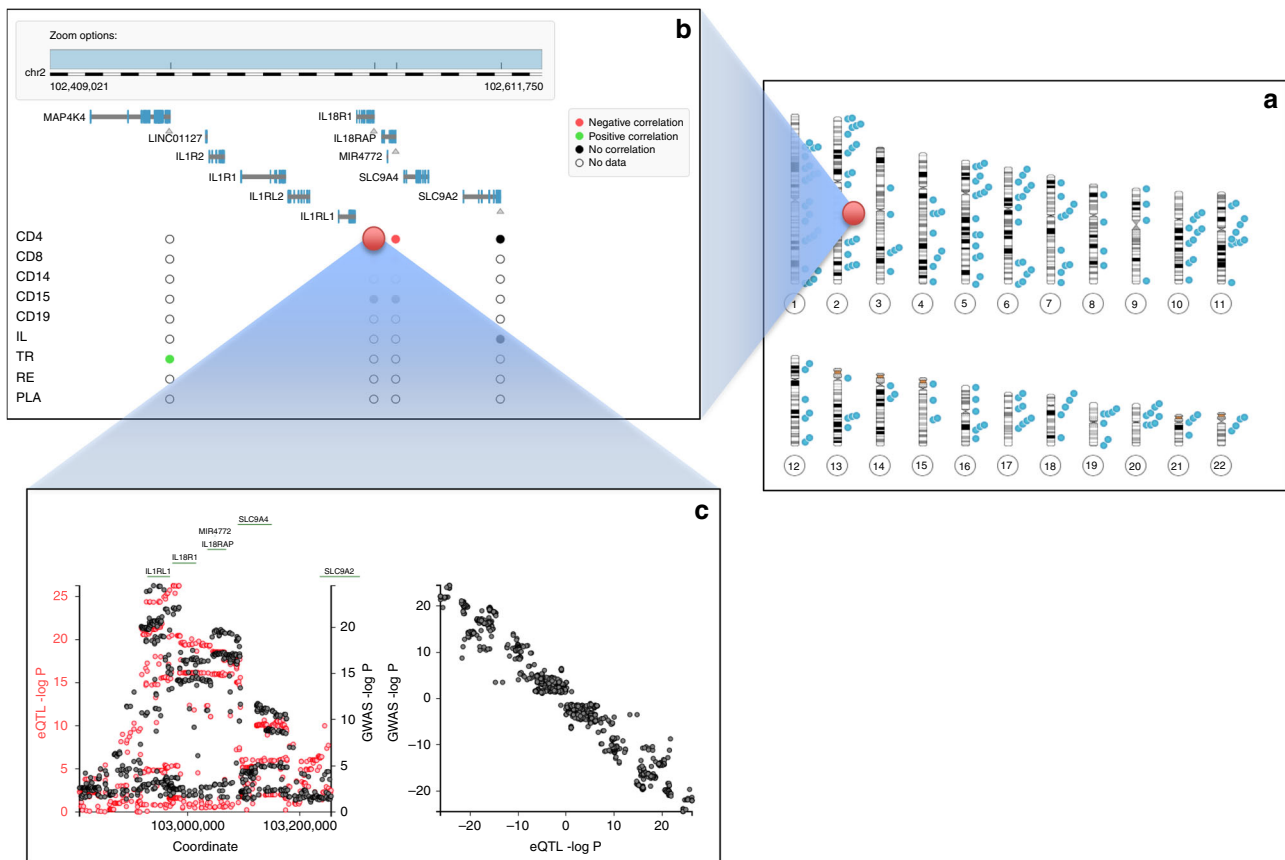
In conclusion, we hereby provide to the scientific community a collection of ~24,000 cis-eQTL in nine cell types that are highly relevant for the study of inflammatory and immune-mediated diseases, particularly of the intestinal tract. The CEDAR dataset advantageously complements existing eQTL datasets including GTEx<sup>10,33</sup>. We propose a paradigm to rationally organize cis-eQTL effects in co-regulated clusters or regulatory modules. We identify ~100 candidate causative genes in 63 out of 200 analyzed risk loci, on the basis of correlated DAP and EAP. We have developed a web-based browser to share the ensuing results with the scientific community (<http://cedar-web.giga.ulg.ac.be>). The

CEDAR website will imminently be extended to accommodate additional common complex disease for which GWAS data are publicly available. We show that the corresponding candidate genes are enriched in causative genes, however, that case-control cohorts larger than those used in this study (12,000 individuals) are required to formally demonstrate causality by means of presently available burden tests.

### Methods

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

**SNP genotyping and imputation.** Total DNA was extracted from EDTA-collected peripheral blood using the MagAttract DNA blood Midi M48 Kit on a QIAcube robot (Qiagen). DNA concentrations were measured using the Quant-iT Picogreen ds DNA Reagents (Invitrogen). Individuals were genotyped for >700 K SNPs using Illumina's Human OmniExpress BeadChips, an iScan system and the Genome Studio software following the guidelines of the manufacturer. We eliminated variants with call rate  $\leq 0.95$ , deviating from Hardy-Weinberg equilibrium ( $p \leq 10^{-4}$ ), or which were monomorphic. We confirmed European ancestry of all individuals by PCA using the HapMap population as reference. Using the real genotypes of 629,570 quality-controlled autosomal SNPs as anchors, we used the Sanger Imputation Services with the UK10K + 1000 Genomes Phase 3 Haplotype panels (<https://imputation.sanger.ac.uk>)<sup>35-37</sup> to impute genotypes at autosomal variants in our population. We eliminated indels, SNPs with MAF  $\leq 0.05$ , deviating from



**Fig. 7** Screen shots of the CEDAR website, showing **a** known CD risk loci on the human karyotype, **b** a zoom in the HD35 risk locus showing the Refseq gene content and summarizing local CEDAR *cis*-eQTL data (white: no expression data, gray: expression data but no evidence for *cis*-eQTL, black: significant *cis*-eQTL but no correlation with DAP, red: significant *cis*-eQTL similar to DAP ( $\theta < -0.60$ ), green: significant *cis*-eQTL similar to DAP ( $\theta > 0.60$ )), and **c** a zoom in the DAP for Crohn's disease (black) and EAP for *IL18R1* (red), as well as the signed correlation between DAP and EAP

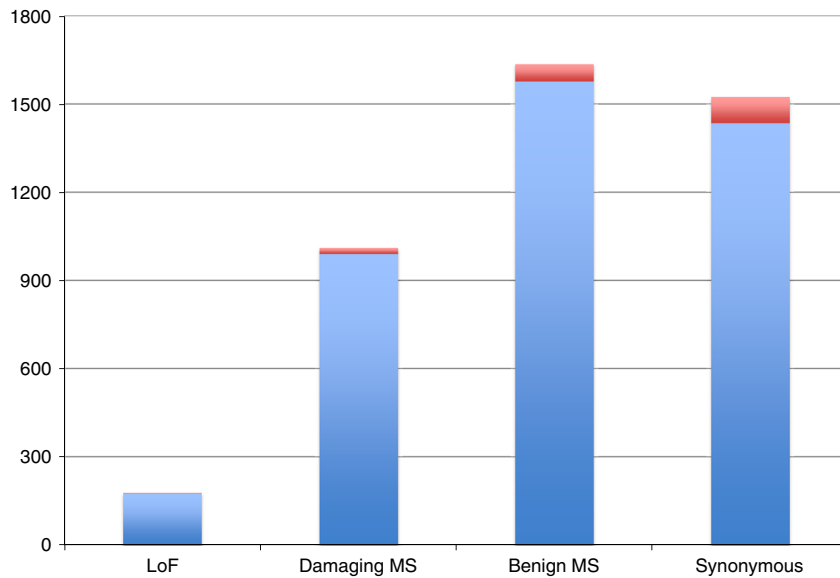
Hardy-Weinberg equilibrium ( $p \leq 10^{-3}$ ), and with low imputation quality ( $\text{INFO} \leq 0.4$ ), leaving 6,019,462 high quality SNPs for eQTL analysis.

**Transcriptome analysis.** Blood samples were kept on ice and treated within 1 h after collection as follows. EDTA-collected blood was layered on Ficoll-Paque PLUS (GE Healthcare) to isolate peripheral blood mononuclear cells by density gradient centrifugation. CD4+ T lymphocytes, CD8+ T lymphocytes, CD19+ B lymphocytes, CD14+ monocytes, and CD15+ granulocytes were isolated by positive selection using the MACS technology (Miltenyi Biotec). To isolate platelets, blood collected on acid-citrate-dextrose (ACD) anticoagulant was centrifuged at 150 g for 10 min. The platelet rich plasma (PRP) was collected, diluted twofold in ACD buffer and centrifuged at  $800 \times g$  for 10 min. The platelet pellet was resuspended in MACS buffer (Miltenyi Biotec) and platelets purified by negative selection using CD45 microbeads (Miltenyi Biotec). Intestinal biopsies were flash frozen in liquid nitrogen immediately after collection and kept at  $-80^\circ\text{C}$  until RNA extraction. Total RNA was extracted from the purified leucocyte populations and intestinal biopsies using the AllPrep Micro Kit and a QIAcube robot (Qiagen). For platelets, total RNA was extracted manually with the RNeasy Mini Kit (Qiagen). Whole genome expression data were generated using HT-12 Expression Beadchips following the instructions of the manufacturer (Illumina). Technical outliers were removed using controls recommended by Illumina and the Lumi package<sup>38</sup>. We kept 29,464/47,323 autosomal probes (corresponding to 19,731 genes) mapped by Re-Annotator<sup>39</sup> to a single gene body with  $\leq 2$  mismatches and not spanning known variants with  $\text{MAF} > 0.05$ . Within cell types, we only considered probes (i.e., “usable” probes) with detection  $p$  value  $\leq 0.05$  in  $\geq 25\%$  of the samples. Fluorescence intensities were  $\text{Log}_2$  transformed and Robust Spline Normalized (RSN) with Lumi<sup>38</sup>. Normalized expression data were corrected for sex, age, smoking status and Sentrix Id using ComBat from the SVA R library<sup>40</sup>. We further corrected the ensuing residuals within tissue for the number of Principal Components (PC) that maximized the number of *cis*-eQTL with  $p \leq 10^{-6}$ <sup>41</sup>. Supplementary Table 2 summarizes the number of usable samples, probes and PC for each tissue type.

***cis*-eQTL analysis.** *cis*-eQTL analyses were conducted with PLINK and using the expression levels precorrected for fixed effects and PC as described above (<http://pnu.gmh.harvard.edu/purcell/plink/>)<sup>42</sup>. Analyses were conducted under an

additive model, i.e., assuming that the average expression level of heterozygotes is at the midpoint between alternate homozygotes. To identify *cis*-eQTL we tested all SNPs in a 2 Mb window centered around the probe (if “usable”).  $P$  values for individual SNPs were corrected for the multiple testing within the window by permutation (10,000 permutations). For each probe–tissue combination we kept the best (corrected)  $p$  value. Within each individual cell type, the ensuing list of corrected  $p$  values was used to compute the corresponding false discovery rates (FDR or  $q$  value). Supplementary Table 3 reports the number of *cis*-eQTL found in the nine analyzed cell types for different FDR thresholds (see also Supplementary Fig. 9).

**Comparing EAP with  $\theta$  to identify *cis*-regulatory modules.** If the transcript levels of a given gene are influenced by the same regulatory variants (one or several) in two tissues, the corresponding EAP (i.e., the  $-\log(p)$  values of association for the SNPs surrounding the gene) are expected to be similar. Likewise, if the transcript levels of different genes are influenced by the same regulatory variants in the same or in different tissues, the corresponding EAP are expected to be similar (cf. main text, Fig. 1). We devised a metric,  $\theta$ , to quantify the similarity between EAP. If two EAP are similar, one can expect the corresponding  $-\log(p)$  values to be positively correlated. One particularly wants the EAP peaks, i.e., the highest  $-\log(p)$  values, to coincide in order to be convinced that the corresponding *cis*-eQTL are driven by the same regulatory variants. To quantify the similarity between EAP while emphasizing the peaks, we developed a weighted correlation. Imagine two vectors  $\mathbf{X}$  and  $\mathbf{Y}$  of  $-\log(p)$  values for  $n$  SNPs surrounding the gene(s) of interest. Using the same nomenclature as in Fig. 1a,  $\mathbf{X}$  could correspond to gene A in tissue 1, and  $\mathbf{Y}$  to gene A in tissue 2, or  $\mathbf{X}$  could correspond to gene A in tissue 1, and  $\mathbf{Y}$  to gene B in tissue 2. We only consider for analysis, SNPs within 1 Mb of either gene (probe) and for which  $x_i$  and/or  $y_i$  is superior to 1.3 (i.e.,  $p$  value  $< 0.05$ ) hence informative for at least one of the two *cis*-eQTL. Indeed, the majority of variants with  $-\log(p) < 1.3$  ( $p > 0.05$ ) for both EAP are by definition not associated with either trait. There is therefore no reason to expect that they could contribute useful information to the correlation metric: their ranking in terms of  $-\log(p)$  values becomes more and more random as the  $-\log(p)$  decreases. We define the weight to be given to each SNP



**Fig. 8** Variants detected by sequencing the coding exons of 45 candidate genes. Variants are sorted in LoF (loss-of-function, i.e., stop gain, frame-shift, splice site), Damaging MS (missense variants considered as damaging by SIFT<sup>5</sup> and damaging or possibly damaging by Polyphen-2<sup>6</sup>), Benign MS (other missense variants), and Synonymous. Blue: variants with MAF < 0.005, Red: variants with MAF ≥ 0.005

in the correlation as:

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

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

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

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

in which

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

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

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

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

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

$r_w$  ignores an important source of information. If two EAP are driven by the same regulatory variant, there should be consistency in the signs of the effects across SNPs in the region. We will refer to the effect of the “reference” allele of SNP  $i$  on the expression levels for the first and second *cis*-eQTL as  $\beta_i^X$  and  $\beta_i^Y$ . If the reference allele of the regulatory variant increases expression for both *cis*-eQTL, the  $\beta_i^X$  and  $\beta_i^Y$ ’s for a SNPs in LD with the regulatory variant are expected to have the same sign (positive or negative depending on the sign of D for the considered SNP). If the reference allele of the regulatory variant increases expression for one *cis*-eQTL and decreases expression for the other, the  $\beta_i^X$  and  $\beta_i^Y$ ’s for a SNPs in LD with the regulatory variant are expected to have opposite sign. We used this notion to develop a weighted and signed measure of correlation,  $r_{ws}$ . The approach was the same as for  $r_w$ , except that the values of  $y_i$  were multiplied by  $-1$  if the signs of  $\beta_i^X$

and  $\beta_i^Y$  were opposite.  $r_{ws}$  is expected to be positive if the regulatory variant affects the expression of both *cis*-eQTL in the same direction and negative otherwise.

We finally combined  $r_w$  and  $r_{ws}$  in a single score referred to as  $\vartheta$ , as follows:

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

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

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

We then evaluated the distribution of  $\vartheta$  for pairs of EAP driven by distinct regulatory variants in the same chromosomal region as follows. We considered 1207 significant *cis*-eQTL (mapping to the 200 IBD risk loci described above). For each one of these, we generated a set of 100 “matching” *cis*-eQTL effects in silico, sequentially considering 100 randomly selected SNPs (from the same locus) as causal. The in silico *cis*-eQTL were designed such that they would explain the same fraction of expression variance as the corresponding real *cis*-eQTL detected with PLINK (cfr. above). When performing *cis*-eQTL analysis under an additive model, PLINK estimates  $\beta_0$  (i.e., the intercept), and  $\beta_1$  (i.e., the slope of the regression), including for the top SNP. Assume that the expression level of the studied gene,  $Z$ , for individual  $i$  is  $z_i$ . Assume that the sample comprises  $n_T$  individuals in total, of which  $n_{11}$  are of genotype “11”,  $n_{12}$  of genotype “12”, and  $n_{22}$  of genotype “22”, for the top *cis*-eQTL SNP. The total expression variance for gene  $Z$  equals:

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

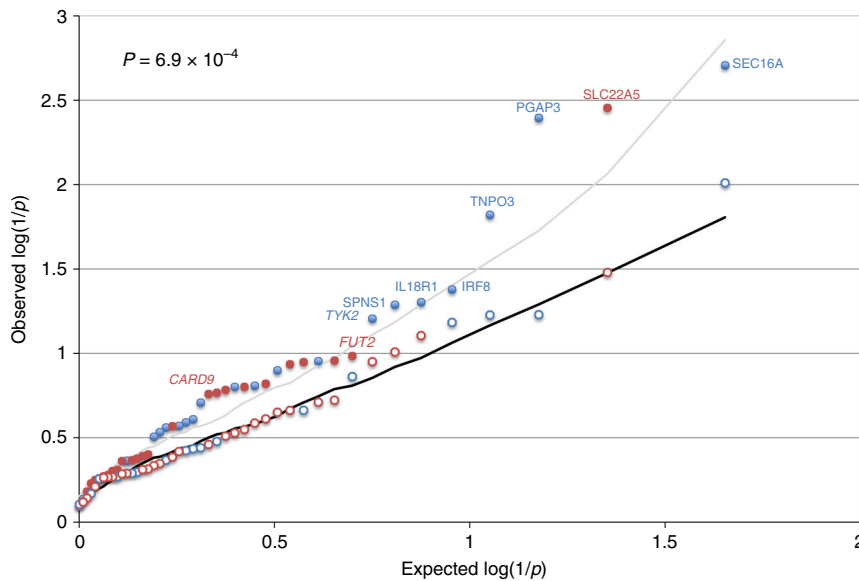
The variance in expression level due to the *cis*-eQTL equals:

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

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

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

To simulate *cis*-eQTL explaining the same  $h_{\text{eQTL}}^2$  as the real eQTL in the CEDAR dataset, we sequentially considered all SNPs in the region. Each one of these SNPs would be characterized by  $n_{11}$  individuals of genotype “11”,  $n_{12}$  of



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

genotype “12”, and  $n_{22}$  of genotype “22”, for a total of  $n_T$  genotyped individuals. We would arbitrarily set  $\bar{z}_{11}$ ,  $\bar{z}_{12}$ , and  $\bar{z}_{22}$  at  $-1$ ,  $0$ , and  $+1$ . As a consequence, the variance due to this *cis*-eQTL equals:

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

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

Knowing  $\sigma_{eQTL}^2$  and  $h_{eQTL}^2$ , and knowing that

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

the residual variance  $\sigma_{RES}^2$  can be computed as

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

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

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

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

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

**Testing for an excess sharing of cRM between cell types.** Assume that cell type 1 is part of  $n_{1T}$  cRM, including  $n_{11}$  private cRM,  $n_{12}$  cRM shared with cell type 2,  $n_{13}$  cRM shared with cell type 3, ..., and  $n_{19}$  cRM shared with cell type 9. Note that  $\sum_{i=1}^9 n_{1i} \geq n_{1T}$ , because cRM may include more than two cell types. Assume that  $n_{1S} = \sum_{i \neq 1}^9 n_{1i}$  is the sum of pair-wise sharing events for cell type 1. We computed, for each cell type  $i \neq 1$ , the probability to observe  $\geq n_{1i}$  sharing events with cell type 1

assuming that the expected number (under the hypothesis of random assortment) is

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

Pair-wise sharing events between tissue 1 and the eight other tissues were generated in silico under this model of random assortment (5000 simulations). The  $p$  value for  $n_{1i}$  was computed as the proportion of simulations that would yield values that would be as large or larger than  $n_{1i}$ . The same approach was used for the nine cell types. Thus, two  $p$  values of enrichment are obtained for each pair of cell types  $i$  and  $j$ , one using  $i$  as reference cell type, and the other using  $j$  as reference cell type. As can be seen from Fig. 5, the corresponding pairs of  $p$  values were always perfectly consistent.

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

**Comparing EAP and DAP using  $\vartheta$ .** The approach used to cluster EAP in cRM was also used to assign DAP for Inflammatory Bowel Disease (IBD) to EAP-defined cRM. We studied 200 IBD risk loci identified in recent GWAS meta-analyses<sup>2,3</sup>. The limits of the corresponding risk loci were as defined in the corresponding publications. We measured the similarity between DAP and EAP using the  $\vartheta$  metric for all *cis*-eQTL mapping to the corresponding intervals (i.e., for all *cis*-eQTL for which the top SNP mapped within the interval). To compute the correlations between DAP and EAP we used all SNPs mapping to the disease interval with  $-\log(p)$  value  $\geq 1.3$  either for DAP, EAP or both.

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

**Evaluating the enrichment of DAP-EAP matching.** To evaluate whether DAP matched EAP more often than expected by chance alone, we analyzed 97 IBD risk loci interrogated by the Immunochip, (i) in order to allow for convenient comparison with Huang et al.<sup>4</sup>, and (ii) because we needed extensively QC genotypes for the IIBDGC data to perform the enrichment analysis with the  $\vartheta$ -based method (see hereafter). Within these 97 IBD risk loci, we focused on 63 regions affecting CD<sup>4</sup>, encompassing at least one significant eQTL, and for which the lead CD-associated SNP had MAF  $> 0.05$ . Indeed, eQTL analyses in the CEDAR dataset



were restricted to SNPs with MAF > 0.05 (see above). We used three methods to evaluate whether the observed number of DAP–EAP matches were higher than expected by chance alone: naïve, frequentist and  $\vartheta$ -based. Analyses were performed separately for the nine cell types.

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

The frequentist approach used the method described by Nica et al.<sup>11</sup>. DAP and EAP were assumed to match if fitting the disease-associated lead SNP in the eQTL analysis caused a larger drop in  $-\log(p)$  than 95% of the SNPs with MAF > 0.05 in the analyzed risk locus. This would yield  $n_F \leq 63$  risk loci for which the DAP would match at least one EAP. To measure the statistical significance of  $n_F$ , we sampled a SNP (MAF > 0.05) at random in each of the 63 risk loci, and counted the number of loci with at least one matching EAP. This “simulation” was repeated 1000 times. The significance of  $n_F$  was measured as the proportion of simulations that would yield  $\geq n_F$  matches.

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

Assume a cohort with  $n_1$  cases and  $n_2$  controls (for instance, the IIBDGC cohort). Assume a SNP with an allelic frequency of  $p$  in the cases + controls, an allelic frequency of  $(p + \delta)$  in cases and  $(p - \delta)$  in controls.

One can easily show that:

$$\delta = -d \frac{n_1}{n_2} \tag{1}$$

The odds ratio (OR) for that SNP equals:

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

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

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

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

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

where

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

$$\beta = 0$$

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

Once we know  $(p + d)$  (i.e., the frequency of the SNP in cases), and hence  $(p - \delta)$  (i.e., the frequency of the SNP in controls), we can use Hardy–Weinberg to determine the frequency of the three genotypes in cases ( $p_{AA}^{IBD}, p_{AB}^{IBD}, p_{BB}^{IBD}$ ) and controls ( $p_{AA}^{CTR}, p_{AB}^{CTR}, p_{BB}^{CTR}$ ). We then create an in silico case–control cohort by sampling (with replacement)  $n_1 \times p_{AA}^{IBD}$  AA cases,  $n_1 \times p_{AB}^{IBD}$  AB cases, ..., and  $n_2 \times p_{BB}^{CTR}$  BB controls from the individuals of the IIBDGC (without discriminating real case and control status). Association analysis of the corresponding data set in the chromosome region of interest generates DAP with max  $-\log(p)$  value similar to the real DAP. This “simulation” was repeated 1000 times. The significance of  $n_\vartheta$  was measured as the proportion of simulations that would yield  $\geq n_\vartheta$  matches.

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

To perform burden tests, we collected DNA samples from 7323 Crohn Disease (CD) cases and 6342 controls of European descent in France (cases: 1899—ctrls: 1731), the Netherlands (2002–1923) and, Belgium (3422–2688). The study protocols were approved by the institutional review board at each center involved with recruitment. Informed consent and permission to share the data were obtained from all subjects, in compliance with the guidelines specified by the recruiting center’s institutional review board.

During the course of this project, we selected 45 genes with high  $|\vartheta|$  values for resequencing (Table 1). We designed primers to amplify all corresponding coding exons plus exon–intron boundaries corresponding to all transcripts reported in the CCDS release 15<sup>43</sup> (Supplementary Data 8). Following Momozawa et al.<sup>24</sup>, the primers were merged in five pools to perform a first round of PCR amplification (25 cycles). We then added 8-bp barcodes and common adapters (for sequencing) to all PCR products by performing a second round of PCR amplification (4 cycles) using primers targeting shared 5’ overhangs introduced during the first PCR. The ensuing libraries were purified, quality controlled and sequenced ( $2 \times 150$ -bp paired-end reads) on a HiSeq 2500 (Illumina) instrument. Sequence reads were sorted by individual using the barcodes, aligned to the human reference sequence (hg19) with the Burrows–Wheeler Aligner (ver. 0.7.12)<sup>44</sup>, and further processed using Genome Analysis Toolkit (GATK, ver. 3.2-2)<sup>45</sup>. We only considered individuals for further analyses if  $\geq 95\%$  of the target regions was covered by  $\geq 20$  sequence reads. Average sequence depth across individuals and target regions was 1060. We called variants for each individual separately using the UnifiedGenotyper and HaplotypeCaller of GATK, as well as VMM (ver. 1.0.2)<sup>46</sup>, and listed all variants detected by either method. Genotypes for all individuals were determined for each variant based on the ratio of reference and alternative alleles amongst sequence reads as determined by Samtools<sup>47</sup>. Individuals were labeled homozygote reference, heterozygote, or homozygote derived when the alternative allele frequency was between 0 and 0.15, between 0.25 and 0.75, and between 0.85 and 1, respectively. If the alternative allele frequency was outside these ranges or a variant position was covered with  $< 20$  sequencing reads, the genotype was considered missing. We excluded variants with call rates  $< 95\%$  or variants that were not in Hardy–Weinberg equilibrium ( $P < 1 \times 10^{-6}$ ). We excluded 281 individuals with  $\geq 2$  minor alleles at 23 variants selected to have a MAF  $\leq 0.01$  in non-Finnish Europeans and  $\geq 0.10$  in Africans or East-Asians in the Exome Aggregation Consortium<sup>27</sup>.

In the end, we used 6597 cases and 5502 controls for further analyses, while 98.5% of the target regions on average was covered with 20 or more sequence reads.

**Gene-based burden test.** We first used SIFT<sup>25</sup> and Polyphen-2<sup>26</sup> to sort the 4175 variants identified by sequencing in four categories: (i) loss-of-function (LoF) or severe, corresponding to stop gain, stop loss, frameshift and splice-site variants, (ii) damaging, corresponding to missense variants predicted by SIFT to be damaging and Polyphen-2 to be possibly or probably damaging, (iii) benign, corresponding to the other missense variants, and (iv) synonymous. We performed the burden test using the LoF plus damaging variants, and used the synonymous variants as controls. We only considered variants with MAF (computed for the entire data set, i.e., cases plus controls)  $\leq 0.005$ . We indeed showed in a previous fine-mapping study that all reported independent effects were driven by variants with MAF  $\geq 0.01$ <sup>4</sup>. By doing so we ensure that the signals of the burden test are independent of previously reported association signals. Thus, 174 LoF, 991 damaging, and 1434 synonymous were ultimately used to perform burden tests.

Burden tests come in two main flavors. In the first, one assumes that disruptive variants will be enriched in either cases (i.e., disruptive variants increase risk) or in controls (i.e., disruptive variance decrease risk). In the second, one assumes that—for a given gene—some disruptive variants will be enriched in cases, while other may be enriched in controls (Supplementary Fig. 11). The first was implemented using CAST<sup>28</sup>. To increase power, we exploited the DAP–EAP information to perform one-sided (rather than two-sided) tests. When  $\vartheta < 0$ , we tested for an enrichment of disruptive variants in cases; when  $\vartheta > 0$ , for an enrichment of disruptive variants in controls.  $P$  values were computed by phenotype permutation, i.e., shuffling case–control status. When applying this test on a gene-by-gene basis using synonymous variants (MAF > 0.005), the distribution of  $p$  values (QQ-plot) indicated that the CAST test was conservative ( $\lambda_{GC} = 0.51$ ) (Supplementary Fig. 12). The second kind of burden test was implemented with SKAT<sup>29</sup>. It is noteworthy that SKAT ignores information from singletons (Supplementary Fig. 11). Just as for CAST,  $p$  values were computed by phenotype permutation, i.e., shuffling case–control status. When applying this test on a gene-by-gene basis using synonymous variants (MAF < 0.005), the distribution of  $p$  values (QQ-plot) indicated that the SKAT test is too permissive ( $\lambda_{GC} = 1.73$ ) (Supplementary Fig. 12). Consequently, gene-based  $p$  values obtained with SKAT were systematically GC corrected using this value of  $\lambda_{GC}$ . We performed the two kinds of



analyses for each gene, as one doesn't a priori know what hypothesis will match the reality best for a given gene.

We also extracted information from the distribution of  $p$  values (or  $-\log(p)$  values) across the 45 analyzed genes. Even if individual genes do not yield  $-\log(p)$  values that exceed the significance threshold (accounting for the number of analyzed genes and tests performed), the distribution of  $-\log(p)$  values may significantly depart from expectations, indicating that the analyzed genes include at least some causative genes. This was done by taking for each gene, the best  $p$  value (whether obtained with CAST or SKAT) and then rank the genes by corresponding  $-\log(p)$  value. The same was done for  $10^5$  phenotype permutations, allowing us to examine the distribution of  $-\log(p)$  values for given ranks and compute the corresponding medians and limits of the 95% confidence band, as well as to compute the probability that  $-2 \sum_{i=1}^{45} \ln(p_i)$  (Fisher's equation to combine  $p$  values) equals or exceeds the observed. Our results show that there is a significant departure from expectation when analyzing the damaging variants ( $p = 6.9 \times 10^{-4}$ ) but not when analyzing the synonymous variants ( $p = 0.66$ ) supporting the presence of genuine causative genes amongst the analyzed list.

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

**Orthogonal tests for age-of-onset and familiarity.** It is commonly assumed that the heritability for common complex diseases is higher in familial and early onset cases<sup>31</sup>. To extract the corresponding information from our data in a manner that would be orthogonal to the gene- and module-based tests described above (i.e., the information about age-of-onset and familiarity would be independent of these burden tests), we devised the following approach.

For age-of-onset, we summed the age-of-onset of the  $n_C$  cases carrying rare disruptive variants for the gene of interest. We then computed the probability that the sum of the age-of-onset of  $n_C$  randomly chosen cases was as different from the mean of age-of-onset as the observed one, yielding a gene-specific two-sided  $p_{SKAT}$  value. In addition, we used the eQTL information to generate gene-specific one-sided  $p_{CAST}$  values, corresponding to the probability that the sum of the age-of-onset of  $n_C$  randomly chosen cases was as low or lower than the observed one (for genes for which decrease in expression level as associated with increased risk), or to the probability that the sum of the age-of-onset of  $n_C$  randomly chosen cases was as high or higher than the observed one (for genes for which increase in expression level as associated with increased risk). These age-of-onset  $p$  values were then combined with the corresponding  $p$  values from the burden test (CAST with CAST, SKAT with SKAT) using Fisher's method.

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

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

Received: 1 September 2017 Accepted: 24 April 2018

Published online: 21 June 2018

## References

- MacArthur, J. et al. The new NHGRI-EBI catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Res.* **45**, D896–D901 (2017).
- Jostins, L. et al. Host-microbe interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature* **491**, 119–124 (2012).
- Liu, J. Z. et al. Association analyses identify 38 susceptibility loci for IBD and highlight shared genetic risk across populations. *Nat. Genet.* **47**, 979–986 (2015).
- Huang, H. et al. Association mapping of IBD loci to single variant resolution. *Nature* **547**, 173–178 (2017).
- Claussnitzer, M. et al. FTO obesity variant circuitry and adipocyte browning in humans. *N. Engl. J. Med.* **373**, 895–907 (2015).
- Hugot, J. P. et al. Association of NOD2 leucine-rich repeat variants with susceptibility to Crohn's disease. *Nature* **411**, 599–603 (2001).
- Hampe, J. et al. A genome-wide association scan of nonsynonymous SNPs identifies a susceptibility variant for Crohn disease in ATG16L1. *Nat. Genet.* **39**, 207–211 (2007).
- Momozawa, Y. et al. Resequencing of positional candidates identifies low frequency IL23R coding variants protecting against inflammatory bowel disease. *Nat. Genet.* **43**, 43–47 (2011).
- Rivas, M. A. et al. Deep resequencing of GWAS loci identifies independent rare variants associated with inflammatory bowel disease. *Nat. Genet.* **43**, 1066–1073 (2011).
- The GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).
- Nica, A. C. et al. Candidate causal regulatory effects by integration of expression QTLs with complex trait genetic associations. *PLoS Genet.* **6**, e1000895 (2010).
- Giambartolomei, C. et al. Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
- Zhu, Z. et al. Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* **48**, 481–487 (2016).
- The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- Nicolae, D. L. Association tests for rare variants. *Annu Rev. Genom. Hum. Genet.* **17**, 117–130 (2016).
- Pritchard, J. K. & Cox, N. J. The allelic architecture of human disease genes: common disease-common variant ... or not? *Hum. Mol. Genet.* **11**, 2417–2423 (2002).
- McGregor, A. P. et al. Morphological evolution through multiple cis-regulatory mutations at a single gene. *Nature* **448**, 587–590 (2007).
- Mackay, T. F. Quantitative trait loci in *Drosophila*. *Nat. Rev. Genet.* **2**, 11–20 (2001).
- Yalcin, B. et al. Genetic dissection of behavioral QTL shows that Rgs2 modulates anxiety in mice. *Nat. Genet.* **36**, 1197–1202 (2004).
- Karim, L. et al. Variants modulating the expression of a chromosome domain encompassing PLAG1 influence bovine stature. *Nat. Genet.* **43**, 405–413 (2011).
- Steinmetz, L. M. et al. Dissecting the architecture of a QTL in yeast. *Nature* **416**, 326–330 (2002).
- Khor, B., Gardet, A. & Xavier, R. Genetics and pathogenesis of inflammatory bowel disease. *Nature* **474**, 307–317 (2011).
- Fuchsberger, C. et al. The genetic architecture of type 2 diabetes. *Nature* **536**, 41–47 (2016).
- Momozawa, Y. et al. Low-frequency coding variants in CETP and CFB are associated with susceptibility of exudative age-related macular degeneration in the Japanese population. *Hum. Mol. Genet.* **25**, 5027–5034 (2016).
- Kumar, P. et al. Predicting the effects of coding non-synonymous variants on protein function using gthe SIFT algorithm. *Nat. Protoc.* **4**, 1073–1081 (2009).
- Adzhubei, I. A. et al. A method and server for predicting damaging missense mutations. *Nat. Methods* **7**, 248–249 (2010).
- Lek, M. et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
- Morgenthaler, S. & Thilly, W. G. A strategy to discover genes that carry multi-allelic or mono-allelic risk for common diseases: a cohort allelic sums test (CAST). *Mutat. Res.* **615**, 28–56 (2007).
- Wu, M. C. et al. Rare-variant association testing for sequencing data with the sequence kernel association test. *Am. J. Hum. Genet.* **89**, 82–93 (2011).
- Richardson, T. G. et al. A pathway-centric approach to rare variant association analysis. *Eur. J. Hum. Genet.* **25**, 123–129 (2017).
- Imielinski, M. et al. Common variants at five new loci associated with early-onset inflammatory bowel disease. *Nat. Genet.* **41**, 1335–1340 (2009).
- Chun, S. et al. Limited statistical evidence for shared genetic effects of eQTLs and autoimmune-disease-associated loci in three major immune-cell types. *Nat. Genet.* **4**, 600–605 (2017).
- The GTEx Consortium. The genotype-tissue expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
- Boyle, E. A. et al. An expanded view of complex traits: from polygenic to omnigenic. *Cell* **169**, 1177–1186 (2017).
- The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 668–674 (2015).

36. Huang, J. et al. Improved imputation of low-frequency and rare variants using the UK10K haplotype reference panel. *Nat. Commun.* **6**, 8111 (2015).
37. McCarthy et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).
38. Du, P. et al. Lumi: a pipeline for processing illumine microarray. *Bioinformatics* **24**, 1547–1548 (2008).
39. Arloth, J. et al. Re-Annotator: annotation pipeline for microarray probe sequences. *PLoS ONE* **10**, e0139516 (2015).
40. Johnson, W. E. et al. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* **8**, 118–127 (2007).
41. Fairfax, B. P. et al. Innate immune activity conditions the effect of regulatory variants upon monocyte gene expression. *Science* **343**, 1246949 (2014).
42. Purcell, S. et al. PLINK: a toolset for whole-genome association and population-based linkage analysis. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
43. Farrell, C. M. et al. Current status and new features of the Consensus Coding Sequence database. *Nucleic Acids Res.* **42**, D865–72 (2014).
44. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–60 (2009).
45. DePristo, M. A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–8 (2011).
46. Shigemizu, D. et al. A practical method to detect SNVs and indels from whole genome and exome sequencing data. *Sci. Rep.* **3**, 2161 (2013).
47. Li, H. et al. The sequence alignment/map (SAM) format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
48. Whitehead Pavlides, J. M. et al. Predicting targets from integrative analyses of summary data from GWAS and eQTL studies for 28 human complex traits. *Genome Med.* **8**, 84–90 (2016).
49. Gamazon, E. R. et al. A gene-based association method for mapping traits using reference transcriptome data. *Nat. Genet.* **47**, 1091–1098 (2015).
50. Hulur, I. et al. Enrichment of inflammatory bowel disease and colorectal cancer risk variants in colon expression quantitative trait loci. *BMC Genom.* **16**, 138–153 (2015).
51. Libioulle, C. et al. Novel Crohn disease locus identified by genome-wide association maps to a gene desert on 5p13.1 and modulates expression of PTGER4. *PLoS Genet* **3**, e58 (2007).
52. Peltekova, V. D. et al. Functional variants of OCTN cation transporter genes are associated with Crohn disease. *Nat. Genet* **39**, 311–318 (2004).
53. McCarroll, S. A. et al. Deletion polymorphism upstream of IRGM expression and Crohn's disease. *Nat. Genet* **40**, 1107–112 (2008).
54. De Lange, K. M. et al. Genome-wide association study implicates immune activation of multiple integrin genes in inflammatory bowel disease. *Nat. Genet.* **49**, 256–261 (2017).

## Acknowledgements

This work was supported by grants to Michel Georges from WELBIO (CAUSIBD), BELSPO (BeMGI), and Horizon 2020 (SYSCID). Computational resources at ULg have been provided by GIGA and the Consortium des Équipements de Calcul Intensif (CÉCI), funded by the Fonds de la Recherche Scientifique de Belgique (F.R.S.-FNRS) under Grant No. 2.5020.11. This work was conducted as part of the BioBank Japan Project supported by the Japan Agency for Medical Research and Development and by the Ministry of Education, Culture, Sports, Sciences and Technology of the Japanese government. The

work of D.A. and I.A. was supported by Russian Ministry of Science and Education under 5-100 Excellence Programme. R.K. Weersma is supported by a VIDI grant (016.136.308) from the Netherlands Organisation for Scientific Research (NWO). DNA samples from the Dutch IBD cohort have been collected within the Parelnoer Institute Project. This nationwide Parelnoer Institute project is part of and funded by the Netherlands Federation of University Medical Centres and has received initial funding from the Dutch Government (from 2007 to 2011). The Parelnoer Institute currently facilitates the uniform nationwide collection of information on and biomaterials of thirteen other diseases. We are grateful to N. Hakozaki, H. Iijima, N. Maki, and other staff of the Laboratory for Genotyping Development, RIKEN Center for the Integrative Medical Sciences. We thank Wouter Coppieters and the other members of the GIGA genomics platform for their support.

## Author contributions

Y.M., J.D., and M.G. conceived experiments, generated data, analyzed data and wrote the manuscript. E.T., V.D., S.R., B.C., F.C., E.D., M.E., A.-S.G., C.L., R.M., M.M., and C.O. generated and analyzed data. I.A., D.A., Y.A., and M.G. conceived and generated the CEDAR website. L.A., G.B., F.H., M.L., B.O., M.J.P., A.E.v.d.M.-d.J., J.J.v.d.W., M. C.V., M.L., J.-P.H., R.K.W., M.D.V., D.F., S.V., M.K., and E.L. collected and provided samples. The IIBDGC provided association *p*-values in the 201 IBD risk loci (DAPs).

## Additional information

**Supplementary Information** accompanies this paper at <https://doi.org/10.1038/s41467-018-04365-8>.

**Competing interests:** The authors declare no competing interests.

**Reprints and permission** information is available online at <http://npg.nature.com/reprintsandpermissions/>

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018

Yukihide Momozawa<sup>1,2</sup>, Julia Dmitrieva<sup>1</sup>, Emilie Théâtre<sup>1</sup>, Valérie Deffontaine<sup>1</sup>, Souad Rahmouni<sup>1</sup>, Benoît Charlotiaux<sup>1</sup>, François Crins<sup>1</sup>, Elisa Docampo<sup>1</sup>, Mahmoud Elansary<sup>1</sup>, Ann-Stephan Gori<sup>1</sup>, Christelle Lecut<sup>3</sup>, Rob Mariman<sup>1</sup>, Myriam Mni<sup>1</sup>, Cécile Oury<sup>3</sup>, Ilya Altukhov<sup>4</sup>, Dmitry Alexeev<sup>5</sup>, Yuri Aulchenko<sup>6,7,8</sup>, Leila Amininejad<sup>9</sup>, Gerd Bouma<sup>10</sup>, Frank Hoentjen<sup>11</sup>, Mark Löwenberg<sup>12</sup>, Bas Oldenburg<sup>13</sup>, Marieke J. Pierik<sup>14</sup>, Andrea E. vander Meulen-de Jong<sup>15</sup>, C. Janneke van der Woude<sup>16</sup>, Marijn C. Visschedijk<sup>17</sup>, The International IBD Genetics Consortium, Mark Lathrop<sup>18</sup>, Jean-Pierre Hugot<sup>19</sup>, Rinse K. Weersma<sup>17</sup>, Martine De Vos<sup>20</sup>, Denis Franchimont<sup>9</sup>, Severine Vermeire<sup>21</sup>, Michiaki Kubo<sup>2</sup>, Edouard Louis<sup>22</sup> & Michel Georges<sup>1</sup>

<sup>1</sup>Unit of Animal Genomics, WELBIO, GIGA-R & Faculty of Veterinary Medicine, University of Liège (B34), 1 Avenue de l'Hôpital, Liège 4000, Belgium. <sup>2</sup>Laboratory for Genotyping Development, RIKEN Center for Integrative Medical Science, 1-7-22, Suehiro-cho, Tsurumi-ku, Yokohama, Kanagawa 230-0045, Japan. <sup>3</sup>Laboratory of Thrombosis and Hemostasis, GIGA-R, University of Liège (B34), 1 Avenue de l'Hôpital, 4000 Liège, Belgium. <sup>4</sup>Moscow Institute of Physics and Technology, Institutskiy Pereulok 9, Dolgoprudny 141700, Russian Federation. <sup>5</sup>Novosibirsk State University, Pirogova ave. 2, Novosibirsk 630090, Russian Federation. <sup>6</sup>PolyOmica, Het Vlaggeschip 61, 's-Hertogenbosch 5237 PA, The Netherlands. <sup>7</sup>Institute of Cytology and Genetics SD RAS, Lavrentyeva ave. 10, 630090 Novosibirsk, Russia. <sup>8</sup>Centre for Global Health Research,

Usher Institute of Population Health Sciences and Informatics, University of Edinburgh, Teviot Place, Edinburgh EH8 9AG, UK. <sup>9</sup>Gastroentérologie Médicale, Faculté de Médecine, Université Libre de Bruxelles, Route de Lennik 808, Anderlecht 1070, Belgium. <sup>10</sup>Department of Gastroenterology and Hepatology, VU University Medical Centre, Amsterdam 1081 HV, The Netherlands. <sup>11</sup>Department of Gastroenterology and Hepatology, University Medical Centre St. Radboud, Nijmegen 6525 GA, The Netherlands. <sup>12</sup>Department of Gastroenterology and Hepatology, Amsterdam Medical Centre, Amsterdam 1105 AZ, The Netherlands. <sup>13</sup>Department of Gastroenterology and Hepatology, University Medical Centre Utrecht, 3584 cXUtrecht, The Netherlands. <sup>14</sup>Department of Gastroenterology and Hepatology, University Medical Centre Maastricht, Maastricht 6229 HX, The Netherlands. <sup>15</sup>Department of Gastroenterology and Hepatology, Leiden University Medical Centre, Leiden 2333 ZA, The Netherlands. <sup>16</sup>Department of Gastroenterology and Hepatology, Erasmus Medical Centre, Rotterdam 3015 CE, The Netherlands. <sup>17</sup>Department of Gastroenterology and Hepatology, University of Groningen and University Medical Center Groningen, Hanzeplein 1, Groningen 9713 GZ, The Netherlands. <sup>18</sup>McGill University Centre for Molecular and Computational Genomics, 740 Dr. Penfield Avenue, Montreal H3A 0G1 QC, Canada. <sup>19</sup>UMR 1149 INSERM/Université Paris-Diderot Sorbonne Paris-Cité, Assistance Publique Hôpitaux de Paris, 48 Bd Sérurier, Paris 75019, France. <sup>20</sup>Department of Gastroenterology, University Hospital, De Pintelaan 185, Gent 9000, Belgium. <sup>21</sup>Translational Research in Gastrointestinal Disorders, Department of Clinical and Experimental Medicine, KU Leuven, UZ Herestraat 49, Leuven 3000, Belgium. <sup>22</sup>CHU-Liège and Unit of Gastroenterology, GIGA-R & Faculty of Medicine, University of Liège, 1 Avenue de l'Hôpital, Liège 4000, Belgium. These authors contributed equally: Yukihide Momozawa, Julia Dmitrieva. <sup>†</sup>A list of The International IBD Genetics Consortium members is provided below.

## The International IBD Genetics Consortium

Clara Abraham<sup>23</sup>, Jean-Paul Achkar<sup>24,25</sup>, Tariq Ahmad<sup>26</sup>, Ashwin N. Ananthakrishnan<sup>27,28</sup>, Vibeke Andersen<sup>29,30,31</sup>, Carl A. Anderson<sup>32</sup>, Jane M. Andrews<sup>33</sup>, Vito Annese<sup>34,35</sup>, Guy Aumais<sup>36,37</sup>, Leonard Baidoo<sup>38</sup>, Robert N. Baldassano<sup>39</sup>, Peter A. Bampton<sup>40</sup>, Murray Barclay<sup>41</sup>, Jeffrey C. Barrett<sup>32</sup>, Theodore M. Bayless<sup>42</sup>, Johannes Bethge<sup>43</sup>, Alain Bitton<sup>44</sup>, Gabrielle Boucher<sup>45</sup>, Stephan Brand<sup>46</sup>, Berenice Brandt<sup>43</sup>, Steven R. Brant<sup>42</sup>, Carsten Büning<sup>47</sup>, Angela Chew<sup>48,49</sup>, Judy H. Cho<sup>50</sup>, Isabelle Cleynen<sup>21</sup>, Ariella Cohain<sup>51</sup>, Anthony Croft<sup>52</sup>, Mark J. Daly<sup>53,54</sup>, Mauro D'Amato<sup>55,56,57</sup>, Silvio Danese<sup>58</sup>, Dirk De Jong<sup>11</sup>, Goda Denapiene<sup>59</sup>, Lee A. Denson<sup>60</sup>, Kathy L. Devaney<sup>27</sup>, Olivier Dewit<sup>61</sup>, Renata D'Inca<sup>62</sup>, Marla Dubinsky<sup>63</sup>, Richard H. Duerr<sup>38,64</sup>, Cathryn Edwards<sup>65</sup>, David Ellinghaus<sup>66</sup>, Jonah Essers<sup>67,68</sup>, Lynnette R. Ferguson<sup>69</sup>, Eleonora A. Festen<sup>17</sup>, Philip Fleshner<sup>70</sup>, Tim Florin<sup>71</sup>, Andre Franke<sup>66</sup>, Karin Fransen<sup>72</sup>, Richard Geary<sup>41,73</sup>, Christian Gieger<sup>74</sup>, Jürgen Glas<sup>46,75</sup>, Philippe Goyette<sup>45</sup>, Todd Green<sup>54,67</sup>, Anne M. Griffiths<sup>76</sup>, Stephen L. Guthery<sup>77</sup>, Hakon Hakonarson<sup>78</sup>, Jonas Halfvarson<sup>78</sup>, Katherine Hanigan<sup>52</sup>, Talin Haritunians<sup>70</sup>, Ailsa Hart<sup>79</sup>, Chris Hawkey<sup>80</sup>, Nicholas K. Hayward<sup>81</sup>, Matija Hedl<sup>23</sup>, Paul Henderson<sup>82,83</sup>, Xinli Hu<sup>84</sup>, Hailiang Huang<sup>53,54</sup>, Ken Y. Hui<sup>50</sup>, Marcin Imielinski<sup>39</sup>, Andrew Ippoliti<sup>70</sup>, Laimas Jonaitis<sup>85</sup>, Luke Jostins<sup>86,87</sup>, Tom H. Karlsen<sup>88,89,90</sup>, Nicholas A. Kennedy<sup>91</sup>, Mohammed Azam Khan<sup>92,93</sup>, Gediminas Kiudelis<sup>85</sup>, Krupa Krishnaprasad<sup>94</sup>, Subra Kugathasan<sup>95</sup>, Limas Kupcinskis<sup>96</sup>, Anna Latiano<sup>34</sup>, Debby Laukens<sup>20</sup>, Ian C. Lawrance<sup>48,97</sup>, James C. Lee<sup>98</sup>, Charlie W. Lees<sup>91</sup>, Marcis Leja<sup>99</sup>, Johan Van Limbergen<sup>76</sup>, Paolo Lionetti<sup>100</sup>, Jimmy Z. Liu<sup>32</sup>, Gillian Mahy<sup>101</sup>, John Mansfield<sup>102</sup>, Dunecan Massey<sup>98</sup>, Christopher G. Mathew<sup>103,104</sup>, Dermot P.B. McGovern<sup>70</sup>, Raquel Milgrom<sup>105</sup>, Mitja Mitrovic<sup>72,106</sup>, Grant W. Montgomery<sup>81</sup>, Craig Mowat<sup>107</sup>, William Newman<sup>92,93</sup>, Aylwin Ng<sup>27,108</sup>, Siew C. Ng<sup>109</sup>, Sok Meng Evelyn Ng<sup>23</sup>, Susanna Nikolaus<sup>43</sup>, Kaida Ning<sup>23</sup>, Markus Nöthen<sup>110</sup>, Ioannis Oikonomou<sup>23</sup>, Orazio Palmieri<sup>34</sup>, Miles Parkes<sup>98</sup>, Anne Phillips<sup>107</sup>, Cyriel Y. Ponsioen<sup>12</sup>, Urös Potocnik<sup>106,111</sup>, Natalie J. Prescott<sup>103</sup>, Deborah D. Proctor<sup>23</sup>, Graham Radford-Smith<sup>52,112</sup>, Jean-Francois Rahier<sup>113</sup>, Soumya Raychaudhuri<sup>84</sup>, Miguel Regueiro<sup>38</sup>, Florian Rieder<sup>24</sup>, John D. Rioux<sup>36,45</sup>, Stephan Ripke<sup>53,54</sup>, Rebecca Roberts<sup>41</sup>, Richard K. Russell<sup>82</sup>, Jeremy D. Sanderson<sup>114</sup>, Miquel Sans<sup>115</sup>, Jack Satsangi<sup>91</sup>, Eric E. Schadt<sup>51</sup>, Stefan Schreiber<sup>43,66</sup>, Dominik Schulte<sup>43</sup>, L. Philip Schumm<sup>116</sup>, Regan Scott<sup>38</sup>, Mark Seielstad<sup>117,118</sup>, Yashoda Sharma<sup>23</sup>, Mark S. Silverberg<sup>105</sup>, Lisa A. Simms<sup>52</sup>, Jurgita Skieceviciene<sup>85</sup>, Sarah L. Spain<sup>32,119</sup>, A. Hillary Steinhart<sup>105</sup>, Joanne M. Stempak<sup>105</sup>, Laura Stronati<sup>120</sup>, Jurgita Sventoraityte<sup>94</sup>, Stephan R. Targan<sup>70</sup>, Kirstin M. Taylor<sup>114</sup>, Anje ter Velde<sup>12</sup>, Leif Torkvist<sup>121</sup>, Mark Tremelling<sup>122</sup>, Suzanne van Sommeren<sup>17</sup>, Eric Vasiliauskas<sup>70</sup>, Hein W. Verspaget<sup>15</sup>, Thomas Walters<sup>76,123</sup>, Kai Wang<sup>39</sup>, Ming-Hsi Wang<sup>24,42</sup>, Zhi Wei<sup>124</sup>, David Whiteman<sup>81</sup>, Cisca Wijmenga<sup>72</sup>, David C. Wilson<sup>82,83</sup>, Juliane Winkelmann<sup>125,126</sup>, Ramnik J. Xavier<sup>27,54</sup>, Bin Zhang<sup>51</sup>, Clarence K. Zhang<sup>127</sup>, Hu Zhang<sup>128,129</sup>, Wei Zhang<sup>23</sup>, Hongyu Zhao<sup>127</sup> & Zhen Z. Zhao<sup>81</sup>



<sup>23</sup>Section of Digestive Diseases, Department of Internal Medicine, Yale School of Medicine, New Haven, CT, USA. <sup>24</sup>Department of Gastroenterology and Hepatology, Digestive Disease Institute, Cleveland Clinic, Cleveland, OH, USA. <sup>25</sup>Department of Pathobiology, Lerner Research Institute, Cleveland Clinic, Cleveland, OH, USA. <sup>26</sup>Peninsula College of Medicine and Dentistry, Exeter, UK. <sup>27</sup>Gastroenterology Unit, Massachusetts General Hospital, Harvard Medical School, Boston, MA 02114, USA. <sup>28</sup>Division of Medical Sciences, Harvard Medical School, Boston, MA, USA. <sup>29</sup>Focused Research Unit for Molecular Diagnostic and Clinical Research (MOK), IRS-Center Sonderjylland, Hospital of Southern Jutland, Åbenrå 6200, Denmark. <sup>30</sup>Institute of Molecular Medicine, University of Southern Denmark, Odense 5000, Denmark. <sup>31</sup>Institute of Regional Health Research, University of Southern Denmark, Odense, Denmark. <sup>32</sup>Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10 1SA, UK. <sup>33</sup>Inflammatory Bowel Disease Service, Department of Gastroenterology and Hepatology, Royal Adelaide Hospital, Adelaide, Australia. <sup>34</sup>Unit of Gastroenterology, Istituto di Ricovero e Cura a Carattere Scientifico-Casa Sollievo della Sofferenza (IRCCS-CSS) Hospital, San Giovanni Rotondo, Italy. <sup>35</sup>Strutture Organizzative Dipartimentali (SOD) Gastroenterologia 2, Azienda Ospedaliero Universitaria (AOU) Careggi, Florence, Italy. <sup>36</sup>Facult de Médecine, Université de Montréal, Montréal, QC H3C 3J7, Canada. <sup>37</sup>Department of Gastroenterology, Hôpital Maisonneuve-Rosemont, Montréal, QC, Canada. <sup>38</sup>Division of Gastroenterology, Hepatology and Nutrition, Department of Medicine, University of Pittsburgh School of Medicine, Pittsburgh, PA 15213, USA. <sup>39</sup>Center for Applied Genomics, The Children's Hospital of Philadelphia, Philadelphia, PA, USA. <sup>40</sup>Department of Gastroenterology and Hepatology, Flinders Medical Centre and School of Medicine, Flinders University, Adelaide, Australia. <sup>41</sup>Department of Medicine, University of Otago, Christchurch, New Zealand. <sup>42</sup>Meyerhoff Inflammatory Bowel Disease Center, Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA. <sup>43</sup>Department for General Internal Medicine, Christian-Albrechts-University, Kiel, Germany. <sup>44</sup>Division of Gastroenterology, Royal Victoria Hospital, Montréal, QC, Canada. <sup>45</sup>Research Center, Montreal Heart Institute, Montréal, QC H1T 1C8, Canada. <sup>46</sup>Department of Medicine II, Ludwig-Maximilians-University Hospital Munich-Grosshadern, Munich, Germany. <sup>47</sup>Department of Gastroenterology, Campus Charité Mitte, Universitätsmedizin Berlin, Berlin, Germany. <sup>48</sup>Harry Perkins Institute for Medical Research, School of Medicine and Pharmacology, University of Western Australia, Murdoch, WA 6150, Australia. <sup>49</sup>IBD Unit, Fremantle Hospital, Fremantle, Australia. <sup>50</sup>Department of Genetics, Yale School of Medicine, New Haven, CT 06510, USA. <sup>51</sup>Department of Genetics and Genomic Sciences, Mount Sinai School of Medicine, New York, NY, USA. <sup>52</sup>Inflammatory Bowel Diseases, Genetics and Computational Biology, Queensland Institute of Medical Research, Brisbane, Australia. <sup>53</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Harvard Medical School, Boston, MA 02114, USA. <sup>54</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02141, USA. <sup>55</sup>Clinical Epidemiology Unit, Department of Medicine Solna, Karolinska Institutet, Stockholm 17176, Sweden. <sup>56</sup>Department of Gastrointestinal and Liver Diseases, BioDonostia Health Research Institute, San Sebastián 20014, Spain. <sup>57</sup>IKERBASQUE, Basque Foundation for Science, Bilbao 48013, Spain. <sup>58</sup>IBD Center, Department of Gastroenterology, Istituto Clinico Humanitas, Milan, Italy. <sup>59</sup>Center of Hepatology, Gastroenterology and Dietetics, Vilnius University, Vilnius, Lithuania. <sup>60</sup>Pediatric Gastroenterology, Cincinnati Children's Hospital Medical Center, Cincinnati, OH, USA. <sup>61</sup>Department of Gastroenterology, Université Catholique de Louvain (UCL) Cliniques Universitaires Saint-Luc, Brussels, Belgium. <sup>62</sup>Division of Gastroenterology, University Hospital Padua, Padua, Italy. <sup>63</sup>Department of Pediatrics, Cedars Sinai Medical Center, Los Angeles, CA, USA. <sup>64</sup>Department of Human Genetics, University of Pittsburgh Graduate School of Public Health, Pittsburgh, PA 15261, USA. <sup>65</sup>Department of Gastroenterology, Torbay Hospital, Torbay, Devon, UK. <sup>66</sup>Institute of Clinical Molecular Biology, Christian-Albrechts-University of Kiel, Kiel 24118, Germany. <sup>67</sup>Center for Human Genetic Research, Massachusetts General Hospital, Harvard Medical School, Boston, MA, USA. <sup>68</sup>Pediatrics, Harvard Medical School, Boston, MA, USA. <sup>69</sup>Faculty of Medical & Health Sciences, School of Medical Sciences, The University of Auckland, Auckland, New Zealand. <sup>70</sup>F. Widjaja Foundation Inflammatory Bowel and Immunobiology Research Institute, Cedars-Sinai Medical Center, Los Angeles, CA 90048, USA. <sup>71</sup>Department of Gastroenterology, Mater Health Services, Brisbane, Australia. <sup>72</sup>Department of Genetics, University Medical Center Groningen, Groningen, The Netherlands. <sup>73</sup>Department of Gastroenterology, Christchurch Hospital, Christchurch, New Zealand. <sup>74</sup>Institute of Genetic Epidemiology, Helmholtz Zentrum München—German Research Center for Environmental Health, Neuherberg, Germany. <sup>75</sup>Department of Preventive Dentistry and Periodontology, Ludwig-Maximilians-University Hospital Munich-Grosshadern, Munich, Germany. <sup>76</sup>Division of Pediatric Gastroenterology, Hepatology and Nutrition, Hospital for Sick Children, Toronto, ON, Canada. <sup>77</sup>Department of Pediatrics, University of Utah School of Medicine, Salt Lake City, UT, USA. <sup>78</sup>Department of Gastroenterology, Faculty of Medicine and Health, Örebro University, SE-70182 Örebro, Sweden. <sup>79</sup>Department of Medicine, St. Mark's Hospital, Harrow, Middlesex, UK. <sup>80</sup>Nottingham Digestive Diseases Centre, Queens Medical Centre, Nottingham, UK. <sup>81</sup>Molecular Epidemiology, Genetics and Computational Biology, Queensland Institute of Medical Research, Brisbane, Australia. <sup>82</sup>Paediatric Gastroenterology and Nutrition, Royal Hospital for Sick Children, Edinburgh, UK. <sup>83</sup>Child Life and Health, University of Edinburgh, Edinburgh, Scotland, UK. <sup>84</sup>Division of Rheumatology Immunology and Allergy, Brigham and Women's Hospital, Boston, MA, USA. <sup>85</sup>Academy of Medicine, Lithuanian University of Health Sciences, Kaunas, Lithuania. <sup>86</sup>Wellcome Trust Centre for Human Genetics, University of Oxford, Headington OX3 7BN, UK. <sup>87</sup>Christ Church, University of Oxford, St Aldates OX1 1DP, UK. <sup>88</sup>Research Institute of Internal Medicine, Department of Transplantation Medicine, Division of Cancer, Surgery and Transplantation, Oslo University Hospital Rikshospitalet, Oslo, Norway. <sup>89</sup>Norwegian PSC Research Center, Department of Transplantation Medicine, Division of Cancer, Surgery and Transplantation, Oslo University Hospital Rikshospitalet, Oslo, Norway. <sup>90</sup>K.G. Jebsen Inflammation Research Centre, Institute of Clinical Medicine, University of Oslo, Oslo, Norway. <sup>91</sup>Gastrointestinal Unit, Western General Hospital University of Edinburgh, Edinburgh, UK. <sup>92</sup>Genetic Medicine, Manchester Academic Health Science Centre, Manchester, UK. <sup>93</sup>The Manchester Centre for Genomic Medicine, University of Manchester, Manchester, UK. <sup>94</sup>QIMR Berghofer Medical Research Institute, Royal Brisbane Hospital, Brisbane, Australia. <sup>95</sup>Department of Pediatrics, Emory University School of Medicine, Atlanta, GA, USA. <sup>96</sup>Department of Gastroenterology, Kaunas University of Medicine, Kaunas, Lithuania. <sup>97</sup>Centre for Inflammatory Bowel Diseases, Saint John of God Hospital, Subiaco, WA 6008, Australia. <sup>98</sup>Inflammatory Bowel Disease Research Group, Addenbrooke's Hospital, Cambridge CB2 0QQ, UK. <sup>99</sup>Faculty of medicine, University of Latvia, Riga, Latvia. <sup>100</sup>Dipartimento di Neuroscienze, Psicologia, Area del Farmaco e Salute del Bambino, Università di Firenze Strutture Organizzative Dipartimentali (SOD) Gastroenterologia e Nutrizione Ospedale Pediatrico Meyer, Firenze, Italy. <sup>101</sup>Department of Gastroenterology, The Townsville Hospital, Townsville, Australia. <sup>102</sup>Institute of Human Genetics, Newcastle University, Newcastle upon Tyne, UK. <sup>103</sup>Department of Medical and Molecular Genetics, King's College London, London SE1 9RT, UK. <sup>104</sup>Sydney Brenner Institute for Molecular Bioscience, University of the Witwatersrand, Johannesburg 2193, South Africa. <sup>105</sup>Inflammatory Bowel Disease Centre, Mount Sinai Hospital, Toronto, ON, Canada. <sup>106</sup>Center for Human Molecular Genetics and Pharmacogenomics, Faculty of Medicine, University of Maribor, Maribor, Slovenia. <sup>107</sup>Department of Medicine, Ninewells Hospital and Medical School, Dundee, UK. <sup>108</sup>Center for Computational and Integrative Biology, Massachusetts General Hospital, Harvard Medical School, Boston, MA, USA. <sup>109</sup>Department of Medicine and Therapeutics, Institute of Digestive Disease, Chinese University of Hong Kong, Hong Kong, Hong Kong. <sup>110</sup>Department of Genomics Life & Brain Center, University Hospital Bonn, Bonn, Germany. <sup>111</sup>Faculty for Chemistry and Chemical Engineering, University of Maribor, Maribor, Slovenia. <sup>112</sup>Department of Gastroenterology, Royal Brisbane and Womens Hospital, Brisbane, Australia. <sup>113</sup>Department of Gastroenterology, Université Catholique de Louvain (UCL) Centre Hospitalier Universitaire (CHU) Mont-Godinne, Mont-Godinne, Belgium. <sup>114</sup>Department of Gastroenterology, Guy's & St. Thomas' NHS Foundation Trust, St.-Thomas Hospital, London, UK. <sup>115</sup>Department of Digestive Diseases, Hospital Quiron Teknon, Barcelona, Spain. <sup>116</sup>Department of Public Health Sciences, University of Chicago, Chicago, IL, USA.

<sup>117</sup>Human Genetics, Genome Institute of Singapore, Singapore, Singapore. <sup>118</sup>Institute for Human Genetics, University of California, San Francisco, CA, USA. <sup>119</sup>Open Targets, Wellcome Trust Genome Campus, Hinxton, Cambridgeshire CB10 1SD, UK. <sup>120</sup>Department of Biology of Radiations and Human Health, Agenzia Nazionale per le Nuove Tecnologie l'energia e lo Sviluppo Economico Sostenibile (ENEA), Rome, Italy. <sup>121</sup>Department of Clinical Science Intervention and Technology, Karolinska Institutet, Stockholm, Sweden. <sup>122</sup>Gastroenterology & General Medicine, Norfolk and Norwich University Hospital, Norwich, UK. <sup>123</sup>Faculty of Medicine, University of Toronto, Toronto, ON, Canada. <sup>124</sup>Department of Computer Science, New Jersey Institute of Technology, Newark, NJ, USA. <sup>125</sup>Institute of Human Genetics, Technische Universität München, Munich, Germany. <sup>126</sup>Department of Neurology, Technische Universität München, Munich, Germany. <sup>127</sup>Department of Biostatistics, School of Public Health, Yale University, New Haven, CT, USA. <sup>128</sup>Department of Gastroenterology, West China Hospital, Chengdu, Sichuan, China. <sup>129</sup>State Key Laboratory of Biotherapy, Sichuan University West China University of Medical Sciences (WCUMS), Chengdu, Sichuan, China.