

1 An artificial neural network approach integrating plasma proteomics and
2 genetic data identifies *PLXNA4* as a new susceptibility locus for
3 pulmonary embolism.

4

5

6 Misbah Razzaq^{1,2}, Maria Jesus Iglesias^{3,4}, Manal Ibrahim-Kosta^{5,6}, Louisa
7 Goumidi^{5,6}, Omar Soukarieh^{1,2}, Carole Proust^{1,2}, Maguelonne Roux², Pierre
8 Suchon^{5,6}, Anne Boland^{2,7}, Delphine Daiain^{2,7}, Robert Olasso^{2,7}, Lynn Butler^{3,4,8}, Jean-
9 François Deleuze^{2,7,9}, Jacob Odeberg^{3,4}, Pierre-Emmanuel Morange^{5,6*}, David-
10 Alexandre Trégouët^{1,2*}

11

12 ¹ Univ. Bordeaux, INSERM, BPH, U1219, F-33000 Bordeaux, France

13 ² Laboratory of Excellence GENMED (Medical Genomics)

14 ³ Science for Life Laboratory, Department of Protein Science, CBH, KTH Royal Institute of
15 Technology, Stockholm, Sweden.

16 ⁴ Department of Clinical Medicine, Faculty of Health Science, the Arctic University of
17 Tromsø, Norway

18 ⁵ Aix Marseille Univ, INSERM, INRAE, C2VN, Marseille, France

19 ⁶ Hematology Laboratory, La Timone University Hospital of Marseille, Marseille, France

20 ⁷ Université Paris-Saclay, CEA, Centre National de Recherche en Génomique Humaine,
21 91057, Evry, France

22 ⁸ Clinical Chemistry and Blood Coagulation Research, Department of Molecular Medicine
23 and Surgery, Karolinska Institute, SE-171 76 Stockholm, Sweden

24 ⁹ Centre d'Etude du Polymorphisme Humain, Fondation Jean Dausset, Paris, France

25

26 * These two authors equally contributed to the work.

27

28 Corresponding authors:

29 Trégouët David-Alexandre: david-alexandre.tregouet@u-bordeaux.fr; INSERM U1219,
30 VINTAGE, case 11, Bordeaux Population Health research center, 146 rue Léo Saignat,
31 33076 Bordeaux, France. Tel: +33 5 47 30 42 54

32 Razzaq Misbah: misbah.razzaq@inserm.fr; INSERM U219, VINTAGE, case 11, Bordeaux
33 Population Health research center, 146 rue Léo Saignat, 33076 Bordeaux, France

34

35

36

37 **Abstract**

38 Venous thromboembolism is the third common cardiovascular disease and is composed of
39 two entities, deep vein thrombosis (DVT) and its fatal form, pulmonary embolism (PE).
40 While PE is observed in ~40% of patients with documented DVT, there is limited biomarkers
41 that can help identifying patients at high PE risk.

42 To fill this need, we implemented a two hidden-layers artificial neural networks (ANN) on
43 376 antibodies and 19 biological traits measured in the plasma of 1388 DVT patients, with or
44 without PE, of the MARTHA study. We used the LIME algorithm to obtain a linear
45 approximate of the resulting ANN prediction model. As MARTHA patients were typed for
46 genotyping DNA arrays, a genome wide association study (GWAS) was conducted on the
47 LIME estimate. Detected single nucleotide polymorphisms (SNPs) were tested for association
48 with PE risk in MARTHA. Main findings were replicated in the EOVT study composed of
49 143 PE patients and 196 DVT only patients.

50 The derived ANN model for PE achieved an accuracy of 0.89 and 0.79 in our training and
51 testing sets, respectively. A GWAS on the LIME approximate identified a strong statistical
52 association peak ($p = 5.3 \times 10^{-7}$) at the *PLXNA4* locus, with lead SNP rs1424597 at which the
53 minor A allele was further shown to associate with an increased risk of PE (OR = 1.49 [1.12 –
54 1.98], $p = 6.1 \times 10^{-3}$). Further association analysis in EOVT revealed that, in the combined
55 MARTHA and EOVT samples, the rs1424597-A allele was associated with increased PE risk
56 (OR = 1.74 [1.27 – 2.38, $p = 5.42 \times 10^{-4}$]) in patients over 37 years of age but not in younger
57 patients (OR = 0.96 [0.65 – 1.41], $p = 0.848$).

58 Using an original integrated proteomics and genetics strategy, we identified *PLXNA4* as a new
59 susceptibility gene for PE whose exact role now needs to be further elucidated.

60

61 **Author Summary**

62 Pulmonary embolism is a severe and potentially fatal condition characterized by the presence
63 of a blood clot (or thrombus) in the pulmonary artery. Pulmonary embolism is often the
64 consequence of the migration of a thrombus from a deep vein to the lung. Together with deep
65 vein thrombosis, pulmonary embolism forms the so-called venous thromboembolism, the
66 third most common cardiovascular disease, and its prevalence strongly increases with age.
67 While pulmonary embolism is observed in ~40% of patients with deep vein thrombosis, there
68 is currently limited biomarkers that can help predicting which patients with deep vein
69 thrombosis are at risk of pulmonary embolism. We here deployed an Artificial Intelligence
70 based methodology integrating both plasma proteomics and genetics data to identify novel
71 biomarkers for PE. We thus identified the *PLXNA4* gene as a novel molecular player involved
72 in the pathophysiology of pulmonary embolism. In particular, using two independent cohorts
73 totalling 1,881 patients with venous thromboembolism among which 467 experienced
74 pulmonary embolism, we identified a genetic polymorphism in the *PLXNA4* gene that
75 associates with ~2 fold increased risk of pulmonary embolism in patients aged more than ~40
76 years.

77

78

79

80

81 **Introduction**

82 Deep vein thrombosis (DVT) and Pulmonary Embolism (PE) are often considered as two
83 sides of the same coin, venous thromboembolism (VTE), the third most common
84 cardiovascular disease. VTE is a complex disease resulting from the interplay of various
85 factors including (epi-)genetics and environmental sources. VTE incidence is estimated
86 at 1 per 1000 patient-years, and its fatal form, PE, is associated with a mortality rate of
87 6% in the acute phase and 20% after one year [1]. PE generally results from the
88 migration of a blood clot from a deep vein to the lung and is observed in ~40% of
89 patients with documented DVT [2]. However, isolated PE without any trace of DVT can
90 also be observed either when the clot has completely migrated to the lung or when it is a
91 pulmonary clot in situ as recently highlighted in COVID-19 patients [3]. Even though
92 some specific risk factors for PE have been identified in DVT patients such as obesity,
93 sickle cell disease [4] as well as some genetic variations in *F5* [4] and *GRK5* [5] genes, the
94 exact, likely multifactorial, biological mechanisms that lead to PE are still not fully
95 characterized. Besides, there is still limited biomarkers that can help discriminating
96 patients that will develop PE from those who won't, the former being then at higher risk
97 of death. Thus, there is clearly a need for novel PE-associated molecular markers to be
98 identified.

99 Plasma is an ideal potential source for VTE biomarkers; the intravascular
100 compartment itself is the site of disease manifestation and tests are relatively non-
101 invasive, quick and cheap. Several types of molecular determinants can be assessed in
102 plasma samples including microRNAs, metabolites and proteins, and all of them have
103 been investigated in the context of VTE. For example, plasma microRNAs have been
104 assessed in relation to VTE recurrence [6,7]. Plasma proteomics has been employed to
105 discover novel proteins associated with VTE risk [8,9] and plasma metabolomics used to
106 identify novel mechanisms involved in VTE etiology [10,11]. Only one study has so far
107 adopted an exploratory plasma proteomics strategy to identify novel proteins associated
108 with high-risk versus low-risk of PE in humans. This study [12] was based on a relatively
109 small sample size and compared 6 patients with high risk of PE to 6 patients at low PE
110 risk, risk being classified based on clinical presentations and symptoms, with plasma
111 samples profiled by matrix-assisted laser desorption/ionization–time-of-flight/time-of-
112 flight mass spectrometry (MALDI-TOF/TOF MS).

113 In this work, we aim at identifying novel molecular phenotypes that could help in
114 better characterizing the biological mechanisms involved in the development of PE in
115 VTE patients. For this, 234 plasma proteins targeted with 376 protein specific
116 antibodies, with the major part derived from the Human Protein Atlas (HPA) repository
117 [13] were profiled in 1388 VTE patients selected from the MARTHA study [14,15] and
118 from whom 283 had experienced a symptomatic PE event. To explore far beyond the
119 search for linear associations between protein levels and PE risk and to identify more
120 complex relationships that could serve as integrative markers of upstream/downstream
121 mechanisms involving molecular determinants that have not necessarily been
122 measured, we deployed a sequential procedure implementing several methodologies
123 selected from the deep-learning domain. Briefly, and as summarized in Figure 1 and
124 more detailed thereafter, the first step consists in applying an under-sampling
125 algorithm (edited nearest neighbors) [16] to remove individuals with strong data
126 heterogeneity that would hamper the efficiency of the downstream analyses, leaving to
127 subsample of 592 VTE patients (497 DVT and 95 PE). This subsample was then used in
128 an Artificial Neural Network (ANN) learning framework in order to predict PE from
129 proteomics data. We then used the Local Interpretable Model-agnostic Explanations
130 (LIME) algorithm [17] to derive a linear approximate of the ANN based predictor for PE
131 risk which would, in addition, have a more meaningful biological interpretation. As
132 MARTHA patients have been previously typed for genome-wide genotype data, we then
133 conducted a genome wide association study of the LIME predictor of PE in order to
134 detect single nucleotide polymorphisms (SNPs) associated with the predictor with the
135 hope that the integration of genetic and proteomic data could provide additional insights
136 into the pathophysiology underlying the identified predictor [18,19]. SNPs with strong
137 statistical association with the LIME predictor were tested for association with PE risk in
138 the whole original MARTHA dataset and significant associations were further tested for
139 replication in an independent study of 339 VTE patients including 143 with PE.
140 Sequencing data were also scrutinized in some patients with observed VTE outcomes
141 poorly predicted by our ANN/LIME prediction models in order to identify rare variants
142 that could be responsible for the observed phenotypes.

143

144 **Figure 1** Analysis workflow of the present study

145

146 Results

147 *Data description* - The MARTHA proteomics substudy was composed of 1,388 VTE
148 patients among which 1,105 were diagnosed for DVT, 95 with isolated PE and 188 with
149 both DVT and PE (Table 1). Patients were phenotyped for 19 quantitative traits known
150 to be involved in thrombotic biological processes (Supplementary Table 1) and for 234
151 different proteins using targeted affinity proteomics with 376 protein specific antibody
152 reagents using multiplexed suspension bead array technology. These proteins were
153 selected for 1) their known roles in the coagulation/fibrinolysis cascade and/or
154 intermediate traits of relevance to thrombosis, 2) their specific expression in endothelial
155 cells (a key cell type involved in thrombosis physiopathology) or 3) encoded by genes
156 identified in pangenomic studies as associated with several cardiovascular disease-
157 linked biological pathways (e.g platelet function, renal function, inflammation). The list
158 of antibody reagents with their target proteins is given in Supplementary Table 2.
159

160 **Table 1 Characteristics of the MARTHA proteomics study**

161

	DVT	PE	DVT+PE
N	1105	95	188
Age at sampling	46.67 (14.90)	48.63 (15.26)	51.57 (16.99)
Age at first VTE	40.89 (15.28)	41.64 (15.02)	44.22 (17.56)
Female sex	716 (65%)	78 (82%)	112 (60 %)
Women under oral contraceptives at VTE event	286 (26%)	35 (37%)	45 (24%)
FV Leiden (rs6025) heterozygotes	255 (23%)	17 (18%)	39 (21%)
Anticoagulant therapy at plasma sampling	303 (27%)	29 (31%)	76 (40%)
Smokers	209 (19 %)	18 (19 %)	24 (13 %)
BMI	25.14 (4.57)	25.20(4.39)	26.43(4.62)

162 DVT : Deep Vein Thrombosis ; PE : Pulmonary Embolism ; BMI : Body Mass Index
163 Data shown correspond to mean (standard deviation) and count (percentage) for
164 continuous categorical variables, respectively
165

166 Exploration of this dataset using high-dimensional visualization techniques including
167 principal component analysis, t-SNE [20] and UMAP [21] did not reveal any specific
168 stratification in the data nor outliers (Figure 2) but rather illustrates that the three class
169 of patients (DVT, PE, DVT+PE) could not be easily separated.

170

171 **Figure 2** Graphical representation of the HPAs and biological MARTHA data projected
172 on the first two principal components derived from standard principal components
173 analysis (a), t-SNE (b) and UMAP (c) techniques.
174

175

176 *Artificial Neural Network for PE* - As the accuracy/efficiency of any ANN strongly
177 depend on the quality/homogeneity of the input data, we first applied the edited nearest
178 neighbors algorithm [16] to perform under sampling of the majority class (DVT) and
179 obtain a more homogeneous set of DVT patients, and further discarded the DVT+PE
180 class to avoid adding noise in discriminating between PE and non PE patients. This
181 strategy led to the selection of a subsample (referred thereafter to as the ANN sample 2)
182 of 592 patients (497 DVT and 95 PE) whose proteomics/biological entered the ANN
183 analysis.

184 A two hidden-layers ANN was then built from the ANN dataset with a training set of
185 576 patients (487 DVT and 89 PE) and a testing set of 16 patients (10 DVT and 6 PE).
186 This allocation was chosen so that the number of PE cases used for training was
187 sufficiently large. Because the training set presented with a strong imbalance with
188 respect to the DVT/PE classes with ~5 times more DVT than PE patients, the ANN was
189 trained iteratively as described in the Materials and Methods section. By completion of
190 the iterative algorithm, the final ANN obtained an area under the operative curve (AUC)
191 of 0.89. Of more interest are the performances of the ANN in the testing set. Indeed, our
192 ANN got F1-scores of 0.82 and 0.60 for the DVT and PE classes, respectively, and a global
193 AUC of 0.79 in the testing set.

194 We then used the LIME algorithm to obtain a local linear approximate of the ANN
195 predictions. In the testing set, the LIME prediction achieved an overall AUC of 0.77
196 instead of 0.79 for ANN. For each of the 16 patients in the testing set, we compared the
197 individual predictions of their observed VTE event provided by the ANN and LIME
198 methods (Table 2). In general, ANN and LIME predictions were rather consistent even if
199 the ANN predictions seem to be more accurate in predicting DVT while LIME appears
200 slightly more accurate in predicting PE. The average prediction in correctly classifying
201 DVT patients was 0.872 by ANN compared to 0.748 by LIME. Note that one DVT patient
202 (individual 10) was wrongly predicted to be PE by the ANN predictor, but not by the
203 LIME predictor. Conversely, the average prediction in correctly classifying PE patients

204 was 0.498 by ANN compared to 0.578 by LIME. Two PE patients (individuals 11 & 12)
205 presented low predictions of being PE, using both ANN and LIME predictors.

206

207 **Table 2 Individual predictions of VT event provided by ANN and LIME in the 16**
208 **patients of the testing set**

209

Individual	Observed clinical class	ANN prediction for class PE	Local Prediction for class PE
1	DVT	0.04	0.31
2	DVT	0.00	0.18
3	DVT	0.03	0.24
4	DVT	0.02	0.17
5	DVT	0.00	0.23
6	DVT	0.02	0.32
7	DVT	0.00	0.25
8	DVT	0.04	0.22
9	DVT	0.25	0.34
10	DVT	0.88	0.26
11	PE	0.00	0.30
12	PE	0.20	0.31
13	PE	0.98	0.94
14	PE	1.0	1.0
15	PE	0.01	0.15
16	PE	0.80	0.77

210

211 We then assessed the correlation of the LIME predictor with the available biological
212 phenotypes. No strong correlation was observed (Supplementary Table 3). However, the
213 LIME predictor showed marginal positive correlation with fibrinogen ($\rho = 0.12$, $p = 5.7 \times$
214 10^{-3}) and factor VIII ($\rho = 0.16$, $p = 0.013$) plasma levels, and marginal negative
215 correlation with prothombin time ($\rho = -0.10$, $p = 0.029$) and protein S ($\rho = -0.10$, $p =$
216 0.021) plasma levels. To go further into the biological interpretation of the LIME
217 predictor, we sought to identify which proteins contribute the most to the definition of
218 the LIME predictor. Figure 3 display the top 20 most contributing antibodies/proteins.
219 Of note, 5 proteins tended to have substantial more importance than the remaining ones,
220 among which three include proteins that had been selected because their gene
221 expression (COX4I2, VCL, VWF) was found to be specifically enriched in endothelial cells
222 [22].

223

224 **Figure 3** List of the top 20 antibodies contributing the most to the prediction model for
225 PE

226

227 *Genetics of the LIME predictor* To get additional information about the biological
 228 mechanisms that could underly the linear LIME predictor, we conducted a GWAS on this
 229 predictor considered as a quantitative linear trait in a sample of 574 individuals of the
 230 ANN subsample with GWAS data. While no SNP reached genome-wide significance, we
 231 observed a peak of strong suggestive statistical association on chromosome 7 at the
 232 *PLXNA4* locus (Supplementary Figure1 – Supplementary Table 4). The sentinel SNP ($p =$
 233 5.33×10^{-7}) was rs1424597 whose minor A allele with frequency of 0.09 was associated
 234 with an increase of $+0.169 \pm 0.034$ in LIME predictor values. In this subsample, the
 235 rs1424597-A allele was slightly more frequent in patients with PE than in patients with
 236 DVT only (0.15 vs 0.08, $p = 4.1 \times 10^{-3}$) (Table 3). The association of rs1424597 with PE
 237 risk was then assessed in the remaining MARTHA samples (738 DVT and 230 PE
 238 (DVT+PE or isolated PE) patients) with available GWAS data and that had not been used
 239 for building our ANN model. In this subsample, we observed a trend for an higher
 240 frequency of the rs1424597-A allele in PE patients compared to non PE patients (0.10 vs
 241 0.08), even if the association did not reach significance ($p = 0.20$).

242

243 **Table 3 Association of rs1424597 with PE risk in the MARTHA and EOVT studies**

	MARTHA				EOVT	
	Participants included in the ANN analysis		Participants outside the ANN analysis			
	DVT N = 480	PE N = 94	DVT N=738	PE N=230	DVT N=196	PE N=143
GG	404 (84%)	71 (75%)	624 (85%)	187 (79%)	149 (76%)	110 (77%)
GA	74 (15%)	18 (19%)	111 (14%)	41 (20%)	47 (24%)	28 (20%)
AA	2 (<1%)	5 (5%)	3 (<1%)	2 (~1%)	0 (-)	5 (3%)
MAF ¹	0.081	0.149	0.079	0.098	0.119	0.133
OR ²	1.98 [1.24 - 3.14] p = 0.0041		1.26 [0.88 - 1.81] p = 0.204		1.12 [0.71 - 1.78] p = 0.613	

244

MAF : Minor Allele Frequency

245

OR : Allelic Odds Ratio [95%CI] and p-value of the Cochran-Armitage trend test for association

246

247

248

249

We further investigated the association of rs1424597 with PE in the EOVT study composed of 143 PE patients 196 DVT patient. In EOVT, the rs1424597-A allele frequency was similar between EOVT patients with PE and with DVT (0.13 vs 0.12, $p =$
 250 0.61) (Table 3). Because by design the EOVT study is enriched with early onset VTE
 251 patients, we assessed whether the discrepancy between MARTHA and EOVT results
 252

252

253 could be due to patient selection criteria (i.e according to age.) We thus split the EOVT
 254 samples according to the median of age of VTE onset, that was 37yrs. As shown in Table
 255 4, the pattern of association of rs1424597 with PE slightly differed according to age. In
 256 EOVT patients younger than 37yrs, its allele frequency tend to be slightly lower in PE
 257 than in DVT patients (0.11s 0.13) while the inverse was observed in patients aged more
 258 than 37 yrs (0.15 vs0.11). Interestingly, the same observations hold in MARTHA when
 259 patients were stratified according to the same age threshold (Table 4). In the combined
 260 MARTHA and EOVT samples, the rs1424597-A allele was associated with an increased
 261 odds ratio (OR) for PE of 1.74 [1.27 – 2.38] ($p = 5.42 \times 10^{-4}$) in patients over 37 years of
 262 age while no association (OR = 0.96 [0.65 – 1.41], $p = 0.848$) was observed in younger
 263 patients. Similar ORs were obtained, OR = 1.73 [1.22 – 2.45] ($p = 1.99 \times 10^{-3}$) and OR =
 264 1.08 [0.77 – 1.53] ($p = 0.628$), respectively, if a more standard age threshold of 40 yrs
 265 [23] had been used.

266

267 **Table 4 Association of rs1424597 with PE risk according to age of onset of venous**
 268 **thrombosis**

	Age of onset <37				Age of onset ≥37			
	MARTHA		EOVT		MARTHA		EOVT	
	DVT N=550	DVT+PE/PE N =147	DVT N=96	PE N=64	DVT N=668	DVT+PE/PE N =177	DVT N=100	PE N=79
GG	464	123	70	51	564	135	79	59
GA	83	23	26	12	102	36	21	16
AA	3	1	0	1	2	6	0	4
MAF ¹	0.081	0.085	0.135	0.109	0.079	0.136	0.105	0.152
OR ²	1.056 [0.664 – 1.678] p = 0.817		0.784 [0.392 – 1.566] p = 0.470		1.820 [1.266 – 2.617] p = 1.16 10 ⁻³		1.53 [0.82 – 2.86] p = 0.196	

269 MAF : Minor Allele Frequency

270 OR : Allelic Odds Ratio [95%CI] and p-value of the Cochran-Armitage trend test for
 271 association

272

273

274 *Genetics of inconsistent LIME predictions*

275 As shown in Table 2, our ANN/LIME models failed to correctly predict the true VTE
 276 outcome in four individuals from the testing set (individuals 10, 11, 12 and 15). First, it
 277 is worthy of note that these 4 individuals were all females. Second, the 3 female PE
 278 patients wrongly predicted to be DVT (individuals 11, 12 and 15) were all under oral
 279 contraceptives (OC) at the time of the PE event (age 45, 35 and 53, respectively), but not
 280 individual 10 incorrectly predicted to be PE. While we cannot rule out the possibility

281 that our ANN/LIME models poorly behave in women under OC, we nevertheless sought
282 to investigate whether discordant predictions could be due to genomic outlier
283 individuals harboring very rare disease causing mutations that could make the global
284 ANN/LIME predictions inaccurate, inline with the idea that the discrepancy between
285 (deep learning derived) predicted and observed phenotypes could be a heritable trait
286 [24]. Among these 4 individuals, only two (Individuals 11 and 15) have been sequenced
287 for their whole genome. Sequence data of these two individuals were then scrutinized
288 for candidate rare variants that could explain the VTE phenotype.

289 Individual 11 is a woman that experienced PE under oral contraceptives (OC) at age
290 45. Of note, her ten closest neighbors inferred from HPA data were all DVT patients
291 which would likely explain why the derived ANN predicted her a DVT outcome instead
292 of PE. She was not found to harbour any candidate variation in known VTE genes but
293 presented in her genome with 61 very rare coding variants with strong predicted
294 deleteriousness that could be good candidates responsible for the PE event.

295 Individual 15 is a woman that had experienced PE at age 53 also under OC. Nine out
296 of 10 of her closest proteomics based neighbors were DVT patients which may also
297 explain why this PE patient was incorrectly predicted to be DVT. This patient was found
298 to carry a very rare nonsynonymous variation (rs121918154;
299 PROC:NM_000312:exon9:c.C814T:p.R272C) in the VTE-associated *PROC* gene. This
300 variation has a minor allele frequency of 0.005% in public database
301 (<https://www.ncbi.nlm.nih.gov/snp/rs121918154>), is predicted to be deleterious by
302 several bioinformatic tools and have been previously reported in VTE patients with
303 protein C deficiency [25,26]. This variation is located in the last exon of the gene and is
304 predicted to alter splicing regulatory elements [27–29], which could lead to a deletion of
305 a part of the peptidase S1 domain that is responsible for the clivage activity of the
306 protein. Of note, this patient exhibited moderately low plasma Protein C levels, 63%,
307 slightly lower than the 65% threshold adopted to declare moderate protein C deficiency
308 [30].

309

310 **Discussion**

311 This work is original in at least three main aspects. First, it is the largest plasma
312 proteomic study with respect to pulmonary embolism in VTE patients. Second, it is to
313 our knowledge the first attempt to deploy ANN methodologies on proteomic data with

314 the aim at identifying new molecular thrombotic players. And finally, the integration of
315 proteomic and genomics data identified *PLXNA4* as a new candidate gene for PE.
316 This work started with the implementation of an ANN methodology on antibody based
317 affinity proteomics data in relation to PE risk. This ANN was not developed as a tool to
318 be used in clinic for predicting PE risk as 1/one is not 100% certain about the identity of
319 the identified tagged proteins [31] (further experimental validation would be needed to
320 assess this) and 2/ plasma protein levels determined with the antibody suspension bead
321 array are not absolute but relative values depending on the current set of studied
322 antibodies. Rather, our ANN based predictor for PE was aimed at serving as an
323 intermediate surrogate biomarker that could generate new knowledge about the
324 (genetics) mechanisms involved in PE. By conducting a GWAS on the derived PE
325 predictor and capitalizing on two case-control samples totalling 467 patients with PE
326 and 1414 patients with DVT, we observed that the *PLXNA4* rs1424597-A allele was
327 associated with a ~2-fold increased risk of PE in VTE patients aged more than ~40yrs.
328 *PLXNA4* codes for Plexin A4, which is part of a receptor complex involved in signal
329 transduction of semaphorin 3A signals linked to cytoskeletal rearrangement, inhibiting
330 integrin adhesion [32,33]. It has a role in axone guidance in nervous system
331 development, and genetic variants in *PLXNA4* have been linked to risk of Alzheimer
332 disease [34,35]. Based on RNA seq data from HPA, FANTOM and GTEx datasets, *PLXNA4*
333 is expressed at medium/high levels in central nervous system, adipose, breast and
334 female reproductive tract tissues, and low levels in a broad range of other tissues
335 (<https://www.proteinatlas.org/ENSG00000221866-PLXNA4/tissue>), indicating roles outside
336 the nervous system. RNA seq data from blood cell populations show expression in
337 plasma cytotoid dendritic cells, NK cells and some T-cell populations, and research based
338 on animal studies suggest a role in immunity and immune function. It has been shown to
339 be a negative regulator of T cell activation [36]. Besides, Wen et al [37] found it to be
340 highly expressed in myeloid cells, where *PLXNA4* had an important function in
341 stimulating TNF-alpha and IL-6 production in macrophages, where knock out mice were
342 protected against lethal dose LPS induced cytokine storms, suggesting it having a critical
343 role in mediating pro-inflammatory cytokine production. The ligand of *PLXNA4*, SEMA3,
344 has also been described with a role in endothelial cell function in an autocrine loop,
345 promoting processes involved in vascular remodeling [38], and also in negatively
346 regulating platelet aggregation [39]. While *PLXNA4* thus has been described with a role

347 in processes/pathways of relevance for thrombosis, little is known about *PLXNA4* in
348 pulmonary embolism. In addition, we did not identify strong elements supporting a
349 functional role of the intronic rs1424597 polymorphisms or of any other
350 polymorphisms in strong linkage disequilibrium with it. Nevertheless, rs1424597 has
351 recently been observed in the FinnGen study (<http://r3.finngen.fi/>) to be marginally
352 associated ($p = 4.5 \cdot 10^{-3}$) with pleural conditions that are inflammatory disorders of the
353 lung. Consistent with this observation, we observed a positive correlation between the
354 rs1424597-associated PE predictor and fibrinogen, a well known inflammatory marker.
355 Additional *PLXNA4* polymorphisms have also been reported to demonstrate strong
356 statistical evidence for association with various lung function markers [40,41]. Besides,
357 our group have previously suggested that *PLXNA4* polymorphisms could interact with
358 polymorphisms at *UQCC1* to modulate the risk of VTE in the general population [42],
359 *UQCC1* being a locus that have also be shown to be involved in lung function [43].
360 Altogether, these observations strongly support for a role of *PLXNA4* in lung function
361 and its precise role in the etiology of pulmonary embolism deserve further investigation.
362 Which polymorphisms could be truly responsible for the observed association with PE
363 risk also merits further works as the rs1424597 is likely tagging for functional
364 variant(s)/haplotypes yet to be characterized. Finally, further studies would be needed
365 to investigate whether the previously suggested *PLXNA4* x *UQCC1* interaction on the risk
366 of VTE (combining both DVT and PE) could be more specific to patients with PE.
367 In addition to searching for common polymorphisms that could associate with our ANN
368 based predictor and with PE risk, we also looked for rare variants that could explain the
369 discrepancy between predicted and observed VTE outcome in our testing set. Two out of
370 four patients with discordant predictions in the testing set have been sequenced for
371 their whole genome. Both were females patients that experienced PE under OC. In one of
372 them, we were able to identify a rare VTE causing mutation in *PROC*. It is not our
373 intention to conclude to any general rule about the relevance of searching of rare
374 variants responsible for any discordancy between ANN predictions and observed
375 outcomes. Especially as we observed that the three PE patients wrongly predicted to be
376 DVT were women who developed PE under OC. These observations could suggest that
377 our plasma proteomics ANN derived predictions may not be valid in such subgroups of
378 VTE patients and highlight the challenge to identify general prediction models for
379 complex diseases. Several additional limitations must be addressed. First, no plasma

380 antibody targeting *PLXNA4* was available when the screening phase of this work was
381 initiated preventing us from validating further its association with PE. Second, no
382 proteomic data was available in the EOVT study to formally replicate the association of
383 our ANN and LIME predictors with PE risk. Third, our GWAS analysis on the ANN
384 derived predictor was performed only in 574 samples which has likely hampered our
385 power to identify genome-wide significant SNPs. We may have then missed additional
386 polymorphisms that could be truly associated with the predictor and could have then
387 help us to better disentangle its underlying molecular biology. Finally, the moderate
388 sample size of the EOVT study has also likely hampered our power for statistically
389 replicating the association of the lead *PLXNA4* polymorphism with PE. In addition, no
390 information was available in the EOVT study to distinguish isolated PE From DVT+PE
391 which prevented us from further testing whether the association of *PLXNA4* with PE risk
392 was mainly restricted to isolated PE as suggested from the MARTHA results. However,
393 the very consistent pattern of association observed according to age strata between the
394 two studies is a strong argument in favor of *PLXNA4* as a new candidate in PE biology.
395 In conclusion, by implementing an original artificial neural network methodology
396 integrating plasma proteomics and genetic data, we identified *PLXNA4* as a new
397 candidate susceptibility gene for PE in VTE patients whose precise role in PE etiology
398 deserve further investigations

399

400 **Materials and Methods**

401

402 Ethical approval

403 Each individual study on which the work is based was approved by its institutional
404 ethics committee and informed written consent was obtained in accordance with the
405 Declaration of Helsinki. Ethics approval were obtained from the “Departement santé de
406 la direction générale de la recherche et de l’innovation du ministère” (Projects DC: 2008-
407 880 & 09.576) and from the institutional ethics committees of the Kremlin-Bicetre
408 Hospital.

409

410 MARTHA study

411 The MARTHA population is composed of VTE patients recruited from the Thrombophilia
412 center of La Timone hospital (Marseille, France) and free of any chronic conditions and

413 of any well characterized genetic risk factors including antithrombin, protein C or
414 protein S deficiency, homozygosity for FV Leiden or Factor II 20210A, and lupus
415 anticoagulant. Detailed description of the MARTHA population has been provided
416 elsewhere [14,44].

417 *MARTHA proteomics substudy.* A sample of 1,388 MARTHA patients with available
418 plasma samples were profiled for targeted plasma proteomic investigations as described
419 below.

420 *MARTHA genetic substudy.* From the whole MARTHA population, 1592 patients
421 with DNA available were genotyped with high-throughput genotyping arrays (see
422 below).

423

424 Plasma proteomic profiling

425 *Generation of antibody suspension bead array (SBA)*

426 The multiplex antibody suspension bead array (SBA) was created by covalent coupling
427 of 339 Human Protein Atlas (HPA) antibodies, 13 from commercial providers and 25
428 monoclonal BSI antibodies (BioSystems International Kft) targeting 234 unique
429 candidate proteins (Supplemental Table 2). Antibodies were individually coupled to
430 carboxylated magnetic beads (MagPlex-C, Luminex Corp.) generating up to 384 different
431 bead identities (IDs), essentially according to methods previously described [45,46]. The
432 final multiplexed suspension bead array was prepared by combining all 384 antibody
433 coupled beads into a single SBA stock with a concentration of approximately 25-40
434 beads of each antibody bead ID/ul.

435

436 *Plasma labelling and protein profiling assay*

437 Plasma samples were diluted 1:10 in filtered 1xPBS and labelled with biotin (NHS-PEG4-
438 Biotin, Thermo Scientific) for 2h at 4°C. The labelling process was terminated by the
439 addition of 12,5ul of 0.5M HCl pH:8.0 to each sample for 20 min and consecutively
440 storage at -20°C until usage [45]. Labelled plasma samples were diluted 1:50 in PVX
441 casein buffer + 10% (v/v) rabbit IgG (0.1% casein, 0.5% polyvinyl alcohol, 0.8%
442 polyvinylpyrrolidone, prepared in 1xPBS). Diluted samples were heat-induced to
443 achieve epitope retrieval for 30 minutes at 56°C. Five microliters of the SBA were mixed
444 with 45ul of heat-treated samples for 16-18 hours, at RT and constant shake. Unbound
445 complexes were removed by 2 consecutive washes with PBS-T and antibody-bound

446 complexes were cross-linked by resuspending the beads in 0.4% PFA-PBS for 10 min. R-
447 phycoerythrin-conjugated streptavidin (1:750, PBS-T; Invitrogen) was added to all
448 samples for 30 min followed by 2 times washes. Relative amount of each protein
449 complex was expressed as median of fluorescence intensity (MFI) by read out on a
450 FlexMAP3D.

451

452 The Early Onset Venous Thrombosis (EOVT) study

453 This study is composed of 339 VTE patients with documented idiopathic isolated PE or
454 DVT selected according to the same criteria as the MARTHA participants, with the
455 exception that the age of VTE onset was below 50 yrs. Detailed description can be found
456 in [44,47]

457

458 Deep-learning framework for identifying a molecular predictor of PE risk.

459 *Step 1 : Normalization*

460 First, all HPA variables were normalized and scaled to have 0 mean et 1 variance to
461 avoid major artificial influence of variables with large range of variations.

462

463 *Step2 : Edited nearest neighbors*

464 As our aim was to identify new molecular markers associated with PE, we
465 hypothesized that conducting our discovery phase on isolated PE, an expected less
466 heterogeneous class of VTE patients than the class of patients with both DVT and PE, will
467 increase our chance to identify novel relevant molecular players. As a consequence, we
468 decided to built our ANN model only on patients with isolated PE (N = 95) or with DVT
469 (N = 1105). However, due to the imbalance nature of this dataset with ~10 more
470 samples in the DVT class than in the PE class, we applied the edited nearest neighbors
471 (ENN) algorithm, an under sampling method usually used in the field of pattern
472 recognition or classification in presence of unbalanced samples [16]. This method relies
473 on under sampling unit of analysis, in our case individuals, from the majority class by
474 removing the most heterogenous units. It consists in computing the Euclidian distance
475 between each pair of individuals from their proteomics data and to remove samples
476 whose clinical phenotype (here DVT) is not consistent with that of his/her k nearest
477 neighbors (k=3 in this work). This led us to the selection of the so called ANN dataset
478 composed of N = 497 DVT and N = 95 PE patients for building our ANN model.

479

480 *Step3: Derivation of an ANN model for PE prediction.*

481 To build our ANN model, the ANN dataset was divided into a training set composed of
482 576 patients (487 DVT and 89 PE) and a testing set of 16 patients (10 DVT and 6 PE), the
483 latter being used for testing the accuracy of the ANN model derived from the former.
484 This allocation was chosen so that the number of PE cases used for training was
485 sufficiently large.

486 Because the application of a standard ANN methodology to our training set would
487 lead to instable network for predicting PE due to the imbalance nature of the input data
488 with ~5 times more DVT than PE patients, an interactive ANN framework was adopted:

489 At each iteration i ,

490 - A random sample of 30 PE patients and 100 DVT patients is selected from the
491 training set and a sample of 70 synthetic PE samples are generated using the ADASYN
492 algorithm [48]. ADASYN is an adaptive synthetic data generation method where new
493 samples are generated based on the weighted distribution for minority class samples
494 with two main advantages, resolving data imbalance and forcing classifiers to be more
495 sensitive to the minority class. This strategy led to a balanced dataset D_i of 100 PE and
496 100 DVT (synthetic) patients on which a ANN is built.

497 - Using the D_i dataset further splitted randomly into 90%/10% training/testing
498 subsamples, a two hidden-layers feed forward neural networks was implemented. The
499 first hidden layer has 395 neurons corresponding to the number of input (proteomic &
500 biological) variables, the second layer 128 while the output layer consisted in 2 neurons,
501 representing the DVT and PE classes respectively. The number of neurons were selected
502 by trial and error approach under the constraint that the number of neurons shall be
503 smaller than the number of input variables and higher than the number of output classes

504 The Rectified Linear Unit (ReLU) function [49] was used to activate hidden layers
505 while the softmax activation function [50] was used to generate class probabilities in the
506 output layer.

507 After fixing the number of nodes, layer and activation function, the process of training
508 the neural network can start. Starting from random weights, forward propagation is
509 used to generate the output of all nodes at all layers while moving from the input to the
510 output layers. The generated final output is compared to the observed class phenotype
511 and an error is calculated using the cross-entropy function [51]. Iteratively, this error

512 was then backpropagated using a gradient descent algorithm [52] (with learning of 0.01
513 and batch size of 32) to update weights according to their contribution to the error. In
514 order to reduce over-fitting and obtain the best performing model, the callback feature
515 proposed by the Keras open-source library (<https://keras.io/>) that was employed for
516 developing this ANN framework was used.

517

518 *Step4: Local Interpretable Model-Agnostic Explanations (LIME)*

519 As a neural network is often considered as a black box without telling much about
520 which, and how, input variables contribute the most to the predictive model, the LIME
521 methodology [17] was applied to the final ANN model obtained at Step 3 in order to
522 inform about which input variables (i.e plasma proteoin levels) contribute top PE risk
523 prediction and what are the relative weights using a linear approximation of the ANN
524 model.

525

526 Genome Wide Genotyping

527 As previously described [15,44], both MARTHA and EOVT participants have been
528 genotyped with high-density genotyping Illumina arrays and imputed for single
529 nucleotide polymorphisms (SNPs) from the 1000G Phase I Integrated Release Version 2
530 Haplotypes using MACH (v1.0.18.c) and Minimac (release 2011-10-27) imputation
531 software.

532

533 Genome-Wide Association analysis (GWAS)

534 Imputed SNPs with imputation quality r^2 greater than 0.5 and with minor allele
535 frequency (MAF) greater than 0.01. were tested for association with the LIME predictor
536 derived in 574 MARTHA participants. Associations with statistical p-value $< 5 \times 10^{-8}$
537 were considered as genome-wide significant.

538

539 Genetic Association Analysis with PE risk

540 The candidate SNP identified from the GWAS on the LIME predictor was tested for its
541 association with PE risk, both in MARTHA and EOVT participants. For this, we employed
542 the Cochran-Armitage trend for association applied to the best guessed genotypes
543 inferred from the imputed allele dosage at the SNP of interest. Meta-analysis of the
544 association results observed in MARTHA and EOVT was conducted using a fixed-effects

545 model based on the inverse-variance weighting and heterogeneity of association
546 between the two studies was assessed by the Cochran-Mantel-Haenszel test statisticsn-
547 Mantel-Haenszel test statistic [53]

548

549 Whole Genome Sequencing

550 From the whole MARTHA study, 200 patients had been selected for whole genome
551 sequencing. These patients were selected to have experienced unprovoked VTE. Besides,
552 these patients should have family history of VTE or multiple unprovoked VTE events, such
553 clinical patterns being compatible with the existence of an underlying VTE causing genetic
554 defect. Genomic DNA was extracted from peripheral blood, using the BioRobot EZ1
555 workstation. The DNA concentration was determined using the Qubit assay kit
556 (Thermofisher). Whole genome sequencing was performed at the Centre National de
557 Recherche en Génomique Humaine (CNRGH, Institut de Biologie François Jacob, Evry,
558 FRANCE). After a complete quality control, 1µg of genomic DNA was used for each sample
559 to prepare a library for whole genome sequencing, using the Illumina TruSeq DNA PCR-Free
560 Library Preparation Kit, according to the manufacturer's instructions. After normalisation and
561 quality control, qualified libraries were sequenced on a HiSeqX5 instrument from Illumina
562 (Illumina Inc., CA, USA) using a paired-end 150 bp reads strategy. One lane of HiSeqX5
563 flow cell was used per sample specific library in order to reach an average sequencing depth
564 of 30x for each sequenced individual. Sequence quality parameters have been assessed
565 throughout the sequencing run and standard bioinformatics analysis of sequencing data was
566 based on the Illumina pipeline to generate FASTQ file for each sample. FastQ sequences were
567 aligned on human genome hg37 using the BWA-mem program [54]. Variant calling was
568 performed using the GATK HaplotypeCaller (GenomeAnalysisTK-v3.3-0,
569 <https://software.broadinstitute.org/gatk/documentation/article.php?id=4148>). Single-sample
570 gVCFs files were then aggregated using GATK CombineGVCFs and joint genotyping calling
571 performed by GATK GenotypeGVCFs. Recalibration was then conducted on the whole gVCF
572 following GATK guidelines. Following GATK VQSR, we retained single nucleotide variants
573 in the 99.5% tranche sensitivity threshold and indels in the 99% tranche sensitivity threshold
574 for further analysis and annotated them using Annovar [55].

575 As a strategy to identify candidate variants that could explain the VTE phenotype in
576 individuals with discordant class prediction, we first prioritized variants that were likely
577 functional (stop loss/stop gain, frameshift, non-synonymous and splicing variants), located in
578 known VTE associated genes (ABO, ARID4A, C4BPB, EIF5A, F2, F3, F5, F8, F9, F13A1,

579 FGG, GRK5, MPHOSPH9, MAST2, NUGCC, OSMR, PLAT, PLCG2, PLEK1, PROC,
580 PROS1, SCARA5, SERPINC1, SLC44A2, STAB2, STX10, STXBP5, THBD, TSPAN15,
581 VWF) [56–58], that have not been reported or at a low frequency (<1%) in public genomic
582 data repositories (dbSNP, GnomAD) and that was present in only one of the 200 sequenced
583 patients. If no candidate variants was identified in known VTE genes, we extended our search
584 to whole coding genes and also took into account the predicted deleteriousness of selected
585 candidates using *in silico* tools such as SIFT, PolyPhen and CADD-v1.2 [59] to further reduce
586 the number of candidates.

587

588

589 **Acknowledgments**

590 Mi.R; O.S. and Ma.R and the production of the MARTHA genomics data were financially
591 supported by the GENMED Laboratory of Excellence on Medical Genomics [ANR-10-
592 LABX-0013], a research program managed by the National Research Agency (ANR) as part
593 of the French Investment for the Future. This work benefited from the financial support from
594 the «EPIDEMIOLOGIE-VTE» Senior Chair from the Initiative of Excellence of the University of
595 Bordeaux. Bioinformatics and statistical analyses benefit from the CBiB computing centre of
596 the University of Bordeaux. The proteomics screening was financed by a grant from
597 Stockholm County Council (SLL 2017-0842) and from Familjen Erling Perssons Foundation.

598

599 **References**

- 600 1. White RH. The epidemiology of venous thromboembolism. *Circulation*. 2003;107: I4-8.
601 doi:10.1161/01.CIR.0000078468.11849.66
- 602 2. Konstantinides SV, Torbicki A, Agnelli G, Danchin N, Fitzmaurice D, Galiè N, et al.
603 2014 ESC guidelines on the diagnosis and management of acute pulmonary embolism.
604 *Eur Heart J*. 2014;35: 3033–3069, 3069a–3069k. doi:10.1093/eurheartj/ehu283
- 605 3. Contou D, Pajot O, Cally R, Logre E, Fraissé M, Mentec H, et al. Pulmonary embolism
606 or thrombosis in ARDS COVID-19 patients: A French monocenter retrospective study.
607 *PLoS One*. 15: e0238413.
- 608 4. van Langevelde K, Flinterman LE, van Hylckama Vlieg A, Rosendaal FR, Cannegieter
609 SC. Broadening the factor V Leiden paradox: pulmonary embolism and deep-vein
610 thrombosis as 2 sides of the spectrum. *Blood*. 2012;120: 933–946. doi:10.1182/blood-
611 2012-02-407551

- 612 5. Rodriguez BAT, Bhan A, Beswick A, Elwood PC, Niiranen TJ, Salomaa V, et al. A
613 Platelet Function Modulator of Thrombin Activation Is Causally Linked to
614 Cardiovascular Disease and Affects PAR4 Receptor Signaling. *Am J Hum Genet.* 2020.
615 doi:10.1016/j.ajhg.2020.06.008
- 616 6. Wang X, Sundquist K, Svensson PJ, Rastkhani H, Palmér K, Memon AA, et al.
617 Association of recurrent venous thromboembolism and circulating microRNAs. *Clin*
618 *Epigenetics.* 2019;11: 28. doi:10.1186/s13148-019-0627-z
- 619 7. Thibord F, Munsch G, Perret C, Suchon P, Roux M, Ibrahim-Kosta M, et al. Bayesian
620 network analysis of plasma microRNA sequencing data in patients with venous
621 thrombosis. *Eur Heart J Suppl.* 2019. Available: [https://www.hal.inserm.fr/inserm-](https://www.hal.inserm.fr/inserm-02310241)
622 [02310241](https://www.hal.inserm.fr/inserm-02310241)
- 623 8. Bruzelius M, Iglesias MJ, Hong M-G, Sanchez-Rivera L, Gyorgy B, Souto JC, et al.
624 PDGFB, a new candidate plasma biomarker for venous thromboembolism: results from
625 the VEREMA affinity proteomics study. *Blood.* 2016;128: e59–e66. doi:10.1182/blood-
626 2016-05-711846
- 627 9. Jensen SB, Hindberg K, Solomon T, Smith EN, Lapek JD, Gonzalez DJ, et al. Discovery
628 of novel plasma biomarkers for future incident venous thromboembolism by untargeted
629 synchronous precursor selection mass spectrometry proteomics. *J Thromb Haemost.*
630 2018;16: 1763–1774. doi:10.1111/jth.14220
- 631 10. Fraser K, Roy N, Goumidi, Louisa, Verdu A, Suchon P, Leal-Valentim F, et al. Plasma
632 Biomarkers and Identification of Resilient Metabolic Disruptions in Patients With
633 Venous Thromboembolism Using a Metabolic Systems Approach. In: *Arteriosclerosis,*
634 *thrombosis, and vascular biology* [Internet]. *Arterioscler Thromb Vasc Biol;* 8 Jun 2020
635 [cited 16 Sep 2020]. doi:10.1161/ATVBAHA.120.314480
- 636 11. Zeleznik OA, Poole EM, Lindstrom S, Kraft P, Van Hylekama Vlieg A, Lasky-Su JA, et
637 al. Metabolomic analysis of 92 pulmonary embolism patients from a nested case-control
638 study identifies metabolites associated with adverse clinical outcomes. *J Thromb*
639 *Haemost.* 2018;16: 500–507. doi:10.1111/jth.13937
- 640 12. Insenser M, Montes-Nieto R, Martínez-García MÁ, Durán EF, Santiuste C, Gómez V, et
641 al. Identification of reduced circulating haptoglobin concentration as a biomarker of the
642 severity of pulmonary embolism: a nontargeted proteomic study. *PLoS ONE.* 2014;9:
643 e100902. doi:10.1371/journal.pone.0100902
- 644 13. Uhlén M, Fagerberg L, Hallström BM, Lindskog C, Oksvold P, Mardinoglu A, et al.
645 Proteomics. Tissue-based map of the human proteome. *Science.* 2015;347: 1260419.
646 doi:10.1126/science.1260419
- 647 14. Oudot-Mellakh T, Cohen W, Germain M, Saut N, Kallel C, Zelenika D, et al. Genome
648 wide association study for plasma levels of natural anticoagulant inhibitors and protein C
649 anticoagulant pathway: the MARTHA project. *Br J Haematol.* 2012;157: 230–239.
650 doi:10.1111/j.1365-2141.2011.09025.x

- 651 15. Germain M, Chasman DI, de Haan H, Tang W, Lindström S, Weng L-C, et al. Meta-
652 analysis of 65,734 individuals identifies TSPAN15 and SLC44A2 as two susceptibility
653 loci for venous thromboembolism. *Am J Hum Genet.* 2015;96: 532–542.
654 doi:10.1016/j.ajhg.2015.01.019
- 655 16. Wilson D L. Asymptotic properties of nearest neighbor rules using edited data. *IEEE*
656 *Transactions on Systems, Man, and Cybernetics.* 1972;3: 408–421.
- 657 17. Ribeiro MT, Singh S, Guestrin C. Why should i trust you?” explaining the predictions of
658 any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on*
659 *Knowledge Discovery and Data Mining.* 2016; 1135–1144.
- 660 18. Sun BB, Maranville JC, Peters JE, Stacey D, Staley JR, Blackshaw J, et al. Genomic
661 atlas of the human plasma proteome. *Nature.* 2018;558: 73–79. doi:10.1038/s41586-018-
662 0175-2
- 663 19. Suhre K, Arnold M, Bhagwat AM, Cotton RJ, Engelke R, Raffler J, et al. Connecting
664 genetic risk to disease end points through the human blood plasma proteome. *Nat*
665 *Commun.* 2017;8: 14357. doi:10.1038/ncomms14357
- 666 20. Maaten L, Hinton G. Visualizing data using t-SNE. *ournal of machine learning research.*
667 2008;9: 2579–2605.
- 668 21. McInnes L, Healy J. UMAP: Uniform manifold approximation and projection for
669 dimension reduction. *arXiv preprint arXiv:180203426.* 2018.
- 670 22. Butler LM, Hallström BM, Fagerberg L, Pontén F, Uhlén M, Renné T, et al. Analysis of
671 Body-wide Unfractionated Tissue Data to Identify a Core Human Endothelial
672 Transcriptome. *Cell Syst.* 2016;3: 287-301.e3. doi:10.1016/j.cels.2016.08.001
- 673 23. Anderson FA, Spencer FA. Risk factors for venous thromboembolism. *Circulation.*
674 2003;107: 19-16. doi:10.1161/01.CIR.0000078469.07362.E6
- 675 24. Jonsson BA, Bjornsdottir G, Thorgeirsson TE, Ellingsen LM, Walters GB, Gudbjartsson
676 DF, et al. Brain age prediction using deep learning uncovers associated sequence
677 variants. *Nature Communications.* 2019;10: 5409. doi:10.1038/s41467-019-13163-9
- 678 25. Allaart CF, Poort SR, Rosendaal FR, Reitsma PH, Bertina RM, Briët E. Increased risk of
679 venous thrombosis in carriers of hereditary protein C deficiency defect. *Lancet.*
680 1993;341: 134–138. doi:10.1016/0140-6736(93)90003-y
- 681 26. Reitsma PH, Poort SR, Allaart CF, Briët E, Bertina RM. The spectrum of genetic defects
682 in a panel of 40 Dutch families with symptomatic protein C deficiency type I:
683 heterogeneity and founder effects. *Blood.* 1991;78: 890–894.
- 684 27. Erkelenz S, Theiss S, Otte M, Widera M, Peter JO, Schaal H. Genomic HEXploring
685 allows landscaping of novel potential splicing regulatory elements. *Nucleic Acids Res.*
686 2014;42: 10681–10697. doi:10.1093/nar/gku736

- 687 28. Ke S, Shang S, Kalachikov SM, Morozova I, Yu L, Russo JJ, et al. Quantitative
688 evaluation of all hexamers as exonic splicing elements. In: *Genome research* [Internet].
689 *Genome Res*; Aug 2011 [cited 21 Sep 2020]. doi:10.1101/gr.119628.110
- 690 29. Tubeuf H, Charbonnier C, Soukarieh O, Blavier A, Lefebvre A, Dauchel H, et al. Large-
691 scale comparative evaluation of user-friendly tools for predicting variant-induced
692 alterations of splicing regulatory elements. *Hum Mutat*. 2020. doi:10.1002/humu.24091
- 693 30. Lijfering WM, Brouwer J-L P, Veeger, Nic J. G. M., Bank I, Coppens M, Middeldorp,
694 Saskia, et al. Selective testing for thrombophilia in patients with first venous thrombosis:
695 results from a retrospective family cohort study on absolute thrombotic risk for currently
696 known thrombophilic defects in 2479 relatives. In: *Blood* [Internet]. *Blood*; 21 May
697 2009 [cited 16 Sep 2020]. doi:10.1182/blood-2008-10-184879
- 698 31. Fredolini C, Byström S, Sanchez-Rivera L, Ioannou M, Tamburro D, Pontén F, et al.
699 Systematic assessment of antibody selectivity in plasma based on a resource of
700 enrichment profiles. *Sci Rep*. 2019;9: 8324. doi:10.1038/s41598-019-43552-5
- 701 32. Hu S, Zhu L. Semaphorins and Their Receptors: From Axonal Guidance to
702 Atherosclerosis. *Front Physiol*. 2018;9: 1236. doi:10.3389/fphys.2018.01236
- 703 33. Fard D, Tamagnone L. Semaphorins in health and disease. *Cytokine Growth Factor Rev*.
704 2020. doi:10.1016/j.cytogfr.2020.05.006
- 705 34. Han Q, Sun Y-A, Zong Y, Chen C, Wang H-F, Tan L, et al. Common Variants in
706 PLXNA4 and Correlation to CSF-related Phenotypes in Alzheimer's Disease. *Front*
707 *Neurosci*. 2018;12: 946. doi:10.3389/fnins.2018.00946
- 708 35. Jun G, Asai H, Zeldich E, Drapeau E, Chen C, Chung J, et al. PLXNA4 is associated
709 with Alzheimer disease and modulates tau phosphorylation. *Ann Neurol*. 2014;76: 379–
710 392. doi:10.1002/ana.24219
- 711 36. Yamamoto M, Suzuki K, Okuno T, Ogata T, Takegahara N, Takamatsu H, et al. Plexin-
712 A4 negatively regulates T lymphocyte responses. *Int Immunol*. 2008;20: 413–420.
713 doi:10.1093/intimm/dxn006
- 714 37. Wen H, Lei Y, Eun S-Y, Ting JP-Y. Plexin-A4-semaphorin 3A signaling is required for
715 Toll-like receptor- and sepsis-induced cytokine storm. *J Exp Med*. 2010;207: 2943–
716 2957. doi:10.1084/jem.20101138
- 717 38. Bussolino F, Valdembri D, Caccavari F, Serini G. Semaphoring vascular morphogenesis.
718 *Endothelium*. 2006;13: 81–91. doi:10.1080/10623320600698003
- 719 39. Kashiwagi H, Shiraga M, Kato H, Kamae T, Yamamoto N, Tadokoro S, et al. Negative
720 regulation of platelet function by a secreted cell repulsive protein, semaphorin 3A.
721 *Blood*. 2005;106: 913–921. doi:10.1182/blood-2004-10-4092
- 722 40. Hardin M, Cho MH, McDonald M-L, Wan E, Lomas DA, Coxson HO, et al. A genome-
723 wide analysis of the response to inhaled β_2 -agonists in chronic obstructive pulmonary
724 disease. *Pharmacogenomics J*. 2016;16: 326–335. doi:10.1038/tpj.2015.65

- 725 41. Imboden M, Bouzigon E, Curjuric I, Ramasamy A, Kumar A, Hancock DB, et al.
726 Genome-wide association study of lung function decline in adults with and without
727 asthma. *J Allergy Clin Immunol*. 2012;129: 1218–1228. doi:10.1016/j.jaci.2012.01.074
- 728 42. Greliche N, Germain M, Lambert J-C, Cohen W, Bertrand M, Dupuis A-M, et al. A
729 genome-wide search for common SNP x SNP interactions on the risk of venous
730 thrombosis. *BMC Med Genet*. 2013;14: 36. doi:10.1186/1471-2350-14-36
- 731 43. Jackson VE, Latourelle JC, Wain LV, Smith AV, Grove ML, Bartz TM, et al. Meta-
732 analysis of exome array data identifies six novel genetic loci for lung function.
733 *Wellcome Open Res*. 2018;3: 4. doi:10.12688/wellcomeopenres.12583.3
- 734 44. Germain M, Saut N, Greliche N, Dina C, Lambert J-C, Perret C, et al. Genetics of
735 venous thrombosis: insights from a new genome wide association study. *PLoS ONE*.
736 2011;6: e25581. doi:10.1371/journal.pone.0025581
- 737 45. Drobin K, Nilsson P, Schwenk JM. Highly multiplexed antibody suspension bead arrays
738 for plasma protein profiling. *Methods Mol Biol*. 2013;1023: 137–145. doi:10.1007/978-
739 1-4614-7209-4_8
- 740 46. Bruzelius M, Iglesias MJ, Hong M-G, Sanchez-Rivera L, Gyorgy B, Souto JC, et al.
741 PDGFB, a new candidate plasma biomarker for venous thromboembolism: results from
742 the VEREMA affinity proteomics study. *Blood*. 2016;128: e59–e66. doi:10.1182/blood-
743 2016-05-711846
- 744 47. Tréguët D-A, Heath S, Saut N, Biron-Andreani C, Schved J-F, Pernod G, et al.
745 Common susceptibility alleles are unlikely to contribute as strongly as the FV and ABO
746 loci to VTE risk: results from a GWAS approach. *Blood*. 2009;113: 5298–5303.
747 doi:10.1182/blood-2008-11-190389
- 748 48. He H, Bai Y, Garcia EA, Li S. Adaptive synthetic sampling approach for imbalanced
749 learning. *IEEE international joint conference on neural networks*. 2008; 1322–1328.
- 750 49. Hahnloser RH, Sarpeshkar R, Mahowald MA, Douglas RJ, Seung HS. Digital selection
751 and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature*. 2000;405:
752 947–951. doi:10.1038/35016072
- 753 50. Bridle JS. Probabilistic Interpretation of Feedforward Classification Network Outputs,
754 with Relationships to Statistical Pattern Recognition. In: Soulié FF, Héroult J (eds)
755 *Neurocomputing NATO ASI Series (Series F: Computer and Systems Sciences)*, vol 68
756 Springer, Berlin, Heidelberg. 1990. doi:doi.org/10.1007/978-3-642-76153-9_28
- 757 51. Hinton GE, Dayan P, Frey BJ, Neal RM. The “wake-sleep” algorithm for unsupervised
758 neural networks. *Science*. 1995;268: 1158–1161. doi:10.1126/science.7761831
- 759 52. Curry HB. The method of steepest descent for non-linear minimization problems. *Quart*
760 *Appl Math*. 1944;2: 258–261.
- 761 53. Mantel N, Haenszel W. Statistical aspects of the analysis of data from retrospective
762 studies of disease. *J Natl Cancer Inst*. 1959;22: 719–748.

- 763 54. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
764 transform. *Bioinformatics*. 2009;25: 1754–1760. doi:10.1093/bioinformatics/btp324
- 765 55. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants
766 from high-throughput sequencing data. *Nucleic Acids Res*. 2010;38: e164.
767 doi:10.1093/nar/gkq603
- 768 56. Desch KC, Ozel AB, Halvorsen M, Jacobi PM, Golden K, Underwood M, et al. Whole-
769 exome sequencing identifies rare variants in STAB2 associated with venous
770 thromboembolic disease. *Blood*. 2020;136: 533–541. doi:10.1182/blood.2019004161
- 771 57. Lindstrom S, Wang L, Smith EN, Gordon W, van Hylckama Vlieg A, de Andrade M, et
772 al. Genomic and Transcriptomic Association Studies Identify 16 Novel Susceptibility
773 Loci for Venous Thromboembolism. *Blood*. 2019. doi:10.1182/blood.2019000435
- 774 58. Trégouët D-A, Morange P-E. What is currently known about the genetics of venous
775 thromboembolism at the dawn of next generation sequencing technologies. *Br J*
776 *Haematol*. 2018;180: 335–345. doi:10.1111/bjh.15004
- 777 59. Kircher M, Witten DM, Jain P, O’Roak BJ, Cooper GM, Shendure J. A general
778 framework for estimating the relative pathogenicity of human genetic variants. *Nat*
779 *Genet*. 2014;46: 310–315. doi:10.1038/ng.2892

780

781 **Supporting information captions**

782

783 **Supplementary Figure 1** Manhattan plot describing the results of the genome wide
784 association study on the LIME predictor

785

786 **Supplementary Table 1** List of biological phenotypes that have been measured in
787 MARTHA patients

788

789 **Supplementary Table 2** List of Human Protein Atlas antibodies measured in MARTHA
790 patients

791

792 **Supplementary Table 3** Correlation between LIME PE predictor and biological traits
793 available in MARTHA participants used for building the ANN

794

795 **Supplementary Table 4** Main statistical associations ($p < 1 \times 10^{-5}$) observed in the
796 GWAS on LIME predictor

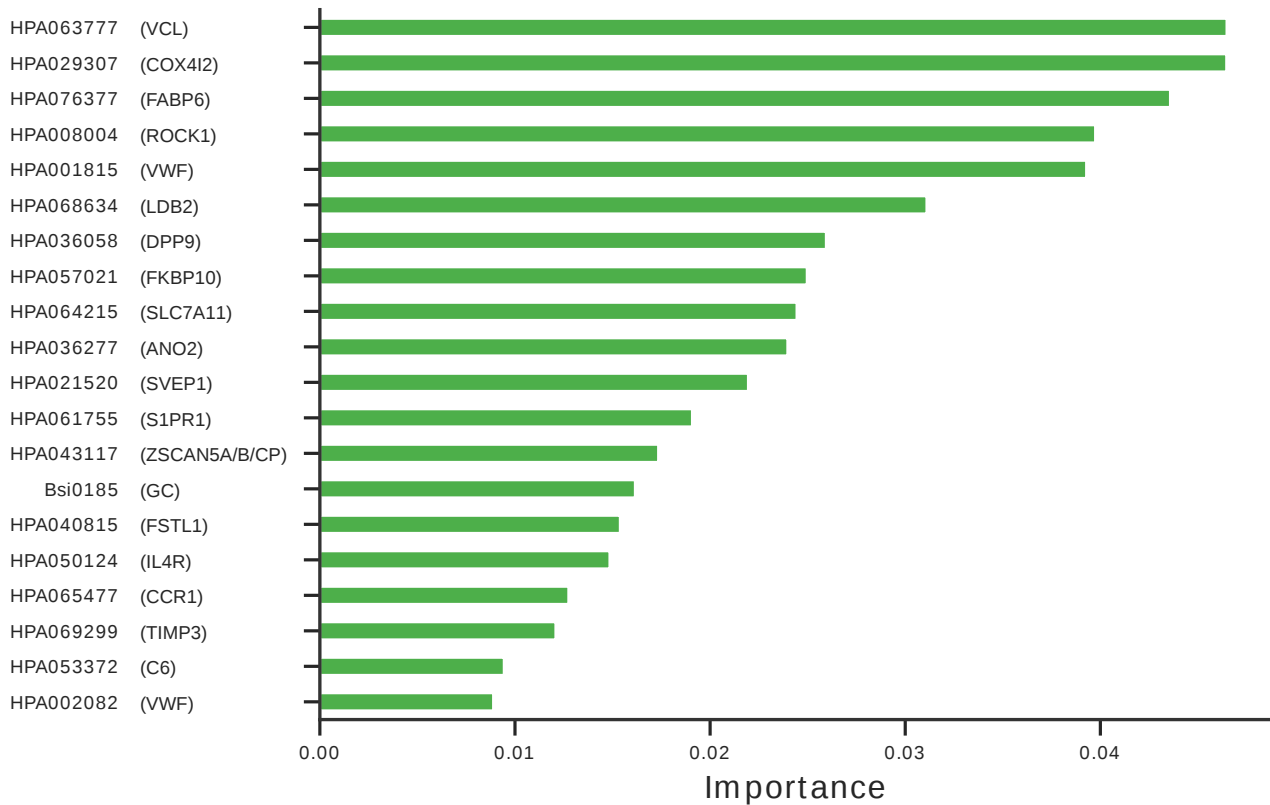
797

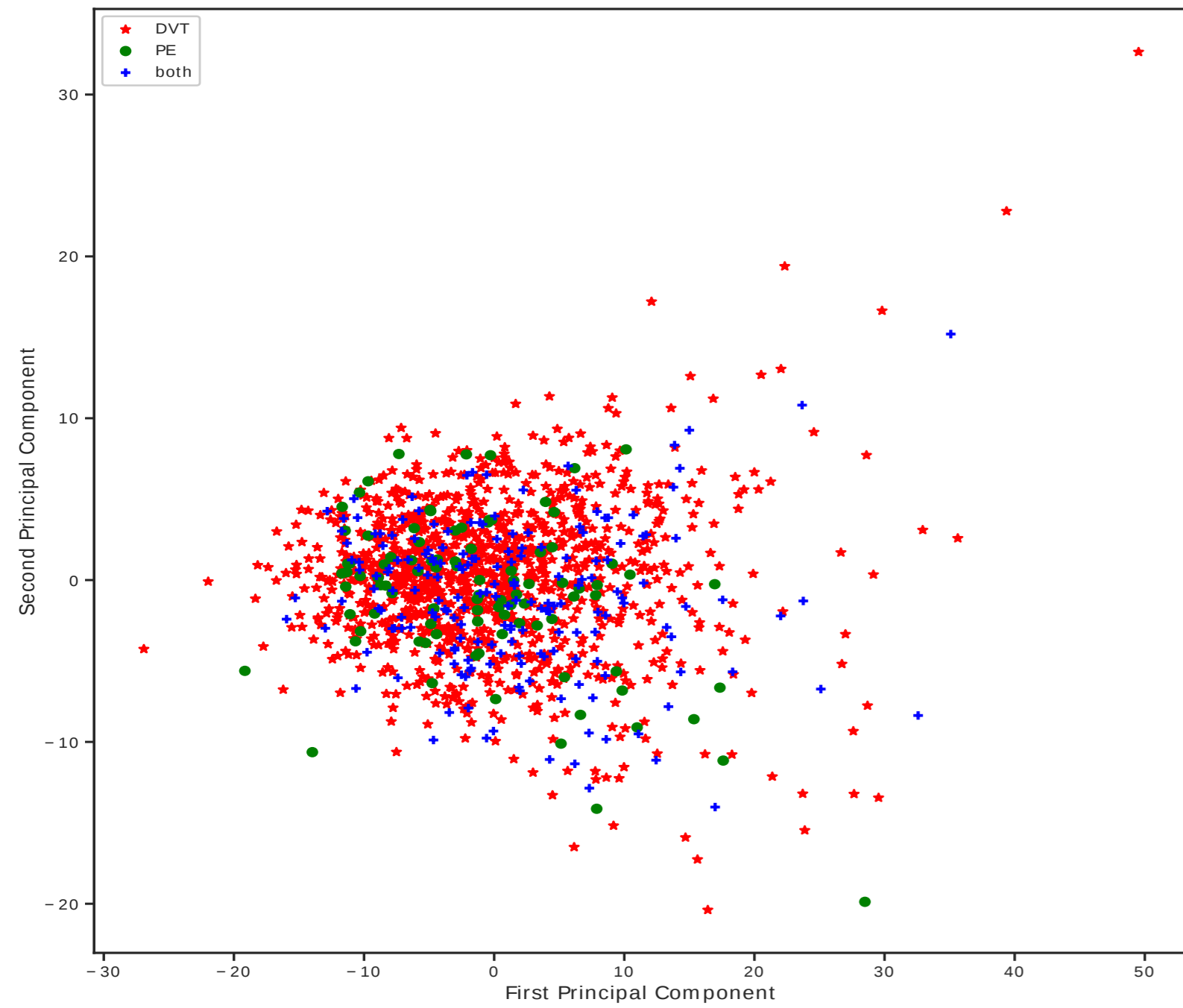
798 **Supplementary Table 5** Rare coding variants identified by whole genome
799 sequencing in individual 11.

800

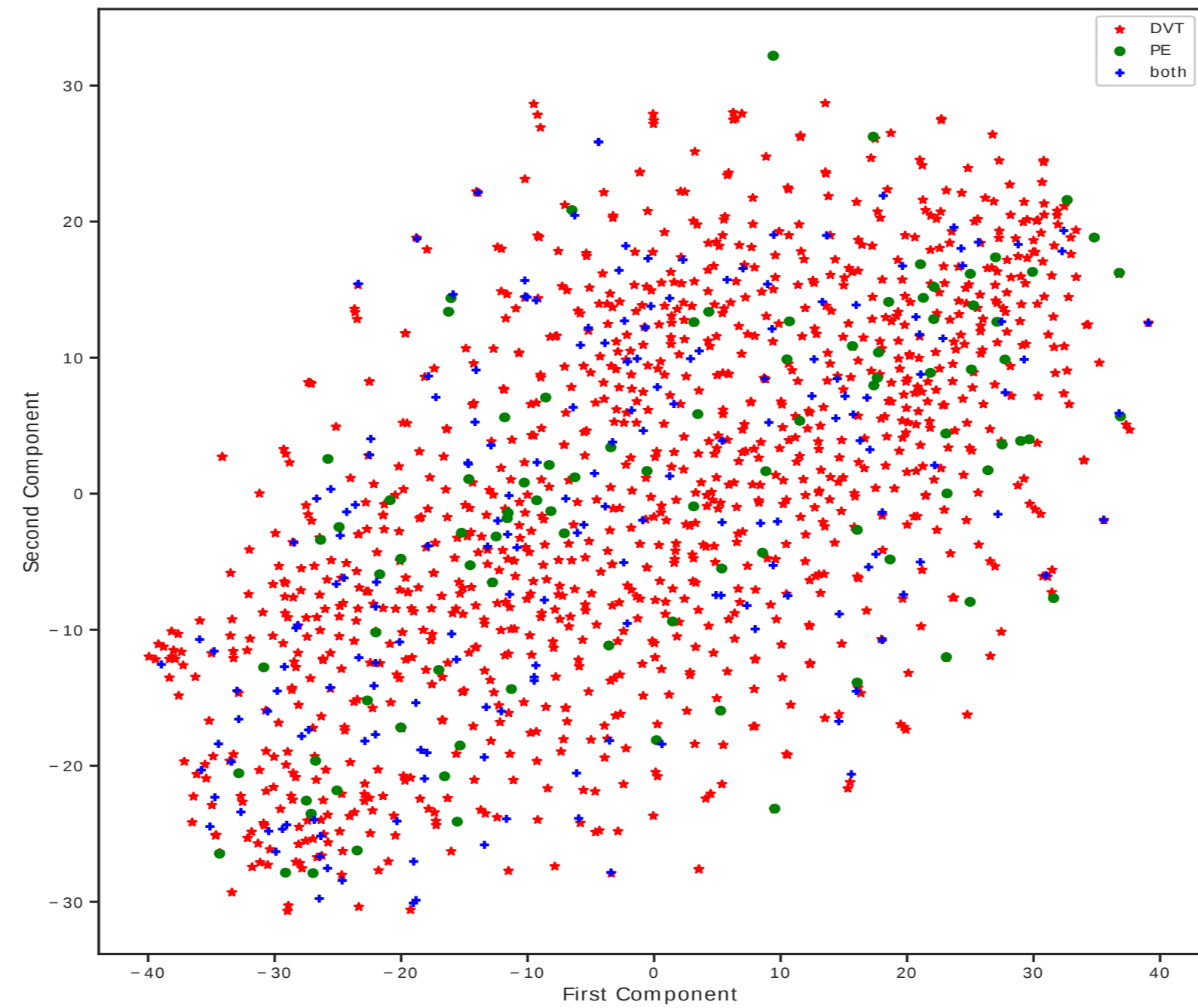
Class PE

Features

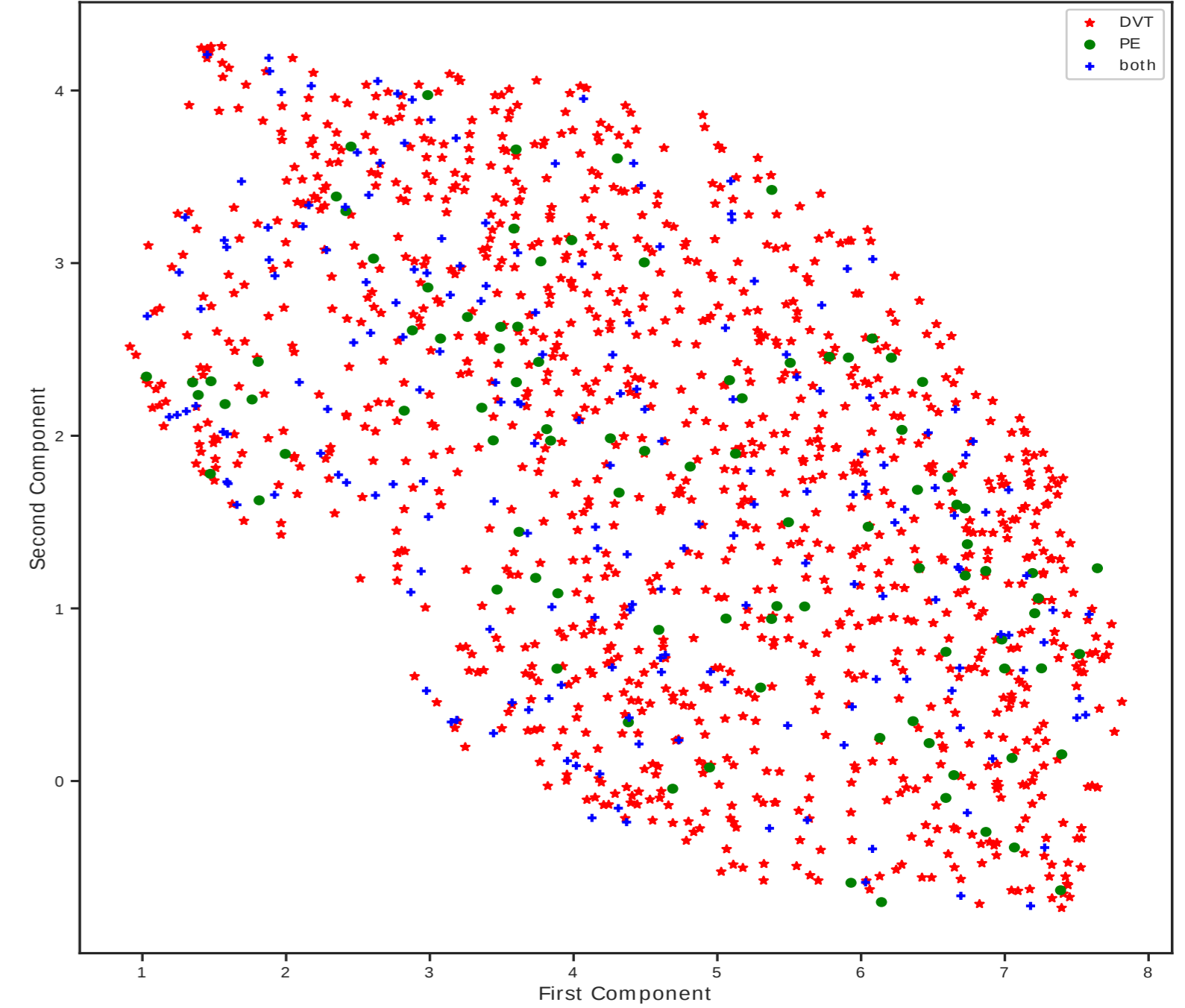




(a) PCA



(b) T-SNE



(c) UMAP

MARTHA

Samples: 1388
Antibodies: 376
Phenotypes: 19
Classes: {DVT:1105, PE:95, DVT+PE:188}

Edited Nearest
Neighbors

Samples: 592
Features: 395
Classes: {DVT:497, PE:95}

Training data

Samples: 576
Features: 395
Classes: {DVT:487, PE:89}

Undersampling

Synthetic Data
ADASYN

ANN training

Iterate

ANN framework

Testing data

Samples: 16
Features: 395
Classes: {DVT:10, PE:6}

LIME explanation
Classes: {DVT:497, PE:95}

Genome wide association study
on LIME predictor
Samples: 574
Classes: {DVT:481, PE:93}

Association of lead SNP on
PE risk in MARTHA
Classes: {DVT:1218, PE:324}

Replication in EOVT
Classes: {DVT:196, PE:143}

Integration of genetics data

