

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

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

23 **Abstract**

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

45 biological mechanisms for HEM.

46 **Keywords**

47 Haemorrhoidal disease, Mendelian randomisation, ERLEC1, ASPN, myxoedema

48

49 Human haemorrhoidal disease (HEM) is a common anorectal disorder. Recently,
50 Zhang *et al.* reported the first and largest genome-wide association study (GWAS)
51 with haemorrhoidal disease (HEM), and these data offered us a resource for
52 understanding the genetic risk factors for HEM.¹ However, being one of the diseases
53 that affect a wide range of people, the etiology of HEM, as well as its molecular
54 mechanism, remains primarily unclear.² In addition, the identification of genes with
55 therapeutic effects needs to be conducted. Here, using a two-sample bidirectional
56 Mendelian randomisation (MR) analysis, we estimated the causal effects of 4907
57 plasma proteins on HEM outcomes, and investigated the effects of plasma proteins
58 that may mediate the impact of risk factors on HEM in order to identify potential
59 therapeutic targets for HEM.

60 In recent years, by incorporating protein quantitative trait loci (pQTLs) into MR
61 analysis, such an approach has been successfully used to prioritize therapy targets.³⁻⁵
62 As stated in the **Supplementary Methods**, 4907 proteins (*cis*-pQTLs) were used as
63 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.⁶⁻⁹ 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
71 of support for colocalization evidence was discovered between two proteins (ERLEC1
72 and ASPN) and HEM ($H4>0.9$) (**Figure 1C**). Furthermore, after controlling for
73 gender and BMI, gene expression analysis of haemorrhoidal tissue and normal
74 specimens revealed that ERLEC1 but not ASPN were differentially expressed (**Figure**
75 **1D**), further supporting that a high ERLEC1 expression level was associated with an
76 increased risk of HEM. Following that, we investigated the tissues in which ERLEC1
77 is expressed in bulk tissues using GTEx v8 (<https://gtexportal.org/>), and found that
78 ERLEC1 was considerably expressed in multiple tissues, including the small intestine
79 and colon, as compared to the whole blood ($p<0.001$) (**Figure 1E**). To further
80 understand the origin of ERLEC1, single-cell ERLEC1 expression was assessed in
81 human rectum tissues, and ERLEC1 was found to be highly expressed in
82 transient-amplifying (TA) cells ($p<0.05$) (**Figure 1F**).¹⁰

83 In order to investigate whether the causal protein mediates the effect of risk
84 factors on HEM, the causal risk factors for HEM were first identified. 5 clinical traits
85 that genetically correlated with HEM were selected (**Supplementary Methods**), with
86 instrumental variables generated from GWASs confined to European populations. It
87 was discovered that a genetically greater risk of myxoedema was linked to an elevated
88 risk of HEM ($p<0.05$) (**Figure 1G**). Although genetic correlations with HEM were

89 reported,¹ there was no evidence that dorsalgia, hernia, diverticular disease, and
90 ankylosing spondylitis were causally associated ($p>0.05$). In order to identify the
91 protein related to HEM risk factors, we conducted MR analysis again on 2 plasma
92 proteins impacting HEM with myxoedema. After filtering, there was a lack of
93 evidence that myxoedema had a causal relationship with these two plasma proteins
94 **(Figure 1H)**.

95 Overall, this study identified some causal associations of circulating proteins and
96 risk factors with HEM by integrating the largest-to-date plasma proteome and GWAS
97 of HEM. ERLEC1 in particular was discovered to be connected with an elevated risk
98 of HEM. In-depth research is needed to investigate the mechanisms by which putative
99 risk factors affect HEM **(Figure 1I)**. Overall, our study could provide further insight
100 into developing potential targets for HEM.

101 **Competing interests**

102 None declared.

103 **Contributors**

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

106 **Acknowledgments**

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

109 **References**

110 1 Zheng T, Ellinghaus D, Juzenas S, *et al*. Genome-wide analysis of 944 133

- 111 individuals provides insights into the etiology of haemorrhoidal disease. *Gut*
112 2021;70:1538-49.
- 113 2 EAM Festen & RK Weersma. Large-scale genetic analyses in an understudied
114 disease: haemorrhoidal disease. *Gut* 2021;70:1429-1430.
- 115 3 Reis G, Moreira Silva EAS, Medeiros Silva DC, *et al.* Early Treatment with
116 Pegylated Interferon Lambda for Covid-19. *N Engl J Med* 2023;388:518-28.
- 117 4 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 5 Dewey, F. E. *et al.* Genetic and Pharmacologic Inactivation of ANGPTL3 and
121 Cardiovascular Disease. *N Engl J Med* 2017;377:211-21.
- 122 6 Zheng J, Haberland V, Baird D, *et al.* Phenome-wide Mendelian randomisation
123 mapping the influence of the plasma proteome on complex diseases. *Nat Genet*
124 2020;52:1122-31.
- 125 7 Chen L, Peters JE, Prins B, *et al.* Systematic Mendelian randomisation using the
126 human plasma proteome to discover potential therapeutic targets for stroke. *Nat*
127 *Commun* 2022;13:1-14.
- 128 8 Yoshiji S, Butler-Laporte G, Lu T, *et al.* Proteome-wide Mendelian randomisation
129 implicates nephronectin as an actionable mediator of the effect of obesity on
130 COVID-19 severity. *Nat Metab* 2023;5:248-64.
- 131 9 Chen J, Xu F, Ruan X, Sun J, *et al.* Therapeutic targets for inflammatory bowel
132 disease: proteome-wide Mendelian randomisation and colocalization analyses.

133 *EBioMedicine* 2023;89:104494.

134 10 Wang Y, Song W, Wang J, *et al.* Single-cell transcriptome analysis reveals
135 differential nutrient absorption functions in human intestine. *J Exp Med*
136 2020;217(2):e20191130.

137

138 **Figure Legends**

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

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

153 **Supplementary Tables1** The significant MR summary statistics obtained in this
154 study.

155

156

157

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 AS (access ID: finn-b-M13_ANKYLOSPON) have a sample size of 1,462 cases
191 and 164,682 controls. The GWAS for myxoedema (access ID: ieu-b-4877) has a
192 sample size of 311,629 cases and 321,173 controls. The GWAS for dorsalgia (access
193 ID: finn-b-M13_DORSALGIA) included 193467 individuals, with 28,785 cases and
194 164,682 controls. A total of 218792 individuals were reported with GWAS of hernia
195 (access ID: finn-b-K11_HERNIA), including 28,235 cases and 190,557 controls.

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

197 To conduct proteome-wide Mendelian randomisation (MR), we first obtained
198 genetic instrumental variables using the protein quantitative trait loci (pQTL) data

199 generated by Ferkingstad *et al.*² The largest-to-date pQTL analysis on plasma
200 proteome (a total of 4907 proteins) in 35,559 Icelanders was performed in their study,
201 and an amount of 18,084 pQTL associations between genetic variation and protein
202 levels in plasma were identified. A total of 4907 pQTLs were successfully
203 downloaded from the deCODE study using aria2c.³ To minimize the risk of horizontal
204 pleiotropy, instrumental variables to *cis*-pQTLs (SNPs located within a 1,000 kb
205 window from the target gene body) of protein were selected for the following
206 analysis.

207 **Mendelian randomisation analysis**

208 MR analysis is an analytical method that uses genetic variation as an
209 instrumental variable (IV) to estimate causal effects. It overcomes the limitations of
210 measurement error and confounding factors that are common in observational studies
211 and is widely used to assess causal relationships.⁴ In this study, the TwoSampleMR
212 package (v0.5.6, <https://mrcieu.github.io/TwoSampleMR/>) was used for MR
213 analysis.⁵ The instrumental variables that determined the exposure in each MR study
214 were specified as genome-wide significant ($p \leq 5e-08$) SNPs. SNPs in the human
215 major histocompatibility complex (MHC) region at chromosome 6:
216 28,477,797-33,448,354 (GRCh37) were excluded from the analysis due to its complex
217 linkage disequilibrium (LD) structure. Using the 1000 Genomes Project European
218 reference panel and an LD threshold of $r^2 < 0.001$ with a clumping window of 10,000
219 kb, PLINK v.1.9 (<http://pngu.mgh.harvard.edu/purcell/plink/>) was employed to derive
220 instrumental variables.⁶⁻⁷ F-statistics were used to determine the strength of each

221 SNP's association with exposure, and F-statistics of more than 10 were considered
222 strong. For the main MR analysis, the inverse variance weighted approach for proteins
223 with two or more instrumental variables and the wald ratio method for proteins with a
224 single instrumental variable was used for evaluating the causal influence of exposure
225 on outcome. In addition, in the case of more instrumental variables used in MR
226 analysis, four additional MR methods (weighted median, simple mode, weighted
227 mode, and MR-Egger method) were used to assess the reliability of the primary
228 results. For exposures with multiple IVs, we additionally investigated heterogeneity
229 across variant-level MR estimations with the "mr_heterogeneity()" function in the
230 TwoSampleMR package (Cochrane's Q test). In addition, a pleiotropy test was
231 performed using MR Egger analysis to determine whether there is horizontal
232 pleiotropy among IVs.

233 Finally, in the event there were more than two IVs in exposure, a leave-one-out
234 analysis was performed, and the MR findings of the remaining IVs were calculated by
235 deleting the IVs one by one to ensure the robustness of the MR data. To acquire robust
236 evidence for the casual estimation, MR findings that meet all of the following criteria
237 were chosen as described by Yoshiji and others: (1) no pleiotropy was found using
238 MR-Egger regression ($p > 0.05$); (2) results with an $I^2 < 50\%$ (no substantial
239 heterogeneity); (3) leave-one-out analysis MR $p < 0.05$ after removing outliers; and (4)
240 reverse MR $p > 0.05$.⁸ The same procedure as mentioned above was utilized to explore
241 the causal effect of the given exposure and associated outcome in the reverse MR
242 analysis. p -values less than a Bonferroni adjusting ($p = 1.01 \times 10^{-5}$ (0.05/4,907)) are

243 deemed significant for multiple testing.

244 **Colocalization analysis**

245 The coloc R package was employed to investigate whether the reported
246 relationships between proteins and HEM were driven by linkage disequilibrium.¹¹ The
247 analysis offers posterior probability for each hypothesis tested: no association in
248 either group (PP0), one GWAS only (PP1), the other GWAS only (PP2), associations
249 with both GWAS but by separate causal signals (PP3), and associations with both
250 GWAS but by the same signals (PP4).¹² A higher PP4 ($PP4 > 0.8$) was considered as
251 strong evidence for colocalization, implying a shared variation between the two
252 phenotypes.^{11,12}

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

254 The GSE154650 dataset was downloaded from NCBI Gene Expression Omnibus
255 (GEO) and analyzed using the R program.¹³ The RPM value of ERLEC1 and ASPN
256 were further subjected to linear model analysis to investigate the differential gene
257 expression in HEM and healthy individuals after correcting for the effects of gender
258 and BMI. The expression data of ERLEC1 from 39 tissues across 838 individuals
259 were obtained from the GTEx v8 (<https://gtexportal.org/>).¹⁴ Mann-Whitney U test was
260 performed to determine the significance of ERLEC1 expression differences between
261 the two groups, and $p < 0.01$ was declared significant.

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

263 For processing scRNA data (GSE125970), the raw data of the gene expression
264 matrix was first downloaded from NCBI Gene Expression Omnibus (GEO) and

265 converted into a Seurat object using the R Seurat package.^{15,16} Low-quality cells were
266 eliminated if they met any of the following requirements: (1) 3000 UMIs; (2) 200
267 genes; and (3) >50% of UMIs derived from the mitochondrial genome. UMI counts
268 were normalized using the NormalizeData function, and the top 2000 features with
269 the greatest cell-to-cell variation were calculated using the FindVariableFeatures
270 function. To correct the batch effects among samples, the "FindIntegrationAnchors"
271 and "IntegrateData" functions were employed. Following that, the ScaleData function
272 was used to scale and center features in the datasets, and the RunPCA function with
273 default parameters was used to reduce dimensionality. The data were then used for
274 nonlinear dimensional reduction with the RunUMAP function and cluster analysis
275 with the FindNeighbors and FindClusters functions. The FindAllMarkers function
276 was used to identify differentially expressed genes (DEG) for a given cluster. The
277 clusters were labeled in the same way that Wang *et al.* did in their study.¹⁶

278 **References**

- 279 1 Zheng T, Ellinghaus D, Juzenas S, *et al.* Genome-wide analysis of 944 133
280 individuals provides insights into the etiology of haemorrhoidal disease. *Gut*
281 2021;70:1538-49.
- 282 2 Ferkingstad E, Sulem P, Atlason BA, *et al.* Large-scale integration of the plasma
283 proteome with genetics and disease. *Nat Genet* 2021;53:1712-21.
- 284 3 Aria2c Multi-source Download Utilily. Available: <http://aria2.sourceforge.net/>
- 285 4 Skrivankova VW, Richmond RC, Woolf BAR, *et al.* Strengthening the reporting
286 of observational studies in epidemiology using mendelian randomisation

- 287 (STROBE-MR): Explanation and elaboration. *BMJ* 2021;375. doi:10.1136/bmj.n2233
- 288 5 Hemani G, Zheng J, Elsworth B, et al. The MR-base platform supports systematic
289 causal inference across the human phenome. *Elife* 2018;7:1-29.
- 290 6 Auton A, Abecasis GR, Altshuler DM, et al. A global reference for human genetic
291 variation. *Nature* 2015;526:68-74.
- 292 7 Purcell S, Neale B, Todd-Brown K, et al. PLINK: A tool set for whole-genome
293 association and population-based linkage analyses. *Am J Hum Genet* 2007;81:559-75.
- 294 8 Yoshiji S, Butler-Laporte G, Lu T, et al. Proteome-wide Mendelian randomisation
295 implicates nephronectin as an actionable mediator of the effect of obesity on
296 COVID-19 severity. *Nat Metab* 2023;5:248-64.
- 297 9 Burgess S, Daniel RM, Butterworth AS, et al. Network Mendelian randomisation:
298 Using genetic variants as instrumental variables to investigate mediation in causal
299 pathways. *Int J Epidemiol* 2015;44:484-95.
- 300 10 Carter AR, Gill D, Davies NM, et al. Understanding the consequences of
301 education inequality on cardiovascular disease: Mendelian randomisation study. *BMJ*
302 2019;365:1-12.
- 303 11. Giambartolomei C, Vukcevic D, Schadt EE, et al. Bayesian test for colocalisation
304 between pairs of genetic association studies using summary statistics. *PLoS Genet*
305 2014;10:e1004383.
- 306 12. Foley CN, Staley JR, Breen PG, et al. A fast and efficient colocalization
307 algorithm for identifying shared genetic risk factors across multiple traits. *Nat*
308 *Commun* 2021;12.

- 309 13. Zheng T, Ellinghaus D, Juzenas S, et al. Genome-wide analysis of 944 133
310 individuals provides insights into the etiology of haemorrhoidal disease. *Gut*
311 2021;70:1538-49.
- 312 14. Carithers LJ, Moore HM. The Genotype-Tissue Expression (GTEx) Project.
313 *Biopreserv Biobank* 2015;13:307-8.
- 314 15. Hao Y, Hao S, Andersen-Nissen E & Mauck WM. Integrated analysis of
315 multimodal single-cell data. Preprint at bioRxiv
316 <https://doi.org/10.1101/2020.10.12.335331>.
- 317 16. Wang Y, Song W, Wang J, et al. Single-cell transcriptome analysis reveals
318 differential nutrient absorption functions in human intestine. *J Exp Med*
319 2020;217(2):e20191130.

