

1 **Research Letter: Therapeutic targets for haemorrhoidal disease: proteome-wide**

2 **Mendelian randomisation and colocalization analyses**

3 Shifang Li^{#*}, Meijiao Gong[#]

4 Laboratory of Immunology and Vaccinology, FARAH, ULiège, Liège 4000, Belgium.

5 [#]Shifang Li and Meijiao Gong contributed equally to this work

6 *Correspondence:

7 Shifang Li, fruceslee@gmail.com

8 Laboratory of Immunology and Vaccinology, FARAH, ULiège, Liège 4000, Belgium

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23 **Abstract**

24 Human haemorrhoidal disease (HEM) is a common anorectal pathology.
25 However, the etiology of HEM, as well as its molecular mechanism, remains largely
26 unclear. In this study, we applied a two-sample bi-direction Mendelian randomisation
27 (MR) analysis to estimate the causal effects of 4907 plasma proteins on HEM
28 outcomes and investigated the mediating impacts of plasma proteins on HEM risk
29 factors to uncover potential HEM treatment targets by integrating GWASs statistics of
30 HEM and plasma protein levels. Following MR analysis, our study identified 5
31 probable causal proteins associated with HEM. ERLEC1 and ASPN levels were
32 genetically predicted to be positively and inversely associated with HEM risk,
33 respectively, with strong evidence of colocalization ($H4 > 0.9$). The findings of an
34 independent cohort corroborate the causal relationship between these two proteins and
35 HEM. Furthermore, gene expression analysis of haemorrhoidal tissue and normal
36 specimens revealed that ERLEC1 but not ASPN were differentially expressed. By
37 analyzing single-cell ERLEC1 expression in human rectum tissues, ERLEC1 was
38 found to be highly expressed in transient-amplifying cells. Interestingly, a genetically
39 greater risk of myxoedema was linked to an elevated risk of HEM. However, there
40 was no evidence that dorsalgia, hernia, diverticular disease, and ankylosing
41 spondylitis were causally associated with HEM. Furthermore, no association was
42 found between myxoedema and the genetically predicted ERLEC1 and ASPN levels.
43 Overall, this study identified some causal associations of circulating proteins and risk
44 factors with HEM by integrating the largest-to-date plasma proteome and GWASs of

45 HEM. The findings could provide further insight into understanding biological
46 mechanisms for HEM.

47 **Keywords**

48 Haemorrhoidal disease, Mendelian randomisation, ERLEC1, myxoedema

49

50 Human haemorrhoidal disease (HEM) is a common anorectal disorder. Recently,
51 Zhang *et al.* reported the first and largest genome-wide association study (GWAS)
52 with haemorrhoidal disease (HEM), and these data offered us a resource for
53 understanding the genetic risk factors for HEM.¹ However, the etiology of HEM, as
54 well as its molecular mechanism, remains primarily unclear.² In addition, the
55 identification of genes with therapeutic effects needs to be conducted. In recent years,
56 by incorporating protein quantitative trait loci (pQTLs) into MR analysis, such an
57 approach has been successfully used to prioritize therapy targets.^{3,4} Here, using a
58 two-sample bidirectional Mendelian randomisation (MR) analysis, we estimated the
59 causal effects of 4907 plasma proteins on HEM outcomes, and investigated the effects
60 of plasma proteins that may mediate the impact of risk factors on HEM in order to
61 identify potential therapeutic targets for HEM.

62 As stated in the **supplementary methods**, 4907 proteins (*cis*-pQTLs) were used
63 as instrumental variables for exposure and HEM as the outcome to estimate the causal
64 effect of plasma protein levels on HEM in a proteome-wide context using MR
65 analysis.⁵⁻⁸ Our study revealed 5 potential causative proteins at the
66 Bonferroni-corrected threshold of $p < 1.01 \times 10^{-5}$, including 3 negative and 2 positive

67 associations (**figure 1A-1B**). MR analysis, for example, revealed that genetically
68 predicted ERLEC1 levels were linked to an increased risk of HEM ($p=5.18e-07$). To
69 determine whether the identified relationships of the circulating protein with HEM
70 shared causative variations, colocalization analysis was carried out and a high level of
71 support for colocalization evidence was discovered between two proteins (ERLEC1
72 and ASPN) and HEM ($H4>0.9$) (**figure 1C**). The findings of the INTERVAL cohort
73 corroborated the causal relationship between these two proteins and HEM (**figure**
74 **1D**).⁹ Interestingly, the deCODE study's lead cis-pQTL for the ERLEC1 (rs2542580)
75 but not ASPN (rs10992273) were not found to be associated with all available
76 secondary traits (**supplementary table1**). Gene expression analysis of haemorrhoidal
77 tissue and normal specimens revealed that ERLEC1 but not ASPN were differentially
78 expressed after controlling for gender and BMI (**figure 1E**), further supporting that a
79 high ERLEC1 expression level was associated with an increased risk of HEM.
80 Following that, we investigated the tissues in which ERLEC1 is expressed in bulk
81 tissues using GTEx v8 (<https://gtexportal.org/>), and found that ERLEC1 was
82 considerably expressed in multiple tissues, including the small intestine and colon, as
83 compared to the whole blood ($p<0.001$) (**figure 1F**). To further understand the origin
84 of ERLEC1, single-cell ERLEC1 expression was assessed in human rectum tissues,
85 and ERLEC1 was found to be highly expressed in transient-amplifying (TA) cells ($p<$
86 0.05) (**figure 1G**).¹⁰

87 In order to investigate whether the causal protein mediates the effect of risk
88 factors on HEM, the causal risk factors for HEM were first identified. 5 clinical traits

89 that genetically correlated with HEM were selected (**supplementary methods**), with
90 instrumental variables generated from GWASs confined to European populations. It
91 was discovered that a genetically greater risk of myxoedema was linked to an elevated
92 risk of HEM ($p < 0.05$) (**figure 1H**). Although genetic correlations with HEM were
93 reported,¹ there was no evidence that dorsalgia, hernia, diverticular disease, and
94 ankylosing spondylitis were causally associated ($p > 0.05$). In order to identify the
95 protein related to HEM risk factors, we conducted MR analysis again on 2 plasma
96 proteins impacting HEM with myxoedema. After filtering, there was a lack of
97 evidence that myxoedema had a causal relationship with these two plasma proteins
98 (**figure 1I**).

99 Overall, by integrating the largest-to-date plasma proteome and GWAS of HEM,
100 we discovered that ERLEC1 could serve as prospective protein therapeutic targets for
101 HEM. In-depth research is needed to investigate the mechanisms by which putative
102 risk factors affect HEM (**figure 1J**).

103 **Competing interests**

104 None declared.

105 **Contributors**

106 SF was involved in conceptualization. SF and MJ were involved in the formal
107 analysis. SF was involved in writing, reviewing, and editing.

108 **Acknowledgments**

109 The authors would like to thank all of the researchers who contributed to the GWAS
110 datasets used in this study for making them available for research purposes.

111 **References**

- 112 1 Zheng T, Ellinghaus D, Juzenas S, *et al.* Genome-wide analysis of 944 133
113 individuals provides insights into the etiology of haemorrhoidal disease. *Gut*
114 2021;70:1538-49.
- 115 2 EAM Festen & RK Weersma. Large-scale genetic analyses in an understudied
116 disease: haemorrhoidal disease. *Gut* 2021;70:1429-1430.
- 117 3 Bovijn J, Lindgren CM & Holmes MV. Genetic variants mimicking therapeutic
118 inhibition of IL-6 receptor signaling and risk of COVID-19. *The Lancet*
119 *Rheumatology* 2020;2:e658-9.
- 120 4 Dewey, F. E. *et al.* Genetic and Pharmacologic Inactivation of ANGPTL3 and
121 Cardiovascular Disease. *N Engl J Med* 2017;377:211-21.
- 122 5 Ferkingstad E, Sulem P, Atlason BA, *et al.* Large-scale integration of the plasma
123 proteome with genetics and disease. *Nat Genet* 2021;53:1712-21.
- 124 6 Zheng J, Haberland V, Baird D, *et al.* Phenome-wide Mendelian randomisation
125 mapping the influence of the plasma proteome on complex diseases. *Nat Genet*
126 2020;52:1122-31.
- 127 7 Chen L, Peters JE, Prins B, *et al.* Systematic Mendelian randomisation using the
128 human plasma proteome to discover potential therapeutic targets for stroke. *Nat*
129 *Commun* 2022;13:1-14.
- 130 8 Yoshiji S, Butler-Laporte G, Lu T, *et al.* Proteome-wide Mendelian randomisation
131 implicates nephronectin as an actionable mediator of the effect of obesity on
132 COVID-19 severity. *Nat Metab* 2023;5:248-64.

133 9 Sun BB, Maranville JC, Peters JE, Stacey D, Staley JR, Blackshaw J, Burgess S,
134 Jiang T, Paige E, Surendran P, *et al.* Genomic atlas of the human plasma proteome.
135 *Nature* 2018;558, 73-79.

136 10 Wang Y, Song W, Wang J, *et al.* Single-cell transcriptome analysis reveals
137 differential nutrient absorption functions in human intestine. *J Exp Med*
138 2020;217(2):e20191130.

139

140 **Figure Legends**

141 **Figure 1 Mendelian randomisation results.** (A) The effect of plasma protein levels
142 on HEM. Volcano plot indicating the effect of plasma protein on HEM using MR
143 analysis. (B) Forest plots shows the effect of plasma ERLEC1 and ASPN levels on
144 HEM. (C) Colocalization analysis of ERLEC1 levels (Up) and ASPN (Down). (D)
145 Forest plots shows the effect of plasma ERLEC1 and ASPN levels on HEM using
146 INTERVAL cohort. (E) Boxplot shows differentially expressed genes in HEM
147 patients when compared to healthy individuals. *p*-values were corrected the effect of
148 gender and BMI using linear model. (F) The violin plot depicts ERLEC1 gene
149 expression across multiple bulk tissues. (G) Data visualization of cell populations in
150 human rectum tissues using UMAP (left) and gene expression of ERLEC1 in different
151 cell types (right). (H) Forest plots showing the causal effect of chosen risk factors on
152 HEM. (I) Forest plots for the effect of myxoedema on plasma ERLEC1 and ASPN
153 levels. (J) Schematic illustration of the proposed model in the study. HEM,
154 haemorrhoidal disease.

155 **Supplementary Methods** The statistics method used in the study.

156 **Supplementary Tables1** The significant MR summary statistics obtained in this
157 study.

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177 **Supplementary Methods**

178 **GWASs of haemorrhoidal disease and risk factors**

179 We used recently published large-scale genome-wide associations (GWASs) for
180 haemorrhoidal disease (HEM).¹ This GWAS summary statistics were derived from
181 944,133 European ancestry individuals (Ncase = 218,920 and Ncontrol = 725,213)
182 from 5 cohorts and downloaded from the GWAS Catalog
183 (<https://www.ebi.ac.uk/gwas/>, access ID: EFO_0009552). Diverticular disease of the
184 intestine, ankylosing spondylitis (AS), dorsalgia, hernia, and myxoedema were
185 evaluated as potential causal risk factors associated with HEM in order to determine
186 the probable causal risk factors. All GWASs for the five risk factors were obtained
187 from the ieu open gwas project (<https://gwas.mrcieu.ac.uk/datasets/>). The summary
188 statistics of the large GWAS (14,357 cases and 182,423 controls) were used for
189 diverticular disease of the intestine (access ID: finn-b-K11_DIVERTIC). The GWAS
190 for ankylosing spondylitis (access ID: finn-b-M13_ANKYLOSPON) have a sample
191 size of 1,462 cases and 164,682 controls. The GWAS for myxoedema (access ID:
192 ieu-b-4877) has a sample size of 311,629 cases and 321,173 controls. The GWAS for
193 dorsalgia (access ID: finn-b-M13_DORSALGIA) included 193467 individuals, with
194 28,785 cases and 164,682 controls. A total of 218792 individuals were reported with
195 GWAS of hernia (access ID: finn-b-K11_HERNIA), including 28,235 cases and
196 190,557 controls.

197 **Plasma protein quantitative trait loci (pQTL) data**

198 To conduct proteome-wide Mendelian randomisation (MR), we first obtained

199 genetic instrumental variables using the protein quantitative trait loci (pQTL) data
200 generated by Ferkingstad *et al.*² The largest-to-date pQTL analysis on plasma
201 proteome (a total of 4907 proteins) in 35,559 Icelanders was performed in their study,
202 and an amount of 18,084 pQTL associations between genetic variation and protein
203 levels in plasma were identified. A total of 4907 pQTLs were successfully
204 downloaded from the deCODE study using aria2c.³ To minimize the risk of horizontal
205 pleiotropy, instrumental variables to *cis*-pQTLs (SNPs located within a 500 kb
206 window from the target gene body) of protein were selected for the following analysis.
207 In order to validate the MR results for ASPN and ERLEC1 using independent study,
208 the *cis*-pQTLs of ASPN and ERLEC1 were obtained from the INTERVAL cohort and
209 used for the following analysis.⁴

210 **Mendelian randomisation analysis**

211 MR analysis is an analytical method that uses genetic variation as an
212 instrumental variable (IV) to estimate causal effects. It overcomes the limitations of
213 measurement error and confounding factors that are common in observational studies
214 and is widely used to assess causal relationships.⁵ In this study, the TwoSampleMR
215 package (v0.5.6, <https://mrcieu.github.io/TwoSampleMR/>) was used for MR
216 analysis.⁶ The instrumental variables that determined the exposure in each MR study
217 were specified as genome-wide significant ($p \leq 5e-08$) SNPs. SNPs in the human
218 major histocompatibility complex (MHC) region at chromosome 6:
219 28,477,797-33,448,354 (GRCh37) were excluded from the analysis due to its complex
220 linkage disequilibrium (LD) structure. Using the 1000 Genomes Project European

221 reference panel and an LD threshold of $r^2 < 0.001$ with a clumping window of 10,000
222 kb, PLINK v.1.9 (<http://pngu.mgh.harvard.edu/purcell/plink/>) was employed to derive
223 instrumental variables.⁷⁻⁸ F-statistics were used to determine the strength of each
224 SNP's association with exposure, and F-statistics of more than 10 were considered
225 strong. For the main MR analysis, the inverse variance weighted approach for proteins
226 with two or more instrumental variables and the wald ratio method for proteins with a
227 single instrumental variable was used for evaluating the causal influence of exposure
228 on outcome. In addition, in the case of more instrumental variables used in MR
229 analysis, four additional MR methods (weighted median, simple mode, weighted
230 mode, and MR-Egger method) were used to assess the reliability of the primary
231 results. For exposures with multiple IVs, we additionally investigated heterogeneity
232 across variant-level MR estimations with the "mr_heterogeneity()" function in the
233 TwoSampleMR package (Cochrane's Q test). In addition, a pleiotropy test was
234 performed using MR Egger analysis to determine whether there is horizontal
235 pleiotropy among IVs. Meanwhile, "phenoscanner"
236 (<https://github.com/phenoscanner/phenoscanner>) was to be utilized to determine any
237 pleiotropy of SNPs used in the MR analysis. An SNP was regarded to be pleiotropic
238 when the reported SNP-traits association was genome-wide significant ($p \leq 5e-08$) in
239 the European population.

240 Finally, in the event there were more than two IVs in exposure, a leave-one-out
241 analysis was performed, and the MR findings of the remaining IVs were calculated by
242 deleting the IVs one by one to ensure the robustness of the MR data. To acquire robust

243 evidence for the casual estimation, MR findings that meet all of the following criteria
244 were chosen as described by Yoshiji and others: (1) no pleiotropy was found using
245 MR-Egger regression ($p>0.05$); (2) results with an $I^2 < 50\%$ (no substantial
246 heterogeneity); (3) leave-one-out analysis MR $p<0.05$ after removing outliers; and (4)
247 reverse MR $p>0.05$.⁹ The same procedure as mentioned above was utilized to explore
248 the causal effect of the given exposure and associated outcome in the reverse MR
249 analysis. p -values less than a Bonferroni adjusting ($p=1.01\times 10^{-5}$ (0.05/4,907)) are
250 deemed significant for multiple testing.

251 **Colocalization analysis**

252 The coloc R package was employed to investigate whether the reported
253 relationships between proteins and HEM were driven by linkage disequilibrium.¹⁰ The
254 analysis offers posterior probability for each hypothesis tested: no association in
255 either group (H0), pQTL only (H1), the GWAS of HEM only (H2), associations with
256 both GWAS but by separate causal signals (H3), and associations with both GWAS
257 but by the same signals (H4).¹¹ A higher H4 ($H4>0.8$) was considered as strong
258 evidence for colocalization, implying a shared variation between the two
259 phenotypes.^{10,11}

260 **Differentially expressed genes analysis in bulk tissues**

261 The GSE154650 dataset was downloaded from NCBI Gene Expression Omnibus
262 (GEO) and analyzed using the R program.¹² The RPM value of ERLEC1 and ASPN
263 were further subjected to linear model analysis to investigate the differential gene
264 expression in HEM and healthy individuals after correcting for the effects of gender

265 and BMI. The expression data of ERLEC1 from 39 tissues across 838 individuals
266 were obtained from the GTEx v8 (<https://gtexportal.org/>).¹³ Mann-Whitney U test was
267 performed to determine the significance of ERLEC1 expression differences between
268 the two groups, and $p < 0.01$ was declared significant.

269 **scRNA-sequencing analysis of human rectum tissues**

270 For processing scRNA data (GSE125970), the raw data of the gene expression
271 matrix was first downloaded from NCBI Gene Expression Omnibus (GEO) and
272 converted into a Seurat object using the R Seurat package.^{14,15} Low-quality cells were
273 eliminated if they met any of the following requirements: (1) 3000 UMIs; (2) 200
274 genes; and (3) >50% of UMIs derived from the mitochondrial genome. UMI counts
275 were normalized using the NormalizeData function, and the top 2000 features with
276 the greatest cell-to-cell variation were calculated using the FindVariableFeatures
277 function. To correct the batch effects among samples, the "FindIntegrationAnchors"
278 and "IntegrateData" functions were employed. Following that, the ScaleData function
279 was used to scale and center features in the datasets, and the RunPCA function with
280 default parameters was used to reduce dimensionality. The data were then used for
281 nonlinear dimensional reduction with the RunUMAP function and cluster analysis
282 with the FindNeighbors and FindClusters functions. The FindAllMarkers function
283 was used to identify differentially expressed genes (DEG) for a given cluster. The
284 clusters were labeled in the same way that Wang *et al.* did in their study.¹⁵

285 **References**

286 1 Zheng T, Ellinghaus D, Juzenas S, *et al.* Genome-wide analysis of 944 133

287 individuals provides insights into the etiology of haemorrhoidal disease. *Gut*
288 2021;70:1538-49.

289 2 Ferkingstad E, Sulem P, Atlason BA, *et al.* Large-scale integration of the plasma
290 proteome with genetics and disease. *Nat Genet* 2021;53:1712-21.

291 3 Aria2c Multi-source Download Utility. Available: <http://aria2.sourceforge.net/>

292 4 Sun BB, Maranville JC, Peters JE, Stacey D, Staley JR, Blackshaw J, Burgess S,
293 Jiang T, Paige E, Surendran P, *et al.* Genomic atlas of the human plasma proteome.
294 *Nature* 2018;558, 73-79.

295 5 Skrivankova VW, Richmond RC, Woolf BAR, *et al.* Strengthening the reporting of
296 observational studies in epidemiology using mendelian randomisation
297 (STROBE-MR): Explanation and elaboration. *BMJ* 2021;375. doi:10.1136/bmj.n2233

298 6 Hemani G, Zheng J, Elsworth B, *et al.* The MR-base platform supports systematic
299 causal inference across the human phenome. *Elife* 2018;7:1-29.

300 7 Auton A, Abecasis GR, Altshuler DM, *et al.* A global reference for human genetic
301 variation. *Nature* 2015;526:68-74.

302 8 Purcell S, Neale B, Todd-Brown K, *et al.* PLINK: A tool set for whole-genome
303 association and population-based linkage analyses. *Am J Hum Genet* 2007;81:559-75.

304 9 Yoshiji S, Butler-Laporte G, Lu T, *et al.* Proteome-wide Mendelian randomisation
305 implicates nephronectin as an actionable mediator of the effect of obesity on
306 COVID-19 severity. *Nat Metab* 2023;5:248-64.

307 10 Giambartolomei C, Vukcevic D, Schadt EE, *et al.* Bayesian test for colocalisation
308 between pairs of genetic association studies using summary statistics. *PLoS Genet*

309 2014;10:e1004383.

310 11 Foley CN, Staley JR, Breen PG, *et al.* A fast and efficient colocalization algorithm
311 for identifying shared genetic risk factors across multiple traits. *Nat Commun*
312 2021;12.

313 12 Zheng T, Ellinghaus D, Juzenas S, *et al.* Genome-wide analysis of 944 133
314 individuals provides insights into the etiology of haemorrhoidal disease. *Gut*
315 2021;70:1538-49.

316 13 Carithers LJ, Moore HM. The Genotype-Tissue Expression (GTEx) Project.
317 *Biopreserv Biobank* 2015;13:307-8.

318 14 Hao Y, Hao S, Andersen-Nissen E & Mauck WM. Integrated analysis of
319 multimodal single-cell data. Preprint at bioRxiv
320 <https://doi.org/10.1101/2020.10.12.335331>.

321 15 Wang Y, Song W, Wang J, *et al.* Single-cell transcriptome analysis reveals
322 differential nutrient absorption functions in human intestine. *J Exp Med*
323 2020;217(2):e20191130.

