



OPEN ACCESS

Original research

# Genome-wide analysis of 944 133 individuals provides insights into the etiology of haemorrhoidal disease

Tenghao Zheng <sup>1,2</sup> David Ellinghaus <sup>3,4</sup> Simonas Juzenas <sup>3,5</sup> François Cossais,<sup>6</sup> Greta Burmeister,<sup>7</sup> Gabriele Mayr,<sup>3</sup> Isabella Friis Jørgensen,<sup>4</sup> Maris Teder-Laving,<sup>8</sup> Anne Heidi Skogholt,<sup>9</sup> Sisi Chen,<sup>10</sup> Peter R Stregé,<sup>10</sup> Go Ito,<sup>3,11</sup> Karina Banasik,<sup>4</sup> Thomas Becker,<sup>12</sup> Frank Bokelmann,<sup>13</sup> Søren Brunak,<sup>4</sup> Stephan Buch,<sup>14</sup> Hartmut Clausnitzer,<sup>15</sup> Christian Datz <sup>16</sup> DBDS Consortium, Frauke Degenhardt,<sup>3</sup> Marek Doniec,<sup>17</sup> Christian Erikstrup,<sup>18</sup> Tõnu Esko,<sup>8</sup> Michael Forster <sup>3</sup> Norbert Frey,<sup>19,20,21</sup> Lars G Fritsche,<sup>22</sup> Maiken Elvestad Gabrielsen,<sup>9</sup> Tobias Gräßle,<sup>23,24</sup> Andrea Gsur <sup>25</sup> Justus Gross,<sup>7</sup> Jochen Hampe <sup>14,26</sup> Alexander Hendricks,<sup>7</sup> Sebastian Hinz,<sup>7</sup> Kristian Hveem,<sup>9</sup> Johannes Jongen,<sup>27,28</sup> Ralf Junker,<sup>15</sup> Tom Hemming Karlsen,<sup>29</sup> Georg Hemmrich-Stanisak,<sup>3</sup> Wolfgang Kruis,<sup>30</sup> Juozas Kupcinskas,<sup>31</sup> Tilman Laubert,<sup>27,28,32</sup> Philip C Rosenstiel,<sup>3,33</sup> Christoph Röcken <sup>34</sup> Matthias Laudes,<sup>35</sup> Fabian H Leendertz,<sup>23,24</sup> Wolfgang Lieb,<sup>36</sup> Verena Limperger,<sup>15</sup> Nikolaos Margetis,<sup>37</sup> Kerstin Mätz-Rensing,<sup>38</sup> Christopher Georg Németh,<sup>12,39</sup> Eivind Ness-Jensen <sup>9,40,41</sup> Ulrike Nowak-Göttl,<sup>15</sup> Anita Pandit,<sup>22</sup> Ole Birger Pedersen,<sup>42</sup> Hans Günter Peleikis,<sup>27,28</sup> Kenneth Peuker,<sup>14,26</sup> Cristina Leal Rodriguez,<sup>4</sup> Malte Christoph Rühlemann <sup>3</sup> Bodo Schniewind,<sup>43</sup> Martin Schulzky,<sup>3</sup> Jurgita Skieceviciene,<sup>31</sup> Jürgen Tepel,<sup>44</sup> Laurent Thomas,<sup>9,45,46,47</sup> Florian Uellendahl-Werth,<sup>3</sup> Henrik Ullum,<sup>48</sup> Ilka Vogel,<sup>49</sup> Henry Volzke,<sup>50</sup> Lorenzo von Fersen,<sup>51</sup> Witigo von Schönfels,<sup>12</sup> Brett Vanderwerff,<sup>22</sup> Julia Wilking,<sup>7</sup> Michael Wittig,<sup>3</sup> Sebastian Zeissig,<sup>14,26</sup> Myrko Zobel,<sup>52</sup> Matthew Zawistowski,<sup>22</sup> Vladimir Vacic,<sup>53</sup> Olga Sazonova,<sup>53</sup> Elizabeth S Noblin,<sup>53</sup> The 23andMe Research Team, Gianrico Farrugia <sup>10</sup> Arthur Beyder,<sup>10</sup> Thilo Wedel,<sup>6</sup> Volker Kahlke,<sup>27,28,54</sup> Clemens Schafmayer,<sup>7</sup> Mauro D'Amato <sup>1,2,55,56</sup> Andre Franke <sup>3,33</sup>

► Additional supplemental material is published online only. To view, please visit the journal online (<http://dx.doi.org/10.1136/gutjnl-2020-323868>).

For numbered affiliations see end of article.

## Correspondence to

Prof Mauro D'Amato, Gastrointestinal Genetics Lab, CIC bioGUNE - BRTA, Bizkaia Science and Technology Park, Building 801A, 48160 Derio, Spain; [mdamato@cicbiogune.es](mailto:mdamato@cicbiogune.es) and Prof Andre Franke, Institute of Clinical Molecular Biology (IKMB) & University Hospital Schleswig-Holstein (UKSH), Christian-Albrechts-University of Kiel (CAU), Rosalind-Franklin-Str. 12, D-24105 Kiel, Germany; [a.franke@mucosa.de](mailto:a.franke@mucosa.de)

TZ, DE and SJ contributed equally.  
CS, MD and AF contributed equally.

Received 18 December 2020  
Revised 19 March 2021  
Accepted 30 March 2021  
Published Online First  
22 April 2021



► <http://dx.doi.org/10.1136/gutjnl-2021-324817>



© Author(s) (or their employer(s)) 2021. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

**To cite:** Zheng T, Ellinghaus D, Juzenas S, *et al.* *Gut* 2021;**70**:1538–1549.

## ABSTRACT

**Objective** Haemorrhoidal disease (HEM) affects a large and silently suffering fraction of the population but its aetiology, including suspected genetic predisposition, is poorly understood. We report the first genome-wide association study (GWAS) meta-analysis to identify genetic risk factors for HEM to date.

**Design** We conducted a GWAS meta-analysis of 218 920 patients with HEM and 725 213 controls of European ancestry. Using GWAS summary statistics, we performed multiple genetic correlation analyses between HEM and other traits as well as calculated HEM polygenic risk scores (PRS) and evaluated their translational potential in independent datasets. Using functional annotation of GWAS results, we identified HEM candidate genes, which differential expression and coexpression in HEM tissues were evaluated employing RNA-seq analyses. The localisation of expressed proteins at selected loci was investigated by immunohistochemistry.

**Results** We demonstrate modest heritability and genetic correlation of HEM with several other diseases from the GI, neuroaffective and cardiovascular domains. HEM PRS validated in 180 435 individuals from independent datasets allowed the identification of those at risk and correlated with younger age of onset and recurrent surgery. We identified 102 independent HEM risk loci harbouring genes whose expression is enriched in blood vessels and GI tissues, and in pathways associated with smooth muscles, epithelial and endothelial development and morphogenesis. Network transcriptomic analyses highlighted HEM gene coexpression modules that are relevant to the development and integrity of the musculoskeletal and epidermal systems, and the organisation of the extracellular matrix.

**Conclusion** HEM has a genetic component that predisposes to smooth muscle, epithelial and connective tissue dysfunction.

## Significance of this study

## What is already known about this subject?

- ▶ Human haemorrhoidal disease (HEM) is a prevalent anorectal pathology characterised by the symptomatic enlargement and distal displacement of anal cushions.
- ▶ As a consequence of scarce research and of being a taboo topic in society, only a few HEM non-genetic risk factors have been suggested, thus the aetiology of the disease still remains unclear.
- ▶ Genetic susceptibility to HEM has been suspected but was never systematically investigated.

## What are the new findings?

- ▶ Here, we report the first well-powered genome-wide association study (GWAS) meta-analysis with a sample size of 944 133 individuals to identify genetic risk factors for HEM.
- ▶ We describe 102 novel independent HEM risk loci that are functionally linked to pathways associated with smooth muscles, epithelial and endothelial development and morphogenesis.
- ▶ We show genetic correlations of HEM with several other diseases that are classified as GI, neuroaffective and cardiovascular disorders.
- ▶ We report significance of computed HEM polygenic risk scores, which were validated in independent population-based cohorts.
- ▶ We show significant enrichment of HEM genes in tissue coexpression modules responsible for the development of musculoskeletal and epidermal systems, and the organisation of the extracellular matrix.
- ▶ Based on our data, we outline HEM as a disorder of impaired neuromuscular motility, smooth muscle contraction and extracellular matrix organisation.

## How might it impact on clinical practice in the foreseeable future?

- ▶ The results from this GWAS provide new insights on HEM genetic predisposition to smooth muscle, epithelial and connective tissue dysfunction.
- ▶ HEM polygenic risk scores identify individuals at risk and correlate with a more severe phenotype.

## INTRODUCTION

Haemorrhoids are normal anal vascular cushions filled with blood at the junction of the rectum and the anus. It is assumed that their main role in humans is to maintain continence<sup>1</sup> but other functions such as sensing fullness, pressure and perceiving anal contents have been suggested given the sensory innervation.<sup>2</sup> Haemorrhoidal disease (hereafter referred to as HEM) occurs when haemorrhoids enlarge and become symptomatic (sometimes associated with rectal bleeding and itching/soiling) due to the deterioration or prolapse of the anchoring connective tissue, the dilation of the haemorrhoidal plexus or the formation of blood clots. Severe forms of HEM often require surgical treatment and the removal of abnormally enlarged and/or thrombosed haemorrhoids.<sup>1</sup> HEM prevalence increases with age and shows staggering figures worldwide (up to 86% prevalence in some reports),<sup>3</sup> whereby a large proportion of cases remain undetected as asymptomatic or mild enough to be self-treated with over-the-counter treatment. HEM represents a considerable medical and socioeconomic burden with an estimated annual cost of US\$800 million in the USA alone, mainly

related to the large number of haemorrhoidectomies performed every year.<sup>4</sup>

A number of HEM risk factors have been suggested, including human erect position. The tight anal sealing provided by the elaborated haemorrhoidal plexus may have developed during human evolution co-occurring with permanent bipedalism, as shown by our histology comparison of four different mammals (human, gorilla, baboon, mouse; online supplemental figure 1, online supplemental material 1). Other suggested risk factors are a sedentary lifestyle, obesity, reduced dietary fibre intake, spending excess time on the toilet, straining during defecation, strenuous lifting, constipation, diarrhoea, pelvic floor dysfunction, pregnancy and giving natural birth, with several being controversially reported. The hypothetical model shown in online supplemental figure 2 summarises the contemporary concepts regarding the pathophysiology of HEM development.<sup>5</sup> Until today, HEM aetiopathogenesis is poorly investigated, and neither the exact molecular mechanisms nor the reason(s) why only some people develop HEM are known. Genetic susceptibility may play a role in HEM development, but no large-scale, genome-wide association study (GWAS) for HEM has ever been conducted. To evaluate the contribution of genetic variation to the genetic architecture of HEM, we carried out a GWAS meta-analysis in 218 920 affected individuals and 725 213 population controls of European ancestry.

## METHODS

Detailed methods are provided in the online supplemental material 1 and summarised in the online supplemental figure 3.

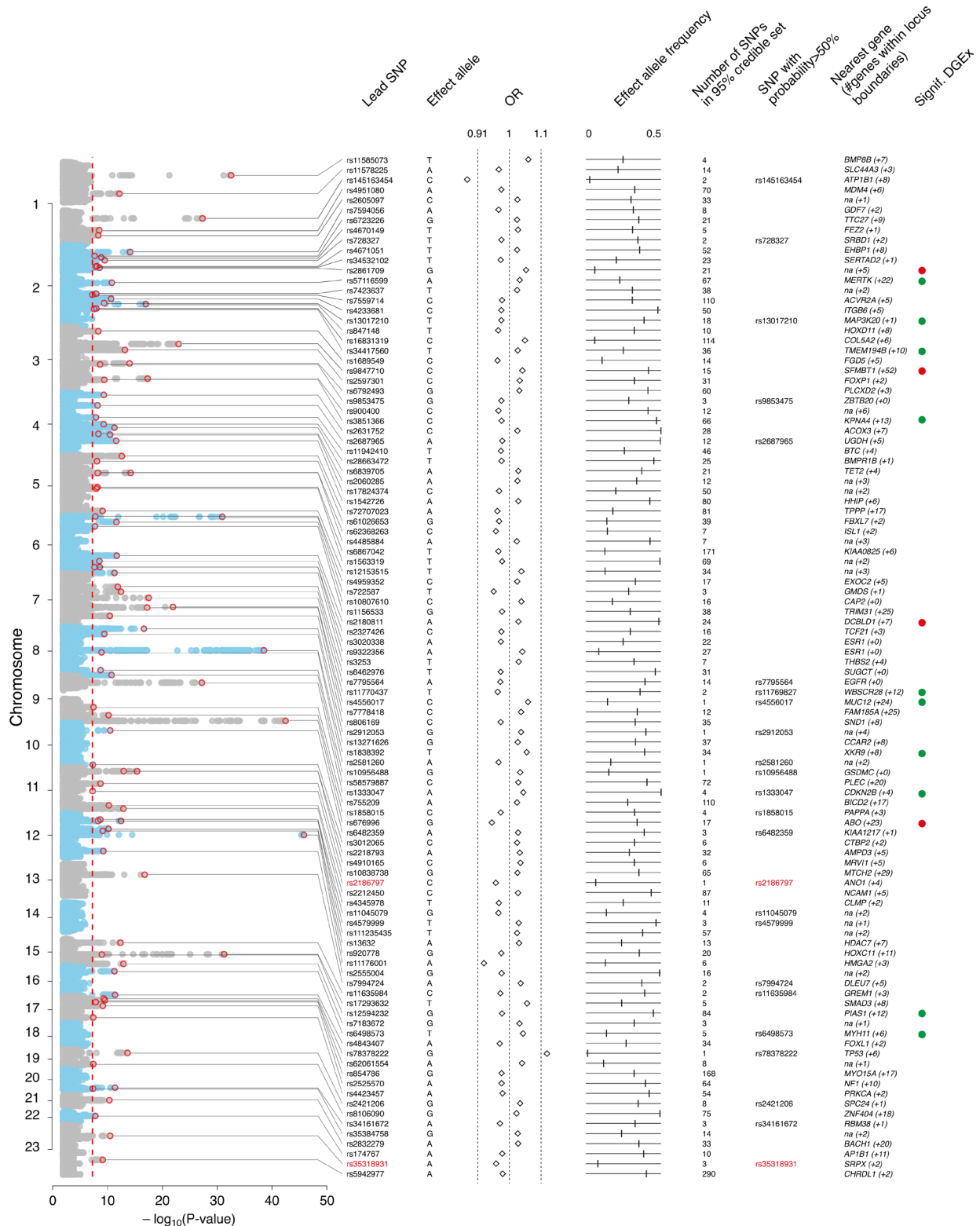
## RESULTS

## GWAS meta-analysis and fine-mapping genomic regions

We conducted a GWAS meta-analysis in 944 133 individuals of European ancestry from five large population-based cohorts (23andMe,<sup>6</sup> UK Biobank (UKBB),<sup>7</sup> Estonian Genome Centre at the University of Tartu,<sup>8</sup> Michigan Genomics Initiative,<sup>9</sup> and Genetic Epidemiology Research on Aging (GERA)<sup>10</sup>) including 218 920 HEM cases and 725 213 controls (online supplemental table 1). HEM cases had higher body mass index (BMI), were significantly older and more often women compared with non-HEM controls; hence, age, sex, BMI (if available) and top principal components from principal component analysis (PCA) were included as covariates in individual GWAS (see 'Methods' section in online supplemental material 1, hereafter referred to as online methods).

After data harmonisation and quality control, 8 494 288 high-quality, common (minor allele frequency >0.01) single nucleotide polymorphisms (SNPs) were included in a fixed-effect inverse variance meta-analysis using the software METAL<sup>11</sup> (online methods). We identified 5480 genome-wide significant associations ( $P_{\text{Meta}} < 5 \times 10^{-8}$ ), which were mapped to 102 independent genomic regions using FUMA<sup>12</sup> (online methods). A summary of the association results for 102 lead SNPs is provided in figure 1 and online supplemental table 2. Although genomic inflation was observed ( $\lambda = 1.3$ ; online supplemental figure 4), this was likely due to polygenicity rather than population stratification, as determined via linkage disequilibrium score regression analysis (LDSC, intercept = 1.06; online methods) and based on a normalised  $\lambda_{1000} = 1.001$ . The heritability of HEM was estimated at 5% (SNP-based heritability,  $h^2_{\text{SNP}}$  computed with LDSC), whereby the newly identified 102 risk variants explained about 0.9% of the variance (online methods).





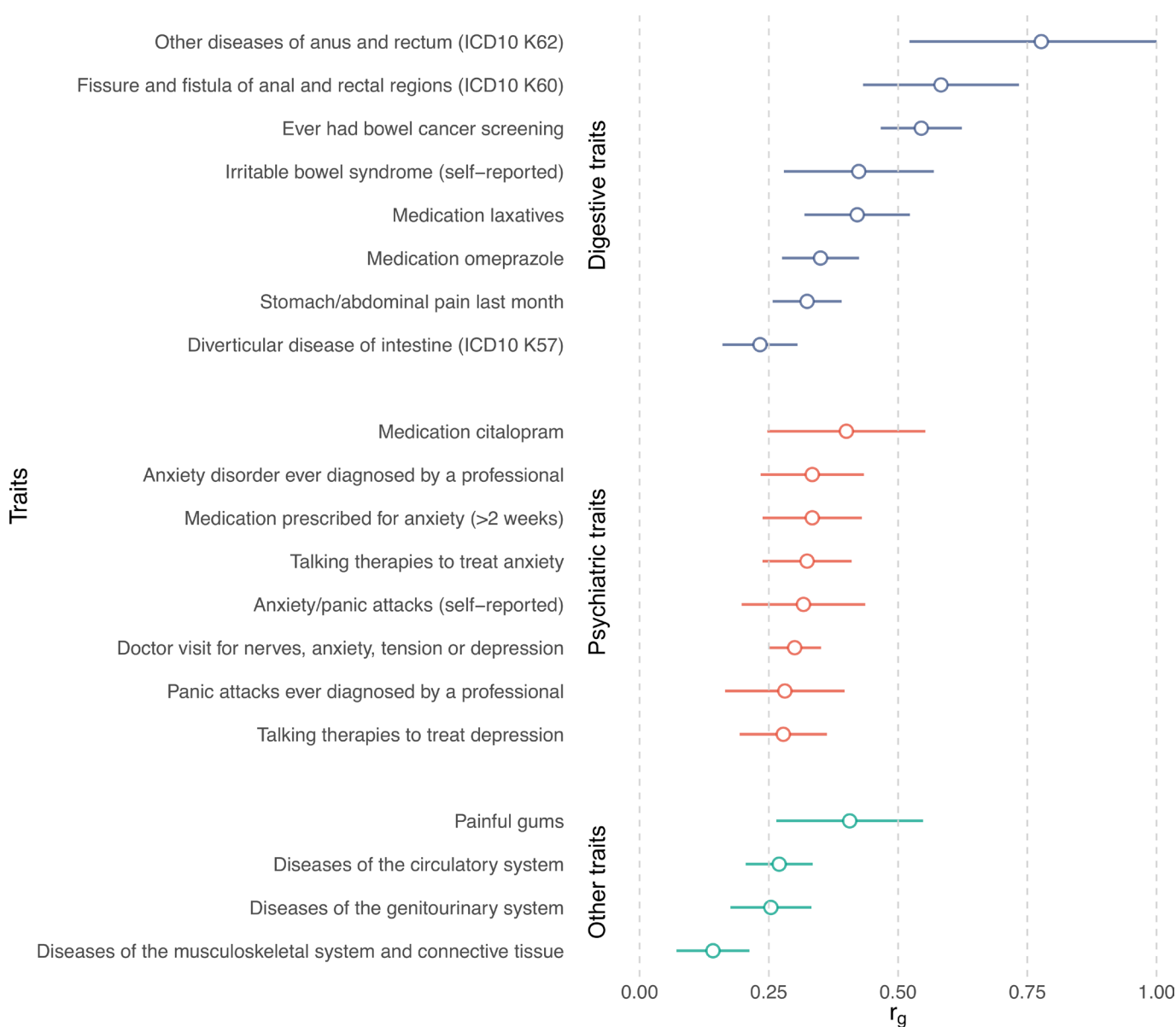
**Figure 1** Annotation of 102 haemorrhoidal disease (HEM) genome-wide association study (GWAS) risk loci. From left to right: Manhattan plot of GWAS meta-analysis results, (genome-wide significance level— $P_{\text{Meta}} < 5 \times 10^{-8}$ —indicated with vertical dotted red line); Lead single nucleotide polymorphism (SNP)—marker associated with the strongest association signal from each locus (also annotated with a red circle in the Manhattan plot); Effect allele—allele associated with reported genetic risk effects (OR), also always the minor allele; OR with respect to the effect allele; Effect allele frequency—frequency of the effect allele in the discovery dataset; Number of SNPs in 95% credible set—the minimum set of variants from Bayesian fine-mapping analysis that is >95% likely to contain the causal variant; SNP with probability >50%—single variant (if detected) with >50% probability of being causal (coding SNPs highlighted in red); Nearest gene (#genes within locus boundaries)—gene closest to the lead SNP (if within 100 kb distance, otherwise 'na') and number of additional genes positionally mapped to the locus using FUMA (online supplemental table 2 and online methods). Signif. DGEx—locus containing HEM genes differentially expressed in RNA Combo-Seq analysis of HEM affected tissue, detected at higher (green) and/or lower (red) level of expression (see online methods).

Bayesian fine-mapping analysis with FINEMAP<sup>13</sup> identified a total of 3323 SNPs that belong to the 95% credible sets of variants most likely to be causal at each locus (online methods, online supplemental figure 5 and online supplemental table 3). For six loci, we pinpointed the association signal to a single causal variant with greater than 95% certainty, including two missense variants rs2186797 (*ANO1*) and rs35318931 (*SRPX*). For another 24 loci, there was evidence that the lead variant is causal with >50% certainty (figure 1).

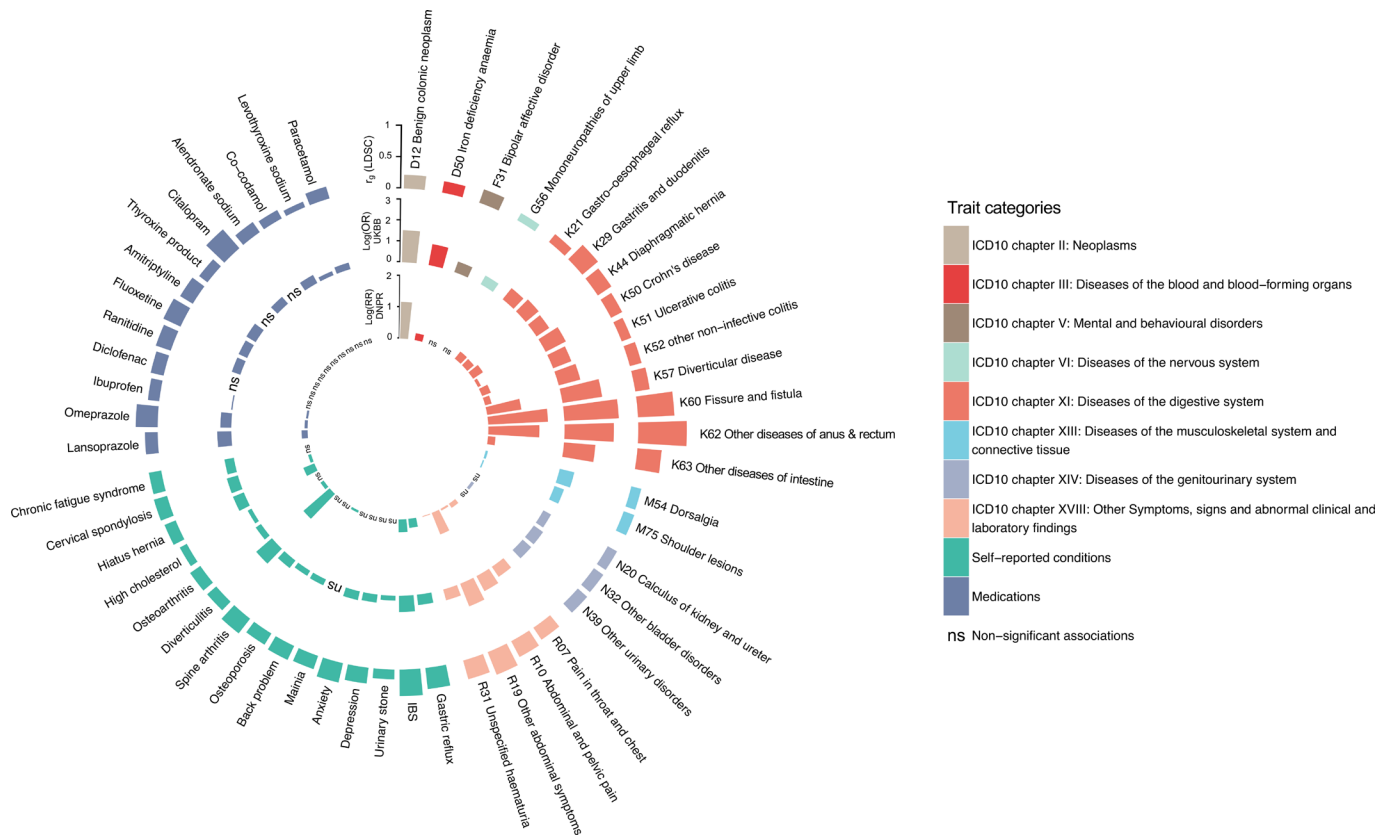
### Cross-trait analyses

A lookup of HEM association signals in previous GWAS studies retrieved from GWAS Atlas,<sup>14</sup> GWAS Catalog<sup>15</sup> and via PhenoScanner v2<sup>16</sup> (online methods) revealed that 76/102 loci had been previously associated with diseases and traits across the metabolic, cardiovascular, digestive, psychiatric, environmental and other domains (online supplemental figure 6 and online supplemental table 4).

We then investigated genetic correlations with 1387 other human traits and conditions using LDSC as implemented in CTG-VL<sup>17</sup> (online methods). The strongest correlations ( $r_g$ ) were observed with diseases and traits from the GI domain, including ‘other diseases of anus and rectum’ ( $r_g=0.78$ ,  $P_{FDR}=4.94\times 10^{-8}$ ), ‘fissure and fistula of anal and rectal regions’ ( $r_g=0.58$ ,  $P_{FDR}=2.70\times 10^{-12}$ ), ‘self-reported IBS’ ( $r_g=0.42$ ,  $P_{FDR}=1.87\times 10^{-7}$ ), use of ‘laxatives’ ( $r_g=0.42$ ,  $P_{FDR}=2.35\times 10^{-14}$ ), and ‘diverticular disease’ ( $r_g=0.23$ ,  $P_{FDR}=6.68\times 10^{-9}$ ) among others (figure 2). Given these similarities, we performed a genome-wide pleiotropy analysis for diverticular disease, IBS and HEM, revealing 44 independent genomic regions shared by at least two phenotypes from the group of diverticular disease, IBS and HEM phenotypes (online supplemental table 5). Other notable correlations were detected for psychiatric and neuro-affective disorders (anxiety, depression and neuroticism), pain-related traits (including abdominal pain and painful gums), diseases of the circulatory system, and diseases of the musculoskeletal system and connective tissue (figure 2 and online supplemental table 6).



**Figure 2** Genetic correlation between haemorrhoidal disease and other traits estimated by linkage disequilibrium score regression analysis. Genetic correlations ( $r_g$  +se) are shown for selected traits, grouped by domain. Only correlations significant after Bonferroni correction were considered (full list available in online supplemental table 6). ICD, International Classification of Diseases.



**Figure 3** Analysis of haemorrhoidal disease (HEM) genetically correlated traits in UK Biobank (UKBB) and the Danish National Patient Registry (DNPR). Traits and conditions identified in linkage disequilibrium score regression analyses of genetic correlation with HEM (outer ring in the circo plot, see also figure 2 and online supplemental table 6) were studied for their differential prevalence in UKBB and DNPR, based on data extracted from participants' healthcare records. Significant results are reported, respectively, as ORs (log(OR), UKBB, middle ring) and relative risk (log(RR), DNPR, inner ring) or 'ns' (for non-significant findings). Diseases and traits are categorised according to ICD10 diagnostic codes or self-reported conditions and use of medications from questionnaire data (see online methods). Self-reported traits in UKBB (dark blue colour) were manually mapped to ICD10-codes in DNPR.

To gain further insight into potential cross-trait overlaps and to validate the LDSC results, we assessed whether genetically correlated traits and conditions also occur more frequently in patients with HEM by analysing data on diagnoses and medications from UKBB. The results were highly consistent with those obtained via LDSC (figure 3, online supplemental table 6). Compared with controls, patients with HEM additionally suffered more often from diverticular disease, IBS and other functional GI disorders (FGIDs), abdominal pain, hypertension, ischaemic heart disease, depression, anxiety, and diseases of the musculoskeletal system and connective tissue among others. To further consolidate these findings, we analysed an independent population-scale healthcare record comprising 8 172 531 individuals from the Danish National Patient Registry (online methods), obtaining similar results (see inner circle of figure 3). Therefore, based on these largely overlapping observations at the genetic and epidemiological level, HEM appears to be strongly associated with other diseases of the digestive, neuroaffective and cardiovascular domains.

### Polygenic risk scores (PRS)

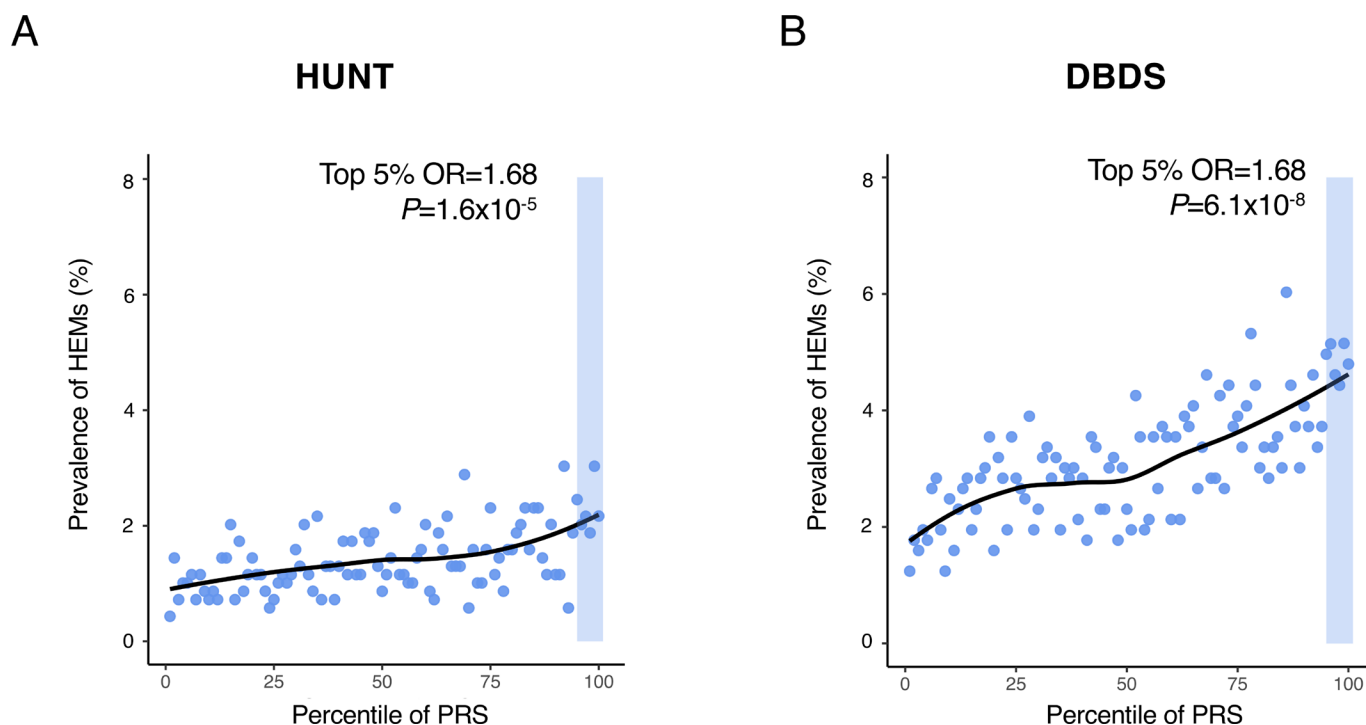
We next exploited meta-analysis summary statistics to compute HEM PRS and evaluate their relevance and translational potential in independent datasets. HEM PRS were calculated with PRSice-2<sup>18</sup> (online methods) and their performance tested in three independent datasets: (1) the Norwegian Trøndelag Health

Study (HUNT<sup>19</sup>; n=69 291; 977 cases vs 68 314 controls), (2) the Danish Blood Donor Study (DBDS<sup>20</sup>; n=56 397; 1754 cases vs 54 643 controls), and (3) 1144 HEM cases from gastroenterology clinics compared with 2740 cross-sectional population controls from Germany<sup>21</sup> (online supplemental table 1). In all three datasets, patients with HEM showed significantly higher PRS values compared with controls (OR=1.24,  $p=1.98 \times 10^{-11}$ ; OR=1.28,  $p=3.52 \times 10^{-22}$  and OR=1.36,  $p=7.64 \times 10^{-15}$ , respectively for the two population-based cohorts HUNT and DBDS, and the German case-control cohort). In HUNT and DBDS, HEM prevalence increased across PRS percentile distributions, with individuals from the top 5% tail being exposed to higher HEM risk compared with the rest of the population (OR=1.68,  $p=1.55 \times 10^{-5}$ , and OR=1.68,  $p=6.11 \times 10^{-8}$ , respectively for HUNT and DBDS, figure 4). Higher HEM PRS were also associated with a more severe phenotype as defined by the need for recurrent invasive procedures (OR=1.03,  $p=8.63 \times 10^{-3}$  in German patients) and a younger age of onset ( $p=1.90 \times 10^{-3}$  in DBDS;  $p=4.01 \times 10^{-3}$  in German patients).

### Functional annotation of GWAS loci, tissue and pathway enrichment analyses

In order to identify most likely candidate genes and relevant molecular pathways, we used independent computational pipelines for the functional annotation of GWAS results. This yielded





**Figure 4** Risk haemorrhoidal disease (HEM) prevalence across polygenic risk score (PRS) percentile distributions. PRS was derived from the results of the association meta-analysis (see online methods). HEM prevalence (%; Y-axis) is reported on a scatter plot in relation to PRS percentile distribution (X-axis) in the Norwegian Trøndelag Health Study (HUNT) (A) and the Danish Blood Donor Study (DBDS) (B) population cohorts. The top 5% of the distribution is highlighted with a shaded area in both cohorts, and the results of testing HEM prevalence in this group versus the rest of the population are also reported (p value and OR from logistic regression; online methods).

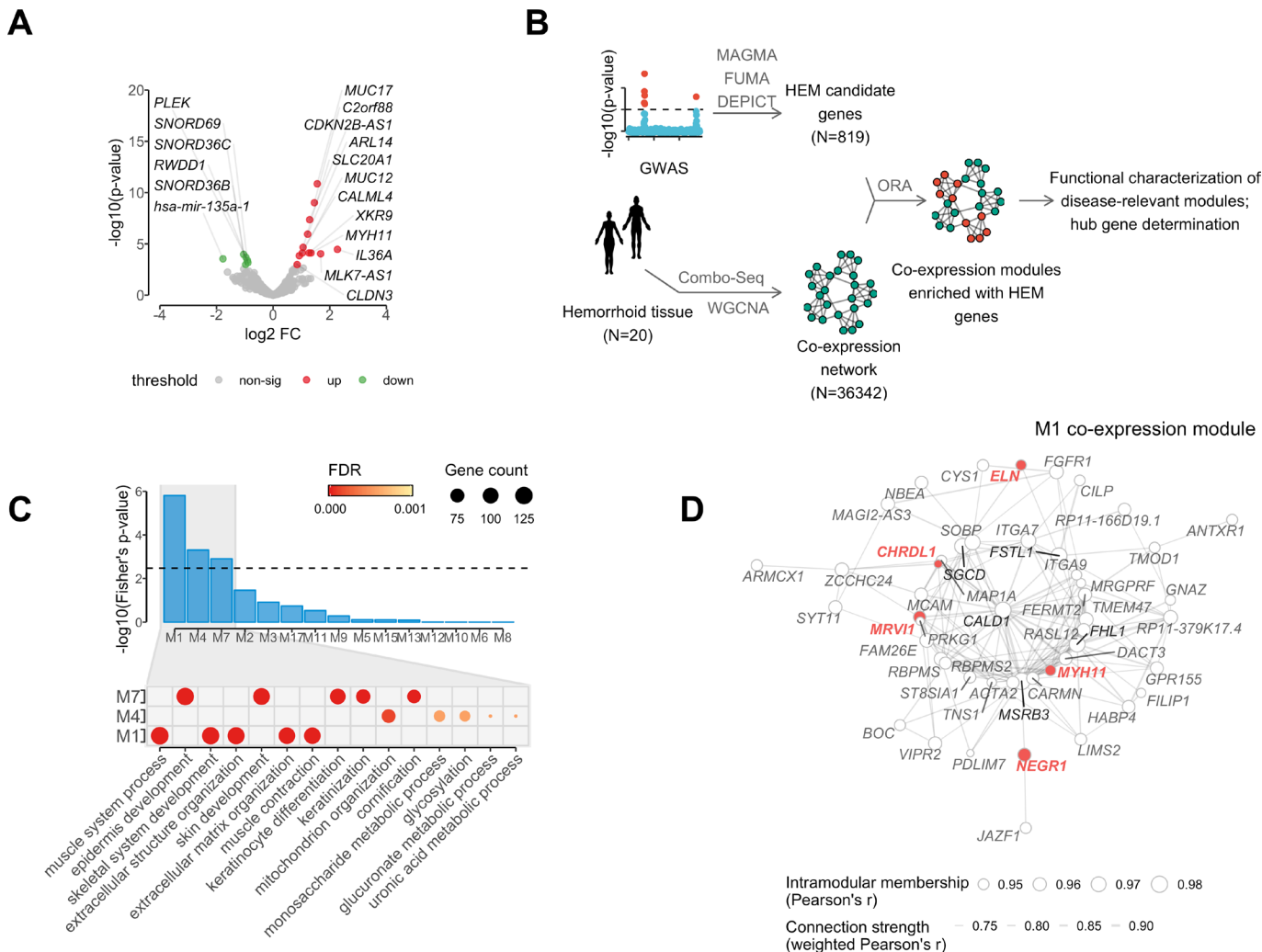
a total of 819 non-redundant HEM-associated transcripts (hereafter referred to as HEM genes; online supplemental table 7) derived from alternative positional (N=540 total from FUMA,<sup>12</sup> DEPICT<sup>22</sup> and MAGMA<sup>23</sup>) and expression quantitative trait (eQTL, N=562 from FUMA) mapping efforts (online methods). Tissue-specific enrichment analyses (TSEA) of 540 positional candidates led to very similar results for all three approaches, with enrichment of HEM gene expression in blood vessels, colon and other relevant tissues (online supplemental figure 7 and online supplemental table 8). Similarly, gene-set enrichment analyses (GSEA) highlighted common pathways important in the development of vasculature and the intestinal tissue including the gene ontology (GO) terms ‘tube morphogenesis and development’, ‘artery morphogenesis and development’, ‘epithelium morphogenesis’, ‘smooth muscle tissue morphogenesis’ and others. Additionally, DEPICT detected enrichment for a number of traits from the Mammalian Phenotype Ontology including ‘abnormal intestinal morphology’, ‘rectal prolapse’, ‘abnormal blood vessel (and artery) morphology’, and ‘abnormal smooth muscle morphology and physiology’ (online supplemental table 8). TSEA and GSEA analyses performed on all 819 transcripts including eQTL genes from FUMA gave rise to similar results although these did not reach statistical significance (not shown).

#### Gene expression in HEM tissue and gene-network analyses

The expression of HEM genes was studied in integrated mRNA and microRNA Combo-Seq analysis of enlarged haemorrhoidal tissue from 20 patients with HEM and normal specimens from 18 controls (online methods, online supplemental table 1). HEM genes were examined with regard to their expression status, differential expression between cases and controls, and their connectivity and topology in gene coexpression networks. After

normalisation for cell-type heterogeneity in different tissues (online methods, online supplemental figure 8), 720 out of 819 candidate genes were found to be expressed in haemorrhoidal tissue, with 287 (39.9%) of these being among the most strongly expressed transcripts (upper quartile) (online supplemental table 7). Compared with normal tissue from controls, 18 HEM candidate genes from 14 independent loci showed differential expression in haemorrhoidal tissue ( $P_{\text{FDR}} < 0.05$  and  $|\log_2 \text{fold change}| > 0.5$ ), with 12 genes showing increased and 6 decreased expression (figure 5A).

To obtain further biological insight from transcriptomic profiling, HEM candidate genes were further characterised for their membership in coexpressed gene modules identified via weighted gene correlation network analysis (WGCNA)<sup>24</sup> (figure 5B; online methods). The final network consisted of 36 342 genes partitioned into 41 coexpression modules, and 3 of these (M1, M4 and M7; referred to as HEM modules) were significantly enriched ( $P_{\text{FDR}} < 0.05$ ) for HEM genes (figure 5C). In total, 260 (35.7%) expressed HEM genes were members of modules M1 (N=121), M4 (N=75) or M7 (N=64) (online supplemental table 9). Functional annotation of these modules revealed M1 enrichment for ‘extracellular matrix (ECM) organisation’ ( $P_{\text{FDR}} = 5.17 \times 10^{-18}$ ) and ‘muscle contraction’ ( $P_{\text{FDR}} = 7.90 \times 10^{-17}$ ); M4 for ‘mitochondrion organisation’ ( $P_{\text{FDR}} = 2.35 \times 10^{-4}$ ) and ‘glycosylation’ ( $P_{\text{FDR}} = 5.70 \times 10^{-4}$ ); and M7 for ‘cornification’ ( $P_{\text{FDR}} = 3.97 \times 10^{-52}$ ) and ‘epidermis development’ ( $P_{\text{FDR}} = 2.07 \times 10^{-41}$ ) (figure 5C). The WGCNA analysis also allowed for the identification of the most interconnected genes, or module hub genes (ie, central nodes in the scale-free network). Of note, several HEM candidate genes, namely *NEGR1*, *MRV11*, *MYH11*, *ELN* and *CHRDL1*, were among the top 50 hub genes for module M1 (figure 5D). M1 is the module



**Figure 5** Analysis of mRNA and microRNA (Combo-Seq) data from haemorrhoidal disease (HEM) affected tissue, in relation to HEM genes coexpression networks. (A) Volcano plot reporting HEM genes differentially expressed in haemorrhoidal tissue (significantly upregulated=red, and downregulated=green); (B) schematic representation of the analytical flow for HEM genes coexpression network identification and characterisation; (C) upper panel (barplot): overrepresentation analysis of HEM genes in coexpression network modules, with significant enrichment ( $P_{\text{FDR}} < 0.05$ ) in modules M1, M4 and M7; lower panel (dotplot): top five gene ontology terms (biological process) from gene set enrichment analysis relative to M1, M4 and M7 coexpression modules (gene counts and false discovery rate (FDR)-adjusted significance level are also reported as indicated); (D) coexpression hub network of module M1. The network represents strength of connections (weighted Pearson's correlation  $> 0.7$ ) among the top 50 hub genes with highest values of intramodular membership (size of the node). HEM genes and the top 5 hub genes are highlighted in red and black, respectively.

most significantly enriched for HEM genes ( $P_{\text{FDR}} = 6.40 \times 10^{-5}$ ), and also the largest ( $n = 3975$ ) coexpression module in the haemorrhoidal tissue (online supplemental table 9), thus pointing to the importance of its associated GO terms ‘ECM organisation’ and muscle contraction in HEM.

**Prioritised HEM genes**

In order to identify genes most likely to play a causative role in HEM, we selected candidates based on a scoring approach by prioritising those associated with one or more of the following: (1) linked to a high-confidence fine-mapped variant (posterior probability (PP)  $> 50\%$ ), (2) differentially expressed in enlarged haemorrhoidal tissue, (3) highlighted by pathway and tissue/cell-type enrichment DEPICT analysis (online methods), or (4) predicted hub of a WGCNA coexpression HEM module (M1, M4, M7). This reduced the number of candidates from 819 to 100 prioritised genes associated with 58 independent HEM loci

(online supplemental table 7). Some notable observations were made in relation to a subset of these prioritised genes, whose associated evidence and known biological function(s) make them remarkably good candidates to play a role in HEM risk.

Two genes, *ANO1* and *SRPX*, were both linked to a single coding variant fine-mapped as causal with very high confidence (rs2186797, *ANO1*, p.Phe608Ser with PP=97.0%, and rs35318931, *SRPX*, p.Ser413Phe with PP=87.3%, respectively). *ANO1* encodes the voltage-gated calcium-activated anion channel anoctamin-1 protein, which is highly expressed in the interstitial cells of Cajal (ICCs) throughout the human GI tract, where it contributes to the control of intestinal motility and peristalsis.<sup>25</sup> The *ANO1*:p.Phe608Ser (F608S) variant is predicted to destabilise local protein structure and to disturb ANO1-activating phospholipid interactions (online supplemental figure 9). Indeed, site-directed mutagenesis and electrophysiology experiments in vitro showed that the amino acid 608 Phe to

Ser change leads to an increased voltage-dependent instantaneous  $\text{Cl}^-$  current, and a slowing of activation and deactivation kinetics (especially at high  $\text{Ca}^{2+}$  concentrations) (online supplemental figure 10). The X-linked gene *SRPX* codes for a Sushi repeat-containing protein whose domain composition implies a role in ECM, and is expressed in various ECM tissues including colon and liver.<sup>26</sup> The *SRPX*:p.Ser413Phe variant (rs35318931) may potentially destabilise the C-terminal domain of unknown function (DUF4174) (online supplemental figure 11), which is conserved in various ECM proteins (online methods, section In silico variant protein analysis).

At locus 7q22.1, the signal was fine-mapped with very high confidence to rs4556017 (PP=96.2%), which exerts eQTL effects on *ACHE* and *SRRT*, both showing high levels of expression in enlarged haemorrhoidal tissue, and found in the M4 and M7 coexpression modules (figure 5), respectively. However, while *SRRT* encodes a poorly characterised capped-RNA binding protein, *ACHE* appears a much better candidate as the gene encodes an enzyme that hydrolyses the neurotransmitter acetylcholine at neuromuscular junctions, and corresponds to the Cartwright blood group antigen Yt.<sup>27</sup> One more locus was fine-mapped to single-variant resolution with high confidence, SNP rs10956488 from locus 8q24.21 (PP=0.962), which is linked to an eQTL for *GSDMC*. Gasdermin C (encoded by *GSDMC*) is a poorly characterised member of the gasdermin family of proteins expressed in epithelial cells and in enlarged haemorrhoidal tissue, though the mechanisms of its eventual HEM involvement remain elusive.

Other genes were linked via eQTL to a variant mapped with PP >50% and were prioritised based on additional experimental evidence. Among these, the fine-mapped SNP rs6498573 from locus 16p13.11 (PP=63.2%) is associated with eQTL effects on *MYH11* (encoding muscle myosin heavy chain 11; online supplemental table 10), a gene coding for a smooth muscle myosin heavy chain that shows mRNA upregulation in enlarged haemorrhoidal tissue and constitutes a hub for the M1 coexpression module associated with ECM organisation and muscle contraction. Notable prioritised candidates were also observed at loci that were not fine-mapped, including *ELN* (encoding elastin, a key component of the ECM found in the connective tissue of many organs, highly expressed in enlarged haemorrhoidal tissue and hub gene for the M1 coexpression module), *COL5A2* (encoding type V collagen protein; highly expressed in HEM tissue and belong to M1 coexpression module), *PRDM6* (encoding a putative histone methyltransferase regulating vascular smooth muscle cells contractility, expressed in enlarged haemorrhoidal tissue and hub for the M1 module), and others (online supplemental table 7).

Finally, while no candidate genes could be highlighted from the top GWAS hit region on chr12q14.3 (rs11176001), both the second and third strongest GWAS signals were detected at loci linked to genes involved in the determination of blood groups, namely *ABO* (rs676996) and the Kell Blood Group Complex Subunit-Related Family member *XKR9* (rs1838392). In addition, blood group antigens are encoded at additional loci, including *ACHE* (rs4556017) and *XKR6* (identified by *MAGMA*). Imputation of human ABO blood types from genotype data revealed the O type to be associated with increased and A and B types decreased HEM risk, both in UKBB and GERA datasets (online supplemental figure 12).

## Localisation of selected HEM gene-encoded proteins

A site-specific analysis of selected candidate proteins for their localisation in anorectal tissues underlined the complex multifactorial nature of cellular components potentially involved in the pathogenesis of HEM. Indeed, analysed candidates displayed a broad spectrum of expression in intestinal mucosal, neuromuscular, immune and anodermal tissues (figure 6, online supplemental figure 13, online supplemental tables 11 and 12). They were also directly colocalised with haemorrhoidal blood vessels, suggesting a putative role connected with the haemorrhoidal vasculature itself.

## DISCUSSION

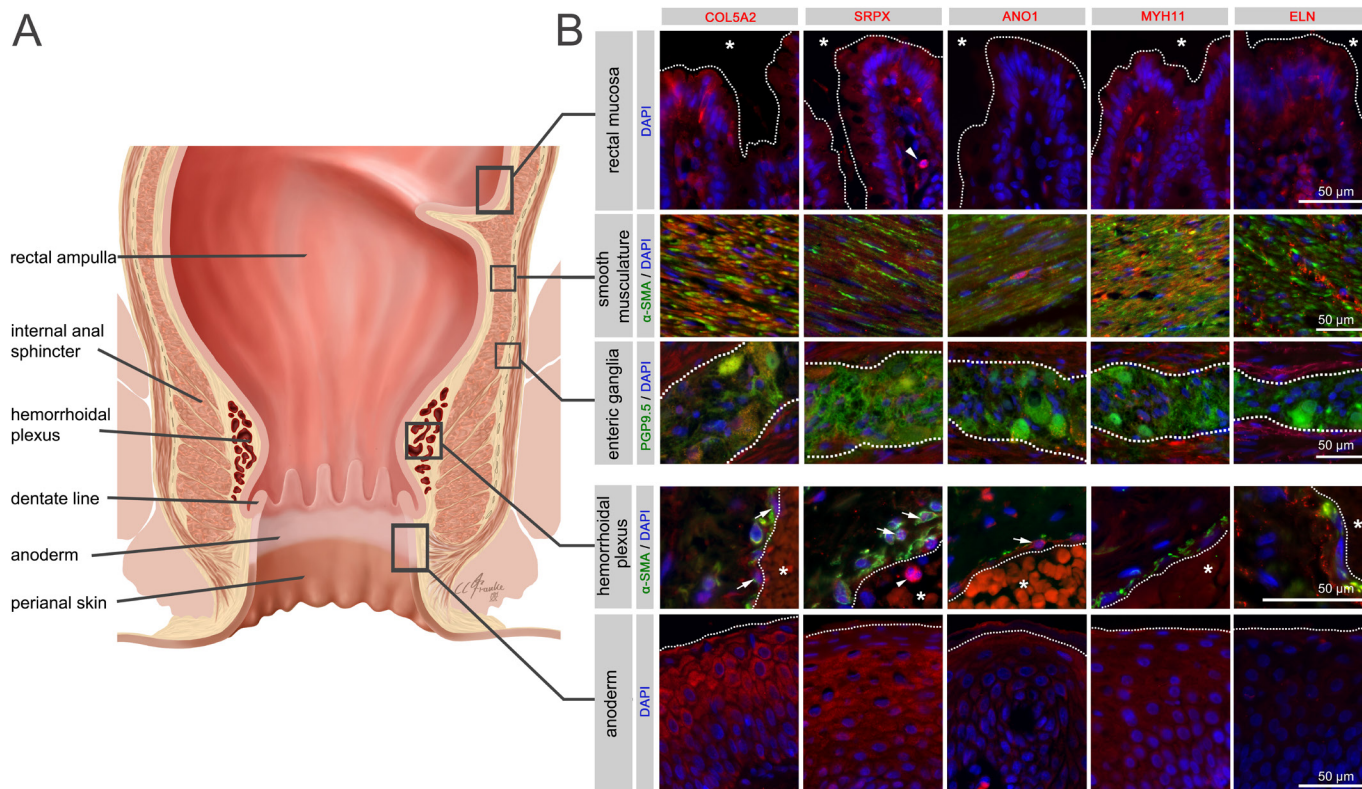
Given the lack of large and systematic epidemiological and molecular studies, and despite its worldwide distribution, HEM still can be regarded as an understudied disease. In this study, we report the largest and most detailed genome-wide analysis of HEM, implemented via a combination of classical GWAS approaches and the use of minimal phenotyping, as recently shown to be effective in boosting sample size for increased statistical power.<sup>28</sup> We demonstrate for the first time that HEM is a partly inherited condition with a weak but detectable heritability estimated at 5% based on SNP data. We identify 102 independent risk loci, which were functionally annotated based on computational predictions and gene expression analysis of diseased and normal tissue. These loci alone explain approximately 0.9% of HEM heritability, which is of similar magnitude with respect to the genetic contribution of other common complex traits.<sup>29</sup> They provide important novel pathophysiological insight, which we discuss below in relation to individual pathways and mechanisms proposed to contribute to HEM aetiology.

## ECM, elasticity of the connective tissue and smooth muscle function

The non-vascular components of the anal cushions consist of the transitional epithelium, connective tissue (elastic and collagenous) and the submucosal anal muscle (muscle of Treitz).<sup>5</sup> Treitz's muscle tightly maintains the anal cushions in their normal position, and its deterioration is considered one of the most important pathogenetic factors in the formation of enlarged and prolapsed haemorrhoids (online supplemental figure 2). Anal cushion fixation is further facilitated by elastic and collagenous connective tissue, whose degeneration, due, for instance, to abnormalities in collagen composition, has been involved in HEM aetiology,<sup>30</sup> although without strong molecular evidence.

The coexpression module M1 identified from HEM tissue is linked to ECM organisation and muscle function and is enriched for HEM gene thus providing novel important evidence for a role of these two interconnected processes in HEM pathogenesis. *ELN* (lead SNP rs11770437) is one of the HEM prioritised genes, and also a main hub gene for the M1 module. *ELN* codes for the elastin protein, a key component of elastic fibres that comprise part of the ECM and confer elasticity to organs and tissues including blood vessels. Mutations in the *ELN* gene have been shown to cause *cutis laxa*, a disease in which dysfunctional elastin interferes with the formation of elastic fibres, thus weakening connective tissue in the skin and blood vessels. Of note, *ELN* has been recently implicated also in diverticular disease,<sup>31</sup> a condition characterised by outpouchings of the colonic wall at sites of relative weakness and/or defective elasticity of the connective and muscle layers, as well as in the common skin condition non-syndromic *striae distensae* (NSD, also known as stretch marks), whose manifestation is due to lost tissue elasticity





**Figure 6** Immunohistochemistry for selected haemorrhoidal disease (HEM) candidate proteins. Illustration of the rectum and anal canal (A) indicating the site-specific localisation of the immunohistochemical panels analysed in (B). Results of fluorescence immunohistochemistry are shown for selected HEM candidate proteins encoded by HEM genes *COL5A2* (rs16831319), *SRPX* (rs35318931), *ANO1* (rs2186797), *MYH11* (rs6498573) and *ELN* (rs11770437) (see also online supplemental table 11 and online supplemental figure 13). Antibody staining was performed on FFPE colorectal tissue specimens from control individuals. Picture layers correspond to the rectal mucosa (top row, epithelial surface delimited by a white dotted line, \*=intestinal lumen), smooth musculature (second row), enteric ganglia (third row, ganglionic boundaries delimited by a white dotted line), haemorrhoidal plexus (fourth row, endothelial surface delimited by a white dotted line, \*=vascular lumen) and the anoderm (bottom row, surface of the anoderm delimited by a white dotted line). Blue: DAPI; green:  $\alpha$ -SMA (anti-alpha smooth muscle actin antibody) for rows 2 and 4 (smooth musculature/haemorrhoidal plexus) and PGP9.5 (member of the ubiquitin hydrolase family of proteins, neuronal marker) for row 3 (enteric ganglia); red: antibody for the respective candidate protein. Arrows point to respective candidate-positive cells within the vascular wall. Arrowheads point to respective candidate-positive nucleated immune cells.

at affected skin sites.<sup>32</sup> Hence, similar mechanisms may underly HEM risk due to genetic variation in *ELN* and, notably, also the *SRPX* (lead SNP rs35318931) and *COL5A2* (lead SNP rs16831319) genes. The *SRPX* lead SNP is also strongly associated with NSD, and is a coding variant potentially impacting the function of an ECM protein. *COL5A2* codes for type V fibril-forming collagen that has regulatory roles during development and growth of type I collagen-positive tissues. Mutations in this gene are known to cause Ehlers-Danlos syndrome, a rare connective tissue disease that affects the skin, joints and blood vessels, and for which a link to HEM has already been postulated.<sup>33</sup> An additional hub gene for the coexpression module M1, *MYH11* (lead SNP rs6498573) encodes a smooth muscle myosin protein that is important for muscle contraction and relaxation, and whose dysfunction has been linked to vascular diseases<sup>34</sup> and GI dysmotility.<sup>35</sup> Functional studies in smooth muscle cells showed that overexpression of *MYH11* led to a paradoxical decrease of protein levels through increased autophagic degradation followed by disruption of contractile signalling.<sup>36</sup> Our transcriptome analysis also showed *MYH11* RNA upregulation in HEM tissue, thus suggesting possible alterations in smooth muscle action that may be relevant, for instance, to the Treitz's muscle function(s). Additional evidence for the involvement of the musculoskeletal system may come from the associations detected

for *GSDMC* (lead SNP rs10956488), an uncharacterised gene also shown to be relevant to lumbar disc herniation and back pain<sup>37</sup> and *PRDM6*, a histone methyltransferase that acts as a transcriptional repressor of smooth muscle gene expression. Finally, besides individual association signals, we detected significant genetic correlation with several other complex diseases of shared aetiology (hernia, dorsalgia), for which connective tissue and/or muscle alterations are described.

### Gut motility

Several lines of evidence link gut motility to the pathophysiology of HEM in this study. In our UKBB analyses, patients with HEM were found to suffer more often from IBS and other dysmotility syndromes than controls, as evidenced also by the increased use of medications including laxatives. These conditions also showed strongest genetic correlation with HEM among all tested traits, indicating similar genetic architecture and predisposing mechanisms. Moreover, given the important role of the gut-brain axis in IBS and other FGIDs,<sup>38</sup> it is possible that the correlations observed for anxiety, depression and other neuroaffective traits may mediate genetic risk effects at least in part via similar mechanisms also involving gut motility. Constipation and prolonged sitting and straining during defecation are associated

with delayed GI transit and reduced peristalsis and are among the proposed HEM risk factors (online supplemental figure 2). Harder stools, due for instance to infrequent defecations, can cause difficulty in bowel emptying and therefore increase pressure and mechanical friction on the haemorrhoidal cushions, leading to excessive engorgement and stretching or tissue damage. The relationship between HEM and gut motility is probably best evidenced by the association signal detected at the *ANO1* locus: its lead SNP rs2186797 corresponds to the missense variant *ANO1*:p.Phe608Ser (F608S) that was fine-mapped with very high confidence and shown to impact anoctamin-1 function. Anoctamin-1 is an ion channel expressed in the ICCs, the pacemakers of the GI tract controlling intestinal peristalsis, has already been implicated in IBS,<sup>39</sup> and is also expressed in the vasculature. Hence, it represents an ideal candidate to also affect HEM risk via genotype-driven modulation of ICC function. An additional interesting candidate is *ACHE*, which shows an eQTL for the fine-mapped lead SNP rs4556017: *ACHE* codes for an enzyme that hydrolyses the neurotransmitter acetylcholine at neuromuscular junctions and is overexpressed in Hirschsprung's disease, a condition in which gut motility is compromised due to the absence of nerve cells (aganglionosis) in the distal or entire segments of the large bowel.<sup>40</sup> Of interest, expression in the enteric ganglia next to the haemorrhoidal plexus was observed for several proteins encoded by HEM genes in our immunofluorescence experiments.

### Vasculature and circulatory system

Previous observations showed that HEM is not varicosities and accelerated blood flow velocities were observed in afferent vessels of patients with HEM.<sup>41</sup> An impaired drainage or filling of the anal cushion may contribute to cushion slippage and may thus be considered as one of many disease-causing factors, as already previously proposed.<sup>41</sup> Our genetic data support the involvement of the vasculature as an important player in HEM pathophysiology. TSEA and GSEA results downstream of HEM GWAS meta-analysis highlight blood vessels and artery morphogenesis among the HEM gene-enriched tissues and GO pathways, respectively. At the same time, moderate genetic correlation is detected for diseases of the circulatory system in the LDSC analyses. We identified a very strong association signal in correspondence of the ABO locus on chromosome 9, which determines the corresponding ABO blood group type (A, B, AB and O). In addition to red blood cells, ABO antigens are expressed on the surface of many cells and tissues, and have been strongly associated with coronary artery disease, thrombosis, haemorrhage,<sup>42</sup> GI bleeding<sup>43</sup> and other conditions related to the circulatory system. Interestingly, increased risk for coronary artery disease has been reported in patients with HEM in at least some studies<sup>44</sup> and replicated in our UKBB cross-disease analyses although amidst other hundreds of associations of similar magnitude of effects. We imputed ABO blood types from genotype data in UKBB, and detected increased HEM risk for carriers of the O type. O type has been reported to be protective for coronary artery disease in UKBB, although also predisposing to hypertension.<sup>45</sup> Hence, the potential mechanism(s) by which variation at this locus impacts HEM risk remain elusive at this stage. Blood antigens are nevertheless likely relevant, as other genes involved in the determination of specific blood groups are also among the 102 HEM GWAS hits (Kell blood group locus *XKR9*, lead SNP rs1838392).

In summary, our data provide important new insight into currently lacking evidence<sup>46</sup> about HEM pathogenesis (online

supplemental figure 14). This sets the stage for more detailed genetic and mechanistic follow-up analyses, the search for therapeutically actionable genes and pathways, and the eventual exploitation for the adoption of preventive measures based on computed individual predisposition.

### Author affiliations

- <sup>1</sup>School of Biological Sciences, Monash University, Clayton, Victoria, Australia
- <sup>2</sup>Unit of Clinical Epidemiology, Department of Medicine Solna, Karolinska Institutet, Stockholm, Sweden
- <sup>3</sup>Institute of Clinical Molecular Biology, Christian-Albrechts-University of Kiel, Kiel, Germany
- <sup>4</sup>Novo Nordisk Foundation Center for Protein Research, Disease Systems Biology, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark
- <sup>5</sup>Institute of Biotechnology, Life Science Centre, Vilnius University, Vilnius, Lithuania
- <sup>6</sup>Institute of Anatomy, Christian-Albrechts-University of Kiel, Kiel, Germany
- <sup>7</sup>Department for General, Visceral, Vascular and Transplantation Surgery, University Medical Center Rostock, Rostock, Germany
- <sup>8</sup>Institute of Genomics, University of Tartu, Tartu, Estonia
- <sup>9</sup>K.G. Jebsen Center for Genetic Epidemiology, Department of Public Health and Nursing, NTNU, Norwegian University of Science and Technology, Trondheim, Norway
- <sup>10</sup>Enteric NeuroScience Program, Division of Gastroenterology and Hepatology, Mayo Clinic, Rochester, MN, USA
- <sup>11</sup>Institute of Advanced Research, Tokyo Medical and Dental University, Tokyo, Japan
- <sup>12</sup>Department of General-, Visceral- Transplant-, Thoracic and Pediatric Surgery, Kiel University, Kiel, Germany
- <sup>13</sup>Medical Office for Surgery Preetz, Preetz, Germany
- <sup>14</sup>Medical Department 1, University Hospital Dresden, Technische Universität Dresden (TU Dresden), Dresden, Germany
- <sup>15</sup>University Hospital Schleswig-Holstein, Institute of Clinical Chemistry, Thrombosis & Hemostasis Treatment Center, Campus Kiel & Lübeck, Kiel, Germany
- <sup>16</sup>Department of Internal Medicine, Hospital Oberndorf, Teaching Hospital of the Paracelsus Private Medical University of Salzburg, Oberndorf, Austria
- <sup>17</sup>Medical office for Colo-Proctology Kiel, Kiel, Germany
- <sup>18</sup>Department of Clinical Immunology, Aarhus University Hospital, Aarhus, Denmark
- <sup>19</sup>Department of Internal Medicine III, University Hospital Schleswig-Holstein, Campus Kiel, Kiel, Germany
- <sup>20</sup>Department of Internal Medicine III, University Hospital Heidelberg, Heidelberg, Germany
- <sup>21</sup>DZHK (German Centre for Cardiovascular Research), Partner Site Hamburg/Kiel/Lübeck, Kiel, Germany
- <sup>22</sup>Department of Biostatistics, University of Michigan, Ann Arbor, MI, USA
- <sup>23</sup>Epidemiology of highly pathogenic microorganisms, Robert Koch-Institute, Berlin, Germany
- <sup>24</sup>Department of Primatology, Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany
- <sup>25</sup>Department of Medicine I, Institute of Cancer Research, Medical University Vienna, Vienna, Austria
- <sup>26</sup>Center for Regenerative Therapies Dresden (CRTD), Technische Universität (TU) Dresden, Dresden, Germany
- <sup>27</sup>Department of Proctological Surgery Park Klinik Kiel, Kiel, Germany
- <sup>28</sup>Proctological Office Kiel, Kiel, Germany
- <sup>29</sup>Research Institute for Internal Medicine, Division of Surgery, Inflammatory Diseases and Transplantation, Oslo University Hospital Rikshospitalet and University of Oslo, Oslo, Norway
- <sup>30</sup>Faculty of Medicine, University of Cologne, Cologne, Germany
- <sup>31</sup>Department of Gastroenterology, Institute for Digestive Research, Lithuanian University of Health Sciences, Kaunas, Lithuania
- <sup>32</sup>University of Lübeck, Lübeck, Germany
- <sup>33</sup>University Hospital of Schleswig-Holstein (UKSH), Kiel Campus, Kiel, Germany
- <sup>34</sup>Department of Pathology, University Medical Center Schleswig-Holstein, Kiel, Germany
- <sup>35</sup>Division of Endocrinology, Diabetes and Clinical Nutrition, Department of Medicine 1, University of Kiel, Kiel, Germany
- <sup>36</sup>Institute of Epidemiology, Christian-Albrechts-University of Kiel, Kiel, Germany
- <sup>37</sup>University of Athens, Athens Euroclinic, Athens, Greece
- <sup>38</sup>Pathology Unit, German Primate Center, Leibniz Institute for Primatology, Göttingen, Germany
- <sup>39</sup>Department of Ophthalmology, Hospital Frankfurt Hoechst, Frankfurt, Germany
- <sup>40</sup>Upper Gastrointestinal Surgery, Department of Molecular Medicine and Surgery, Karolinska Institutet, Stockholm, Sweden
- <sup>41</sup>Department of Medicine, Levanger Hospital, Nord-Trøndelag Hospital Trust, Levanger, Norway
- <sup>42</sup>Department of Clinical Immunology, Naestved Hospital, Naestved, Denmark
- <sup>43</sup>General Hospital Lüneburg, Lüneburg, Germany



- <sup>44</sup>Department of General and Thoracic Surgery, Hospital Osnabrück, Osnabrück, Germany
- <sup>45</sup>Department of Clinical and Molecular Medicine, Norwegian University of Science and Technology, Trondheim, Norway
- <sup>46</sup>BioCore – Bioinformatics Core Facility, Norwegian University of Science and Technology, Trondheim, Norway
- <sup>47</sup>Clinic of Laboratory Medicine, St.Olavs Hospital, Trondheim University Hospital, Trondheim, Norway
- <sup>48</sup>Department of Clinical Immunology, Copenhagen University Hospital, Copenhagen, Denmark
- <sup>49</sup>Department of Surgery, Community Hospital Kiel, Kiel, Germany
- <sup>50</sup>Institute for Community Medicine, University Medicine Greifswald, Greifswald, Germany
- <sup>51</sup>Department Research & Conservation, Zoo Nuremberg, Nuremberg, Germany
- <sup>52</sup>Department of Gastroenterology, Helios Hospital Weißeritztal, Freital, Germany
- <sup>53</sup>23andMe Inc, Sunnyvale, CA, USA
- <sup>54</sup>Christian-Albrechts-University of Kiel, Kiel, Germany
- <sup>55</sup>Gastrointestinal Genetics Lab, CIC bioGUNE - BRTA, Derio, Spain
- <sup>56</sup>Ikerbasque, Basque Foundation for Science, Bilbao, Spain

**Twitter** Malte Christoph Rühlemann @mruehlemann, Arthur Beyder @beyderlab and Mauro D'Amato @damato\_mauro

**Contributors** AF, CS and MD'A were responsible for the concept and the design of the study. CS and GB coordinated the recruitment of the German clinical cohort. VK, JG, IV, BS, HGP, AH, JJ, JW, TB, MD, FB, SH, JK, JS, WK, MZo, CGN, TL, JT and WS recruited and phenotyped patients of the clinical cohorts. VK collected and phenotyped the biopsy samples that were used in the immunohistochemistry and transcriptome analyses. Animal samples for histology were taken and provided by KP, SZ, LF, TG, FHL and KM-R. Immunohistochemistry and histology analyses were performed by FC and TW with support of JW and CR. TZ, DE and FD performed genotype and phenotype data collection. TZ, DE, and FD performed GWAS data quality control and imputation. SJ performed the transcriptome (RNA-Seq) analysis with help from GI and PR. FU-W, SBu, and MF helped with bioinformatic analyses. IFJ, KB, SBr, CLR, CE, OBP, and HU conducted analyses in DBDS and DNPR. HC, VL, RJ, NF, WL, ML, and UN-G contributed German control data. SBu, CD, HV, AG and JH contributed the GWAS data on diverticular disease. Norwegian HUNT data were contributed by AHS, THK, MEG, LT, EN-J and KH. The Estonian dataset was provided by TE and MT-L. Protein models and analyses were performed by GM. The analysis of the MGI data set was performed by MZa, BV, AP and LGF. The 23andMe data set was analyzed and provided by OS, ESN, and VV. MW implemented ABO blood-group inference. SC, PRS, GF and AB performed site-directed ANO1 mutagenesis and whole-cell electrophysiology experiments. TZ, DE implemented statistical models and performed the (meta-)analysis. TZ, DE, and SJ curated and interpreted results. MS designed online supplementary figure S2 with scientific input from NM and AF. The GWAS data browser was implemented by MR and set up by GH-S. AF, MD'A and DE wrote the manuscript draft with substantial contributions from TZ and SJ. All authors reviewed, edited and approved the final manuscript.

**Funding** This project was funded by Andre Franke's and Clemens Schafmayer's DFG grant "Discovery of risk factors for hemorrhoids" (ID: FR 2821/19-1). The study received infrastructure support from the DFG Cluster of Excellence 2167 "Precision Medicine in Chronic Inflammation (PMI)" (DFG Grant: "EXC2167"). The project was supported by grants from the Swedish Research Council to MD (VR 2017-02403), the Novo Nordisk Foundation (grants NNF17OC0027594 and NNF14CC0001) and BigTempHealth (grant 5153-00002B). We are indebted to the valuable assistance by Tanja Wesse (Genotyping of the German samples) and Petra Röthgen (German DZHK control data set). We thank Clemens Franke (Institute of Anatomy, Christian-Albrechts-University of Kiel, Kiel, Germany) for graphic assistance (figure 6, online supplementary figure S1 and S13). This research has been conducted using the UK Biobank Resource under Application Number 31435. The work on cross-trait analysis for diverticular disease presented in this manuscript was supported by the German Research Council (DFG, ID: Ha3091/9-1) and the Austrian Science Fund (FWF, ID: I1542-B13). Data access to the UK biobank data for diverticular disease was granted under project numbers 22691. EGCUT work has also been supported by the European Regional Development Fund and grants SP1G118045T, No. 2014-2020.4.01.15-0012 GENTRANSMED and 2014-2020.4.01.16-0125 This study was also funded by EU H2020 grant 692145, Estonian Research Council Grant PUT1660. Data analyzes with Estonian datasets were carried out in part in the High-Performance Computing Center of University of Tartu. The work on site-directed mutagenesis of ANO1 and whole-cell electrophysiology was funded by NIH DK057061. We would like to thank the research participants and employees of 23andMe for making this work possible. The following members of the 23andMe Research Team contributed to this study: Michelle Agee, Adam Auton, Robert K. Bell, Katarzyna Bryc, Sarah L. Elson, Pierre Fontanillas, Nicholas A. Furlotte, David A. Hinds, Karen E. Huber, Aaron Kleinman, Nadia K. Litterman, Jennifer C. McCreight, Matthew H. McIntyre, Joanna L. Mountain, Elizabeth S. Noblin, Carrie A.M. Northover, Steven J. Pitts, J. Fah Sathirapongsasuti, Olga V. Sazonova, Janie F. Shelton, Suyash Shringarpure, Chao Tian, Joyce Y. Tung, and Vladimir Vacić.

**Competing interests** Vladimir Vacić and Olga V. Sazonova are/were employed by and hold stock or stock options in 23andMe, Inc. All other authors have nothing to declare.

**Patient consent for publication** Not required.

**Provenance and peer review** Not commissioned; externally peer reviewed.

**Data availability statement** Data are available in a public, open access repository. Genome-wide summary statistics of our analyses are publicly available through our web browser (<http://hemorrhoids.online>) and have been submitted to the European Bioinformatics Institute ([www.ebi.ac.uk/gwas](http://www.ebi.ac.uk/gwas)) under the accession number GCST90014033. RNA-seq data have been deposited at NCBI Gene Expression Omnibus (GEO) under the accession number GSE154650.

**Supplemental material** This content has been supplied by the author(s). It has not been vetted by BMJ Publishing Group Limited (BMJ) and may not have been peer-reviewed. Any opinions or recommendations discussed are solely those of the author(s) and are not endorsed by BMJ. BMJ disclaims all liability and responsibility arising from any reliance placed on the content. Where the content includes any translated material, BMJ does not warrant the accuracy and reliability of the translations (including but not limited to local regulations, clinical guidelines, terminology, drug names and drug dosages), and is not responsible for any error and/or omissions arising from translation and adaptation or otherwise.

**Open access** This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

#### ORCID iDs

Tenghao Zheng <http://orcid.org/0000-0003-1587-7481>  
 David Ellinghaus <http://orcid.org/0000-0002-4332-6110>  
 Simonas Juzenas <http://orcid.org/0000-0001-9293-0691>  
 Christian Datz <http://orcid.org/0000-0001-7838-4532>  
 Michael Forster <http://orcid.org/0000-0001-9927-5124>  
 Andrea Gsur <http://orcid.org/0000-0002-9795-1528>  
 Jochen Hampe <http://orcid.org/0000-0002-2421-6127>  
 Christoph Röcken <http://orcid.org/0000-0002-6989-8002>  
 Eivind Ness-Jensen <http://orcid.org/0000-0001-6005-0729>  
 Malte Christoph Rühlemann <http://orcid.org/0000-0002-0685-0052>  
 Gianrico Farrugia <http://orcid.org/0000-0003-3473-5235>  
 Mauro D'Amato <http://orcid.org/0000-0003-2743-5197>  
 Andre Franke <http://orcid.org/0000-0003-1530-5811>

#### REFERENCES

- Jacobs D. Clinical practice. *Hemorrhoids*. *N Engl J Med* 2014;371:944–51.
- Duthie HL, Cairns FW. Sensory nerve-endings and sensation in the anal region of man. *Br J Surg* 1960;47:585–95.
- Haas PA, Haas GP, Schmaltz S, Fox JA, Jr. The prevalence of hemorrhoids. *Dis Colon Rectum* 1983;26:435–9.
- Yang JY, Peery AF, Lund JL, et al. Burden and cost of outpatient hemorrhoids in the United States Employer-Insured population, 2014. *Am J Gastroenterol* 2019;114:798–803.
- Margetis N. Pathophysiology of internal hemorrhoids. *Ann Gastroenterol* 2019;32:264–72.
- Tung JY, Do CB, Hinds DA, et al. Efficient replication of over 180 genetic associations with self-reported medical data. *PLoS One* 2011;6:e23473.
- Bycroft C, Freeman C, Petkova D, et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature* 2018;562:203–9.
- Leitsalu L, Haller T, Esko T, et al. Cohort profile: Estonian Biobank of the Estonian genome center, University of Tartu. *Int J Epidemiol* 2015;44:1137–47.
- Fritsche LG, Gruber SB, Wu Z, et al. Association of polygenic risk scores for multiple cancers in a Phenome-wide study: results from the Michigan genomics initiative. *Am J Hum Genet* 2018;102:1048–61.
- Kvale MN, Hesselton S, Hoffmann TJ, et al. Genotyping informatics and quality control for 100,000 subjects in the genetic epidemiology research on adult health and aging (GERA) cohort. *Genetics* 2015;200:1051–60.
- Willer CJ, Li Y, Abecasis GR. Metal: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 2010;26:2190–1.
- Watanabe K, Taskesen E, van Bochoven A, et al. Functional mapping and annotation of genetic associations with FUMA. *Nat Commun* 2017;8:8.
- Benner C, Spencer CCA, Havulinna AS, et al. FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics* 2016;32:1493–501.
- Watanabe K, Stringer S, Frei O, et al. A global overview of pleiotropy and genetic architecture in complex traits. *Nat Genet* 2019;51:1339–48.



- 15 Buniello A, MacArthur JAL, Cerezo M, *et al.* The NHGRI-EBI GWAS catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res* 2019;47:D1005–12.
- 16 Kamat MA, Blackshaw JA, Young R, *et al.* PhenoScanner V2: an expanded tool for searching human genotype-phenotype associations. *Bioinformatics* 2019;35:4851–3.
- 17 Cuéllar-Partida G, Lundberg M, Kho PF, *et al.* Complex-Traits genetics virtual lab: a community-driven web platform for post-GWAS analyses. *bioRxiv* 2019;518027.
- 18 Choi SW, O'Reilly PF. PRSice-2: polygenic risk score software for biobank-scale data. *Gigascience* 2019;8. doi:10.1093/gigascience/giz082. [Epub ahead of print: 01 07 2019].
- 19 Krokstad S, Langhammer A, Hveem K, *et al.* Cohort profile: the HUNT study, Norway. *Int J Epidemiol* 2013;42:968–77.
- 20 Hansen TF, Banasik K, Erikstrup C, *et al.* Dbds genomic cohort, a prospective and comprehensive resource for integrative and temporal analysis of genetic, environmental and lifestyle factors affecting health of blood donors. *BMI Open* 2019;9:e028401.
- 21 Krawczak M, Nikolaus S, von Eberstein H, *et al.* PopGen: population-based recruitment of patients and controls for the analysis of complex genotype-phenotype relationships. *Community Genet* 2006;9:55–61.
- 22 Pers TH, Karjalainen JM, Chan Y, *et al.* Biological interpretation of genome-wide association studies using predicted gene functions. *Nat Commun* 2015;6:5890.
- 23 de Leeuw CA, Mooij JM, Heskes T, *et al.* MAGMA: generalized gene-set analysis of GWAS data. *PLoS Comput Biol* 2015;11:e1004219.
- 24 Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics* 2008;9:559.
- 25 Malysz J, Gibbons SJ, Saravanaperumal SA, *et al.* Conditional genetic deletion of *Ano1* in interstitial cells of Cajal impairs  $Ca^{2+}$  transients and slow waves in adult mouse small intestine. *Am J Physiol Gastrointest Liver Physiol* 2017;312:G228–45.
- 26 Naba A, Clauser KR, Whittaker CA, *et al.* Extracellular matrix signatures of human primary metastatic colon cancers and their metastases to liver. *BMC Cancer* 2014;14:518.
- 27 Bartels CF, Zelinski T, Lockridge O. Mutation at codon 322 in the human acetylcholinesterase (ACHE) gene accounts for YT blood group polymorphism. *Am J Hum Genet* 1993;52:928–36.
- 28 DeBoever C, Tanigawa Y, Aguirre M, *et al.* Assessing digital phenotyping to enhance genetic studies of human diseases. *Am J Hum Genet* 2020;106:611–22.
- 29 Okada Y, Wu D, Trynka G, *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature* 2014;506:376–81.
- 30 Nasserri YY, Krott E, Van Groningen KM, *et al.* Abnormalities in collagen composition may contribute to the pathogenesis of hemorrhoids: morphometric analysis. *Tech Coloproctol* 2015;19:83–7.
- 31 Schafmayer C, Harrison JW, Buch S, *et al.* Genome-Wide association analysis of diverticular disease points towards neuromuscular, connective tissue and epithelial pathomechanisms. *Gut* 2019;68:854–65.
- 32 Tung JY, Kiefer AK, Mullins M, *et al.* Genome-Wide association analysis implicates elastic microfibrils in the development of nonsyndromic striae distensae. *J Invest Dermatol* 2013;133:2628–31.
- 33 Plackett TP, Kwon E, Gagliano RA, *et al.* Ehlers-Danlos syndrome-hypermobility type and hemorrhoids. *Case Rep Surg* 2014;2014:1–3.
- 34 Zhu L, Vranckx R, Khau Van Kien P, *et al.* Mutations in myosin heavy chain 11 cause a syndrome associating thoracic aortic aneurysm/aortic dissection and patent ductus arteriosus. *Nat Genet* 2006;38:343–9.
- 35 Gilbert MA, Schultz-Rogers L, Rajagopalan R, *et al.* Protein-elongating mutations in MYH11 are implicated in a dominantly inherited smooth muscle dysmotility syndrome with severe esophageal, gastric, and intestinal disease. *Hum Mutat* 2020;41:973–82.
- 36 Kwartler CS, Chen J, Thakur D, *et al.* Overexpression of smooth muscle myosin heavy chain leads to activation of the unfolded protein response and autophagic turnover of thick filament-associated proteins in vascular smooth muscle cells. *J Biol Chem* 2014;289:14075–88.
- 37 Suri P, Palmer MR, Tsepilov YA, *et al.* Genome-Wide meta-analysis of 158,000 individuals of European ancestry identifies three loci associated with chronic back pain. *PLoS Genet* 2018;14:e1007601.
- 38 Holtmann G, Shah A, Morrison M. Pathophysiology of functional gastrointestinal disorders: a holistic overview. *Dig Dis* 2017;35 Suppl 1:5–13.
- 39 Mazzone A, Gibbons SJ, Eisenman ST, *et al.* Direct repression of anoctamin 1 (*ANO1*) gene transcription by Gli proteins. *FASEB J* 2019;33:6632–42.
- 40 Moore SW, Johnson G. Acetylcholinesterase in Hirschsprung's disease. *Pediatr Surg Int* 2005;21:255–63.
- 41 Aigner F, Gruber H, Conrad F, *et al.* Revised morphology and hemodynamics of the anorectal vascular plexus: impact on the course of hemorrhoidal disease. *Int J Colorectal Dis* 2009;24:105–13.
- 42 Franchini M, Lippi G. Relative risks of thrombosis and bleeding in different ABO blood groups. *Semin Thromb Hemost* 2016;42:112–7.
- 43 Huang ES, Khalili H, Strate LL, *et al.* Su1786 ABO blood group and the risk of gastrointestinal bleeding. *Gastroenterology* 2012;142:S-503.
- 44 Chang S-S, Sung F-C, Lin C-L, *et al.* Association between hemorrhoid and risk of coronary heart disease: a nationwide population-based cohort study. *Medicine* 2017;96:e7662.
- 45 Groot HE, Villegas Sierra LE, Said MA, *et al.* Genetically determined ABO blood group and its associations with health and disease. *Arterioscler Thromb Vasc Biol* 2020;40:830–8.
- 46 Sandler RS, Peery AF. Rethinking what we know about hemorrhoids. *Clin Gastroenterol Hepatol* 2019;17:8–15.

## **Supplementary Materials**

**Genome-wide analysis of 944,133 individuals provides insights into the etiology of hemorrhoidal disease**

## Table of Contents

<b>METHODS</b> .....	<b>4</b>
<b>Histology of hemorrhoidal plexus</b> .....	<b>4</b>
<b>Study cohorts and patients' material</b> .....	<b>5</b>
<i>23andMe</i> .....	5
<i>UK Biobank (UKBB)</i> .....	5
<i>Estonian Genome Project of University of Tartu (EGCUT)</i> .....	6
<i>Michigan Genomics Initiative (MGI)</i> .....	6
<i>Genetic Epidemiology Research on Aging</i> .....	6
<i>German case-control cohort</i> .....	7
<i>The Trøndelag Health Study (HUNT)</i> .....	7
<i>Danish Blood Donor Study</i> .....	8
<i>Danish National Patient Registry</i> .....	8
<i>Hemorrhoidal tissue</i> .....	9
<b>Genotyping, quality control and genotype imputation of cohorts included in this study</b> .....	<b>9</b>
<i>23andMe</i> .....	9
<i>UK Biobank</i> .....	11
<i>Estonian Genome Project of University of Tartu (EGCUT)</i> .....	12
<i>Michigan Genomics Initiative (MGI)</i> .....	13
<i>Genetic Epidemiology Research on Aging (GERA)</i> .....	14
<i>German case-control cohort</i> .....	14
<i>The Trøndelag Health Study (HUNT)</i> .....	15
<i>Danish Blood Donor Study</i> .....	17
<b>GWAS association analysis for discovery cohorts</b> .....	<b>18</b>
<i>23andMe</i> .....	18
<i>UK Biobank (UKBB), Estonian Genome Project of University of Tartu (EGCUT), Michigan Genomics Initiative (MGI), Genetic Epidemiology Research on Aging (GERA)</i> .....	19
<b>GWAS meta-analysis across discovery cohorts</b> .....	<b>19</b>
<b>Annotation of HEM GWAS risk loci and gene mapping</b> .....	<b>20</b>
<b>Bayesian fine-mapping analysis</b> .....	<b>20</b>
<b>Heritability analysis via linkage disequilibrium score regression (LDSC)</b> .....	<b>21</b>
<b>Genome-wide pleiotropy analysis</b> .....	<b>21</b>
<b>Tissue and pathway enrichment analyses</b> .....	<b>21</b>
<b>Polygenic risk scores (PRS) analysis</b> .....	<b>22</b>
<b>Phenome-wide association studies (PheWAS)</b> .....	<b>22</b>
<b>Cross-trait analyses</b> .....	<b>23</b>
<b>RNA library preparation and RNA-sequencing</b> .....	<b>23</b>
<b>Mapping and quality assessment of RNA-Seq data</b> .....	<b>24</b>
<b>Gene signature-based determination of anal canal zones</b> .....	<b>24</b>
<b>Differential gene expression analysis</b> .....	<b>25</b>
<b>Identification and characterization of enriched co-expression modules</b> .....	<b>25</b>
<b>ABO blood group analysis</b> .....	<b>26</b>
<b>Fluorescence Immunohistochemistry</b> .....	<b>26</b>



<b><i>In silico</i> variant protein analysis .....</b>	<b>27</b>
<b>Site-directed ANO1 mutagenesis and whole-cell electrophysiology .....</b>	<b>29</b>
<b>SUPPLEMENTARY FIGURES.....</b>	<b>31</b>
Online Supplementary Figure S1. <i>Histological analysis of the anorectum in four different species</i> .....	31
Online Supplementary Figure S2. <i>Suggested integrated model that summarizes the contemporary thinking on the pathophysiology of HEM (figure and legend are mainly taken from Figure 2 in Nikolaos Margetis' review[92])</i> .....	33
Online Supplementary Figure S3. <i>Schematic overview of the study workflow</i> .....	36
Online Supplementary Figure S4. <i>Quantile-quantile (QQ) plot of GWAS meta-analysis results</i> .....	37
Online Supplementary Figure S5. <i>Regional association plots of HEM GWAS risk loci</i> .....	50
Online Supplementary Figure S6. <i>Previously reported associations of HEM risk loci with other traits and diseases, clustered by biological areas</i> .....	51
Online Supplementary Figure S7. <i>Gene set enrichment analyses of HEM genes</i> .....	52
Online Supplementary Figure S8. <i>Gene signature-based determination of anal canal zones</i> .....	53
Online Supplementary Figure S9. <i>ANO1 Alignment of TM4-5 and ANO1 structure</i> .....	54
Online Supplementary Figure S10. <i>F608S mutant of ANO1 has high instantaneous current but slow voltage-dependent activation and deactivation kinetics in vitro</i> .....	57
Online Supplementary Figure S11. <i>Sushi repeat-containing protein (SRPX) structure und alignment</i> .....	59
Online Supplementary Figure S12. <i>ABO blood groups and HEM risk in UKBB and GERA</i> .....	60
Online Supplementary Figure S13. <i>Immunohistochemistry for selected HEM candidate proteins</i> .....	61
Online Supplementary Figure S14. <i>Graphical abstract of the study</i> .....	63
<b>SUPPLEMENTARY REFERENCES.....</b>	<b>64</b>

## METHODS

### Histology of hemorrhoidal plexus

For histologic examination and phylogenetic comparison of the hemorrhoidal plexus, formalin fixed anorectal specimens were obtained from *Homo sapiens*, *Gorilla gorilla gorilla*, baboon (*Papio anubis*), and mouse (10-week old male C57BL/6JRj mouse).

Human tissue was retrieved from a healthy donor (female, 54 years) who was recruited by the body donation program of the Institute of Anatomy, Kiel University. The donors had previously given written consent to the use of their samples for teaching and research purposes; the donors were free from diseases related to the gastrointestinal tract and the anorectum. The gorilla specimen comes from a 43-year-old female western lowland gorilla from Nuremberg Zoo (Germany). The animal had to be euthanized due to a terminal metastatic adenocarcinoma of the uterus. Rectum and surrounding tissue were removed during post-mortem examination four hours after death, cut and fixed in 10% neutral buffered formalin. Rectum samples from the baboon were taken from a ten year old male olive baboon kept at the German Primate Center Göttingen and included in a study authorized by the governmental veterinary authority, i.e. the Lower Saxony State Office for Consumer Protection (Food Safety Ref. No. 33.19-42502-04-18/3036 according to the regulations of the German Welfare Act and the European Directive 2010/63/EU on the protection of animals used for experimental and other scientific purpose). The rectal specimen was collected during routine necropsy following a standardized necropsy protocol and fixed in 10% buffered formalin.

All tissue samples were taken from the anal canal at the level of the hemorrhoidal plexus, dehydrated, embedded in paraffin wax, cut into sections (6 µm) and processed for hematoxylin-eosin and Azan stainings. The findings were evaluated and documented with a Keyence microscope (BZ-X800) using the integrated stitching tool BZ-X800 Analyser software version 1.1.1.8.

## Study cohorts and patients' material

### **23andMe**

The 23andMe study dataset contains participants drawn from the research participant base of the personal genetics company, 23andMe, Inc[1]. Genetic data and comprehensive phenotypic information from health surveys were available for 402,845 unrelated individuals of European ancestry. Study participants were divided into HEM cases and controls based on their self-completed HEM health questionnaires, resulting in 174,785 HEM cases and 228,060 controls in the current 23andMe GWAS. Demographic data of 23andMe samples are reported in **online supplementary table S1**. Participants provided informed consent and participated in the research online, under a protocol approved by the external AAHRPP-accredited IRB, Ethical & Independent Review Services (E&I Review). The full GWAS summary statistics for the 23andMe discovery data set will be made available through 23andMe to qualified researchers under an agreement with 23andMe that protects the privacy of the 23andMe participants. Please visit <https://research.23andme.com/dataset-access> for more information and to apply to access the data.

### **UK Biobank (UKBB)**

The UKBB is a large population-based study in the United Kingdom with extensive phenotypic and genotype data from approximately 500,000 participants[2]. Each individual underwent cognitive, physical assessment and sampling for DNA collection when enrolled, and health-related information was collected including data from their electronic health records (EHRs). The diagnoses in the EHRs are coded in the terminology of the International Statistical Classification of Diseases and Related Health Problems (ICD) terminology. For this GWAS study we included 408,592 individuals of European ancestry (self-reported "white" and of genetic Caucasian descent). Of these, 23,856 samples met our criteria for HEM cases (either ICD10 code I84 or ICD9 code 455 in the medical records). The other part of the cohort (n=384,736) served as study controls. The demographic data of the individuals are reported in **online supplementary table S1**. UKBB received ethical approval from the competent Research Ethics Committee (REC reference 11/NW/0382) and the project's Application ID is 31435.

### ***Estonian Genome Project of University of Tartu (EGCUT)***

The Estonian Biobank is a population-based cohort of the Estonian Genome Center at the University of Tartu (EGCUT), Estonia, with a current size of app. 200,000 participants aged over 18[3]. The whole project is conducted according to the Estonian Gene Research Act and all participants have signed the broad informed consent. Upon recruitment, the biobank participants filled out a detailed questionnaire, covering lifestyle, diet and clinical diagnoses (described by ICD10 codes). In this study, individuals with any entry of the ICD10 code for HEM (I84) were included as HEM cases. Further, we selected 30,441 controls with genome-wide data as study controls, resulting in 6,956 HEM cases and 30,441 population controls. The demographic data of the individuals are reported in **online supplementary table S1**. This study has been reviewed and approved by the Estonian Committee on Bioethics and Human Research.

### ***Michigan Genomics Initiative (MGI)***

The Michigan Genomics Initiative (MGI) is a longitudinal cohort of participants in Michigan Medicine, USA[4]. MGI participants were recruited primarily through surgical procedures at Michigan Medicine and gave consent for link their EHRs and genetic data for research purposes. We used a current data freeze of 40,000 European individuals for GWAS analysis. Of these, 4,539 HEM cases were defined based on a review of EHRs (either ICD10 code I84, ICD10-CM code K64 or ICD9 code 455). The rest of the cohort with genome-wide data was defined as study controls (n=35,338). The demographic data of the individuals are reported in **online supplementary table S1**. This study has been reviewed and approved by the Michigan Institutional Review Board.

### ***Genetic Epidemiology Research on Aging***

The Genetic Epidemiology Research on Aging (GERA) Cohort comprises more than 100,000 adults who are members of the Kaiser Permanente Medical Care Plan, Northern California Region (KPNC), USA. The health status of participants in the GERA cohort was assessed using EHRs collected at Kaiser Permanente's facilities in Northern California from January 1, 1995 to March 15, 2013. HEM cases (n=8,813) were those in which at least two ICD9 code diagnoses of HEM (ICD9 code 455) were recorded on separate days. Their genome-wide data were compared with those of the

remaining cohort (n=46,780) as controls. The demographic data of the individuals are reported in **online supplementary table S1**. The GERA data access was applied for on the dbGaP website (dbGaP Study Accession: phs000674.v3.p3) and the study was approved by the dbGap Access Review Committee.

### ***German case-control cohort***

Initiated by the Department of General and Thoracic Surgery and the biobank PopGen[5] of the Medical Faculty of Kiel University, Kiel, Germany, a cohort of HEM patients with symptomatic hemorrhoids and the need of invasive treatment was newly established. Between January 2016 and December 2017, individuals with a prior diagnosis of high-grade hemorrhoids were identified based on the medical records of five hospitals and practices in the North German region using German procedural codes (OPS-301 by German Institute for medical Documentation and Information). The main inclusion criteria were the need for hemorrhoidectomy or invasive treatment (rubber band ligation, sclerotherapy) on high grade hemorrhoidal disease, verified by DRG-code (Diagnosis related Groups). Patients receiving exclusively conservative treatment were not included in this study as the aim was to recruit patients with a strong phenotype of advanced hemorrhoidal disease. The cohort included 1,007 patients undergoing surgical/invasive treatment of a high grade hemorrhoidal disease. In total, 1,144 cases and 2,740 controls were available for PRS analysis (section **Polygenic risk score (PRS) analysis, Methods**). The demographic data of the individuals are reported in **online supplementary table S1**. The study protocol was approved by the ethics committee (ref: A156/03-1/15) of the Medical Faculty of Kiel University and written informed consent was obtained from all study participants.

### ***The Trøndelag Health Study (HUNT)***

The Trøndelag Health Study (HUNT) is a large population-based cohort from the county Nord-Trøndelag in Norway. All residents in the county, aged 20 years and older, have been invited to participate. Data was collected through three cross-sectional surveys, HUNT1 (1984-1986), HUNT2 (1995-1997) and HUNT3 (2006-2008), and has been described in detail previously[6], with the fourth survey recently completed (HUNT4, 2017-2019). All genotyped participants have signed a written informed consent regarding the use of data from questionnaires, biological samples and linkage to other registries for research purposes. Cases were defined as having



an ICD10 K64 diagnosis and the remainder of the cohort were used as controls. In total, 977 cases and 68,314 controls were available for PRS analysis (section **Polygenic risk score (PRS) analysis, Methods**). The demographic data of the individuals are reported in **online supplementary table S1**.

### ***Danish Blood Donor Study***

The Danish Blood Donor Study (DBDS) is a large prospective cohort of nation-wide Danish blood donors (n=56,397) and comprises both extensive phenotype data as well as genome-wide genotyping data[7, 8]. HEM cases are defined using the ICD-8 code 455 or ICD-10 codes I84 or K64, resulting in 1,754 cases in the DBDS cohort as registered in the National Patient Registry. In total, 1,754 cases and 54,643 controls were available for genome-wide polygenic risk score (PRS) analysis. The demographic data of the individuals are reported in **online supplementary table S1**. This study was approved according to the Danish Blood Donor study protocol (ref: 1700407) as a part of “Genetics of healthy ageing and specific diseases among blood donors”.

### ***Danish National Patient Registry***

The Danish National Patient Registry (DNPR) is a population-wide registry containing all diagnoses made in hospitals in Denmark from 1977 to 2018 and includes more than 8 million patients. The diagnoses in the registry are coded in the terminology of the International Statistical Classification of Diseases and Related Health Problems (ICD) 8th Revision (1997-1993) or 10th Revision (1994-2018) terminology. All patients with a hemorrhoid disease code in the disease registry were identified. In the ICD-8 period patients with ‘Hemorrhoids’ are recognized using the code 455. As the ICD-10 code for hemorrhoids changed in 2013, we combined patients diagnosed with ‘Hemorrhoids’ (ICD-10: I84) from 1994-2012 and patients diagnosed with ‘Hemorrhoids and perianal venous thrombosis’ (K64) from 2013-2018. Information about drugs administered in hospitals is available for more than 1.6 million patients in the period 2006-2016 and is defined using the Anatomical Therapeutic Chemical (ATC) Classification System of the WHO. We integrated data from two different electronic medication modules corresponding to the administrative databases for hospital internal drug consumption from two health regions of Denmark (Capital Region and Region Zealand): OPUS-medicin and Elektronisk patient medicinering[9].

The demographic data of the individuals are reported in **online supplementary table S1**. This DNPR study has been approved by the Danish Data Protection Agency, Copenhagen (ref: FSEID-00003092, FSEID-00003724 and 3-3013-1731/1).

### ***Hemorrhoidal tissue***

A group of 38 individuals undergoing surgery for hemorrhoids (n=20; cases) and anal fissures (n=18; controls) was included in this study. Hemorrhoidal tissue samples were obtained either by Milligan Morgan open hemorrhoidectomy or by stapled hemorrhoidopexy for grade (Goligher) 3 and 4 hemorrhoids. Approximately 1cm<sup>3</sup> of hemorrhoidal tissue was obtained from Milligan-Morgan-specimens just above the dentate line. In hemorrhoidopexy patients, approximately 1cm<sup>3</sup> of biopsies were taken from the “doughnut” tissue at 3 o’clock in the prone position. Healthy hemorrhoid tissue samples were taken from adjacent zones (1-2 cm above the dentate line) of anal fissures. Clinical and demographic data for the sampled individuals are listed in **online supplementary table S1**. This study was approved by the bioethical committee of medical faculty, University Hospital Schleswig-Holstein Kiel, Germany. All participants provided written informed consent.

## **Genotyping, quality control and genotype imputation of cohorts included in this study**

### ***23andMe***

DNA extraction and genotyping were performed on saliva samples by National Genetics Institute (NGI), a CLIA licensed clinical laboratory and a subsidiary of Laboratory Corporation of America. Samples had been genotyped on one of four genotyping platforms. The V1 and V2 platforms were variants of the Illumina HumanHap550+ BeadChip, including about 25,000 custom SNPs selected by 23andMe, with a total of about 560,000 SNPs. The V3 platform was based on the Illumina OmniExpress+ BeadChip, with custom content to improve the overlap with the V2 array, with a total of about 950,000 SNPs. The V4 platform is a fully custom array, including a lower redundancy subset of V2 and V3 SNPs with additional coverage of lower-frequency coding variation, and about 570,000 SNPs. Samples that failed to reach 98.5% call rate were re-analyzed. Individuals who repeatedly failed

analyses were recontacted by the 23andMe customer service to provide additional samples.

For GWAS quality control (QC) analysis, we limited participants to a set of individuals with  $\geq 97\%$  European descent, determined by analysis of local ancestry[10]. In brief, the algorithm initially partitions phased genomic data into short windows of about 100 SNPs. Within each window, a support vector machine (SVM) was used to classify each haplotype into one of 31 reference populations. SVM classifications were translated into a hidden Markov model (HMM) that takes into account switch errors and incorrect assignments and reports probabilities for each reference population in each window. Finally, simulated admixed individuals were used to recalibrate the HMM probabilities so that the reported assignments are consistent with the simulated admixture ratios. Reference population data was derived from public datasets (the Human Genome Diversity Project, HapMap, and 1000 Genomes) and from 23andMe customers who reported having four grandparents from the same country. For each analysis, a maximal set of unrelated individuals was selected using a segmental identity-by-descent (IBD) estimation algorithm[11]. Individuals were identified as related if they shared more than 700 cM IBD, including regions where the two individuals share either one or both genomic segments identical-by-descent. This degree of relatedness (about 20% of the genome) corresponds approximately to the expected minimum proportion between cousins and first-degree cousins in an outbred population. SNPs deviating from Hardy-Weinberg equilibrium ( $P < 10^{-20}$ ), having a call rate  $< 95\%$ , or with large discrepancies in allele frequency compared to the European 1000 Genomes reference data were excluded. SNPs with large differences in allele frequency (chi squared  $P < 10^{-15}$ ) were identified by computing a 2x2 table of allele counts for European 1000 Genomes samples and 2000 randomly sampled 23andMe customers of European ancestry.

Genotype data were imputed using the September 2013 1000 Genomes Phase1 reference haplotypes[12]. Phasing and imputation was performed separately for the data of each genotyping platform. Phasing was performed with a phasing tool, Finch, developed internally by 23andMe, which implements the Beagle haplotype graph-based phasing algorithm[13] and which was modified to separate the steps of constructing the haplotype graph and phasing. Finch extends the Beagle model to

allow genotyping errors and recombination events to handle cases where there are no consistent paths through the haplotype graph for the individual to be phased. From a representative sample of genotyped individuals, haplotype graphs were generated for European and non-European samples for each 23andMe genotyping platform. Subsequently, an out-of-sample phasing of all genotyped individuals against the corresponding graph was performed. In preparation for imputation, the phased chromosomes were divided into segments of no more than 10,000 genotyped SNPs, with overlaps of 200 SNPs. Each phased segment was imputed against all-ethnicity 1000 Genomes haplotypes (excluding monomorphic and singleton sites) using Minimac2[14], using 5 rounds and 200 states for parameter estimation. For the X chromosome, we created separate haplotype graphs for the non-pseudoautosomal region and each pseudoautosomal region, and these regions were separately phased. Then we imputed males and females together using Minimac2, as for the autosomes, and treated males as homozygous pseudo-diploids for the non-pseudoautosomal region. After QC and genotype imputation a total of 7,024,410 SNPs with imputation quality score  $Rsq > 0.8$  and minor allele frequency (MAF)  $> 1\%$  in 174,785 cases and 228,060 controls were available for association analysis.

### **UK Biobank**

DNA samples were genotyped on custom UK Biobank (UKBB) arrays. 408,951 individuals from UKBB were genotyped for 825,927 variants using a custom Affymetrix UK Biobank Axiom Array, and 49,626 individuals were genotyped for 807,411 variants using a custom Affymetrix UK BiLEVE Axiom Array chip from the UK BiLEVE study[15], which is a subset of UKBB.

All SNPs were subjected to quality control (QC): checks, such as deviations from Hardy-Weinberg equilibrium ( $P < 10^{-5}$ ), batch and plate effects, sex effects, and array effects across control replicates. The SNPs that failed call rate  $< 0.95$  were set to missing for all individuals. The QC was performed centrally for each sample tested for heterozygosity and missing rates. Samples with excessive relatedness ( $> 10$  suspected third-degree relatives) were excluded. Full details of the QC of the genetic data performed centrally by UK Biobank are available in the original publication[2]. To identify sample outliers (i.e. subjects of non-Europeans ancestry), we performed



principal component analysis (PCA) with FlashPCA2[16]. PCA revealed no non-European ancestry outliers. Genotypes of 408,592 UKBB participants with European ancestry (self-reported “white” and genetic Caucasian) were used after QC. Of these, 23,856 samples satisfied our criteria for being HEM cases (either ICD10 code I84 or ICD9 code 455 in medical records) and the remainder of the cohort (n=384,736) served as study controls.

Genetic variants were imputed centrally by UKBB using IMPUTE4[2] and a reference panel that merged the UK10K and 1000 Genomes Phase 3 panel as well as the Haplotype Reference Consortium (HRC) panel[2]. After QC and genotype imputation, a total of 9,572,556 SNPs with an imputation quality score INFO>0.8 and MAF >1% in 23,856 cases and 384,736 controls were available for association analysis.

### ***Estonian Genome Project of University of Tartu (EGCUT)***

The Estonian cohort originates from the population-based biobank of the Estonian Genome Project of University of Tartu (EGCUT). The EGCUT project has been conducted according to the Estonian Gene Research Act and all participants have signed the broad informed consent. The current cohort size is about 200,000 aged 18 years and older, which is very close to the age distribution in the adult Estonian population. Subjects were recruited by general practitioners and doctors in hospitals. The persons who visited the general practitioner’s practices or hospitals were selected at random. Each participant completed a computer assisted personal interview during 1-2 hours in a doctor’s office, which included personal data (place of birth, place(s) of living, nationality etc.), genealogical data (family history, three generations), educational and occupational history and lifestyle data (physical activity, dietary habits, smoking habits, alcohol consumption, women’s health, quality of life). Diseases were defined according to the ICD10 coding. Illumina Human CoreExome, OmniExpress, 370CNV BeadChip and Illumina Global Screening Array (GSA) arrays were used for genotyping.

QC included filtering based on sample call rate (<98%), heterozygosity (> mean  $\pm$  3SD), genotype and phenotype sex discordance, cryptic relatedness (IBD >20%) and outliers of European ancestry based on a multidimensional scaling (MDS) analysis

including 210 HapMap reference samples[17]. SNP QC included testing for call rate (<99%), MAF (<1%) and extreme deviation from Hardy–Weinberg equilibrium ( $P < 10^{-4}$ ).

Pre-phasing was performed using SHAPEIT2[18]. Genotype imputation was performed using the Estonian-specific reference panel[19] and IMPUTE2[20] with default parameters. After QC and genotype imputation, a total of 7,462,975 SNPs with imputation quality score INFO>0.8 and minor allele frequency (MAF) >1% in 6,956 cases and 30,441 controls were available for association analysis.

### ***Michigan Genomics Initiative (MGI)***

DNA samples were genotyped on custom Illumina HumanCoreExome v12.1 bead chips. Samples were excluded if they exhibited (1) a calling rate < 99%, (2) an estimated contamination > 2.5% (BAF Regress)[21] or (3) deviating sex information if the derived sex did not match the self-reported gender. Variants were excluded if they (1) deviated from Hardy-Weinberg equilibrium ( $P_{HWE} < 10^{-5}$ ), (2) had a calling rate < 99%. After quality control, 392,323 polymorphic variants were kept in the following analyses. Next, we estimated the pair-wise relationship of the samples using the software KING[22] and we limited the dataset within a subset of individuals without first- or second-degree relationship. The genetic ancestry of the samples were derived by projecting the principal components of the samples onto that of the Human Genome Diversity Project (HGDP) reference panel (938 unrelated individuals)[23]. Principal component analysis was performed using PLINK v1.90[24], including a subset of LD pruned variants ( $r^2 < 0.5$ ) with MAF >1% shared between the HGDP reference and the MGI data. We retained only samples of recent European ancestry (defined as samples that fell into a circle around the center of the reference HGDP populations in the PC1 versus PC2 space).

Genotype imputation was conducted using the Haplotype Reference Consortium (HRC) panel and the Michigan Imputation Server[25]. After QC and genotype imputation, a total of 6,536,218 SNPs with imputation quality score Rsq>0.8 and MAF >1% in 4,539 cases and 35,338 controls were available for association analysis.

### **Genetic Epidemiology Research on Aging (GERA)**

DNA samples were collected from participants of the Genetic Epidemiology Research on Aging (GERA) cohort and genotyped on high-density custom designed Affymetrix Axiom arrays. Genetic variants with >5% of missing data, MAF <1% in either disease sets or in controls or deviating from Hardy-Weinberg equilibrium ( $P < 10^{-5}$ ) were excluded. Samples with >2% missing data and overall increased/decreased heterozygosity rates were removed. For robust duplicate/relatedness testing (IBS/IBD estimation) and population structure analysis, a pruned subset of 144,799 independent SNPs was used. Pair-wise percentage IBD values were computed using PLINK. By definition,  $Z_0: P(\text{IBD}=0)$ ,  $Z_1: P(\text{IBD}=1)$ ,  $Z_2: P(\text{IBD}=2)$ ,  $Z_0+Z_1+Z_2=1$ , and  $\text{PI\_HAT}: P(\text{IBD}=2) + 0.5 * P(\text{IBD}=1)$  (proportion IBD). One individual (the one showing greater missingness) from each pair with  $\text{PI\_HAT} > 0.1875$  was removed. To identify sample outliers (i.e. subjects of non-Europeans ancestry), we performed principal component analysis (PCA) using the smartpca program[26], based on a set of 144,799 “high-performing” markers after exclusion of SNPs that had an  $r^2$  value greater than 0.5, were within 5 MB of each other, within the MHC region, had a call rates lower than 99.5% and that were located in regions with inversions on chromosomes 8p23 and 17q21.

Genotype data were pre-phased with SHAPE-IT v2.5[18], and then imputed with IMPUTE2 v2.3.1[27] using the 1000 Genomes Phase 3 data as a reference panel. After QC and genotype imputation, a total of 6,897,996 SNPs with imputation quality score  $\text{INFO} > 0.8$  and  $\text{MAF} > 1\%$  in 8,813 cases and 46,780 controls were available for association analysis.

### **German case-control cohort**

DNA samples were genotyped using Illumina’s Global Screening Array version 1.0. Patients with a reported “migration background” were excluded. 3,505 eligible patients were contacted by their treating physician by mail. The initial submission rate was 40%. After consent to participate, the Popgen Biobank sent a study kit with a questionnaire on clinical and socio-demographic characteristics and a set of blood tubes, so that a blood sample could be collected at the family doctor's office and returned to the study center. In addition, a subset of study participants were asked to

complete a comprehensive questionnaire on their dietary habits and usual physical activity. Patients were excluded from the study in the absence of informed consent/blood sample or after withdrawal of consent.

Variants that had >2% missing data, a minor allele frequency <0.1% in either of the different disease sets or in controls, had different missing genotype rates in affected and unaffected individuals ( $P_{Fisher} < 10^{-5}$ ) or deviated from Hardy-Weinberg equilibrium ( $P_{HWE} < 10^{-5}$ ) were excluded. Samples that had >2% missing data and overall increased/decreased heterozygosity rates (with an average marker heterozygosity of  $\pm 5$  s.d. away from the sample mean) were removed. For robust duplicate/relatedness testing (IBS/IBD estimation) and population structure analysis, we used a pruned subset of 100,596 independent SNPs (MAF>0.05) SNPs excluding X- and Y-chromosomes, SNPs in LD (leaving no pairs with  $r^2 > 0.2$ ), and 11 high-LD regions as described by Price *et al.*[28]. Pair-wise percentage IBD values were computed using PLINK2. By definition, Z0: P(IBD=0), Z1: P(IBD=1), Z2: P(IBD=2),  $Z0+Z1+Z2=1$ , and PI\_HAT:  $P(\text{IBD}=2) + 0.5 * P(\text{IBD}=1)$  (proportion IBD). One individual (the one showing greater missingness) from each pair with PI\_HAT>0.1875 was removed. To identify sample outliers (i.e.subjects of non-Europeans ancestry), we performed principal component analysis (PCA) with FlashPCA2[16], on the basis of a set of 100,596 independent markers (described above).

Genotype imputation was conducted using the Haplotype Reference Consortium (HRC) panel and the Sanger Imputation Service[25]. After QC and genotype imputation, a total of 7,117,385 SNPs with imputation quality score INFO>0.8 and MAF >1% in 1,144 cases and 2,740 population controls were available for association analysis.

### ***The Trøndelag Health Study (HUNT)***

DNA was extracted from whole blood from HUNT2 and HUNT3. Genotyping was a research collaboration between researchers from the Norwegian University of Science and Technology (NTNU) and the University of Michigan. Each individual with a DNA sample of an appropriate DNA concentration was selected for genotyping. Samples



were taken at random and genotyped in batches. All genotyping was performed at the Genomics-Core Facility (GCF) at NTNU.

Genotype quality control and genotype imputation were conducted by the K.G. Jebsen Center for Genetic Epidemiology, Department of Public health and Nursing, Faculty of Medicine and Health Sciences, NTNU. In total, DNA from 71,860 HUNT samples was genotyped using one of three different Illumina HumanCoreExome arrays: HumanCoreExome12 v1.0, HumanCoreExome12 v1.1 and UM HUNT Biobank v1.0. Samples were excluded if they did not achieve a 99% call rate, had a contamination >2.5% as estimated with BAF Regress[29], had large chromosomal copy number variants, a lower call rate of a technical duplicate pair and a twin pair, gonosomal constellations other than XX and XY, or whose derived sex was inconsistent with the reported sex. Samples that passed quality control were analyzed in a second round of genotype calling following the Genome Studio quality control protocol described elsewhere[30]. Genomic position, strand orientation and the reference allele of genotyped variants were determined by aligning their probe sequences against the human genome (Genome Reference Consortium Human genome build 37 and revised Cambridge Reference Sequence of the human mitochondrial DNA; <http://genome.ucsc.edu>) using BLAT. Variants were excluded if their probe sequences could not be perfectly mapped to the reference genome, cluster separation was <0.3, Gentrain score was <0.15, showed deviations from Hardy-Weinberg equilibrium in unrelated samples of European ancestry with  $P$ -value <0.0001, their call rate was <99%, or another assay with higher call rate genotyped the same variant. Ancestry of all samples was inferred by projecting all genotyped samples onto top principal components of the Human Genome Diversity Project (HGDP) reference panel (938 unrelated individuals; downloaded from <http://csg.sph.umich.edu/chaolong/LASER/>)[23, 31], using PLINK v1.90. Recent European ancestry was defined for samples that fell into an ellipsoid spans European populations of the HGDP panel. The different arrays were harmonized by reducing them to a set of overlapping variants and excluding variants that had frequency differences >15% between data sets, or that were monomorphic in one data set and had a MAF >1% in another data set. The resulting genotype data were phased using Eagle2 v2.3[32].

Imputation was performed on the 69,716 samples of recent European ancestry using Minimac3 (v2.0.1, <http://genome.sph.umich.edu/wiki/Minimac3>)[33] with default settings (2.5 Mb reference based chunking with 500kb windows) and a customized Haplotype Reference consortium release 1.1 (HRC v1.1) for autosomal variants and HRC v1.1 for chromosome X variants[25]. The customized reference panel represented the merged panel of two reciprocally imputed reference panels: (1) 2,201 low-coverage whole-genome sequences samples from the HUNT study and (2) HRC v1.1 with 1,023 HUNT WGS samples removed before merging. After QC and genotype imputation, over 24.9 million SNPs with imputation quality score  $R^2 \geq 0.3$  in 977 cases and 68,314 population controls were available for association analysis.

### **Danish Blood Donor Study**

DNA samples were genotyped at deCode genetics, Iceland, using Illumina's Global Screening Array as described elsewhere[7]. Details on genotype quality control and imputation are available in Hansen et al., 2019[7]. First- and second-degree relatives were excluded from the analysis. The phenotypic data used in this project includes sex, age, self-reported BMI and selected diagnoses from the Danish National Patient Registry. Participant were classified as having HEM using the ICD-8 code 455 or ICD-10 codes I84 or K64 from the National Patient Registry, resulting in the identification of 1,754 HEM cases in the DBDS cohorts.

The DBDS Genomic Consortium is represented by the following scientists: Andersen Steffen, Department of Finance, Copenhagen Business School, Copenhagen, Denmark; Banasik Karina, Novo Nordisk Foundation Center for Protein Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark; Brunak Søren, Novo Nordisk Foundation Center for Protein Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark; Burgdorf Kristoffer, Department of Clinical Immunology, Copenhagen University Hospital, Copenhagen, Denmark; Erikstrup Christian, Department of Clinical Immunology, Aarhus University Hospital, Aarhus, Denmark; Hansen Thomas Folkmann, Danish Headache Center, department of Neurology Rigshospitalet, Glostrup, Denmark; Hjalgrim Henrik, Department of Epidemiology Research, Statens Serum Institut, Copenhagen, Denmark; Jemec Gregor, Department of Clinical Medicine, Sealand University hospital, Roskilde, Denmark;

Jennum Poul, Department of clinical neurophysiology at University of Copenhagen, Copenhagen, Denmark; Johansson Per Ingemar, Department of Clinical Immunology, Copenhagen University Hospital, Copenhagen, Denmark; Nielsen Kasper Rene, Department of Clinical Immunology, Aalborg University Hospital, Aalborg, Denmark; Nyegaard Mette, Department of Biomedicine, Aarhus University, Denmark; Mie Topholm Bruun, Department of Clinical Immunology, Odense University Hospital, Odense, Denmark; Pedersen Ole Birger, Department of Clinical Immunology, Naestved Hospital, Naestved, Denmark; Petersen Mikkel, Department of Clinical Immunology, Aarhus University Hospital, Aarhus, Denmark; Sørensen Erik, Department of Clinical Immunology, Copenhagen University Hospital Copenhagen, Denmark; Ullum Henrik, Department of Clinical Immunology, Copenhagen University Hospital, Copenhagen, Denmark; Werge Thomas, Institute of Biological Psychiatry, Mental Health Centre Sct. Hans, Copenhagen University Hospital, Roskilde, Denmark; Gudbjartsson Daniel, deCODE genetics, Reykjavik, Iceland; Stefansson Kari, deCODE genetics, Reykjavik, Iceland; Stefánsson Hreinn, deCODE genetics, Reykjavik, Iceland; Þorsteinsdóttir Unnur, deCODE genetics, Reykjavik, Iceland.

## **GWAS association analysis for discovery cohorts**

### ***23andMe***

For comparisons between cases and controls, association test results were performed by logistic regression analysis assuming additive allelic effects. For tests using imputed data, imputed allelic dosages were used rather than best-guess genotypes. Age, biological sex, BMI, the top five principal components from principal component analysis (to account for potential residual population structure) as well as indicators for genotype platforms (to account for genotype batch effects) were included as covariates in the regression analysis. The association test *P*-value was computed using a likelihood ratio test. Results for the X chromosome are computed similarly, with male genotypes coded as if they were homozygous diploid for the observed allele. For chromosome X association analysis, haplotypic allele calls in males outside pseudoautosomal regions (PAR) are converted to homozygous calls by doubling the haplotypic allele (assuming inactivation of large parts of one of the two female X chromosomes[34] and sex was used as a covariate for association testing.

Association summary statistics were adjusted for an estimated genomic control inflation factor  $\lambda_{GC}=1.200$ .

***UK Biobank (UKBB), Estonian Genome Project of University of Tartu (EGCUT), Michigan Genomics Initiative (MGI), Genetic Epidemiology Research on Aging (GERA)***

For each individual case-control data set, association testing was performed using a linear mixed model (LMM) under an additive genetic model for all measured and imputed genetic variants in dosage format using BOLT-LMM[35] (UKBB, GERA) or SAIGE[36] (MGI). Within association analysis, we adjusted for the following covariates: sex, age, BMI (available for UKBB and GERA), the top ten principal components from principal component analysis and a binary indicator variable for genotyping platform (e.g. UKBB Axiom Array vs. UK BiLEVE Axiom Array) to account for the different genotyping chips. For GWAS data set from EGCUT, association testing was carried out with EPACTS [<https://github.com/statgen/EPACTS>], adjusting for age, sex, binary indicator variable for genotyping platform and top four principal components from principal component analysis. For chromosome X association analysis, see text above. The genomic control inflation factors for UKBB, EGCUT, MGI and GERA were  $\lambda_{GC}=1.0966$ , 1.0263, 0.9822 and 0.9541, respectively. For GWAS meta-analysis across discovery cohorts (23andme, UKBB, EGCUT, MGI and GERA).

**GWAS meta-analysis across discovery cohorts**

Prior to GWAS meta-analysis, separate GWAS analyses for discovery cohorts were performed either via logistic regression or mixed linear model association analysis using BOLT-LMM[35] or SAIGE[36] including sex, age, BMI (where available), top principle components (PCs) from principal component analysis (PCA; to control for potential residual population stratification) and genotyping array (if relevant) as covariates. File-level QC of the five individual GWAS summary statistics and meta-level QC from discovery cohorts were carried out using the R package “EasyQC” (v9.2)[37]. In short, the QC process verified data integrity and harmonized both SNP marker IDs and allele coding across the datasets. We only included markers with imputation quality metrics (INFO or Rsq)>0.8 and MAF>1% in the meta-analysis.



Markers with deviating allele frequency (difference >20% from the Haplotype Reference Consortium (HRC) genome reference panel v1.1 comprising 32,488 reference individuals of European ancestry[25]) were removed along with indels and multi-allelic markers. The resulting summary statistics of the five discovery cohorts (with a total of 218,920 HEM cases and 725,213 controls) were meta-analyses via fixed-effect meta-analysis based on METAL's inverse-variance weighted approach[38]. We used the generally accepted threshold of  $5 \times 10^{-8}$  for meta-analysis  $P$ -values to define statistical significance ( $P_{\text{Meta}} < 5 \times 10^{-8}$ ). Genome-wide summary statistics of our analyses are publicly available through our web browser (<http://hemorrhoids.online>) and have been submitted to the European Bioinformatics Institute ([www.ebi.ac.uk/gwas](http://www.ebi.ac.uk/gwas)) under accession number GCST90014033.

### **Annotation of HEM GWAS risk loci and gene mapping**

We used independent computational pipelines for the functional annotation of GWAS meta-analysis results, using FUnctional Mapping and Annotation of Genome-Wide Association Studies (FUMA v1.3.5)[39], Data-driven Expression-Prioritized Integration for Complex Traits (DEPICT)[40], and Multi-marker Analysis of GenoMic Annotation (MAGMA, also implemented in FUMA)[41]. The 102 newly identified genome-wide significant risk loci were defined in FUMA (using default parameters and eQTL databases including GTEx v7) as non-overlapping genomic regions that extend a linkage disequilibrium (LD) window ( $r^2 = 0.6$ ) around each lead SNP association signal with  $P_{\text{Meta}} < 5 \times 10^{-8}$ . Annotation of these regions with FUMA resulted in 712 transcripts mapped to risk loci (415 positional and 562 eQTL candidates), while 217 genes were identified using DEPICT, and 255 in MAGMA independent gene-based tests, bringing the total of non-redundant HEM candidate genes to 819 (**online supplementary table S7**). Regional association plots of all 102 risk loci were generated using LocusZoom[42].

### **Bayesian fine-mapping analysis**

A Bayesian fine-mapping analysis was carried out using FINEMAP[43] for the 102 genome-wide significant risk loci in order to calculate the posterior inclusion probability (PIP) for each lead SNP as causal and to determine a credible set for each risk locus, i.e. a minimum set of variants containing all causal variants with certainty  $\geq 0.95\%$ . As

input for fine-mapping we extracted all genetic variants located within the 102 risk loci (as defined by FUMA) and calculated the local LD structure using genotypes from UKBB samples (**online supplementary table S1**) serving as a reference population.

### **Heritability analysis via linkage disequilibrium score regression (LDSC)**

Narrow-sense heritability ( $h^2_{\text{SNP}}$ ) for HEM and the genetic correlation ( $r_g$ ) between HEM and other traits were estimated using LD score regression, as implemented in the online platform CTG-VL[44]. This platform integrates summary statistics of 1,387 traits from multiple resources such as UKBB, the Psychiatric Genomics Consortium (PGC) and the Genetic Investigation of ANthropometric Traits (GIANT) consortium. Significantly correlated pairs of traits were reported after FDR correction for multiple comparisons at  $\alpha=0.05$ .

### **Genome-wide pleiotropy analysis**

We conducted cross-phenotype association analysis based on subsets (ASSET) methodology[45] across association summary statistics from diverticular disease[46], irritable bowel syndrome (IBS)[47] and HEM to identify shared risk loci. The subset-based meta-analysis (SBM) method maintains similar type-I error rates as for standard meta-analysis and identifies the best subset of non-null studies while in parallel accounting for multiple-hypothesis testing and shared individuals. This method offers a substantial power increase (sometimes approaching between 100-500%)[45] compared to standard univariate meta-analysis approaches, where the (heterogeneous) effect of a specific SNP is not exclusively restricted to a single phenotype. Under the assumption that association signals from shared risk loci based on positional overlap are tagging same causal variant for different phenotypes, the SBM approach improves power compared to standard fixed-effects meta-analysis methodology.

### **Tissue and pathway enrichment analyses**

Gene-set and tissue-specific enrichment analyses (respectively GSEA and TSEA) of HEM genes were carried out using integrated default pipelines in FUMA, and DEPICT implemented in the CTG-VL platform[44]-[40]. HEM gene lists were derived from three

alternative approaches including positional and/or eQTL mapping in FUMA, MAGMA gene-based analyses (also implemented in FUMA), and DEPICT functional annotations, and tested against Gene Ontology (GO) terms and 30 GTEx v7 general tissue types. Statistical significance was defined using  $P_{\text{Benjamini-Hochberg}} < 0.05$ .

### **Polygenic risk scores (PRS) analysis**

The analysis of polygenic risk scores (PRS) was performed on the basis of a pruning and thresholding approach, using the  $P$  value and LD-driven clumping procedure as implemented in PRSice-2[48]. Effect estimates and corresponding standard errors from GWAS meta-analysis results were used as the base dataset to generate weights over a range of  $P$  values ( $0.5$  to  $5 \times 10^{-8}$ ) and  $r^2$   $0.1$  LD thresholds, with the most appropriate thresholds selected as those that include SNPs with the highest Nagelkerke's  $R^2$  value. The selected model was then applied to the QCed genetic datasets from the German case-control cohort, HUNT and DBDS, respectively. Logistic regression was used to test HEM PRS distribution in cases and controls, taking into account sex, age, BMI and the top 10 PCs from PCA. For HUNT and DBDS we also studied HEM prevalence across PRS percentile distributions. PRSs were binned into percentiles and HEM prevalence from the top 5% of PRS distribution was compared to the remainder of the population in a logistic regression adjusting for sex, age, BMI and the top 10 PCs. Additional analyses were performed to evaluate the relationship between HEM PRS and age at diagnosis (Spearman's correlation test) and need for invasive treatments (number of surgeries and/or rubber-band ligation; tested with linear regression correcting for sex, age, BMI and the top 10 PCs).

### **Phenome-wide association studies (PheWAS)**

For each of the 102 GWAS risk loci, we queried the lead SNP and its LD proxies ( $r^2 > 0.8$ , from 1000 Genomes Project samples of European descent) using PhenoScanner v2[49], and manually inspecting the GWAS catalog[50] and GWAS ATLAS[51]. Only genome-wide significant associations ( $P < 5 \times 10^{-8}$ ) were taken into account, and those from GWAS ATLAS were collapsed by trait categories and plotted with the R package "ggforce" into an alluvial diagram.

## Cross-trait analyses

Traits genetically correlated with HEM (from LDSC analyses) were tested for their prevalence in HEM patients vs controls in UKBB and DNPR. In UKBB, we derived the ICD10 diagnoses from data-fields “41202” (primary diagnosis) and “41204” (secondary diagnosis), self-reported medical conditions from data-field “20002”, and self-reported medication use from data-field “20003”. Differential prevalence was tested using a logistic regression model adjusted for sex and age, including FDR correction for multiple comparisons. For DNPR, a previously published method[52] was used to identify diseases that significantly co-occur more often with HEM diagnoses. Each combination of pair-wise disease co-occurrences was compared to a comparison group matched by sex, age, type of hospital encounter and week of discharge. The relative risk (RR) is used to evaluate the strength of the correlation between significant disease pairs (disease A followed by disease B and *vice versa*). Here, we have used this method to evaluate the temporal co-occurrence of selected diseases and medications with the HEM diagnosis in the DNPR, including FDR correction for multiple comparisons.

## RNA library preparation and RNA-sequencing.

The RNA-Sequencing (RNA-Seq) libraries were prepared from 20 ng of total RNA from freshly frozen tissue extracted with the mirVana miRNA Isolation Kit according to the manufacturer’s protocol (Ambion). The NEXTFLEX Combo-Seq Kit (Perkin Elmer) was used to generate combined mRNA and microRNA libraries following manufacturer’s instructions. In short, poly-A-tailed RNA species were reverse transcribed to generate DNA:RNA duplexes whose RNA molecules were specifically sheared by RNase H, resulting in RNA fragments containing 5'-monophosphate and a 3'-hydroxyl groups. These mRNA fragments were 3'-polyadenylated together with small RNAs and then 5' 4N adapters were ligated to their 5' ends. Finally, first strand synthesis followed by PCR amplification were used to add sequences required for Illumina sequencing. The generated RNA libraries were quality-controlled using the Agilent 2200 TapeStation (Agilent Technologies), randomized and then deep sequenced (5 samples per lane), 1x50bp using the Illumina HiSeq 4000 platform.

## Mapping and quality assessment of RNA-Seq data

The sequenced reads were demultiplexed and obtained as fastq files for each sample. Data pre-processing, quality control, mapping to genome (build hg38) and transcriptome annotation (miRBase v21, Ensembl 83) were performed using the `exceRpt`[53] pipeline. More precisely, reads were trimmed for 3' adapter sequences, 4N nucleotides at 5' end and low-quality bases (<Q20). The trimmed sequences shorter than 15 bp were discarded and only high-quality reads were then mapped to genome (with minimum sequence match of 15 nucleotides), annotated and quantified. All RNA-Seq libraries were quality controlled for library size (>10M of mapped reads), transcriptome genome ratio (> 0.95) and outliers for number of detected unique genes and microRNAs (< Q1–1.5 IQR). Low abundant gene-level and microRNA arm level counts that were expressed below 0.1 RPM in less than 25% of the samples per trait were removed from downstream analyses. The generated quality-controlled counts and raw sequencing reads have been deposited at NCBI Gene Expression Omnibus (GEO)[54] under the accession number [GSE154650](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE154650).

## Gene signature-based determination of anal canal zones

Histologically, the anal canal can be divided into three zones according to the epithelial lining. The upper part is of the mucosal type (intestinal) and the lower part is of the squamous keratinized (anoderm), while the middle part, where the epithelium varies, is called the anal transitional zone[55, 56]. Due to the gradient nature of the anal canal epithelium, keratinocyte and sebocyte marker gene signatures from the xCell catalog[57] were used to discriminate the different histological zones. More specifically, the quality-controlled gene counts were normalized using the variance stabilizing transformation (VST) implemented in the DESeq2 R package[58]. The normalized gene counts were then ranked according to their expression level using the `rank()` function from the base R package and submitted to single sample gene set enrichment analysis (ssGSEA)[59] implemented in the GSVA R package[60]. The obtained normalized enrichment scores (NES) of keratinocytes and sebocytes were used to cluster samples into 3 groups (in accordance to the number of histological zones) by employing the base R function `kmeans()` with `k=3` and `nstart=20` as parameters. The obtained clusters were assigned to histological zones by the relative abundance of keratinocytes and sebocytes (i.e. sebum-producing epithelial cells), and



further confirmed by the expression levels of previously defined marker genes, including *KRT4*, *KRT8*, *KRT13* and *KRT20*[56, 61] genes (**online supplementary figure S8**). Multidimensional scaling (MDS) analysis using Spearman's rank correlation distance (1-correlation coefficient) was performed on VST normalized expression data and was used to explore the results.

### Differential gene expression analysis

The quality-controlled count data were further analyzed using edgeR[62] workflow for differential expression analysis. Negative binomial generalized log-linear models were fitted to the trimmed mean of M-values (TMM) normalized count data of HEM genes using glmFit() function with trended dispersion estimates and the offsets for GC-content correction generated by EDASeq (default parameters). The glmLRT() function was used to calculate log-likelihood-ratio statistics and *P*-values of differential expression. The models were adjusted for BMI and histological zones of anal canal (see the previous paragraph). The nominal *P*-values were corrected for multiple testing according to Benjamini and Hochberg. Transcripts with an FDR corrected *P*-value < 0.05 and a log<sub>2</sub> fold change > 0.5 (in either direction) were considered to be significantly differentially expressed.

### Identification and characterization of enriched co-expression modules

Weighted gene co-expression network analysis of hemorrhoid-specific tissue was performed using the automated WGCNA[63] pipeline implemented in the CEMiTool[64] R package. The quality-controlled and VST normalized data (36,342 genes in 20 samples) was used to calculate signed scale-free topology overlap matrix, which was subsequently used to define gene co-expression modules in an unsupervised manner. More specifically, Pearson correlation coefficients for each gene-gene comparison (including miRNAs) were used to calculate adjacencies defined as following:  $a_{ij} = |0.5 + 0.5 \times \text{cor}(x_i, x_j)|^\beta$ , where  $x_i$  and  $x_j$  are expression values of  $i^{\text{th}}$  and  $j^{\text{th}}$  genes and where  $\beta$  is a soft threshold power based on scale-free topology, which was identified by employing pickSoftThreshold() function from the WGCNA R package. The generated adjacencies were then used to compute topological overlap measures (TOM) and their dissimilarity measures (1-TOM) were further used for

average linkage hierarchical clustering and dynamic tree cutting (cutoff value of 0.995) to identify gene co-expression modules. Each gene co-expression module contained a minimum of 50 genes and was summarized into eigengene, which is the first principal component of their expression values. Highly similar modules were identified by correlation of their eigengenes ( $>0.7$  Pearson's  $r$ ) and merged together. The intramodular connectivity of each gene was measured by Pearson's correlation of module eigengene and its expression value. The top 10% of genes having the highest connectivity values were considered as being module hub genes (central nodes in the scale-free network). A Fisher's exact test was used to identify modules with significantly ( $P_{FDR} < 0.05$ ) overrepresented in HEM genes. The ClusterProfiler[65] R package was used to identify gene ontology (GO) terms "biological process" pathway enrichments for HEM-significant modules.

### **ABO blood group analysis**

The association between ABO blood types and HEM was tested on 408,592 and 55,593 individuals from UKBB and GERA, respectively. We first imputed ABO blood group information individually based on genotypes at the *ABO locus* on chromosome 9q34.2. We extracted the genotypes of three SNPs: rs8176719, rs41302905 and the adjacent rs8176747 and inferred blood group status based on these SNPs as previously described[66]. Next, the risk of HEM was assessed on samples from each blood groups of the ABO blood group system ("A", "B", "AB" and "O"). An association test based on logistic regression was employed to test for significant HEM association for each of the four blood groups, adjusting for sex, age, BMI and the top 10 PCs from PCA. FDR correction was applied for multiple testing.

### **Fluorescence Immunohistochemistry**

Fluorescence immunohistochemistry was performed as previously described[46, 67]. Briefly, anorectal specimens were fixed in 4% paraformaldehyde for 24 hours. Paraffin-embedded tissue sections were pre-treated with citrate buffer and primary antibodies were incubated overnight. Used primary and secondary antibodies are listed in **online supplementary table S13**. All antibodies were diluted in antibody diluent (ThermoFisher Scientific). Nuclei were counterstained with DAPI (Roche, Mannheim, Germany). Image acquisition was performed on a fluorescence inverted

microscope (Axiovert 200 M, Zeiss, Gottingen, Germany) coupled to an AxioCam MR3 camera (Zeiss) using Axiovision software (version 4.7, Zeiss).

### ***In silico* variant protein analysis**

To construct a first hypothetical model of whether *SRPX* and *ANO1* missense lead variants (shown in red in **figure 1**) are likely to interfere with functionally active domains at the protein level, we conducted protein domain analyses for *SRPX* and *ANO1*.

*ANO1* (also *TMEM16A*) is an anion channel protein that enables the passive flow of Cl anions through the membrane as a result of increased intracellular Ca<sup>2+</sup> levels. The decrease in an anion flow occurs over time after prolonged stimulation eventually leads to complete desensitization to saturated Ca<sup>2+</sup>. In addition to elevated Ca-levels, *ANO1* function is regulated by the PIP2 (Phosphatidylinositol(4,5)-bisphosphate) signal lipid which binds at the cytoplasmic membrane interface[68]. The Interaction with PIP2 has been shown to slow down the *ANO1* regulatory process, probably by hindering the gradual collapse of the ion conduction pore[69].

The *ANO1* protein functions as a homodimer, with each subunit consisting of ten transmembrane helices and its own anion conduction pore (**online supplementary figure S9**) which is composed of helices 3-7 and contains a conserved Ca<sup>2+</sup> binding site[70, 71]. Ion flow through the pore is made possible by local structural rearrangements that open the channel in response to Ca<sup>2+</sup> binding[70].

The variant F608S is located at the beginning of the transmembrane helix 5, i.e. near the cytoplasmic interface (**online supplementary figure S9**). Although helix 5 is part of the ion conduction pore, the sidechain of F608 points in the opposite direction to the dimer interface and is located near the predicted PIP2 binding residues. Adjacent K609 forms a stabilizing salt bridge with E594 in the TM4-TM5 linker which is conserved in all members of the *TMEM16* protein family. Mutation of this salt bridge results in a rapid Ca<sup>2+</sup> desensitization, similar to a direct mutation of the predicted PIP2 binding residues[69].

F608 and its sequential and structural neighbors are conserved among *ANO1* orthologs (**online supplementary figure S9**). The variant causes a change from the aromatic and very hydrophobic phenylalanine to the smaller and polar/hydrophilic serine. All members of the *TMEM16* superfamily conserved a non-polar residue at this

position, suggesting that the polar sidechain of the serine may cause a structural conflict in the region. The variant could interfere with the K609-E594 salt bridge which stabilizes the PIP2 binding. F608S may thus interfere with the PIP2 binding and consequently accelerates ANO1 degradation, similar to the rapid desensitization that was demonstrated by mutational analyses of basic amino acids in the vicinity and the salt bridge[69].

The SRPX (also DRS, ETX1, SRPX1) variant rs35318931 causes a Ser413Phe exchange at the C-terminal domain of unknown function (**online supplementary figure S11**). The protein is further composed of three Sushi domains, and one HYR domain. Sushi domains are components involved in extracellular protein-protein interactions and are often found in complement control proteins[72]. The HYR (hyalin repeat) domain is predicted to contribute to cell adhesion since the domain enables the hyalin protein to bind to the receptor[73]. The SRPX C-terminal domain is a phylogenetically widespread protein domain that is well-conserved in vertebrates (**online supplementary figure S11**) and also in many bacteria, and has been named the DUDES domain (DRO1-URB-DRS-Equarin-SRPX)[74]. Protein structural analyses assign a thioredoxin-like fold to this domain, although the location of potential functional cysteines seem unique for SRPX and SRPX2 proteins[75]. Therefore, the conserved structural core allows fold recognition, but the lack of suitable structural templates including loops and termini complicates in-silico functional prediction for SRPX (**online supplementary figure S11**). SRPX was originally identified as a tumor suppressor[76] and, in this context, to the induction of apoptosis[77] and downregulation of glucose metabolism via Lactate dehydrogenase-B[78]. Proteomics studies found SRPX expression in the extracellular matrix (ECM) of different tissues (lung[79], cartilage[80] and colon and liver[81]) and is upregulated in the ECM during cardiac remodeling[82]. Further, SRPX was also shown to interact with PELO at the actin cytoskeleton[83].

Other members of the DUDES protein family were shown to localize in the extracellular matrix, e.g, SRPX2 in brain[84], equarin in chick lens[85]. CCDC80 is a remote homologous that binds activated JAK2 and is consequently more abundant in the extracellular matrix. JAK2-binding was also detected by the paralog SRPX2, and interaction is therefore also predicted for SRPX[86]. CCDC80 is composed of three DUDES domains, that are independently able to bind JAK2, assuming the SRPX

DUDES domain is responsible for protein association with the ECM. The ECM provides structural integrity for tissues, and involves in cell differentiation, activation and migration. HEM tissue is less stable and show abnormalities in the ECM collagen composition (compared to healthy tissue[87]).

The variant Ser413Phe locates at the beginning of strand 3 in the central beta sheet. The preceding loop is highly variable among homologs[75] but the following strand is one of the best conserved regions within the protein family, including an invariant F414. The change from the polar and small amino acid serine to the larger, aromatic and hydrophobic phenylalanine potentially destabilizes the domain structure due to its location adjacent the conserved hydrophobic core of the protein fold.

The SRPX domain structure was derived from the UniProt database and by search against the NCBI Conserved Domains Database (CDD). SRPX and SRPX2 protein sequences were derived from UniProt, Ensembl and the consensus sequence of pfam13778 from the CDD. Sequence alignments were conducted using Muscle. The sequence alignment was visualized using JalView applying the Clustal coloring scheme. Protein sequence identifiers UniProt or Ensembl: SRPX: human, P78539; mouse, Q9R0M3; cow, F1MQX1; zebrafish, Q58ED3; xenopus tropicalis, ENSXETT00000018780.4. SRPX2: human, O60687; mouse, Q8R054; cow, Q5EA25; zebrafish, E7F8X0, xenopus tropicalis, ENSXETT00000014699.4.

The structure-based alignment for modeling the SRPX C-terminal domain of unknown function (DUF4174/pfam13778, 332-451) is based on secondary structure predictions, structural alignments of two templates (PDBs 3drn/chain A, 3cmi/chain A) and multiple sequence alignment including the consensus sequence of pfam13778. Structural models of SRPX and ANO1 were visualized using PyMOL.

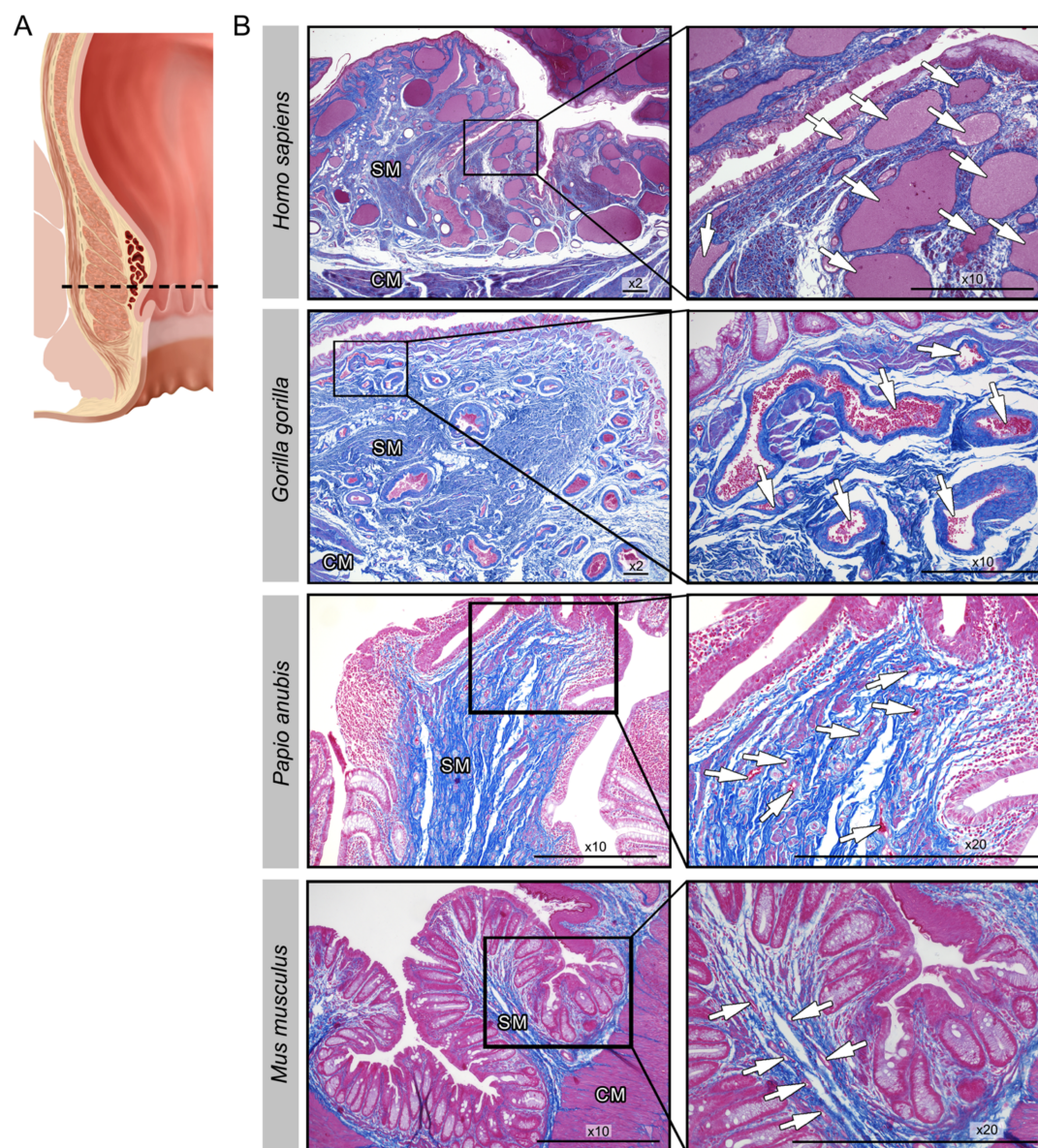
### **Site-directed ANO1 mutagenesis and whole-cell electrophysiology**

F608S (F671S in NM\_018043) variant was introduced into the full-length human ANO1 gene with exon 0 (123 bp[88, 89]) and exon b (66 bp[90]) by a single nucleotide change (c.2012 T→C), using the QuikChange Lightning Site-Directed Mutagenesis Kit (Agilent Technologies, Santa Clara, CA, USA) according to the manufacturer's instructions. The integrity of the construct and presence of the F608S mutation were verified by DNA sequencing. The primer sequences were: (forward) 5'-cttccgcaggaggagta-3' and (reverse) 5'-cagcaggaaagccttgagatcagcctctcctc-3'.



HEK293 cells were co-transfected with pEGFP-C1 plus either wild-type ANO1 or F608S-ANO1 by Lipofectamine 3000 (Thermo Fisher Scientific, MA). Ca<sup>2+</sup>-activated Cl<sup>-</sup> currents were recorded by whole-cell electrophysiology as previously described by Strege et al.[91]

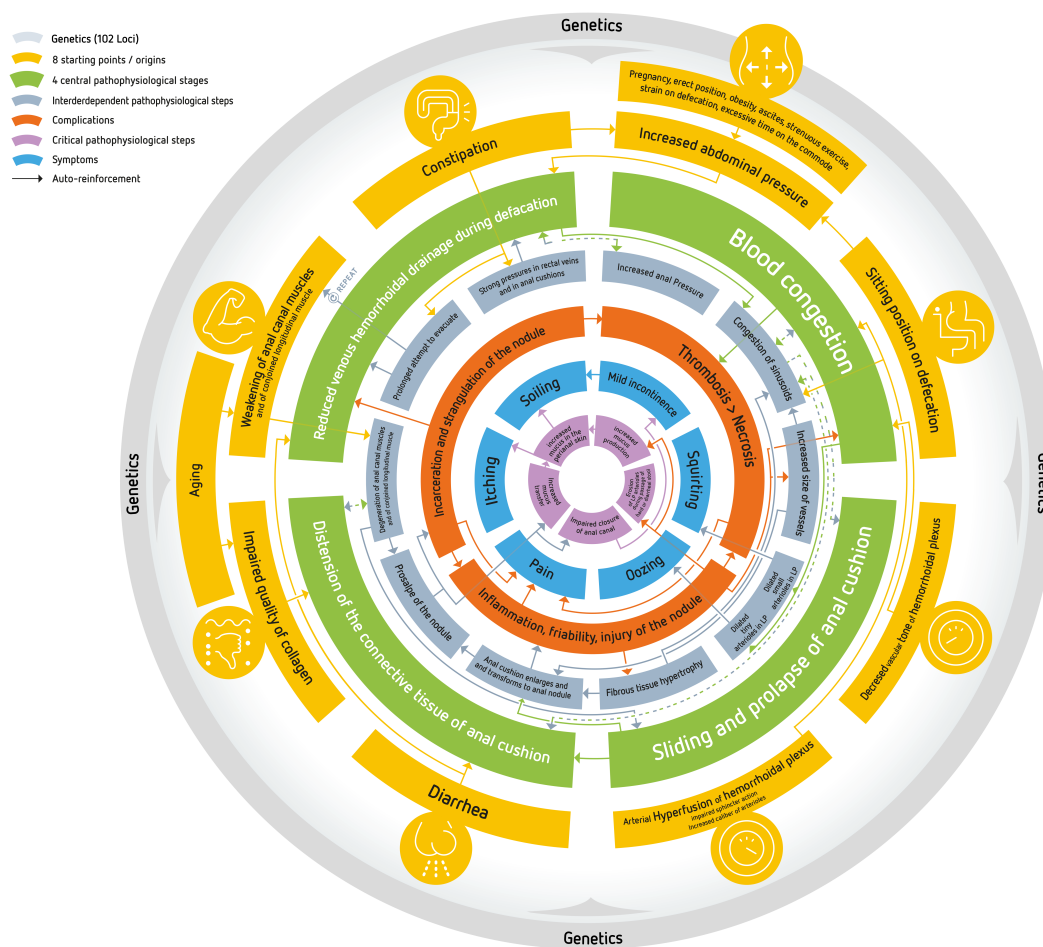
## SUPPLEMENTARY FIGURES



### Online Supplementary Figure S1. *Histological analysis of the anorectum in four different species.*

The left panel (A) shows the section plane of the anal canal at the level of the hemorrhoidal plexus. Panel (B) shows the hemorrhoidal plexus of 4 different species: *Homo sapiens* (top row), *Gorilla gorilla* (second row), baboon (*Papio anubis*; third row), and mouse (10-week-old male C57BL/6JRj mouse; bottom row). While the human anorectum shows a well-developed hemorrhoidal plexus with densely packed blood

vessels of large diameters, the gorilla sample displays a rudimentary hemorrhoidal plexus with fewer and smaller blood vessels. Both baboon and mouse samples exhibit only small-sized and scattered blood vessels which resemble normal vascularization patterns of the regular rectal mucosa. Azan staining with visualization of connective tissue (blue) as well as cell nuclei, erythrocytes and smooth muscle (all purple red). Magnifications for human and gorilla (left 2x, right 10x), for baboon and mouse (left 10x, right 20x). Scale bars: 500  $\mu$ m. White arrows: hemorrhoidal/submucosal blood vessels, SM = submucosa, CM = circular muscle layer/internal anal sphincter.

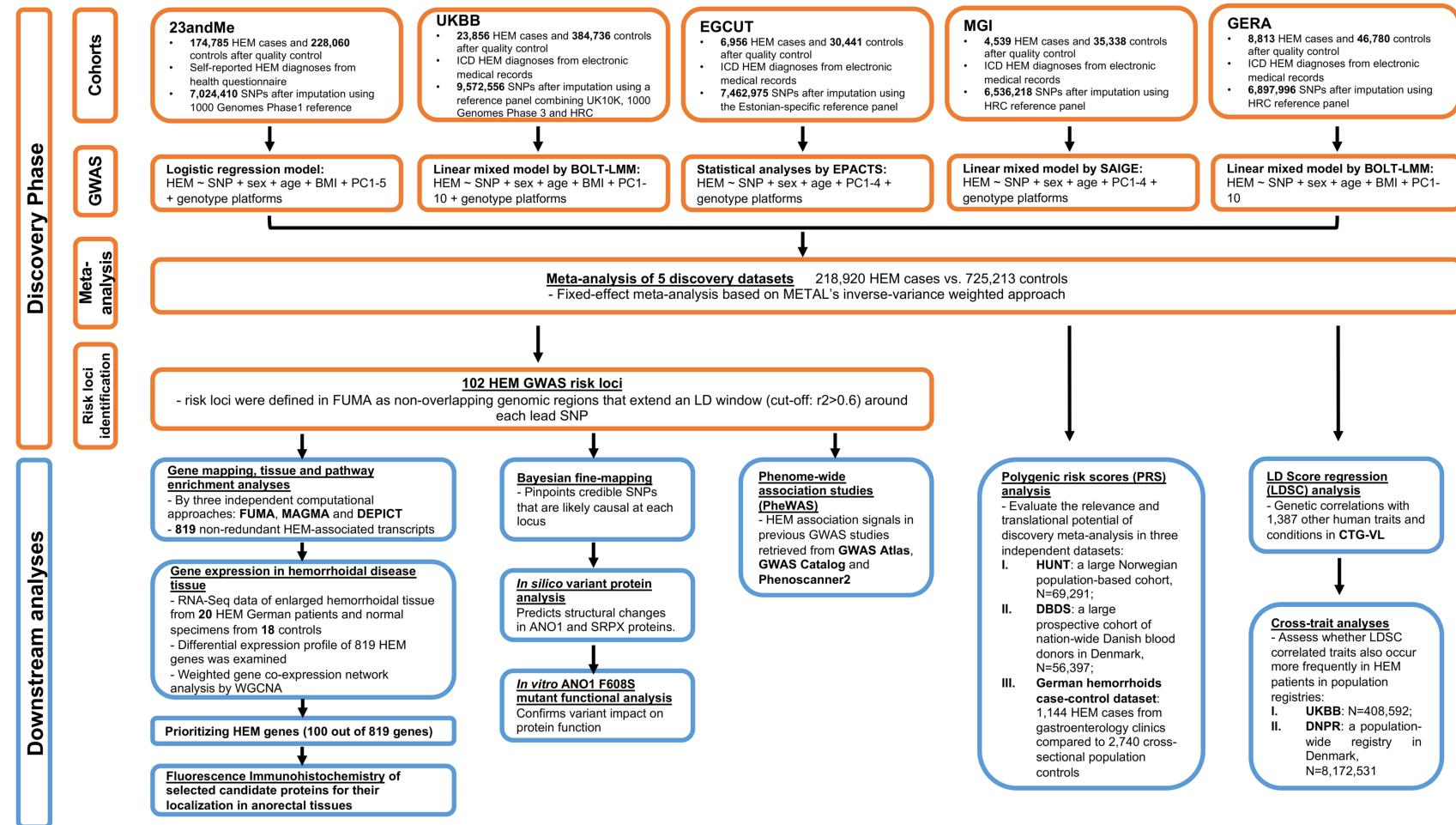


**Online Supplementary Figure S2. Suggested integrated model that summarizes the contemporary thinking on the pathophysiology of HEM (figure and legend are mainly taken from Figure 2 in Nikolaos Margetis' review[92]).**

HEM is a complex and multifactorial disease, most likely resulting from separate origins and combinations thereof. Different origins and causes have been suggested (orange), which force the hemorrhoidal plexus in different abnormal directions and probably converge in four central pathophysiological events (green). Different consecutive pathophysiological stages (grey) connect the primary causes and the 4 central events. These pathophysiological stages are interconnected, interdependent, and mutually reinforcing, creating a vicious cycle. This multidirectional network is continuously auto-reinforced, as shown by the arrows, and over time provides only

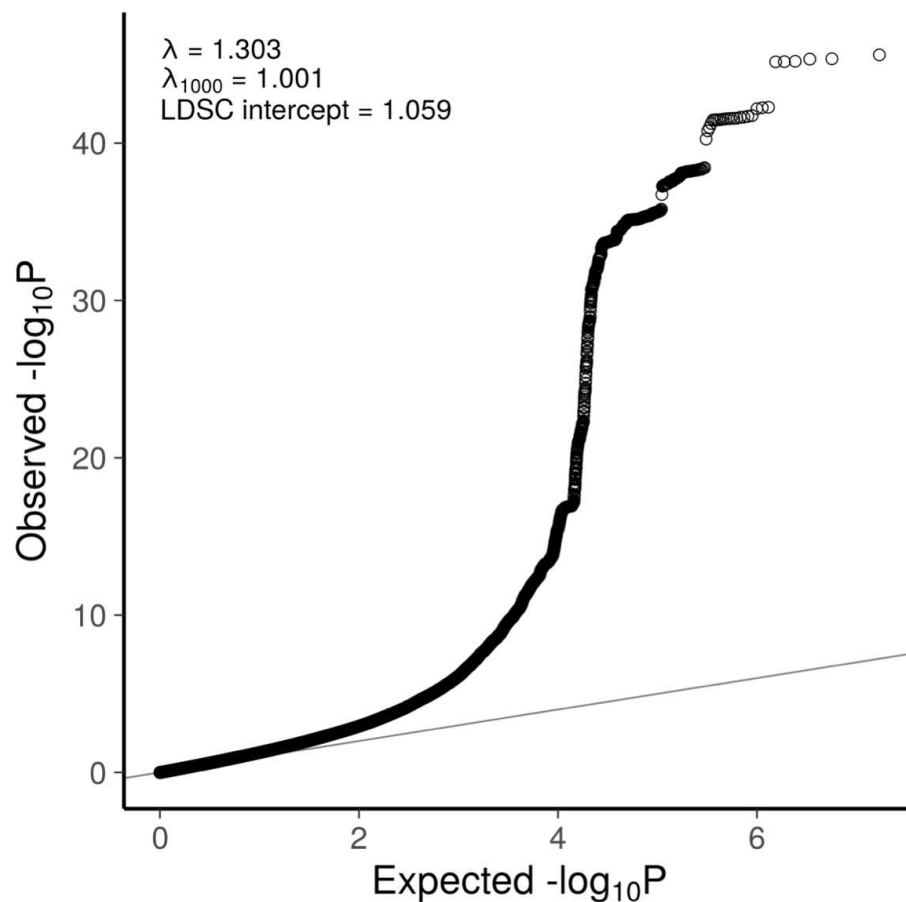
one outcome, with the hemorrhoids deteriorating. Ultimately, symptoms (blue) and complications (red) occur. For further details we refer to Margetis' review[92].





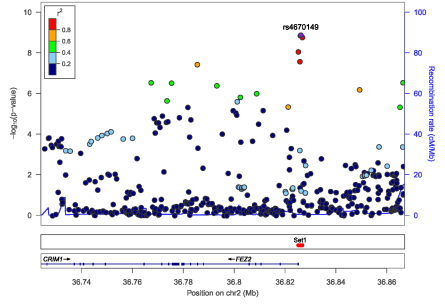
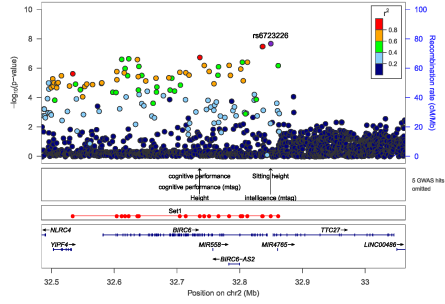
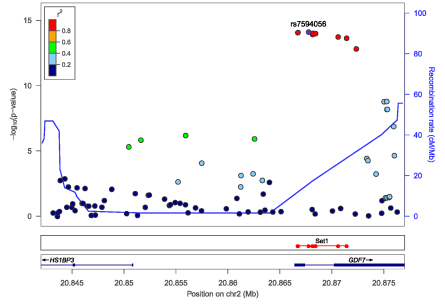
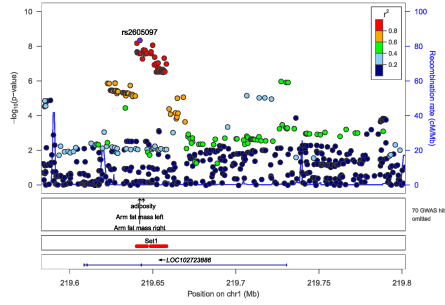
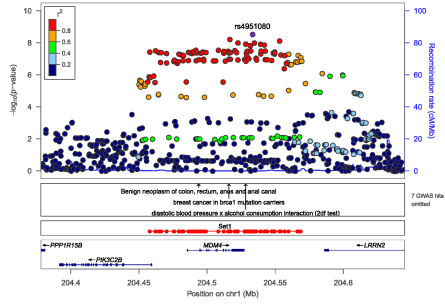
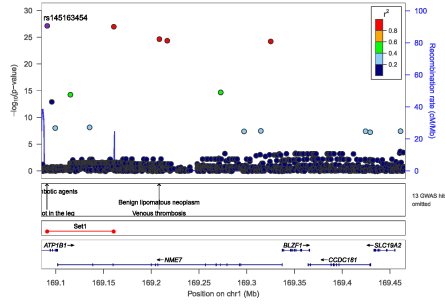
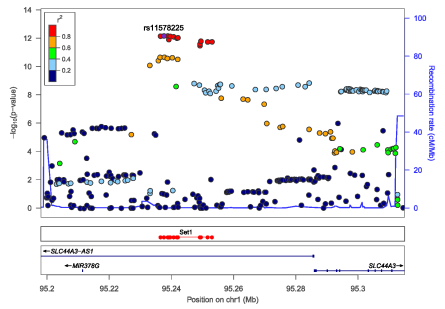
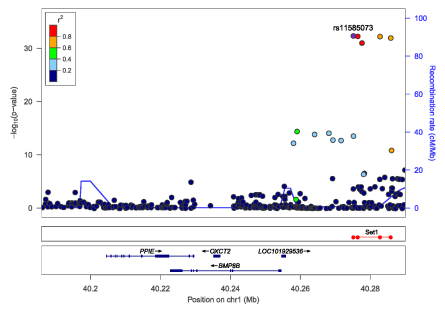
**Online Supplementary Figure S3. Schematic overview of the study workflow.**

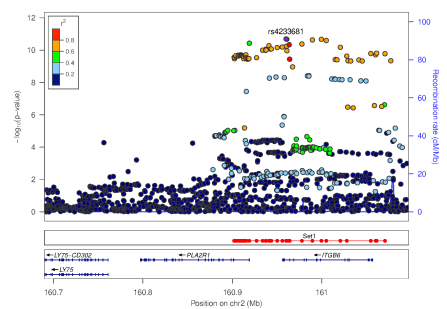
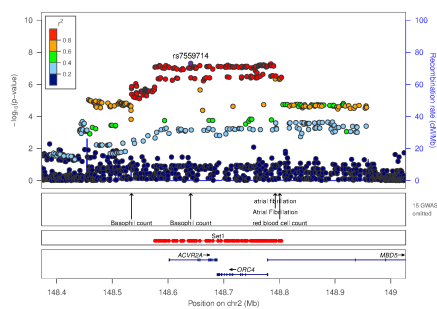
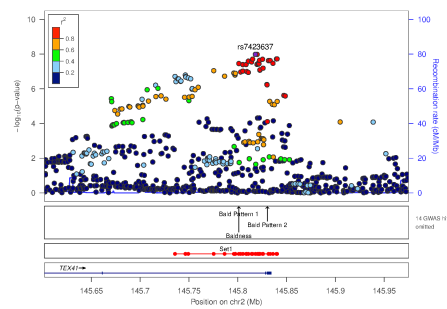
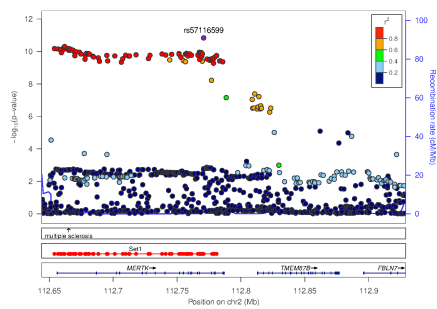
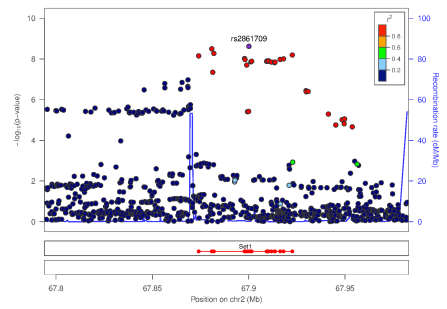
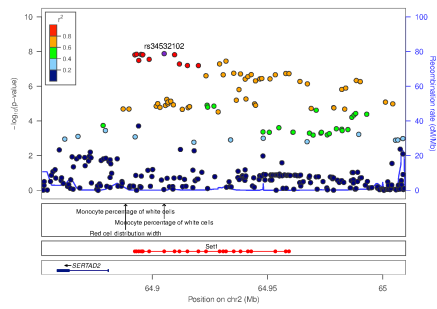
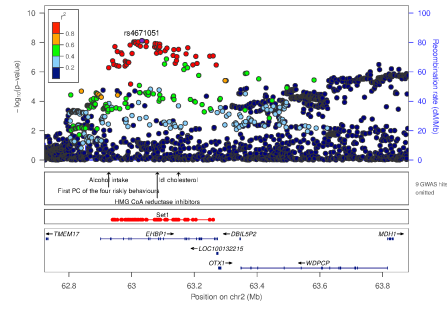
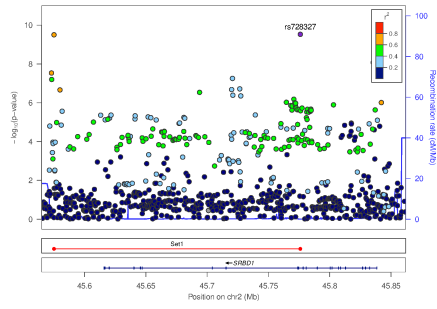
The flowchart shows the study design and analytic strategy of both the discovery phase and the downstream analyses, which includes the study aims, cohorts and numbers of samples of each analytic stage. HEM: hemorrhoids disease. BMI: body mass index, UKBB:UK Biobank, EGCUT: Estonian Genome Center at the University of Tartu, MGI: Michigan Genomics Initiative, GERA: Genetic Epidemiology Research on Aging, HUNT: The Trøndelag Health Study, DBDS: Danish Blood Donor Study, DNPR: Danish National Patient Registry. QC: quality control. IBD: identity by descent. ICD: International Classification of Diseases. Rsq: R square. MAF: minor allele frequency. SNP: Single nucleotide polymorphisms. PC: principal component. LD: linkage disequilibrium. HRC: haplotype reference consortium.



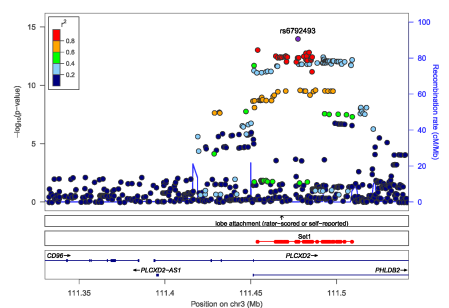
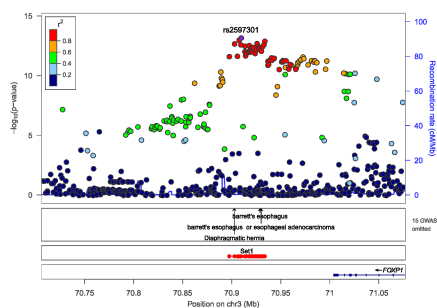
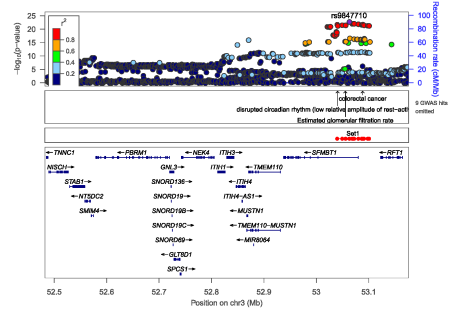
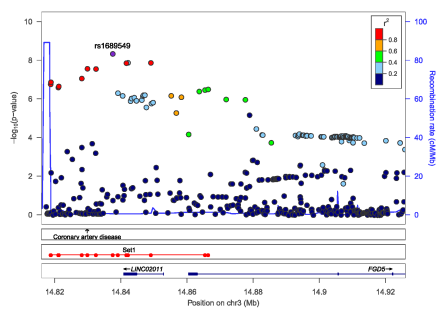
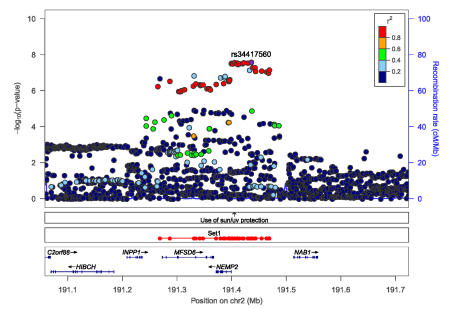
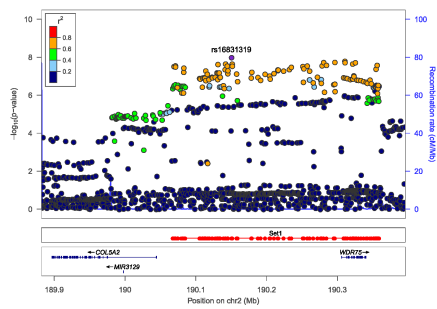
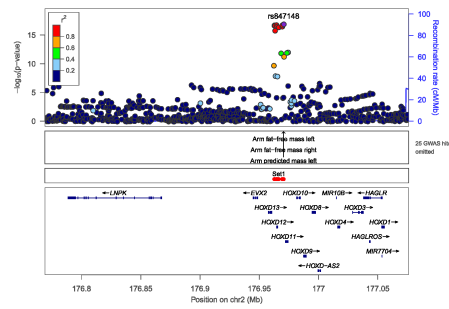
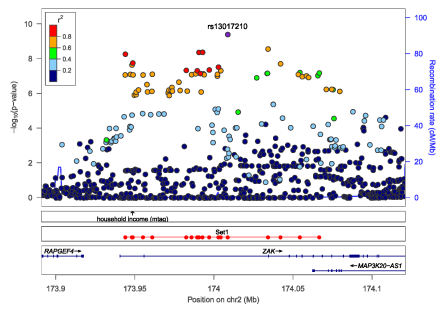
**Online Supplementary Figure S4. Quantile-quantile (QQ) plot of GWAS meta-analysis results.**

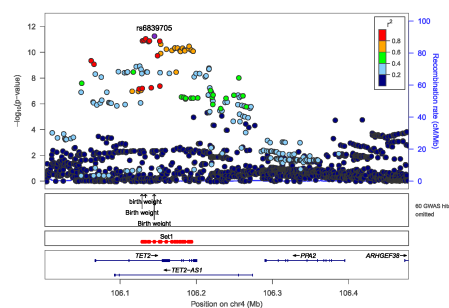
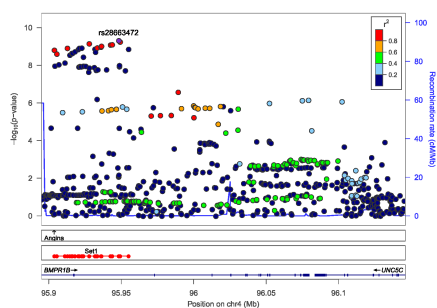
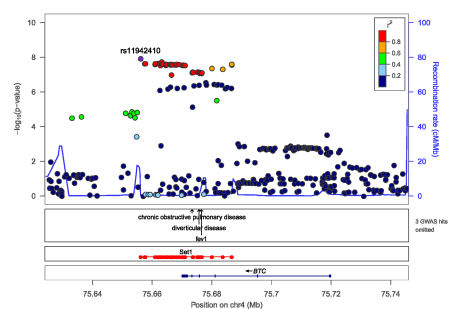
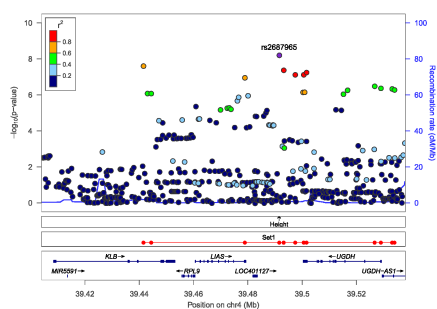
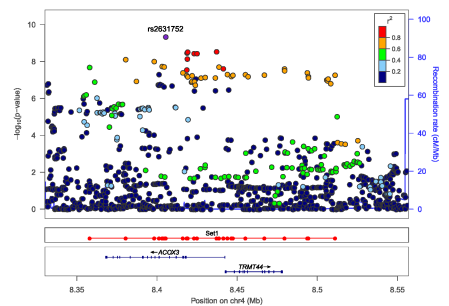
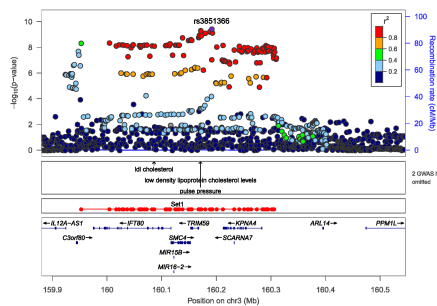
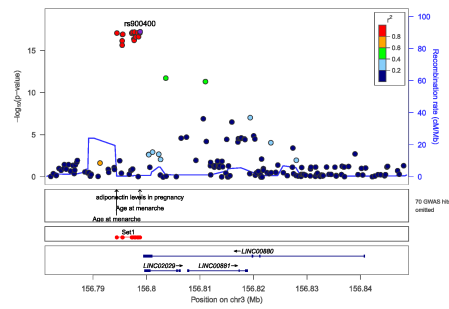
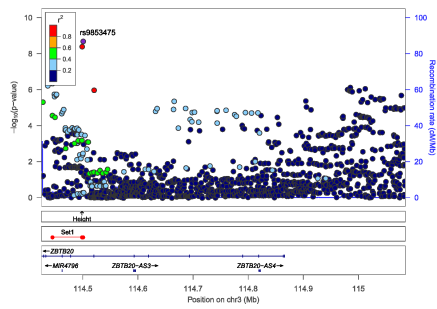
Only markers that passed the imputation quality score  $R^2 > 0.8$  and  $MAF > 1\%$  were used for the plot. The genomic inflation factor  $\lambda$  is defined as the ratio of the medians of the sample  $\chi^2$  test statistics and the 1-d.f.  $\chi^2$  distribution (0.455)[93]. Lambda inflation statistics are influenced by the sample size. To facilitate comparison with other studies,  $\lambda_{1000}$  converts a given lambda from  $n$  cases and  $m$  controls so that the value corresponds to an analysis with 1000 cases and 1000 controls. Although genomic inflation was observed ( $\lambda = 1.303$ ) this was probably due to polygenicity rather than population stratification as determined by linkage disequilibrium score regression analysis (LDSC, intercept=1.059)[94].

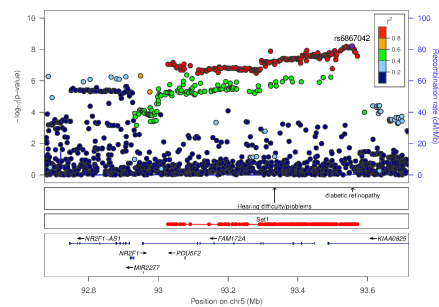
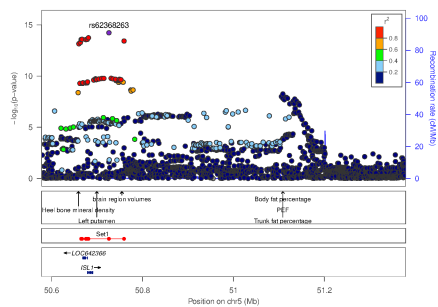
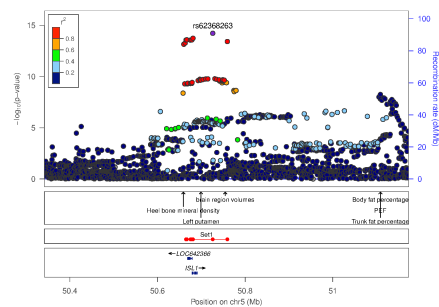
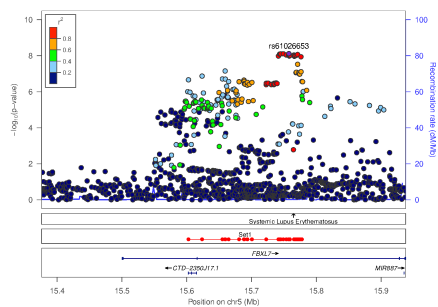
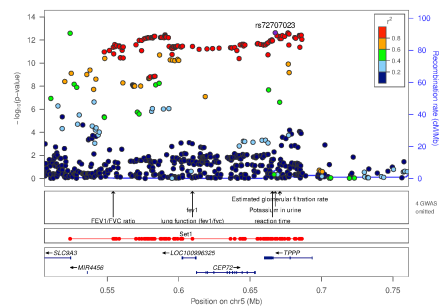
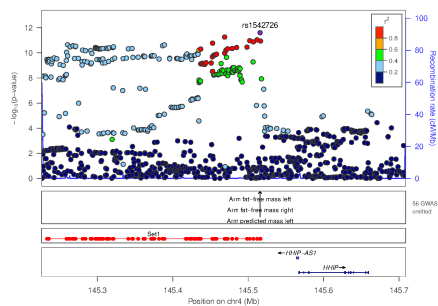
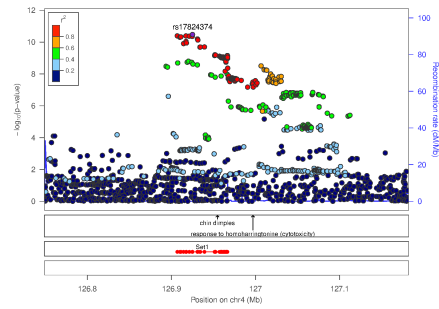
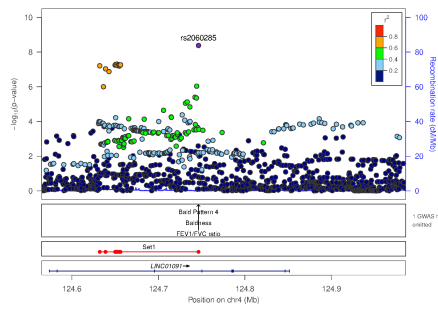


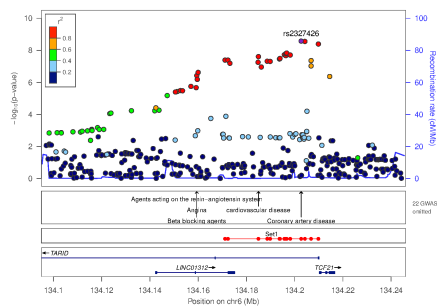
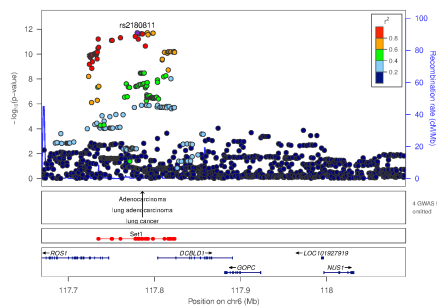
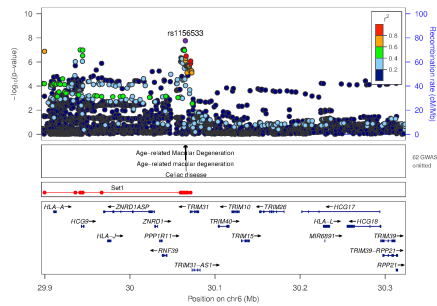
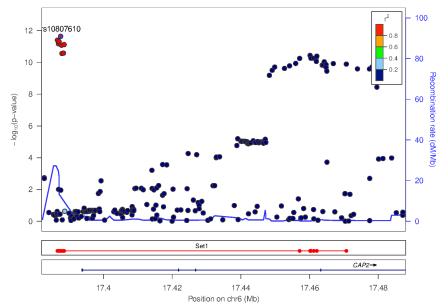
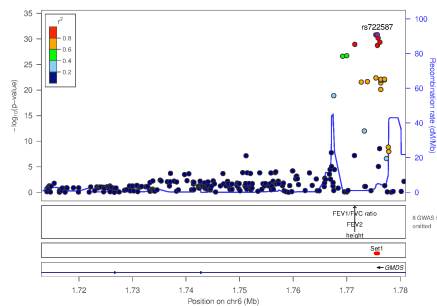
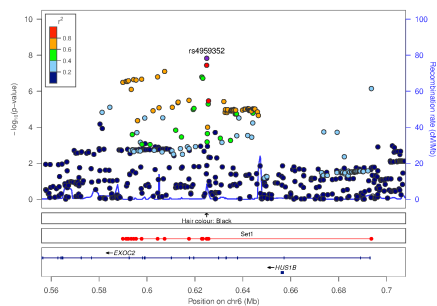
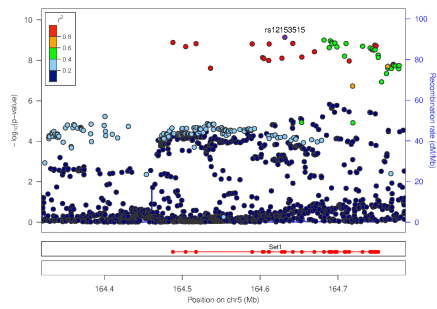
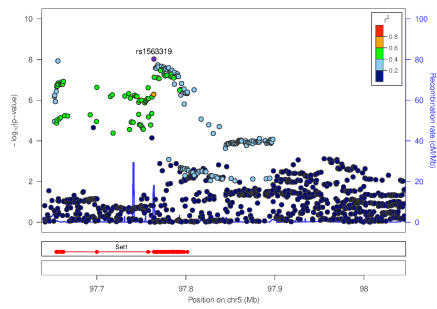


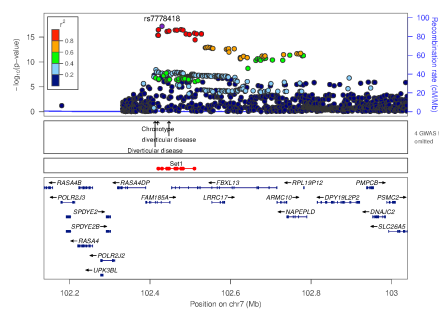
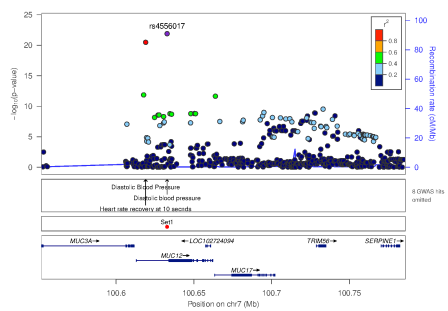
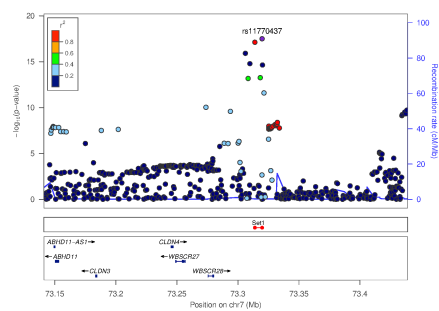
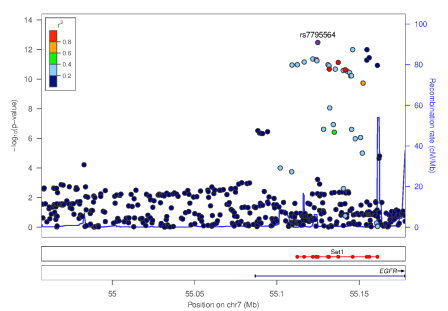
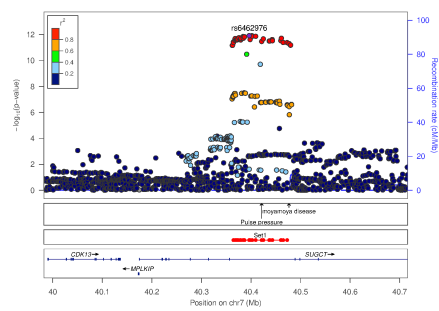
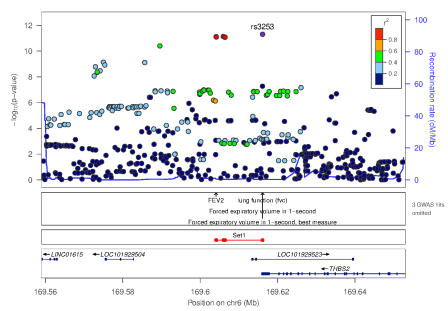
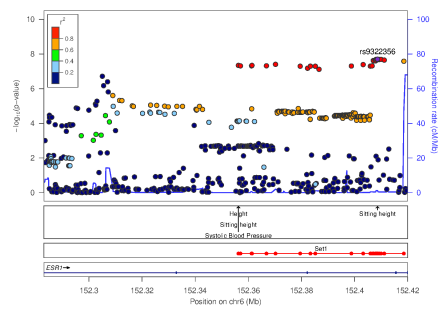
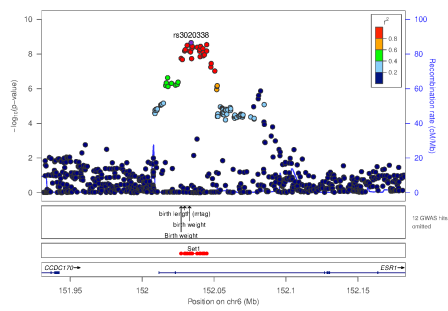




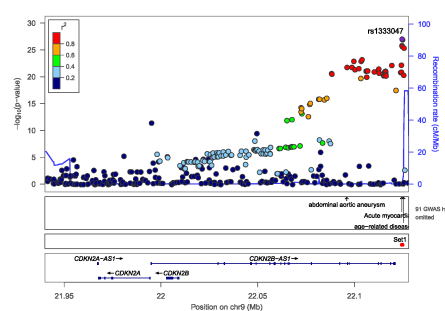
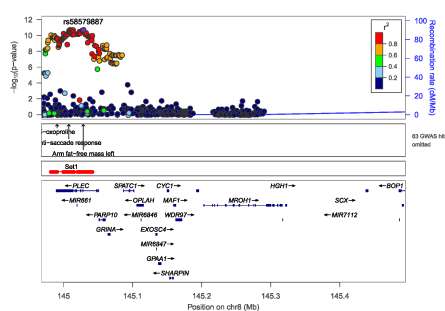
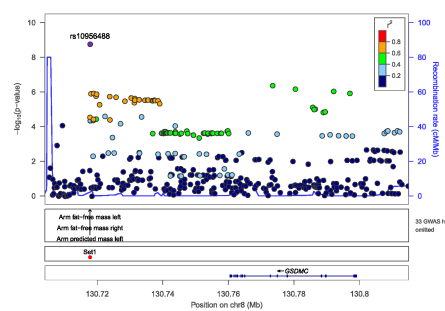
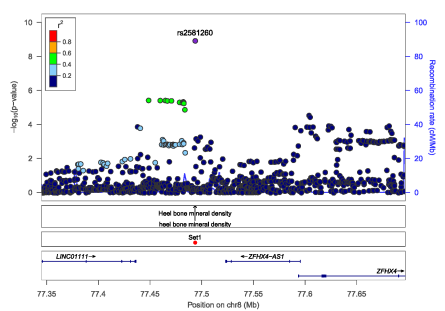
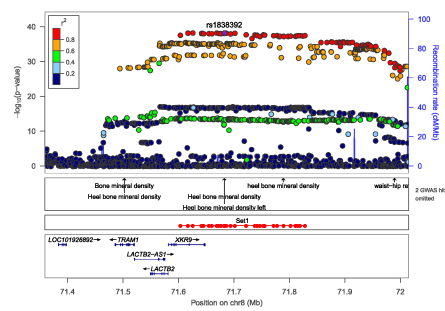
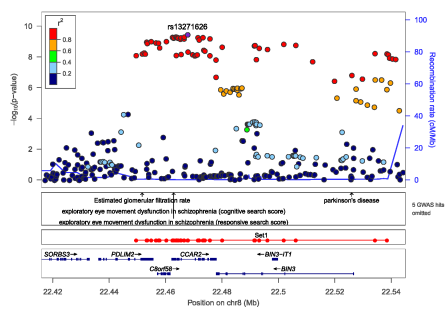
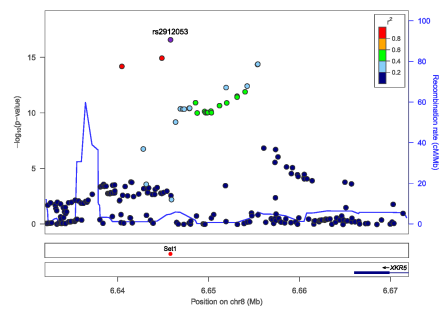
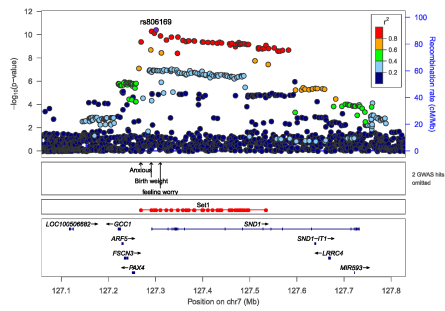


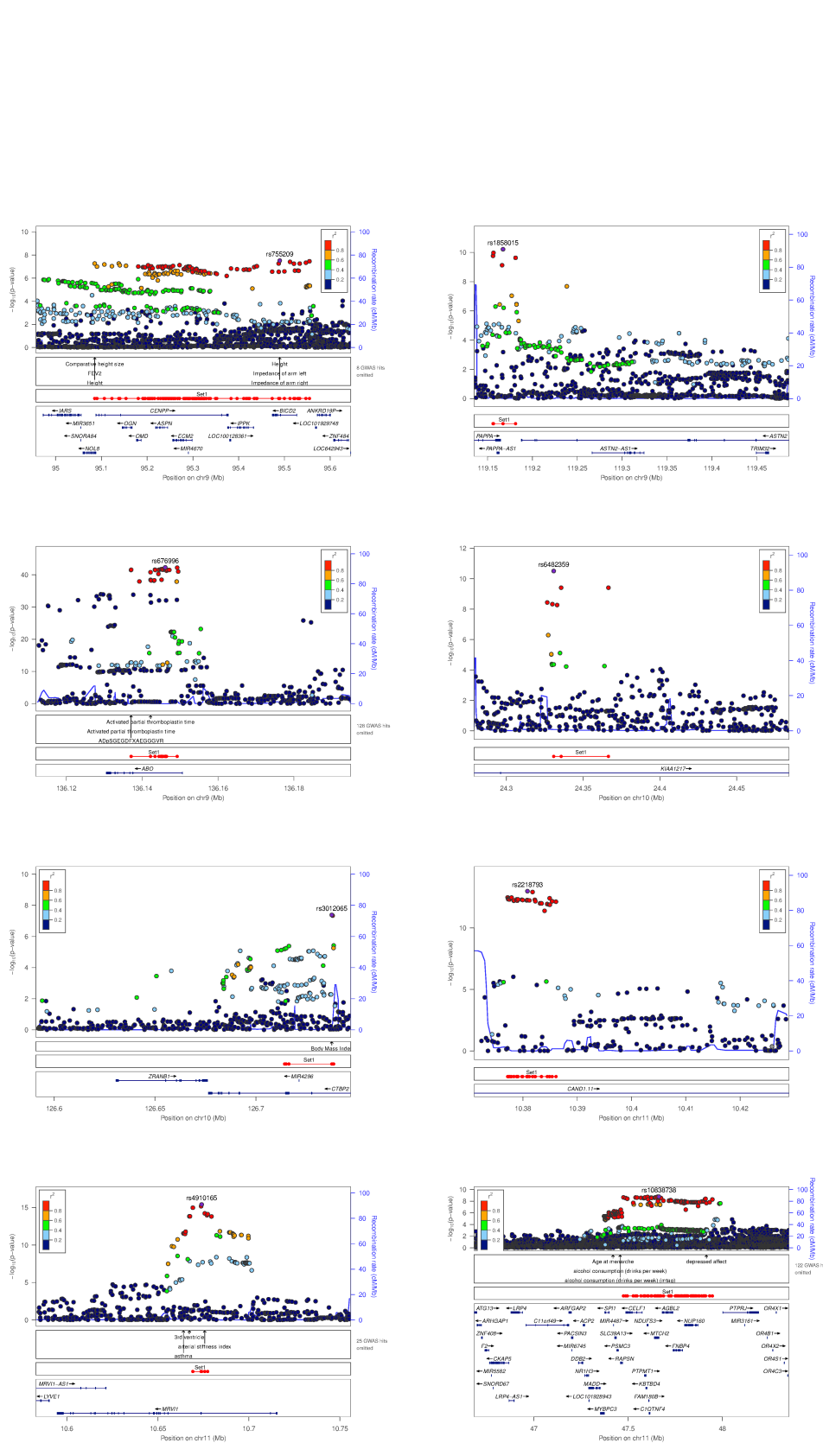


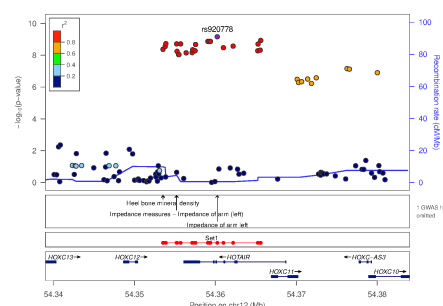
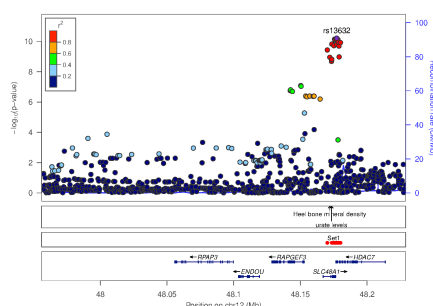
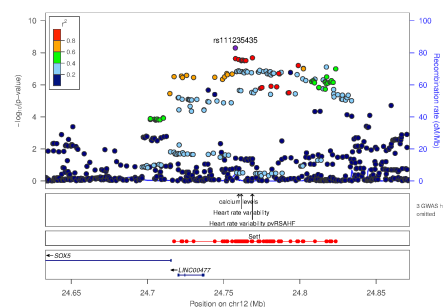
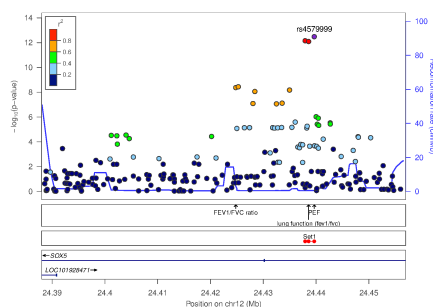
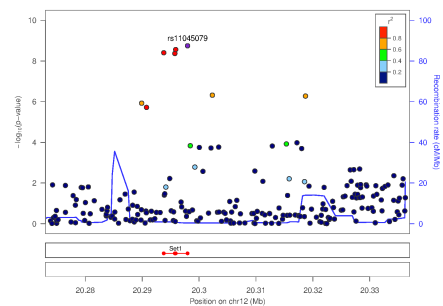
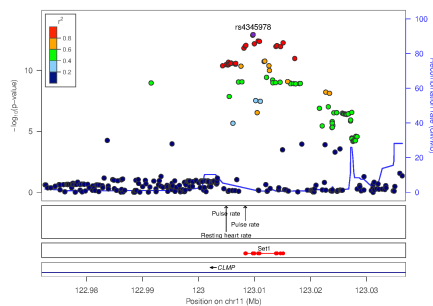
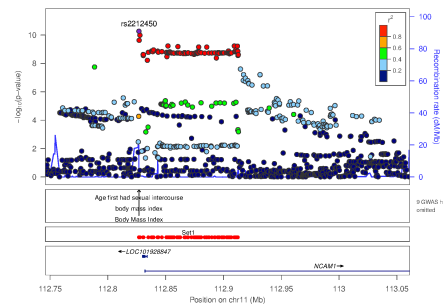
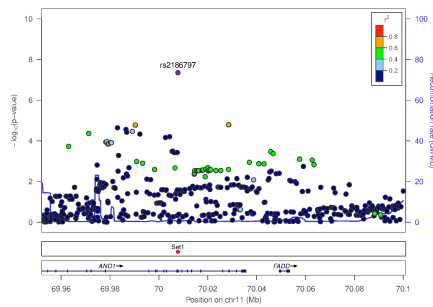


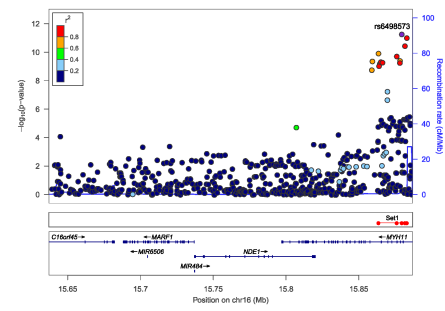
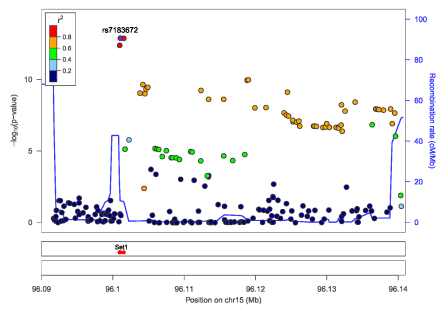
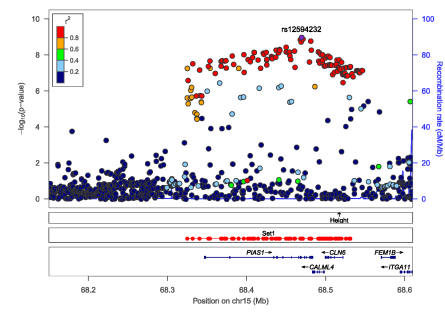
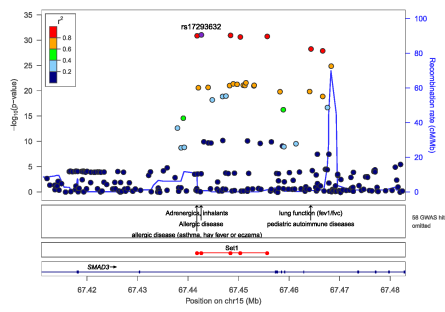
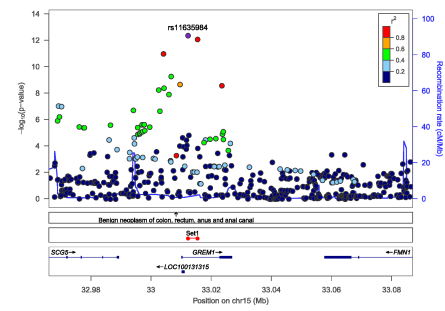
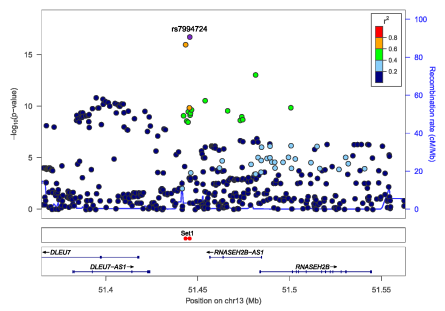
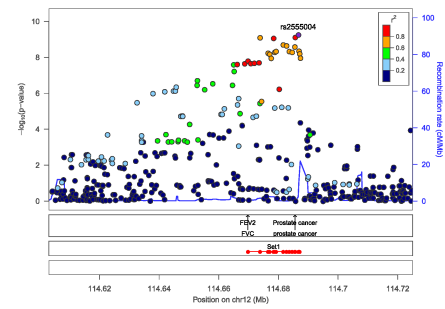
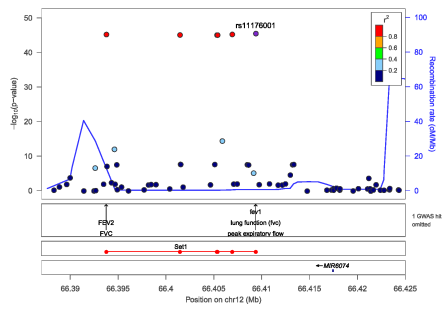


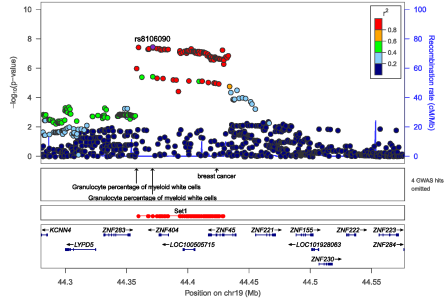
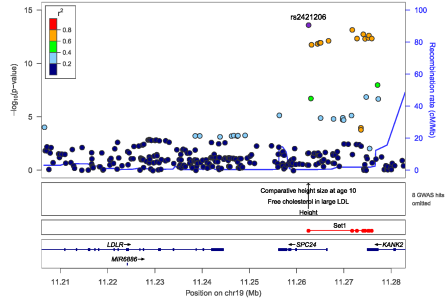
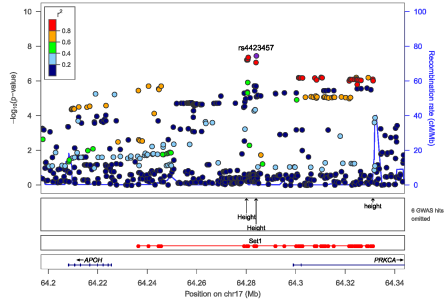
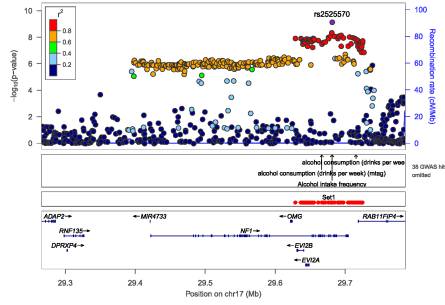
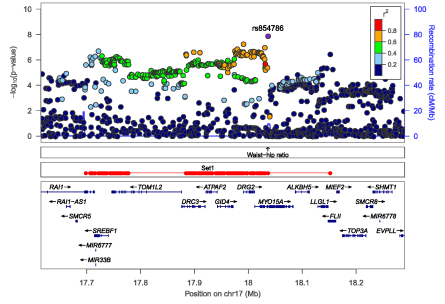
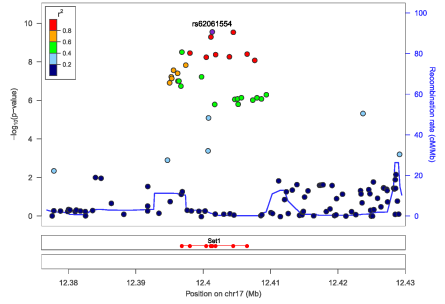
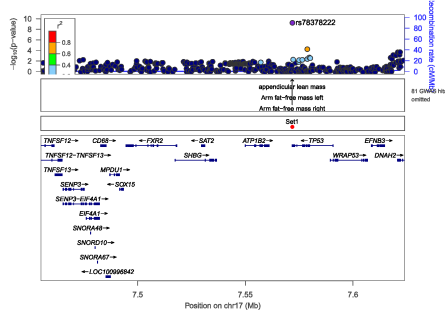
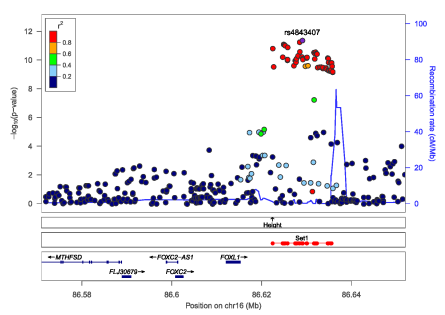




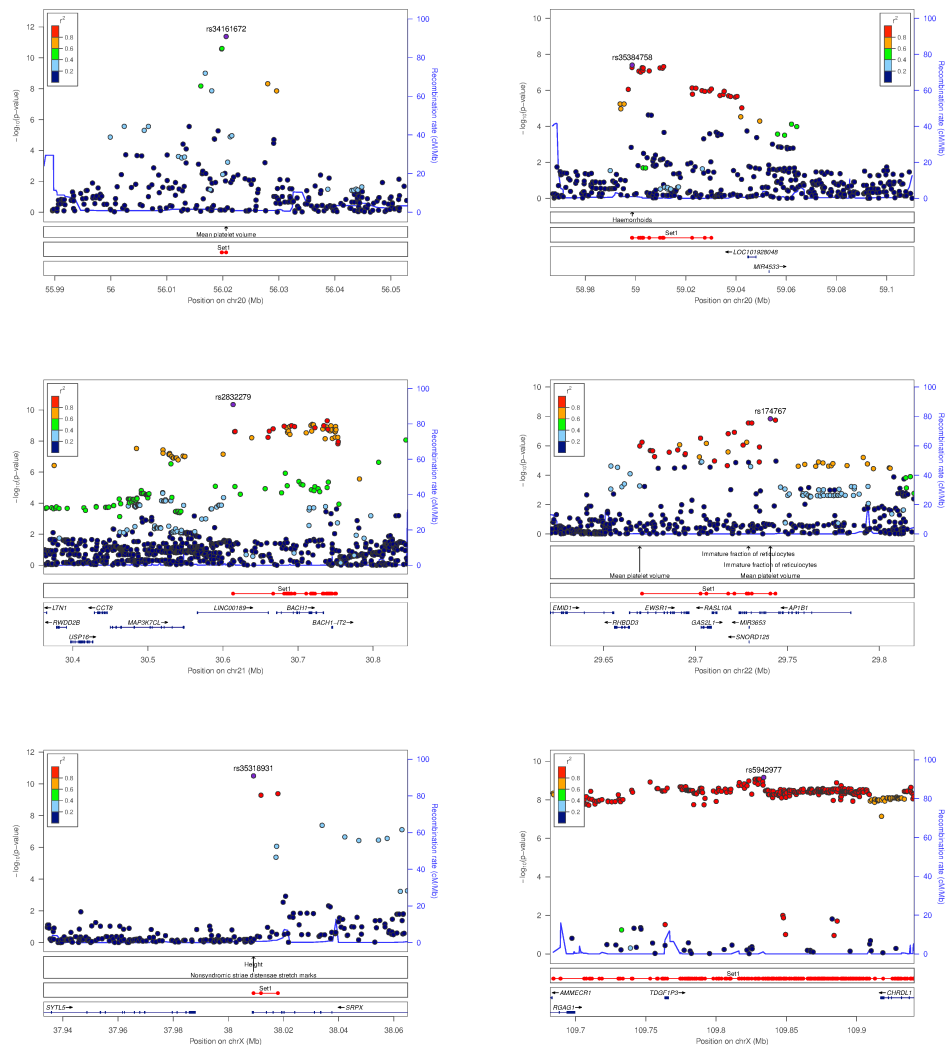






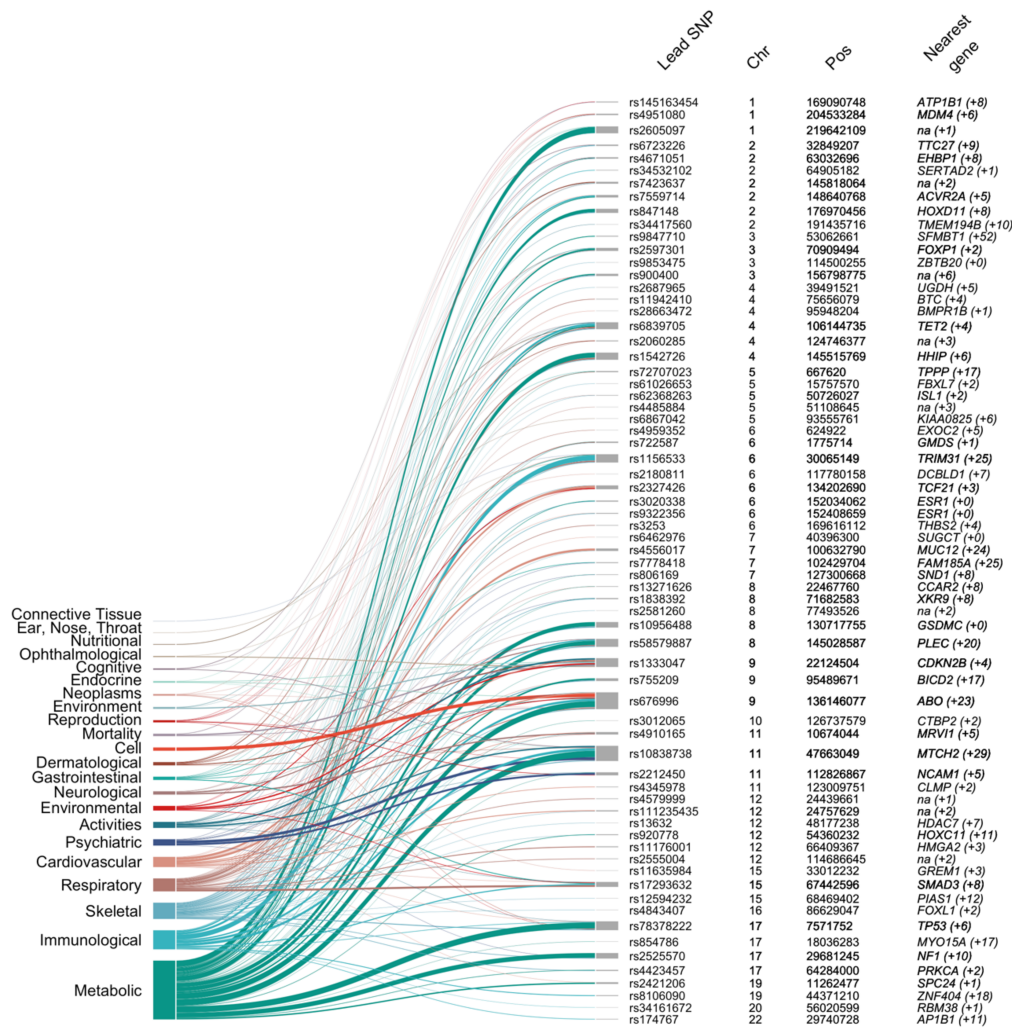






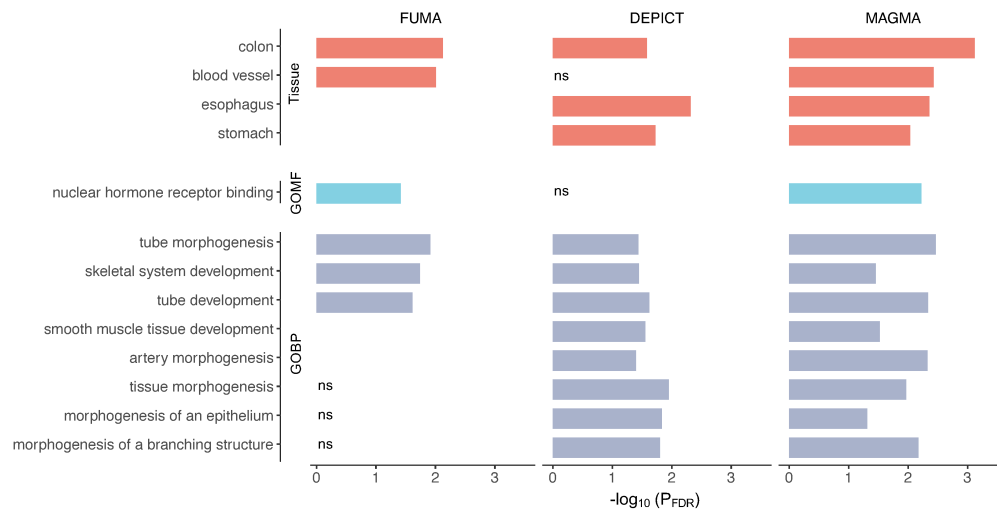
### Online Supplementary Figure S5. Regional association plots of HEM GWAS risk loci.

Shown are the  $-\log_{10} P$ -values from meta-analysis with regard to the physical location of markers and the degree of linkage disequilibrium ( $r^2$ ). Purple circle: lead SNP; line: recombination intensity (cM/Mb). Positions and gene annotations are according to NCBI's build 37 (hg19). Plots were generated using LocusZoom[42], also reporting the 95%-fine mapped credible sets at each locus.



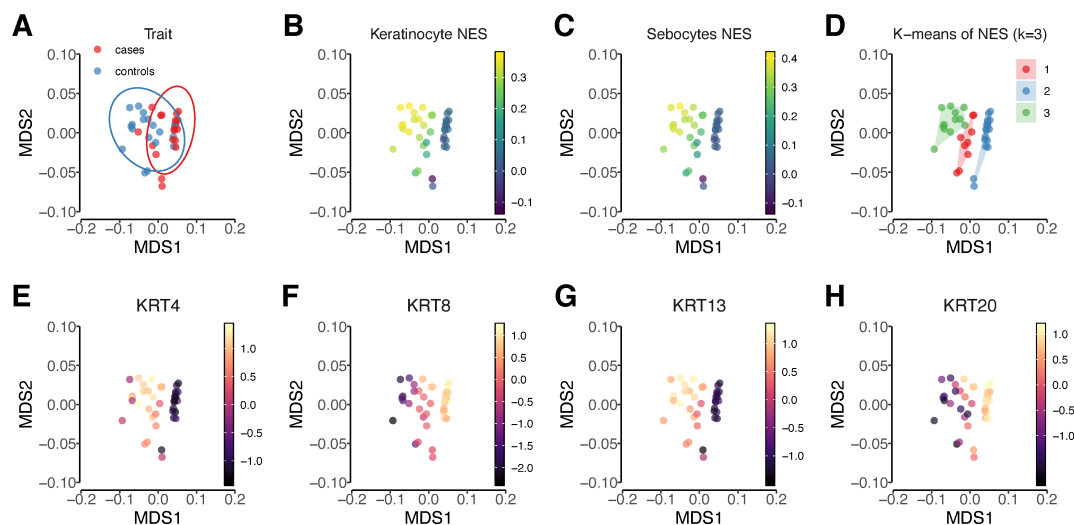
**Online Supplementary Figure S6: Previously reported associations of HEM risk loci with other traits and diseases, clustered by biological areas.**

The plot shows associations with other traits, extracted from the GWAS ATLAS for the 102 lead SNPs (and/or their  $r^2 > 0.8$  LD proxies) ordered by chromosome and chromosomal position. Associations are grouped by domain and represented with different colors. The ribbon size is proportional to the number of traits associated at the genome-wide significance level ( $P_{\text{Meta}} < 5 \times 10^{-8}$ ). Columns from left to right: Lead SNP – marker showing strongest association signal from each locus; Chr – chromosome; Pos – SNP position on chromosome (genome build hg19); Nearest gene (#genes within locus boundaries) – gene closest to the lead SNP (if within 100 kb distance, otherwise “na”) (**Methods**).



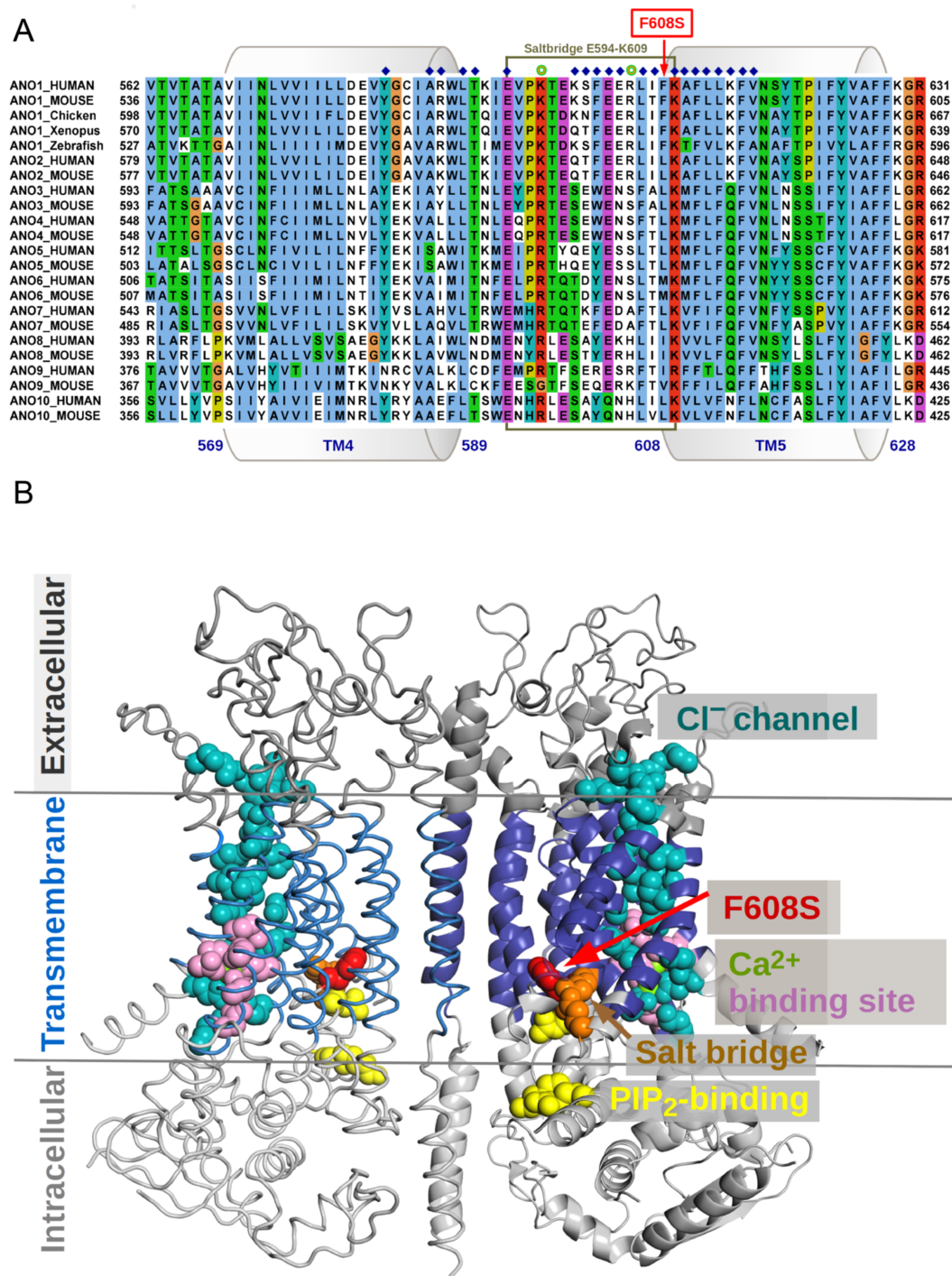
### Online Supplementary Figure S7. Gene set enrichment analyses of HEM genes.

Tissues and pathways are shown, which resulted significantly enriched in at least 2/3 analyses (using FUMA, MAGMA or DEPICT generate HEM gene lists; **Methods**). Gene Ontology Biological Processes (GOBP) and Molecular Function (GOMF) categories are reported; ns=non-significant; some tissues/pathways were not available in all analyses (missing bars).



**Online Supplementary Figure S8. Gene signature-based determination of anal canal zones.**

Multidimensional scaling (MDS) analysis of the transcriptome data using Spearman's correlation distance ( $1 - \text{correlation coefficient}$ ); **(A)** Colored by trait status, where the cases are enlarged hemorrhoidal tissue samples and controls are healthy hemorrhoidal tissue; **(B)** Colored by normalized enrichment score (NES) of keratinocyte cells; **(C)** Colored by NES of sebocytes; **(D)** Colored by clusters obtained by applying the k-means algorithm; **(E-H)** Colored by normalized expression values of anal canal marker genes, including *KRT4*, *KRT8*, *KRT13* and *KRT20*.



**Online Supplementary Figure S9. ANO1 Alignment of TM4-5 and ANO1 structure.**

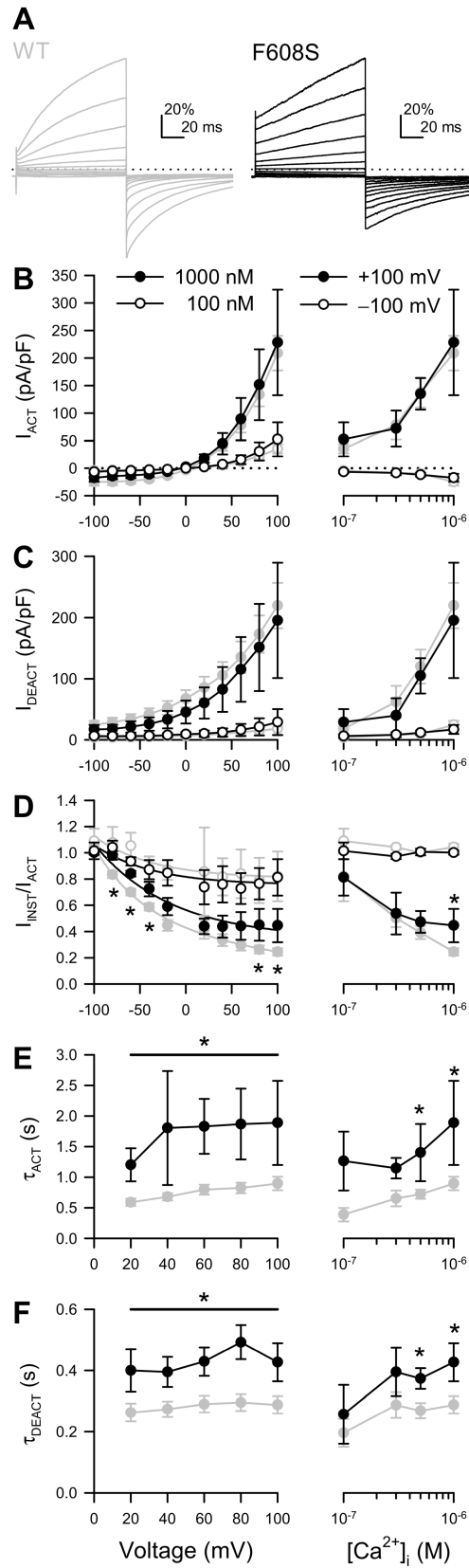
**(A)** Protein sequence alignment of ANO1 transmembrane helices TM4 and TM5 including the intracellular linker connected to TM5 via a salt bridge. Blue diamonds



mark conserved amino acid positions next to and in close proximity to the F608S variant. Green spheres mark PIP2-binding positions K597 and R605, whose mutation has been shown to lead to rapid channel inactivation through increased desensitization to  $\text{Ca}^{2+}$ . The same effect was observed with the mutation of E594 or K609 which form a stabilizing salt bridge[69]. Only hydrophobic amino acids (blue) are conserved at the site of the F608S variant. It is therefore predicted that the mutation to the polar serine destabilizes the local protein structure and affects the integrity of this salt bridge. Accelerated desensitization of the anion channel may result from conformational changes of the putative PIP2 binding site due to a disruption of the salt bridge[69]. **Consequently, F608S may be able to down-regulate ANO1 activity.**

**(B)** Structural model of the ANO1 dimer and localization of the F608S variant. The F608S variant (red spheres) is located at the beginning of transmembrane helix 5 and thus at the membrane-cytosolic interface and a predicted PIP2 interaction site[69] (yellow spheres). The exchange of the hydrophobic sidechain of phenylalanine (F) to a polar serine (S) within a conserved hydrophobic region is expected to destabilize the structure by disrupting the stability conducted by the salt bridge of K609 and E594 (orange spheres), which could accelerate the down-regulation of ANO1 by a faster channel inactivation by desensitization to  $\text{Ca}^{2+}$ . This effect was shown by an alanine mutation of the salt bridge[69]. The ANO1 structural model is based on cryo-electron microscopy of the murine homolog (PDB ID 5oyb[70]) The two monomers are distinguished by representation as ribbons and cartoons, respectively.

Extracellular and intracellular domains are colored dark and light grey, the transmembrane domain is blue. The Cl anion channel is highlighted as teal spheres. The  $\text{Ca}^{2+}$ -binding site is colored pink, the calcium atoms are shown as yellow-green spheres. Predicted PIP2-interacting residues (R481, K597 and R605[69]) are depicted as yellow spheres. Protein sequences were derived from the UniProt sequence database and visualized using JalView[95] with the Clustal X color scheme.



**Online Supplementary Figure S10. F608S mutant of ANO1 has high instantaneous current but slow voltage-dependent activation and deactivation kinetics *in vitro*.**

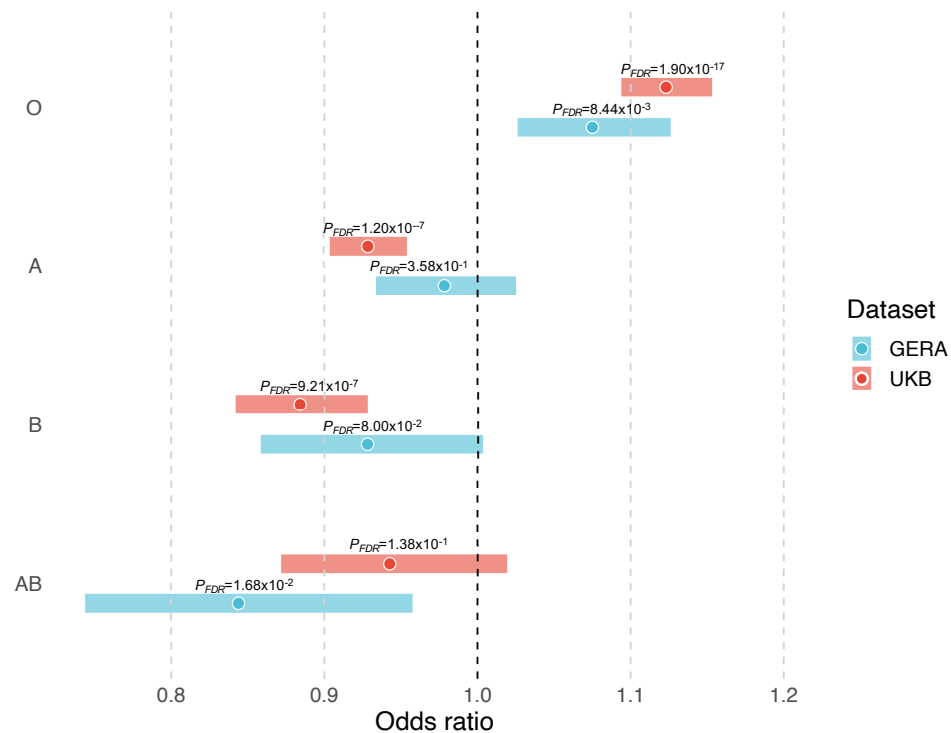
**(A)** Representative Cl<sup>-</sup> currents recorded from HEK293 cells transfected with wild-type ANO1 (left) or F608S-ANO1 (right), elicited by stepping for 1 s from -100 mV holding voltage to 100 through +100 mV. **(B-F)** Left, voltage-dependence at 1000 (●) or 100 nM [Ca<sup>2+</sup>]<sub>i</sub> (○); or right, [Ca<sup>2+</sup>]<sub>i</sub>-dependence at +100 (●) or -100 mV (○); of Cl<sup>-</sup> current parameters from HEK293 cells expressing WT- (gray) or F608S-ANO1 (black): Cl<sup>-</sup> current densities at the 1-s plateau (**B**, I<sub>ACT</sub>), tail currents immediately upon deactivation (**C**, I<sub>DEACT</sub>), ratios of the instantaneous Cl<sup>-</sup> current at 20 ms versus the plateau current at 1 s (**D**, I<sub>INST</sub>/I<sub>ACT</sub>), time constants of Cl<sup>-</sup> current during activation (**E**, T<sub>ACT</sub>) or deactivation (**F**, T<sub>DEACT</sub>) (\*P < 0.05, F608S vs. WT, by unpaired two-tailed t-test; n = 5-27 cells per [Ca<sup>2+</sup>]<sub>i</sub>).

To determine the functional impact of F608S on human ANO1, we recorded whole-cell voltage-dependent Ca<sup>2+</sup>-activated Cl<sup>-</sup> currents from HEK293 cells expressing wild-type or F608S-ANO1 at 100-1000 nM intracellular Ca<sup>2+</sup> ([Ca<sup>2+</sup>]<sub>i</sub>). Ca<sup>2+</sup>-activated Cl<sup>-</sup> current densities of F608S-ANO1 were similar to WT for both activation (**B**) and deactivation (**C**) at all tested voltages and [Ca<sup>2+</sup>]<sub>i</sub> concentrations. However, the kinetics of the two constructs were different. F608S-ANO1 Cl<sup>-</sup> currents had a larger ratio of instantaneous-to-plateau current (I<sub>INST</sub>/I<sub>ACT</sub>) at high [Ca<sup>2+</sup>]<sub>i</sub> (**D**). Moreover, F608S-ANO1 activated and deactivated slower than WT, as reflected in an increase in the time constants of activation (T<sub>ACT</sub>, **E**) and deactivation (T<sub>DEACT</sub>, **F**) at positive voltages (+20 to +100 mV) and high [Ca<sup>2+</sup>]<sub>i</sub> (500-1000 nM).



**Online Supplementary Figure S11. *Sushi repeat-containing protein (SRPX) structure und alignment.***

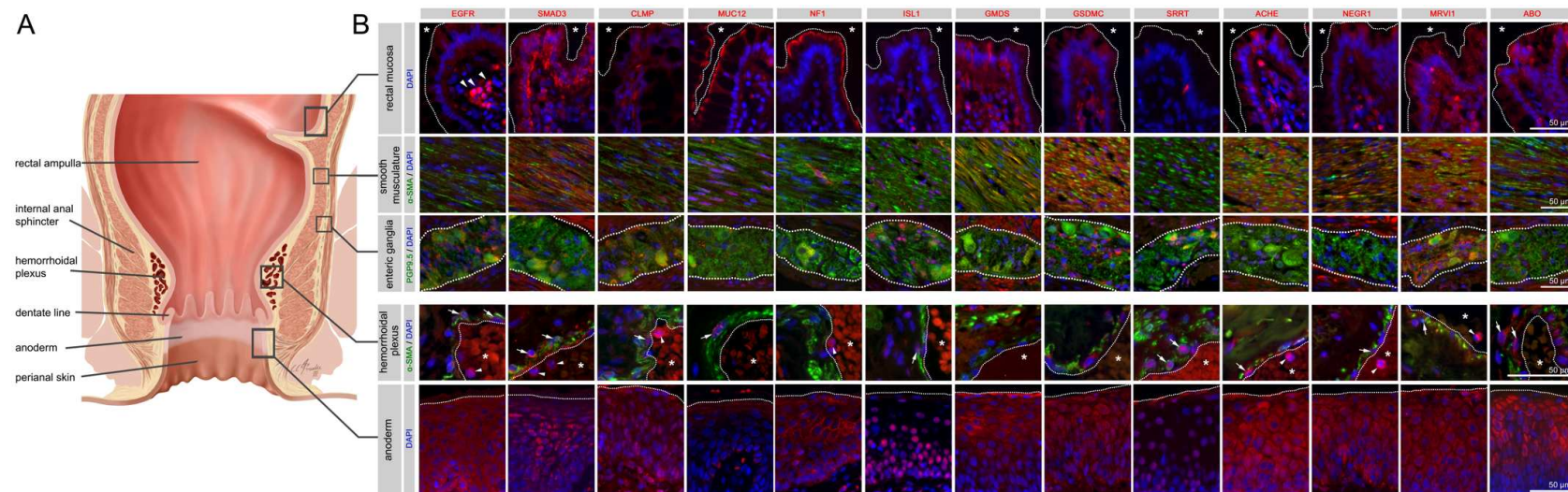
**(A)** SRPX domain structure and the predicted protein fold of the C-terminal domain. The N-terminal signal peptide is shown as a green dashed line. Predicting the 3D location of the Ser413Phe variant is based on a model with lower confidence, with loop and helical structures being less reliable than the central beta sheet. In this model it is predicted that the polar Ser413 stabilizes loops originating from strands 1, 3 and 4, and a mutation to a hydrophobic phenylalanine could interfere with this function **(B)** Multiple sequence alignment with predicted secondary structures. Conserved sequence positions are largely consistent with the pfam13778 family, in particular with the central beta sheet, which enhances the confidence of the core regions in the above structural model. Ser413Phe is located adjacent to the conserved beta strand 3 and the invariant Phe414 which supports an important structural role of the variant. For further details see **Methods**, section ***In silico* variant protein analysis**.



**Online Supplementary Figure S12. ABO blood groups and HEM risk in UKBB and GERA.**

The plot shows odds ratios (and 95% confidence intervals) from testing ABO blood groups vs HEM risk in UKBB and GERA (**Methods**). An association test based on logistic regression is used to test for a significant HEM association for each of the four blood groups, taking into account sex, age, BMI and the top 10 PCs from PCA. FDR correction was applied to correct for multiple testing. FDR: false discovery rate.





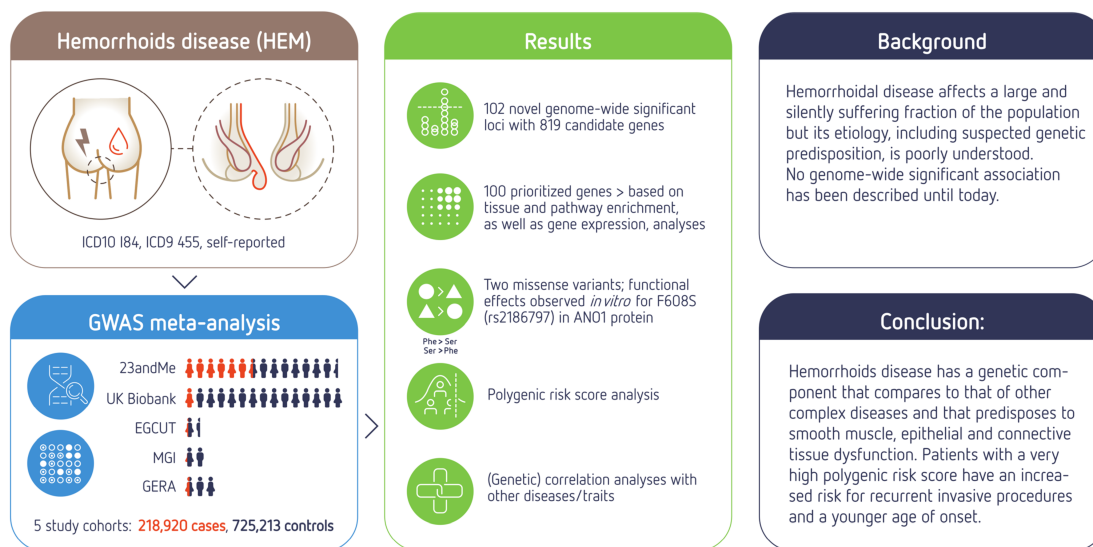
### Online Supplementary Figure S13. Immunohistochemistry for selected HEM candidate proteins.

Illustration of the rectum and anal canal (A) with indication the site-specific localization of the immunohistochemical panels analyzed in (B). Fluorescence immunohistochemistry (B) for selected HEM candidate proteins (see also **online supplementary table S11**), encoded by candidate genes within our 102 identified genome-wide significant loci, are shown. *SRPX* (rs35318931), *ANO1* (rs2186797) and *MYH11* (rs6498573) were determined as prioritized HEM genes in our study. *ANO1* and *SRPX* are interesting HEM candidate genes since the lead SNPs at these loci are (missense) coding variants. *MYH11* is also a main hub gene within the M1 co-expression module of our transcriptome analysis. Given the ABO blood group association observed in our study in HEM patients (**online supplementary figure S12**), we have included ABO as further target for immunohistochemistry.

Antibody staining was performed on colorectal FFPE tissue specimens from control individuals. The rows correspond to the rectal mucosa (top row, epithelial surface delimited by dashed line, \*: intestinal lumen), smooth muscle (second row), enteric ganglia

(third row, ganglionic boundaries delimited by dashed line), hemorrhoidal plexus (fourth row, endothelial surface delimited by dashed line, \*: vascular lumen), and the anoderm (bottom row, border of the anoderm delimited by dashed line). Blue: DAPI; green:  $\alpha$ -SMA (anti-alpha smooth muscle actin antibody) for row 2 and 4 (smooth musculature/hemorrhoidal plexus) and PGP9.5 (member of the ubiquitin hydrolase family of proteins, neuronal marker) for row 3 (enteric ganglia); red: antibody for the respective candidate protein. Arrows point to corresponding candidate-positive cells within the vascular wall. Arrowheads point to corresponding candidate-positive nucleated immune cells.

## Visual Abstract

Online Supplementary Figure S14. *Graphical abstract of the study.*

## SUPPLEMENTARY REFERENCES

- 1 Tung JY, Do CB, Hinds DA, Kiefer AK, Macpherson JM, Chowdry AB, *et al.* Efficient replication of over 180 genetic associations with self-reported medical data. *PLoS One* 2011;**6**:e23473.
- 2 Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* 2018;**562**:203-9.
- 3 Leitsalu L, Haller T, Esko T, Tammesoo ML, Alavere H, Snieder H, *et al.* Cohort Profile: Estonian Biobank of the Estonian Genome Center, University of Tartu. *Int J Epidemiol* 2015;**44**:1137-47.
- 4 Fritsche LG, Gruber SB, Wu Z, Schmidt EM, Zawistowski M, Moser SE, *et al.* Association of Polygenic Risk Scores for Multiple Cancers in a Phenome-wide Study: Results from The Michigan Genomics Initiative. *Am J Hum Genet* 2018;**102**:1048-61.
- 5 Krawczak M, Nikolaus S, von Eberstein H, Croucher PJ, El Mokhtari NE, Schreiber S. PopGen: population-based recruitment of patients and controls for the analysis of complex genotype-phenotype relationships. *Community Genet* 2006;**9**:55-61.
- 6 Krokstad S, Langhammer A, Hveem K, Holmen TL, Midthjell K, Stene TR, *et al.* Cohort Profile: the HUNT Study, Norway. *Int J Epidemiol* 2013;**42**:968-77.
- 7 Hansen TF, Banasik K, Erikstrup C, Pedersen OB, Westergaard D, Chmura PJ, *et al.* DBDS Genomic Cohort, a prospective and comprehensive resource for integrative and temporal analysis of genetic, environmental and lifestyle factors affecting health of blood donors. *BMJ Open* 2019;**9**:e028401.
- 8 Burgdorf KS, Simonsen J, Sundby A, Rostgaard K, Pedersen OB, Sorensen E, *et al.* Socio-demographic characteristics of Danish blood donors. *PLoS One* 2017;**12**:e0169112.
- 9 Jensen TB, Jimenez-Solem E, Cortes R, Betzer C, Boge Breinholt S, Meidahl Petersen K, *et al.* Content and validation of the Electronic Patient Medication module (EPM)-the administrative in-hospital drug use database in the Capital Region of Denmark. *Scand J Public Health* 2020;**48**:43-8.
- 10 Durand EY, Do CB, Mountain JL, Macpherson JM. Ancestry Composition: A Novel, Efficient Pipeline for Ancestry Deconvolution. *bioRxiv* 2014:010512.
- 11 Henn BM, Hon L, Macpherson JM, Eriksson N, Saxonov S, Pe'er I, *et al.* Cryptic distant relatives are common in both isolated and cosmopolitan genetic samples. *PLoS One* 2012;**7**:e34267.
- 12 Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin RM, *et al.* A map of human genome variation from population-scale sequencing. *Nature* 2010;**467**:1061-73.
- 13 Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am J Hum Genet* 2007;**81**:1084-97.
- 14 Fuchsberger C, Abecasis GR, Hinds DA. minimac2: faster genotype imputation. *Bioinformatics* 2015;**31**:782-4.
- 15 Wain LV, Shrine N, Miller S, Jackson VE, Ntalla I, Soler Artigas M, *et al.* Novel insights into the genetics of smoking behaviour, lung function, and chronic obstructive pulmonary disease (UK BiLEVE): a genetic association study in UK Biobank. *Lancet Respir Med* 2015;**3**:769-81.
- 16 Abraham G, Inouye M. Fast principal component analysis of large-scale genome-wide data. *PLoS One* 2014;**9**:e93766.

- 17 International HapMap C, Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve LL, *et al.* A second generation human haplotype map of over 3.1 million SNPs. *Nature* 2007;**449**:851-61.
- 18 Delaneau O, Marchini J, Zagury JF. A linear complexity phasing method for thousands of genomes. *Nat Methods* 2011;**9**:179-81.
- 19 Mitt M, Kals M, Parn K, Gabriel SB, Lander ES, Palotie A, *et al.* Improved imputation accuracy of rare and low-frequency variants using population-specific high-coverage WGS-based imputation reference panel. *Eur J Hum Genet* 2017;**25**:869-76.
- 20 Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet* 2009;**5**:e1000529.
- 21 Jun G, Flickinger M, Hetrick KN, Romm JM, Doheny KF, Abecasis GR, *et al.* Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am J Hum Genet* 2012;**91**:839-48.
- 22 Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. Robust relationship inference in genome-wide association studies. *Bioinformatics* 2010;**26**:2867-73.
- 23 Li JZ, Absher DM, Tang H, Southwick AM, Casto AM, Ramachandran S, *et al.* Worldwide human relationships inferred from genome-wide patterns of variation. *Science* 2008;**319**:1100-4.
- 24 Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 2015;**4**:7.
- 25 McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet* 2016;**48**:1279-83.
- 26 Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet* 2006;**2**:e190.
- 27 Howie B, Fuchsberger C, Stephens M, Marchini J, Abecasis GR. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat Genet* 2012;**44**:955-9.
- 28 Price AL, Weale ME, Patterson N, Myers SR, Need AC, Shianna KV, *et al.* Long-range LD can confound genome scans in admixed populations. *Am J Hum Genet* 2008;**83**:132-5; author reply 5-9.
- 29 Jun TH, Rouf Mian MA, Michel AP. Genetic mapping revealed two loci for soybean aphid resistance in PI 567301B. *Theor Appl Genet* 2012;**124**:13-22.
- 30 Guo Y, He J, Zhao S, Wu H, Zhong X, Sheng Q, *et al.* Illumina human exome genotyping array clustering and quality control. *Nat Protoc* 2014;**9**:2643-62.
- 31 Wang C, Zhan X, Liang L, Abecasis GR, Lin X. Improved ancestry estimation for both genotyping and sequencing data using projection procrustes analysis and genotype imputation. *Am J Hum Genet* 2015;**96**:926-37.
- 32 Loh PR, Palamara PF, Price AL. Fast and accurate long-range phasing in a UK Biobank cohort. *Nat Genet* 2016;**48**:811-6.
- 33 Das S, Forer L, Schonherr S, Sidore C, Locke AE, Kwong A, *et al.* Next-generation genotype imputation service and methods. *Nat Genet* 2016;**48**:1284-7.
- 34 König IR, Loley C, Erdmann J, Ziegler A. How to include chromosome X in your genome-wide association study. *Genet Epidemiol* 2014;**38**:97-103.
- 35 Loh PR, Tucker G, Bulik-Sullivan BK, Vilhjalmsón BJ, Finucane HK, Salem RM, *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet* 2015;**47**:284-90.
- 36 Zhou W, Nielsen JB, Fritsche LG, Dey R, Gabrielsen ME, Wolford BN, *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat Genet* 2018;**50**:1335-41.



- 37 Winkler TW, Day FR, Croteau-Chonka DC, Wood AR, Locke AE, Magi R, *et al.* Quality control and conduct of genome-wide association meta-analyses. *Nat Protoc* 2014;**9**:1192-212.
- 38 Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 2010;**26**:2190-1.
- 39 Watanabe K, Taskesen E, van Bochoven A, Posthuma D. Functional mapping and annotation of genetic associations with FUMA. *Nat Commun* 2017;**8**:1826.
- 40 Pers TH, Karjalainen JM, Chan Y, Westra HJ, Wood AR, Yang J, *et al.* Biological interpretation of genome-wide association studies using predicted gene functions. *Nat Commun* 2015;**6**:5890.
- 41 de Leeuw CA, Mooij JM, Heskes T, Posthuma D. MAGMA: generalized gene-set analysis of GWAS data. *PLoS Comput Biol* 2015;**11**:e1004219.
- 42 Pruim RJ, Welch RP, Sanna S, Teslovich TM, Chines PS, Gliedt TP, *et al.* LocusZoom: regional visualization of genome-wide association scan results. *Bioinformatics* 2010;**26**:2336-7.
- 43 Benner C, Spencer CC, Havulinna AS, Salomaa V, Ripatti S, Pirinen M. FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics* 2016;**32**:1493-501.
- 44 Cuéllar-Partida G, Lundberg M, Kho PF, D'Urso S, Gutiérrez-Mondragón LF, Ngo TT, *et al.* Complex-Traits Genetics Virtual Lab: A community-driven web platform for post-GWAS analyses. *bioRxiv* 2019:518027.
- 45 Bhattacharjee S, Rajaraman P, Jacobs KB, Wheeler WA, Melin BS, Hartge P, *et al.* A subset-based approach improves power and interpretation for the combined analysis of genetic association studies of heterogeneous traits. *Am J Hum Genet* 2012;**90**:821-35.
- 46 Schafmayer C, Harrison JW, Buch S, Lange C, Reichert MC, Hofer P, *et al.* Genome-wide association analysis of diverticular disease points towards neuromuscular, connective tissue and epithelial pathomechanisms. *Gut* 2019;**68**:854-65.
- 47 Bonfiglio F, Zheng T, Garcia-Etxebarria K, Hadizadeh F, Bujanda L, Bresso F, *et al.* Female-Specific Association Between Variants on Chromosome 9 and Self-Reported Diagnosis of Irritable Bowel Syndrome. *Gastroenterology* 2018;**155**:168-79.
- 48 Choi SW, O'Reilly PF. PRSice-2: Polygenic Risk Score software for biobank-scale data. *Gigascience* 2019;**8**.
- 49 Kamat MA, Blackshaw JA, Young R, Surendran P, Burgess S, Danesh J, *et al.* PhenoScanner V2: an expanded tool for searching human genotype-phenotype associations. *Bioinformatics* 2019;**35**:4851-3.
- 50 Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res* 2019;**47**:D1005-D12.
- 51 Watanabe K, Stringer S, Frei O, Umicevic Mirkov M, de Leeuw C, Polderman TJC, *et al.* A global overview of pleiotropy and genetic architecture in complex traits. *Nat Genet* 2019;**51**:1339-48.
- 52 Jensen AB, Moseley PL, Oprea TI, Ellesoe SG, Eriksson R, Schmock H, *et al.* Temporal disease trajectories condensed from population-wide registry data covering 6.2 million patients. *Nat Commun* 2014;**5**:4022.
- 53 Rozowsky J, Kitchen RR, Park JJ, Galeev TR, Diao J, Warrell J, *et al.* exceRpt: A Comprehensive Analytic Platform for Extracellular RNA Profiling. *Cell Syst* 2019;**8**:352-7 e3.
- 54 Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res* 2002;**30**:207-10.
- 55 Fenger C. The anal transitional zone. *Acta Pathol Microbiol Immunol Scand Suppl* 1987;**289**:1-42.



- 56 Iacobuzio-Donahue CA. Inflammatory and Neoplastic Disorders of the Anal Canal. In: Robert Odze JG, ed. *Surgical Pathology of the GI Tract, Liver, Biliary Tract, and Pancreas*: Elsevier Inc., 2009:733-61.
- 57 Aran D, Hu Z, Butte AJ. xCell: digitally portraying the tissue cellular heterogeneity landscape. *Genome Biol* 2017;**18**:220.
- 58 Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* 2014;**15**:550.
- 59 Barbie DA, Tamayo P, Boehm JS, Kim SY, Moody SE, Dunn IF, *et al.* Systematic RNA interference reveals that oncogenic KRAS-driven cancers require TBK1. *Nature* 2009;**462**:108-12.
- 60 Hanzelmann S, Castelo R, Guinney J. GSEA: gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics* 2013;**14**:7.
- 61 Williams GR, Talbot IC, Northover JM, Leigh IM. Keratin expression in the normal anal canal. *Histopathology* 1995;**26**:39-44.
- 62 McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res* 2012;**40**:4288-97.
- 63 Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics* 2008;**9**:559.
- 64 Russo PST, Ferreira GR, Cardozo LE, Burger MC, Arias-Carrasco R, Maruyama SR, *et al.* CEMiTool: a Bioconductor package for performing comprehensive modular co-expression analyses. *BMC Bioinformatics* 2018;**19**:56.
- 65 Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS* 2012;**16**:284-7.
- 66 Severe Covid GG, Ellinghaus D, Degenhardt F, Bujanda L, Buti M, Albillos A, *et al.* Genomewide Association Study of Severe Covid-19 with Respiratory Failure. *N Engl J Med* 2020;**383**:1522-34.
- 67 Cossais F, Lange C, Barrenschee M, Moding M, Ebsen M, Vogel I, *et al.* Altered enteric expression of the homeobox transcription factor Phox2b in patients with diverticular disease. *United European Gastroenterol J* 2019;**7**:349-57.
- 68 De Jesus-Perez JJ, Cruz-Rangel S, Espino-Saldana AE, Martinez-Torres A, Qu Z, Hartzell HC, *et al.* Phosphatidylinositol 4,5-bisphosphate, cholesterol, and fatty acids modulate the calcium-activated chloride channel TMEM16A (ANO1). *Biochim Biophys Acta Mol Cell Biol Lipids* 2018;**1863**:299-312.
- 69 Le SC, Jia Z, Chen J, Yang H. Molecular basis of PIP2-dependent regulation of the Ca(2+)-activated chloride channel TMEM16A. *Nat Commun* 2019;**10**:3769.
- 70 Paulino C, Kalienkova V, Lam AKM, Neldner Y, Dutzler R. Activation mechanism of the calcium-activated chloride channel TMEM16A revealed by cryo-EM. *Nature* 2017;**552**:421-5.
- 71 Dang S, Feng S, Tien J, Peters CJ, Bulkley D, Lolicato M, *et al.* Cryo-EM structures of the TMEM16A calcium-activated chloride channel. *Nature* 2017;**552**:426-9.
- 72 Kirkitadze MD, Barlow PN. Structure and flexibility of the multiple domain proteins that regulate complement activation. *Immunol Rev* 2001;**180**:146-61.
- 73 Callebaut I, Gilges D, Vigon I, Mornon JP. HYR, an extracellular module involved in cellular adhesion and related to the immunoglobulin-like fold. *Protein Sci* 2000;**9**:1382-90.
- 74 Bommer GT, Jager C, Durr EM, Baehs S, Eichhorst ST, Brabletz T, *et al.* DRO1, a gene down-regulated by oncogenes, mediates growth inhibition in colon and pancreatic cancer cells. *J Biol Chem* 2005;**280**:7962-75.

- 75 Pawlowski K, Muszewska A, Lenart A, Szczepinska T, Godzik A, Grynberg M. A widespread peroxiredoxin-like domain present in tumor suppression- and progression-implicated proteins. *BMC Genomics* 2010;**11**:590.
- 76 Inoue H, Pan J, Hakura A. Suppression of v-src transformation by the drs gene. *J Virol* 1998;**72**:2532-7.
- 77 Tambe Y, Yoshioka-Yamashita A, Mukaiho K, Haraguchi S, Chano T, Isono T, *et al.* Tumor prone phenotype of mice deficient in a novel apoptosis-inducing gene, drs. *Carcinogenesis* 2007;**28**:777-84.
- 78 Tambe Y, Hasebe M, Kim CJ, Yamamoto A, Inoue H. The drs tumor suppressor regulates glucose metabolism via lactate dehydrogenase-B. *Mol Carcinog* 2016;**55**:52-63.
- 79 Burgstaller G, Oehrle B, Gerckens M, White ES, Schiller HB, Eickelberg O. The instructive extracellular matrix of the lung: basic composition and alterations in chronic lung disease. *Eur Respir J* 2017;**50**.
- 80 Wilson R, Norris EL, Brachvogel B, Angelucci C, Zivkovic S, Gordon L, *et al.* Changes in the chondrocyte and extracellular matrix proteome during post-natal mouse cartilage development. *Mol Cell Proteomics* 2012;**11**:M111 014159.
- 81 Naba A, Clauser KR, Whittaker CA, Carr SA, Tanabe KK, Hynes RO. Extracellular matrix signatures of human primary metastatic colon cancers and their metastases to liver. *BMC Cancer* 2014;**14**:518.
- 82 Perea-Gil I, Uriarte JJ, Prat-Vidal C, Galvez-Monton C, Roura S, Lucia-Valldeperas A, *et al.* In vitro comparative study of two decellularization protocols in search of an optimal myocardial scaffold for recellularization. *Am J Transl Res* 2015;**7**:558-73.
- 83 Burnicka-Turek O, Kata A, Buyandelger B, Ebermann L, Kramann N, Burfeind P, *et al.* Pelota interacts with HAX1, EIF3G and SRPX and the resulting protein complexes are associated with the actin cytoskeleton. *BMC Cell Biol* 2010;**11**:28.
- 84 Royer-Zemmour B, Ponsolle-Lenfant M, Gara H, Roll P, Leveque C, Massacrier A, *et al.* Epileptic and developmental disorders of the speech cortex: ligand/receptor interaction of wild-type and mutant SRPX2 with the plasminogen activator receptor uPAR. *Hum Mol Genet* 2008;**17**:3617-30.
- 85 Song X, Tanaka H, Ohta K. Multiple roles of Equarin during lens development. *Dev Growth Differ* 2014;**56**:199-205.
- 86 O'Leary EE, Mazurkiewicz-Munoz AM, Argetsinger LS, Maures TJ, Huynh HT, Carter-Su C. Identification of steroid-sensitive gene-1/Ccdc80 as a JAK2-binding protein. *Mol Endocrinol* 2013;**27**:619-34.
- 87 Nasser YY, Krott E, Van Groningen KM, Berho M, Osborne MC, Wollman S, *et al.* Abnormalities in collagen composition may contribute to the pathogenesis of hemorrhoids: morphometric analysis. *Tech Coloproctol* 2015;**19**:83-7.
- 88 Mazzone A, Gibbons SJ, Bernard CE, Nowsheen S, Middha S, Almada LL, *et al.* Identification and characterization of a novel promoter for the human ANO1 gene regulated by the transcription factor signal transducer and activator of transcription 6 (STAT6). *FASEB J* 2015;**29**:152-63.
- 89 Strege PR, Bernard CE, Mazzone A, Linden DR, Beyder A, Gibbons SJ, *et al.* A novel exon in the human Ca<sup>2+</sup>-activated Cl<sup>-</sup> channel Anol1 imparts greater sensitivity to intracellular Ca<sup>2+</sup>. *Am J Physiol Gastrointest Liver Physiol* 2015;**309**:G743-9.
- 90 Ferrera L, Caputo A, Ubby I, Bussani E, Zegarra-Moran O, Ravazzolo R, *et al.* Regulation of TMEM16A chloride channel properties by alternative splicing. *J Biol Chem* 2009;**284**:33360-8.
- 91 Strege PR, Gibbons SJ, Mazzone A, Bernard CE, Beyder A, Farrugia G. EAVK segment "c" sequence confers Ca<sup>(2+)</sup>-dependent changes to the kinetics of full-length human Anol1. *Am J Physiol Gastrointest Liver Physiol* 2017;**312**:G572-G9.

- 92 Margetis N. Pathophysiology of internal hemorrhoids. *Ann Gastroenterol* 2019;**32**:264-72.
- 93 Devlin B, Roeder K. Genomic control for association studies. *Biometrics* 1999;**55**:997-1004.
- 94 Finucane HK, Bulik-Sullivan B, Gusev A, Trynka G, Reshef Y, Loh PR, *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat Genet* 2015;**47**:1228-35.
- 95 Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ. Jalview Version 2--a multiple sequence alignment editor and analysis workbench. *Bioinformatics* 2009;**25**:1189-91.

## Visual Abstract

