

1 Identification of rare sequence variation underlying 2 heritable pulmonary arterial hypertension

3 Stefan Gräf^{1,2,3,*,\$,#}, Matthias Haime^{1,2,3,*}, Marta Bleda^{1,*}, Charaka Hadinnapola^{1,*}, Laura
4 Southgate^{4,5,*}, Wei Li¹, Joshua Hodgson¹, Bin Liu¹, Richard M. Salmon¹, Mark Southwood⁶,
5 Rajiv D. Machado⁷, Jennifer M. Martin^{1,2,3}, Carmen M. Treacy^{1,6}, Katherine Yates^{1,2,3}, Louise
6 C. Daugherty^{2,3}, Olga Shamardina^{2,3}, Deborah Whitehorn^{2,3}, Simon Holden⁸, Micheala
7 Aldred⁹, Harm J. Bogaard¹⁰, Colin Church¹¹, Gerry Coghlan¹², Robin Condliffe¹³, Paul A.
8 Corris¹⁴, Cesare Danesino^{15,16}, Mélanie Eyries¹⁷, Henning Gall¹⁸, Stefano Ghio¹⁶, Hossein-
9 Ardeschir Ghofrani^{18,19}, J. Simon R. Gibbs²⁰, Barbara Girerd²¹, Arjan C. Houweling¹⁰, Luke
10 Howard¹⁹, Marc Humbert²¹, David G. Kiely¹³, Gabor Kovacs^{22,23}, Robert V. MacKenzie
11 Ross²⁴, Shahin Moledina²⁵, David Montani²¹, Michael Newnham¹, Andrea Olschewski²²,
12 Horst Olschewski^{22,23}, Andrew J. Peacock¹¹, Joanna Pepke-Zaba⁶, Inga Prokopenko¹⁹,
13 Christopher J. Rhodes¹⁹, Laura Scelsi¹⁶, Werner Seeger¹⁸, Florent Soubrier¹⁷, Dan F. Stein¹,
14 Jay Suntharalingam²⁴, Emilia M. Swietlik¹, Mark R. Toshner¹, David A. van Heel²⁶, Anton
15 Vonk Noordegraaf¹⁰, Quinten Waisfisz¹⁰, John Wharton¹⁹, Stephen J. Wort^{27,19}, Willem H.
16 Ouwehand^{2,3}, Nicole Soranzo^{2,28}, Allan Lawrie²⁹, Paul D. Upton¹, Martin R. Wilkins¹⁹, Richard
17 C. Trembath⁵, Nicholas W. Morrell^{1,3,\$,#}

18
19 *These authors contributed equally to this work.

20 \$These authors jointly supervised this work.

21 #Corresponding authors.

22 Affiliations

23
24 ¹Department of Medicine, University of Cambridge, Cambridge, United Kingdom.

25 ²Department of Haematology, University of Cambridge, Cambridge, United Kingdom. ³NIHR
26 BioResource - Rare Diseases, Cambridge, United Kingdom. ⁴Molecular and Clinical
27 Sciences Research Institute, St George's, University of London, London, United Kingdom.

28 ⁵Division of Genetics & Molecular Medicine, King's College London, London, United
29 Kingdom. ⁶Papworth Hospital, Papworth, United Kingdom. ⁷Institute of Medical and
30 Biomedical Education, St George's University of London, Lincoln, United Kingdom.

31 ⁸Addenbrooke's Hospital, Cambridge, United Kingdom. ⁹Cleveland Clinic, Cleveland, Ohio,
32 United States. ¹⁰VU University Medical Center, Amsterdam, The Netherlands. ¹¹Golden
33 Jubilee National Hospital, Glasgow, United Kingdom. ¹²Royal Free Hospital, London, United

34 Kingdom. ¹³Sheffield Pulmonary Vascular Disease Unit, Royal Hallamshire Hospital,
35 Sheffield, United Kingdom. ¹⁴University of Newcastle, Newcastle, United Kingdom.

36 ¹⁵Department of Molecular Medicine, University of Pavia, Pavia, Italy. ¹⁶Fondazione IRCCS
37 Policlinico San Matteo, Pavia, Italy. ¹⁷Département de génétique, hôpital Pitié-Salpêtrière,
38 Assistance Publique-Hôpitaux de Paris, and UMR_S 1166-ICAN, INSERM, UPMC Sorbonne

39 Universités, Paris, France. ¹⁸University of Giessen and Marburg Lung Center (UGMLC),
40 member of the German Center for Lung Research (DZL) and of the Excellence Cluster
41 Cardio-Pulmonary System (ECCCP), Giessen, Germany. ¹⁹Imperial College London,
42 London, United Kingdom. ²⁰National Heart & Lung Institute, Imperial College London,
43 London, United Kingdom. ²¹Université Paris-Sud, Faculté de Médecine, Université Paris-

44 Saclay; AP-HP, Service de Pneumologie, Centre de référence de l'hypertension pulmonaire;

45 INSERM UMR_S 999, Hôpital Bicêtre, Le Kremlin-Bicêtre, Paris, France. ²²Ludwig
46 Boltzmann Institute for Lung Vascular Research, Graz, Austria. ²³Medical University of Graz,
47 Graz, Austria. ²⁴Royal United Hospitals Bath NHS Foundation Trust, Bath, United Kingdom.
48 ²⁵Great Ormond Street Hospital, London, United Kingdom. ²⁶Blizard Institute, Queen Mary
49 University of London, London, United Kingdom. ²⁷Royal Brompton Hospital, London, United
50 Kingdom. ²⁸Wellcome Trust Sanger Institute, Hinxton, United Kingdom. ²⁹Department of
51 Infection, Immunity & Cardiovascular Disease, University of Sheffield, Sheffield, United
52 Kingdom.

53

54

55

56 Correspondence and requests for materials should be addressed to N.W.M
57 (email:'nwm23@cam.ac.uk[mailto:nwm23@cam.ac.uk]') or to S.G.
58 (email:'sg550@cam.ac.uk[mailto:sg550@cam.ac.uk]').

59

60 Abstract

61 Pulmonary arterial hypertension (PAH) is a rare disorder with a poor prognosis. Deleterious
62 variation within components of the transforming growth factor- β pathway, particularly the
63 bone morphogenetic protein type 2 receptor (*BMPR2*), underlie most heritable forms of PAH.
64 To identify the missing heritability we perform whole genome sequencing in 1038 PAH index
65 cases and 6385 PAH-negative control subjects. Case-control analyses reveal significant
66 overrepresentation of rare variants in *ATP13A3*, *AQP1* and *SOX17*, and provide
67 independent validation of a critical role for *GDF2* in PAH. We demonstrate familial
68 segregation of mutations in *SOX17* and *AQP1* with PAH. Mutations in *GDF2*, encoding a
69 *BMPR2* ligand, lead to reduced secretion from transfected cells. In addition, we identify
70 pathogenic mutations in the majority of previously reported PAH genes, and provide
71 evidence for further putative genes. Taken together these findings contribute new insights
72 into the molecular basis of PAH and indicate unexplored pathways for therapeutic
73 intervention.

74 Introduction

75 Idiopathic and heritable pulmonary arterial hypertension (PAH) are rare disorders
76 characterised by occlusion of arterioles in the lung¹, leading to marked increases in
77 pulmonary vascular resistance². Life expectancy from diagnosis averages 3-5 years³, with
78 death ensuing from failure of the right ventricle.

79
80 Mutations in the gene encoding the bone morphogenetic protein type 2 receptor (*BMPR2*), a
81 receptor for the transforming growth factor-beta (TGF- β) superfamily^{4, 5} account for over 80%
82 of families with PAH, and approximately 20% of sporadic cases⁶. Mutations have been
83 identified in genes encoding other components of the TGF- β /bone morphogenetic protein
84 (BMP) signalling pathways, including *ACVRL1*⁷ and *ENG*⁸. On endothelial cells, *BMPR2* and
85 *ACVRL1* form a signaling complex, utilizing *ENG* as a co-receptor. Case reports of rare
86 sequence variation in the BMP signalling intermediaries, *SMAD1*, *SMAD4* and *SMAD9*^{9, 10},
87 provide compelling evidence for a central role of dysregulated BMP signalling in PAH
88 pathogenesis.

89
90 Analysis of coding variation in *BMPR2*-negative kindreds revealed heterozygous mutations
91 in genes not directly impacting on the TGF- β /BMP pathway, including *CAV1*¹¹, and the
92 potassium channel, *KCNK3*¹². Deletions and loss of function mutations in *TBX4*, an essential
93 regulator of embryonic development, were identified in childhood onset PAH¹³. A clinically
94 and pathologically distinct form of PAH, known as pulmonary veno-occlusive disease or
95 pulmonary capillary haemangiomatosis (PVOD/PCH), was shown recently to be caused by
96 biallelic recessive mutations in *EIF2AK4*^{14, 15}, a kinase in the integrated stress response.

97
98 The purpose of the present study was to identify additional rare sequence variation
99 contributing to the genetic architecture of PAH, and to assess the relative contribution of rare
100 variants in genes implicated in prior studies. A major finding is that rare likely causal
101 heterozygous variants in several previously unidentified genes (*ATP13A3*, *AQP1* and
102 *SOX17*) were significantly overrepresented in the PAH cohort, and we provide independent
103 validation for *GDF2* as a causal gene.

104 Results

105 Description of the PAH cohort

106 In total, 1048 PAH cases (1038 index cases and 10 related affected individuals) were
107 recruited for WGS. Of these, 908 (86.7%) were diagnosed with idiopathic PAH, 58 (5.5%)
108 gave a family history of PAH and 60 (5.7%) gave a history of drug exposure associated with
109 PAH¹⁶. Twenty two cases (2.1%) held a clinical diagnosis of PVOD/PCH (Figure 1a).
110 Demographic and clinical characteristics of the PAH cohort are provided in Supplementary
111 Table 1. An additional UK family was recruited separately for novel gene identification
112 studies. Briefly, the proband was diagnosed at 12 years with a persistent ductus arteriosus
113 and elevated pulmonary arterial pressure. Explant lung histology following heart-lung
114 transplantation revealed the presence of plexiform lesions. Two of the proband's offspring
115 were also diagnosed with childhood-onset PAH, one of which had an atrial septal defect.

116 The proband's parents, siblings and a third child showed no evidence of cardiovascular
117 disease.

118 Pathogenic variants in previously reported PAH disease genes

119 Our filtering strategy detected rare deleterious variation in previously reported PAH genes in
120 19.9% of the PAH cohort. For *BMP2*, rare heterozygous mutations were identified in 160 of
121 1048 cases (15.3%). The frequency of *BMP2* mutations in familial PAH was 75.9%, in
122 sporadic cases 12.2%, and 8.3% in anorexigen-exposed PAH cases. Forty-eight percent of
123 *BMP2* mutations were reported previously¹⁷, and the remainder were newly identified in
124 this study. Fourteen percent of *BMP2* mutations resulted in the deletion of larger protein-
125 coding regions ranging from 5 kb to 3.8 Mb in size. Supplementary Data 1 provides the
126 breakdown of *BMP2* SNVs and indels, and the larger deletions are shown in Figure 2a-c
127 with a detailed summary in Supplementary Table 2.

128

129 Of the other genes previously reported in PAH we identified deleterious heterozygous rare
130 variants in *ACVRL1* (9 cases, 0.9%), *ENG* (6 cases, 0.6%), *SMAD9* (4 cases, 0.4%), *KCNK3*
131 (4 cases, 0.4%), and *TBX4* (14 cases, 1.3%). We identified one case with highly deleterious
132 variants in both *BMP2* (p.Cys123Arg) and *SMAD9* (p.Arg294Ter). Details of consequence
133 types, deleteriousness and conservation scores, and minor allele frequencies are provided in
134 Supplementary Data 2. Fourteen cases (1.3%) with biallelic *EIF2AK4* mutations were
135 found¹⁹. No pathogenic coding variants in *CAV1*, *SMAD1* or *SMAD4* were identified.–Taken
136 together, rare causal variation in non-*BMP2* disease genes (*TBX4*, *ENG*, *ACVRL1*,
137 *SMAD9*, *KCNK3* and *EIF2AK4*) accounted for 4.7% of the entire PAH cohort. The clinical
138 characteristics of cases with variants in these previously reported genes are shown in
139 Supplementary Table 3.

140

141 In a case-control comparison of the frequencies of deleterious variants confined to the
142 previously reported PAH genes, we observed significant overrepresentation of rare variants
143 in *BMP2*, *TBX4*, *ACVRL1* and biallelic variants in *EIF2AK4* only ($P < 0.05$) (Supplementary
144 Table 4).

145 Identification of novel PAH disease genes

146 The strategy to identify novel causative genes in PAH employed a series of case-control
147 analyses (Figure 1b). To identify signals that might be masked by variants in previously
148 reported PAH genes, we excluded subjects with rare variants and deletions in *BMP2*,
149 *EIF2AK4*, *ENG*, *ACVRL1*, *TBX4*, *SMAD9* and *KCNK3*. A genome-wide comparison of
150 protein truncating variants (PTVs), representative of high impact variants, identified a higher
151 frequency of PTVs in *ATP13A3* (6 cases) ($P_{adj} = 0.0346$). Moreover, we identified additional
152 PTVs in several putative PAH genes, including *EVI5* (5 cases, 1 control) and *KDR* (4 cases,
153 0 controls) (Figure 3a), that require further validation to evaluate their contribution to PAH
154 pathogenesis (Supplementary Table 5).

155

156 We next analysed rare missense variants overrepresented in the PAH cohort, again
157 excluding subjects with variants in the previously reported PAH genes. This revealed
158 significant overrepresentation of rare variants in *GDF2* after correction for multiple testing
159 ($P_{adj} = 0.0023$), followed by *AQP1* (Figure 3b and Supplementary Table 6). Next, in a

160 combined analysis of rare missense variants and PTV, only *GDF2* remained significant ($P =$
161 0.001). Rare variants in additional putative genes occurred at higher frequency in cases
162 compared to controls, including *AQP1*, *ALPPL2*, *ATP13A3*, *OR8U1*, *IFT74*, *FLNA*, *SOX17*,
163 *ATP13A5*, *C3orf20* and *PIWIL1* (uncorrected P value < 0.0005), but were not significant after
164 correction for multiple testing (Figure 3c and Supplementary Table 7).

165
166 In order to increase power to detect rare associations, we deployed SKAT-O on filtered rare
167 PTVs and missense variants. Excluding previously reported genes, this analysis revealed an
168 association with rare variants in *AQP1* ($P_{adj} = 4.28 \times 10^{-6}$) and *SOX17* ($P_{adj} = 6.7 \times 10^{-5}$) (Figure
169 3d). *AQP1* and *SOX17* were both also nominally significant in the combined burden tests,
170 described above. Association was also found with rare variants in *MFRP* ($P_{adj} = 1.3 \times 10^{-5}$).
171 However, we consider *MFRP* a false-positive finding for reasons given in the Discussion.
172 Supplementary Table 8 shows the top 50 most significant genes identified by SKAT-O,
173 providing further candidates to be evaluated in future studies. Details of rare variants in
174 novel PAH genes (*GDF2*, *ATP13A3*, *AQP1*, *SOX17*) identified in cases are provided in
175 Supplementary Data 3.

176
177 Notably, a genome-wide assessment of larger structural variation did not identify any
178 additional large deletions after exclusion of subjects harbouring deletions in *BMP2* (Figure
179 2d-e).

180
181 The proportion of PAH cases with mutations in the new genes was 3.5%. The clinical
182 characteristics of PAH cases with mutations in these genes are provided in Supplementary
183 Table 3b. Of note, cases with mutations in *SOX17* and *AQP1* were significantly younger at
184 diagnosis (32.8 ± 16.2 years [$P = 0.002$] and 36.9 ± 14.3 years [$P = 0.013$], respectively)
185 compared to cases with no mutations in the previously established genes (51.7 ± 16.6
186 years).

187 Non-coding variation around PAH disease genes

188 An initial analysis for enrichment of variants in the non-coding sequence surrounding
189 previously reported and newly identified PAH disease genes, including upstream gene
190 regions, 5' UTRs, intronic sequence, 3' UTRs and downstream gene regions, did not detect
191 an significant overrepresentation in the PAH cohort. Details of the non-coding variants that
192 passed the filtering strategy are provided in Supplementary Data 4.

193 Independent validation and familial segregation analysis

194 To provide further validation of the potentially causal role of mutations in the new genes
195 identified, we examined whole-exome data from an independent UK family with three
196 affected individuals across two generations. Microsatellite genotyping across chromosome
197 2q33 had previously demonstrated non-sharing of haplotypes in affected individuals,
198 consistent with exclusion of linkage to the *BMP2* locus. No pathogenic variants were
199 identified in the protein-coding regions of the *BMP2* gene or other TGF- β pathway genes.
200 Analysis of exome sequence data from individual II-1 identified a novel heterozygous
201 c.411C>G (p.Y137*) PTV in the *SOX17* gene. Segregation analysis in the extended family
202 demonstrated that the mutation had arisen *de novo* in the affected father (II-1) and was

203 transmitted to the affected offspring (III-1). All unaffected family members were confirmed as
204 wild-type (Figure 4a).

205

206 Three HPAH subjects harbouring rare variants in *AQP1*, identified in the NIHR BR-RD WGS
207 study, were also selected for familial co-segregation analysis (Figure 4b-d). No pathogenic
208 variants in any of the previously reported genes were identified in these families. The first
209 pedigree comprised three affected individuals across two generations. Sanger sequencing
210 confirmed the presence of the heterozygous *AQP1* c.583C>T (p.R195W) missense variant
211 in the proband (E011942), the affected father (E011942.f) and the healthy younger paternal
212 uncle (E011942.u1). An additional unaffected uncle did not carry the *AQP1* variant. These
213 results indicate likely incomplete penetrance in the unaffected carrier, as observed in
214 *BMPR2* families²⁰. No additional clinical information was available for the deceased
215 grandparents (Figure 4b). The remaining two families comprised affected parent-offspring
216 individuals. By Sanger sequencing we independently confirmed a heterozygous *AQP1*
217 c.527T>A (p.Val176Glu) missense variant in proband (E012415) and his affected father
218 (Figure 4c), as well as a heterozygous *AQP1* c.583C>T (p.R195W) missense variant in
219 proband (E010634) and her affected father (Figure 4d). These results highlight recurrent
220 *AQP1* variation across unrelated families and demonstrate co-segregation with the
221 phenotype.

222 Predicted functional impact of variants in novel PAH genes

223 To evaluate the potential functional impact of rare variants identified in the likely causative
224 new genes we performed structural analysis of *GDF2*, *ATP13A3*, *AQP1*, and *SOX17*. In
225 addition we undertook a functional analysis of the *GDF2* variants identified.

226

227 Heterozygous mutations in *GDF2* exclusive to PAH cases comprised 1 frameshift variant
228 and 7 missense variants. *GDF2* encodes growth and differentiation factor 2, also known as
229 bone morphogenetic protein 9 (BMP9), the major circulating ligand for the endothelial
230 *BMPR2/ACVRL1* receptor complex²¹. Amino acid substitutions were assessed against the
231 published crystal structure²² of the prodomain bound form of *GDF2* (Figure 5). Variants
232 clustered at the interface between the prodomain and growth factor domain. Since the
233 prodomain is important for the processing of *GDF2*, it is likely that amino acid substitutions
234 reduce the stability of the prodomain-growth factor interface. In keeping with these
235 predictions, HEK293T cells transfected with *GDF2* variants exclusive to PAH cases,
236 demonstrated reduced secretion of mature *GDF2* into the cell supernatants (Figure 5d),
237 compared with wild type *GDF2*.

238

239 We identified 3 heterozygous frameshift variants, 2 stop gained, 2 splice region variants in
240 *ATP13A3*, which are predicted to lead to loss of ATPase catalytic activity (Figure 6a). In
241 addition, we identified 4 heterozygous likely pathogenic missense variants in PAH cases,
242 two near the conserved ATPase catalytic site and predicted to destabilise the conformation
243 of the catalytic domain (Figure 6b-d). The distribution of variants (Figure 6a) suggests that
244 these mutations impact critically on the function of the protein.

245

246 The majority of rare variants identified in *AQP1*, which encodes aquaporin 1, are situated
247 within the critical water channel (Figure 7). In particular the p.Arg195Trp variant, identified in
248 5 PAH cases, locates at the hydrophilic face of the pore. This arginine at position 195 helps

249 define the constriction region of the AQP1 pore structure and is conserved across the water
250 specific aquaporins²⁶. Rare variants in *SOX17*, included 4 nonsense variants (including the
251 PTV identified in the additional UK family) predicted to lead to loss of the beta-catenin
252 binding region, and 6 missense variants predicted to disrupt interactions with Oct4 and beta-
253 catenin^{27, 28} (Figure 8).

254
255 GDF2 is known to be secreted from the liver, but the cellular localization of proteins encoded
256 by the other novel genes is less well characterised. Thus we employed
257 immunohistochemistry to examine localisation in the normal and hypertensive human
258 pulmonary vasculature. Figure 9 shows that AQP1, ATP13A3 and SOX17 are predominantly
259 localised to the pulmonary endothelium in normal human lung and to endothelial cells within
260 plexiform lesions of patients with idiopathic PAH. In addition, we determined the relative
261 mRNA expression levels of *AQP1*, *ATP13A3* and *SOX17* in primary cultures of pulmonary
262 artery smooth muscle cells (PASMCs), pulmonary artery endothelial cells (PAECs) and
263 blood outgrowth endothelial cells (BOECs)³¹. AQP1 was expressed in PASMCs and
264 endothelial cells, with a trend towards higher levels in PASMCs (Figure 10a). ATP13A3 was
265 highly expressed in both cell types (Figure 10b), whereas SOX17 was almost exclusively
266 expressed in endothelial cells (Figure 10c). Although AQP1 and SOX17 are known to play
267 roles in endothelial function, the function of ATP13A3 in vascular cells is entirely unknown.
268 Thus, we determined the impact of ATP13A3 knockdown on proliferation and apoptosis of
269 BOECs. Loss of ATP13A3 led to marked inhibition of serum-stimulated proliferation of
270 BOECs, and increased apoptosis in serum-deprived conditions (Figure 10d-f).

271 Discussion

272 We report a comprehensive analysis of rare genetic variation in a large cohort of index cases
273 with idiopathic and heritable forms of PAH. Whilst we utilised WGS, the main goal was the
274 identification of rare causal variation underlying PAH in the protein coding sequence. The
275 approach involved a rigorous case-control comparison using a tiered search for variants.
276 First, we searched for high impact PTVs overrepresented in cases, having excluded
277 previously established PAH genes. This revealed PTVs in *ATP13A3*, a poorly characterised
278 P-type ATPase of the P5 subfamily³². There is little information regarding the function of the
279 ATPase, *ATP13A3*, which appears widely expressed in mouse tissues³². Although, the
280 precise substrate specificity is unknown, ATP13A3 plays a role in polyamine transport³³.
281 Based on available RNA sequencing data, *ATP13A3* is highly expressed in human
282 pulmonary vascular cells and cardiac tissue (<https://www.encodeproject.org>). We confirmed
283 that *ATP13A3* mRNA is expressed in primary cultured pulmonary artery smooth muscle cells
284 and endothelial cells, and provide preliminary data that loss of *ATP13A3* inhibits proliferation
285 and increases apoptosis of endothelial cells. These findings are consistent with the widely
286 accepted paradigm that endothelial apoptosis is a major trigger for the initiation of PAH^{34, 35}.
287 It will be of considerable interest to determine the role of *ATP13A3* in vascular cells and
288 whether it is functionally associated with BMP signalling, or represents a distinct therapeutic
289 target in PAH.

290
291 Analysis of missense variation, and a combined analysis of all predicted deleterious
292 variation, revealed that mutation at the *GDF2* gene is also significant determinant of
293 predisposition to PAH. Of the new genes identified, *GDF2* provides further evidence for the

294 central role of the BMP signalling pathway in PAH. *GDF2* encodes the major circulating
295 ligand for the endothelial BMPR2/ACVRL1 receptor complex²¹. Taken together, the genetic
296 findings suggest that a deficiency in *GDF2*/BMPR2/ACVRL1 signalling in pulmonary artery
297 endothelial cells is critical in PAH pathobiology. The majority of *GDF2* variants detected in
298 our adult-onset PAH cohort were heterozygous missense variants, in contrast to a previous
299 case report of childhood onset PAH due to a homozygous nonsense mutation³⁶. The finding
300 of causal *GDF2* variants in PAH cases, associated with reduced production of GDF2 from
301 cells, provides further support for investigating replacement of this factor as a therapeutic
302 strategy in PAH³⁷.

303

304 To maximise the assessment of rare variation in a case-control study design, we deployed
305 the SKAT-O test. This approach revealed a significant association of rare variation in the
306 aquaporin gene, *AQP1*, and the transcription factor encoded by *SOX17*. Of note, both *AQP1*
307 and *SOX17* were within the top 8 ranked genes in our combined PTV and missense burden
308 test analysis (Supplementary Table 7), providing further confidence in their causative
309 contribution to PAH.

310

311 Aquaporin 1 belongs to a family of membrane channel proteins that facilitate water transport
312 in response to osmotic gradients²⁶, and *AQP1* is known to promote endothelial cell migration
313 and angiogenesis³⁸. Thus, approaches that maintain or restore pulmonary endothelial
314 function could offer new therapeutic directions in PAH. Conversely, *AQP1* inhibition in
315 pulmonary artery smooth muscle cells ameliorated hypoxia-induced pulmonary hypertension
316 in mice³⁹, suggesting that further studies are required to determine the key cell type
317 impacted by *AQP1* mutations in human PAH, and the functional impact of these *AQP1*
318 variants on water transport. The demonstration of familial segregation of *AQP1* variants with
319 PAH provides further support for the potentially causal role of these mutations in disease.
320 However, we also identified an unaffected *AQP1* variant carrier consistent with reduced
321 penetrance, which is well described for other PAH genes, including *BMPR2*.

322

323 Although functional studies are required to confirm the mechanisms by which mutations in
324 *SOX17* cause PAH, this finding provides additional support for the vascular endothelium as
325 the major initiating cell type in this disorder. *SOX17* encodes the SRY-box containing
326 transcription factor 17, which plays a fundamental role in angiogenesis⁴⁰ and arteriovenous
327 differentiation⁴¹. Moreover, conditional deletion of *SOX17* in mesenchymal progenitors leads
328 to impaired formation of lung microvessels⁴². The demonstration of familial segregation of
329 the *SOX17* p.Y137* PTV with early onset PAH provides additional evidence for a causal role
330 for these variants in PAH. The co-existence of a patent ductus arteriosus in the index case
331 and an atrial septal defect (ASD) in one of the affected offspring is of interest and suggests
332 an association with congenital heart disease. Small ASDs are not uncommon in idiopathic
333 PAH, and a more detailed clinical phenotyping of *SOX17* mutation carriers will be required to
334 determine whether the presence of ASDs and other congenital heart abnormalities are more
335 common in carriers of these mutations.

336

337 Whilst the SKAT-O analysis also provided support for the *MFRP* gene, recessive bi-allelic
338 mutations in *MFRP* cause retinal degeneration and posterior microphthalmos⁴³. The
339 expression of *MFRP* transcripts is largely confined to the central nervous system⁴⁴ and the
340 majority of variants were present in the Genome Aggregation Database (GnomAD,

341 <http://gnomad.broadinstitute.org>). On the basis of these considerations, variants in *MFRP*
342 are unlikely to contribute to PAH aetiology.

343

344 This analysis provides new insights on the frequency and validity of previously reported
345 genes in PAH. We confirmed that mutations in *BMP2* are the most common genetic cause
346 and validated rare causal variants in *ACVRL1*, *ENG*, *SMAD9*, *TBX4*, *KCNK3* and *EIF2AK4*.
347 Although our findings question the validity of *CAV1*, *SMAD1* and *SMAD4* as causal genes,
348 previous reports might represent private mutations occurring in very rare families. The use of
349 WGS in this study allowed closer interrogation of larger deletions around the *BMP2* locus
350 than has been possible previously. Nevertheless, additional analyses are required to
351 determine the full impact of structural variation (inversions, duplications, smaller deletions) at
352 this and other loci.

353

354 The non-PAH cohort used in the case-control comparisons for this study comprised
355 individuals, or relatives of individuals, with other rare diseases recruited to the NIHR
356 Bioresource for Rare Diseases (NIHR BR-RD) in the UK (see Methods). In general, for very
357 rare causal variants, the comparison between PAH cases and non-PAH rare disease
358 controls should not reduce our ability to detect overrepresentation of rare variants in a
359 particular gene in the PAH cohort, if mutations in that gene are specific to PAH. However, if
360 rare variants in a gene were responsible for more than one phenotype, it is possible that this
361 would reduce the power to detect overrepresentation in the PAH cohort. For example, if
362 mutations occurred in different functional domains of the expressed protein, this might lead
363 to PAH if mutations affected one domain, but other phenotypes if they affected another
364 domain. Overcoming this potential limitation will require additional analysis of the functional
365 impact of variants and their distribution within a gene, and more detailed information on the
366 phenotypes of subjects in the non-PAH group.

367

368 Taken together, this study identifies rare sequence variation in new genes underlying
369 heritable forms of PAH, and provides a unique resource for future large-scale discovery
370 efforts in this disorder. Mutations in previously established genes accounted for 19.9% of
371 PAH cases. Including new genes identified in this study (*GDF2*, *ATP13A3*, *AQP1*, *SOX17*),
372 the total proportion of cases explained by mutations increased to 23.5%. It is likely that
373 independent confirmation of the expanded list of putative genes identified in this study will
374 increase further the proportion of cases explained by mutations, but this will require larger
375 international collaborations. The results suggest that the genetic architecture of PAH,
376 beyond mutations in *BMP2*, is characterised by substantial genetic heterogeneity and
377 consists of rare heterozygous coding region mutations shared by small numbers of cases.
378 The contribution of rare variation within non-coding regulatory regions to PAH aetiology
379 remains to be determined. This will require functional annotation of regulatory and other non-
380 coding regions specific for relevant cell types, further case-control analyses of these regions
381 and ultimately functional studies of gene regulation to assess the pathogenicity of non-
382 coding variants. Our findings to date provide support for a central role of the pulmonary
383 vascular endothelium in disease pathogenesis, and suggest new mechanisms that could be
384 exploited therapeutically in this life-limiting disease.

385

386 Methods

387 Ethics and patient selection

388 Cases were recruited from the UK National Pulmonary Hypertension Centres, Universite
389 Sud Paris (France), the VU University Medical Center Amsterdam (The Netherlands), the
390 Universities of Gießen and Marburg (Germany), San Matteo Hospital, Pavia (Italy), and
391 Medical University of Graz (Austria). All cases had a clinical diagnosis of idiopathic PAH,
392 heritable PAH, drug- and toxin-associated PAH, or PVOD/PCH established by their expert
393 centre. The non-PAH cohort for the case-control comparison comprised 6385 unrelated
394 subjects recruited to the NIHR BR-RD study. All PAH and non-PAH patients provided written
395 informed consent (UK Research Ethics Committee: 13/EE/0325), or local forms consenting
396 to genetic testing in deceased patients and non-UK cases. An additional UK family
397 diagnosed with HPAH was ascertained as described previously⁴⁵. Blood and saliva samples
398 were collected under written informed consent of the participants or their parents for use in
399 gene identification studies (UK Research Ethics Committee: 08/H0802/32).

400 Composition of non-PAH control cohort

401 The non-PAH control cohort consisted of subjects with bleeding, thrombotic and platelet
402 disorders (15.5%), cerebral small vessel disease (2.1%), Ehlers-Danlos syndrome (0.3%),
403 subjects recruited to Genomics England Ltd (19.8%), hypertrophic cardiomyopathy (3.6%),
404 intrahepatic cholestasis of pregnancy (4.1%), Leber hereditary optic neuropathy (0.9%),
405 multiple primary tumours (7.8%), neuropathic pain disorder (2.6%), primary immune
406 disorders (15.3%), primary membranoproliferative glomerulonephritis (2.3%), retinal
407 dystrophies/paediatric neurology and metabolic disease (19.8%), stem cell and myeloid
408 disorders (2.1%), steroid resistant nephrotic syndrome (3.6%), and others (0.3%), or their
409 first degree relatives.

410 High-throughput sequencing

411 DNA extracted from venous blood underwent whole genome sequencing using the Illumina
412 TruSeq DNA PCR-Free Sample Preparation kit (Illumina Inc., San Diego, CA, USA) and
413 Illumina HiSeq 2000 or HiSeq X sequencer, generating 100 - 150 bp reads with a minimum
414 coverage of 15X for ~95% of the genome (mean coverage of 35X). Whole-exome
415 sequencing was conducted for individual II-1 (Figure 4a) using genomic DNA extracted from
416 peripheral blood. Paired-end sequence reads were generated on an Illumina HiSeq 2000.

417 Generation of analysis-ready data sets

418 Sequencing reads were pre-processed by Illumina with Isaac Aligner and Variant Caller (v2,
419 Illumina Inc.) using human genome assembly GRCh37 as reference. Variants were
420 normalised, merged into multi-sample VCF files by chromosome using the gVCF
421 aggregation tool agg (<https://github.com/Illumina/agg>) and annotated with Ensembl's Variant
422 Effect Predictor (VEP). Following read alignment to the reference genome (GRCh37), variant
423 calling and annotation of whole-exome data for individual II:1 were performed using GATK
424 UnifiedGenotyper⁴⁶ and ANNOVAR⁴⁷, respectively. Annotations included minor allele
425 frequencies from other control data sets (i.e. ExAC⁴⁸, 1000 Genomes Project⁴⁹ and UK10K⁵⁰)
426 as well as deleteriousness and conservation scores (i.e. CADD⁵¹, SIFT⁵², PolyPhen-2⁵³ and

427 Gerp⁵⁴) enabling further filtering and assessment of the likely pathogenicity of variants. To
428 take forward only high quality calls, the pass frequency (proportion of samples containing
429 alternate alleles that passed the original variant filtering) and call rate (proportion of samples
430 with reference or alternate genotypes) were combined into the overall pass rate (OPR: pass
431 frequency x call rate) and variants with an OPR of 80% or higher were retained.

432 Estimation of ethnicity and relatedness

433 We estimated the population structure and relatedness based on a representative set of
434 SNPs using the R package GENESIS to perform PC-Air⁵⁵ and PC-Relate⁵⁶, respectively.
435 The selected 35,114 autosomal SNPs were present on Illumina genotyping arrays
436 (HumanCoreExome-12v1.1, HumanCoreExome-24v1.0, HumanOmni2.5-8v1.1), do not
437 overlap quality control excluded regions or multiallelic sites in the 1000 Genomes (1000G)
438 Phase 3 dataset⁴⁹, do not have any missing genotypes in NIHR BR-RD, had a MAF of 0.3 or
439 above and LD pruning was performed using PLINK⁵⁷ with a window size of 50 bp, window
440 shift of 5 bp and a variance inflation factor threshold of 2. The 2,110 samples from the
441 1000G Project including the European (EUR), African (AFR), South Asian (SAS) and East
442 Asian (EAS) populations (excluding the admixed American population) were filtered for the
443 selected SNPs and the filtered data were used to perform a principal component analysis
444 (PCA) using PC-Air. We modelled the scores of the leading five principal components as
445 data generated by a population specific multivariate Gaussian distribution and estimated the
446 corresponding mean and covariance parameters. Genotypes from the NIHR BR-RD samples
447 were projected onto the loadings for the leading five principal components from the 1000G
448 PCA and we computed the likelihood that each sample belonged to each subpopulation
449 under a mixture of multivariate Gaussians models. Each sample was allocated to the
450 population with the highest likelihood, unless the highest likelihood was similar to likelihood
451 values for other populations, as might be expected for example under admixed ancestry or if
452 the sample came from a population not included in 1000G. Such ambiguous samples were
453 labeled as “other”. PC-Relate was used to identify related individuals in NIHR BR-RD. We
454 used the first 20 PCs from PC-Air to adjust for relatedness and extracted the pairwise
455 Identity-By-State distances and kinship values. The pairwise information was used by
456 Primus to infer family networks and calculate the maximum set of unrelated samples.

457
458 Of the 9,110 NIHR BR-RD samples, we assigned 80.2% to Non-Finish European (n=7,307),
459 7.2% to South Asian (n=649), 2.3% to African (n=213), 0.08% to East Asian (n=78), 0.02%
460 to Finnish-European (n=19) and 9.2% to Other (n=844) and retrieved a maximum set of
461 7,493 unrelated individuals (UWGS10K), representing 82.2% of the entire NIHR BR-RD
462 cohort.

463 Cohort definition and allele frequency calculation

464 Based on the relatedness analysis, we defined the following sample subsets: (a) the
465 maximum number of unrelated non-PAH controls (UPAHC, n=6385), (b) all affected PAH
466 cases (PAHAFF, n=1048), and (c) all unrelated PAH index cases (PAHIDX, n=1038). These
467 subsets were used to annotate the variants in the multi-sample VCF file with calculated
468 minor allele frequencies using the fill-tags extension of BCFtools⁵⁸.

469 Rare variant filtering

470 Filtering of rare variants was performed as follows: 1) variants with a MAF less than 1 in
471 10,000 in UPAHC subjects, UK10K and ExAC were retained (adjusted for X chromosome
472 variants to 1 in 8,000); 2) variants with a combined annotation dependent depletion
473 deleteriousness (CADD) score of less than 15 were excluded. CADD scores were calculated
474 using the CADD web service (<http://cadd.gs.washington.edu>) for variants lacking a score; 3)
475 premature truncating variants (PTVs) or missense variants of the canonical transcript were
476 retained; 4) missense variants predicted to be both tolerated and benign by SIFT and
477 PolyPhen-2, respectively, were removed.

478
479 To identify likely causative mutations (as reported in Supplementary Table 3), variants in
480 previously reported and putative genes, identified in this study, were examined in more detail
481 to exclude variants that did not segregate in families (where data available). Furthermore,
482 variants shared between cases and non-PAH controls, as well as variants of uncertain
483 significance that co-occurred with previously reported causative mutations or high impact
484 PTVs were also excluded.

485 Burden analysis of protein-truncating and missense variants

486 Filtered variants were grouped per gene and consequence type (predicted PTV / missense)
487 and subjects with at least one variant were counted (no double counting) per group and
488 tested for association with disease. We applied a one-tailed Fisher's exact test with *post hoc*
489 Bonferroni correction to calculate the *P* value for genome-wide significance.

490 Rare variant analysis using SKAT-O

491 To further investigate the aggregated effect that rare variants contribute to PAH aetiology,
492 we applied a Sequence Kernel Association test (SKAT-O). SKAT-O increases the power of
493 discovery under different inheritance models by combining variance-component and burden
494 tests. Variants were filtered based on MAF as specified above, and only PTV and missense
495 variants were included. For the analysis we implemented SKAT-O in RvTests v1.9.9⁵⁹ with
496 default parameters and weights being Beta(1,25), and applying a correction for read length,
497 gender and the first five principal components of the ethnicity PCA. Variants were collapsed
498 considering only the protein-coding region in the canonical transcript of the protein-coding
499 genes in the genome assembly GRCh37.

500 Analysis of large deletions

501 Copy number variation was identified using Canvas⁶⁰ and Manta⁶¹. Deletions called by both
502 Manta and Canvas with a reciprocal overlap of $\geq 20\%$ were retained. Of these, deletions
503 were excluded if both failed standard Illumina quality metrics or overlapped with known
504 benign deletions in healthy cohorts⁶². Deletions with a reciprocal overlap of $\geq 50\%$ between
505 samples were merged and filtered for a frequency of less than 1 in 1,000 in WGS10K and
506 overlapping exonic regions of protein coding genes (GRCh37 genome assembly). The
507 number of subjects with deletions were added up by gene (no double counting of subjects)
508 and tested for association with the disease. We applied a one-tailed (greater) Fisher's exact
509 test with Bonferroni *post hoc* correction for multiple testing to determine the *P* values for
510 genome-wide significance.

511 Confirmation of variants

512 Variant sequencing reads for SNVs, indels and deletions were visualised for validation on
513 Integrative Genomes Viewer (IGV)¹⁸, and were confirmed by diagnostic capture-based high-
514 throughput sequencing, if the IGV inspection was not satisfactory. For the familial
515 segregation analysis, linkage to the *BMPR2* locus was first examined by microsatellite
516 genotyping analysis. Mutation screening of the *BMPR2*, *ACVRL1*, *ENG*, *AQP1* and *SOX17*
517 genes was conducted by capillary sequencing using BigDye Terminator v3.1 chemistry. All
518 DNA fragments were resolved on an ABI Fragment Analyzer (Applied Biosystems). All
519 primer sequences are listed in Supplementary Table 9. The family trees were drawn using
520 the R package FamAgg⁶³.

521 Structural analysis of novel variants

522 The domain structures and the functional groups of the novel PAH genes were plotted
523 according to the entry in UniProtKB. Clustal Omega was used for sequence alignment.
524 Structural data were obtained from RCSB Protein Data Bank and analysed according to
525 published reports. Figures were generated using PyMOL Molecular Graphics System.

526 Production of pGDF2 Wild Type and Variant Proteins

527 The cloning of human wild type pro-GDF2 (pGDF2) in pCEP4 has been described
528 previously⁶⁴. Site-directed mutagenesis was performed according to the manufacturer's
529 instructions (QuickChange Site Directed Mutagenesis Kit, Agilent Technologies). Mutations
530 were confirmed by Sanger sequencing. HEK-EBNA cells were transfected with plasmids
531 containing either wild-type or mutant pGDF2 for 14 hours. The transfecting supernatant was
532 removed and replaced with CDCHO media (Invitrogen) for 5 days to express the proteins.
533 The conditioned media containing GDF2 and the variants were harvested and snap-frozen
534 on dry-ice before being stored at -80°C. For each variant, conditioned media from three
535 independent transfections were collected for further characterisation.

536 GDF2 ELISA

537 High binding 96-well ELISA plates (Greiner, South Lanarkshire, UK) were coated with
538 0.2µg/well of mouse monoclonal anti-human GDF2 antibody (R&D Systems, Oxfordshire,
539 UK) in PBS (0.1M phosphate pH7.4, 0.137M NaCl, 2.7mM KCl, Sigma) overnight at 4°C in a
540 humidified chamber. Plates were washed with PBS containing 0.05% (v/v) Tween-20 (PBS-
541 T), followed by blocking with 1% bovine serum albumin in PBS-T (1% BSA/PBS-T) for 90min
542 at room temperature. Recombinant human GDF2 standards (1-3000pg/ml) or conditioned
543 media samples (100µl/well of 1:30, 1:100, 1:300, 1:1000, 1:3000 and 1:10000 dilutions)
544 were then added and incubated for 2h at room temperature. After washing, plates were then
545 incubated with 0.04µg/well biotinylated goat anti-human GDF2 (R&D Systems) in 1%
546 BSA/PBS-T for 2hr. Plates were washed, then incubated with ExtrAvidin(r)-Alkaline
547 phosphatase (Sigma) diluted 1:400 in 1% BSA/PBS-T for 90 min. Plates were washed with
548 PBS-T followed by water. The ELISA was developed with a colorimetric substrate
549 comprising 1mg/ml 4-Nitrophenyl phosphate disodium salt hexahydrate (Sigma) in 1M
550 Diethanolamine, pH9.8 containing 0.5mM MgCl₂. The assay was developed in the dark at
551 room temperature and the absorbance measured at 405nm.

552 Cell culture and treatments

553 Distal human pulmonary artery smooth muscle cells (PASMCs) were cultured from explants
554 dissected from lung resection specimens. Small pulmonary arterioles (0.5 to 2mm diameter)
555 were dissected and divided into small pieces before plating in T25 flasks. Explants were left
556 to adhere for 2 hours and then incubated in DMEM/20% FBS plus amino acids at 37°C in
557 95% air/5% CO₂ until PASMCs had formed confluent monolayers. Cells were then
558 trypsinized, and for subsequent passages cells were maintained in DMEM supplemented
559 with 10% FBS. The cellular phenotype of PASMCs was confirmed by positive
560 immunofluorescence staining with anti-smooth muscle specific alpha-actin (Clone IA4
561 Sigma-Aldrich; 1:100 dilution). The derivation of human tissues and cells was approved by
562 Papworth Hospital ethical review committee (Ref 08/H0304/56+5) and all subjects provided
563 informed and written consent.

564

565 Human blood outgrowth endothelial cells (BOECs) were derived from 40–80 ml of peripheral
566 venous blood isolated from healthy subjects. The study was approved by the
567 Cambridgeshire 3 Research Ethics Committee (Ref 11/EE/0297), and all subjects provided
568 informed and written consent. BOECs were cultured in 10% FBS supplemented with EGM-
569 2MV (Life Technologies, Carlsbad, CA). Cells were used between passages 4 and 8⁶⁵. The
570 endothelial phenotype of BOECs was determined by flow cytometry for expression of
571 endothelial surface markers, as described previously³¹. Cells were routinely tested to
572 exclude mycoplasma infection.

573

574 Human pulmonary artery endothelial cells (PAECs) were purchased from Lonza (Cat. No.
575 CC-2530; Basel, Switzerland). Cells were maintained in EGM-2 with 2% FBS (Lonza).
576 PAECs were used for experiments between passages 4 and 8. For experiments cells were
577 cultured in the presence of EBM-2 containing Antibiotic-Antimycotic (Invitrogen,
578 Renfrewshire, UK). Cells were routinely tested to exclude mycoplasma contamination.

579 RNA preparation and quantitative reverse transcription-PCR

580 Total RNA was extracted using RNeasy Mini Kit with DNase digestion (Qiagen, West
581 Sussex, UK), according to the manufacturer's instructions. cDNA was prepared from 1 µg of
582 RNA using High Capacity Reverse Transcriptase kit (Applied Biosystems, Foster City, CA).
583 Quantitative PCR reactions employed MicroAmp optical 96-well reaction plates (Applied
584 Biosystems). 50 ng µl⁻¹ cDNA was used with SYBR Green Jumpstart Taq Readymix
585 (Sigma-Aldrich), ROX reference dye (Invitrogen) using custom made sense and anti-sense
586 primers (all 200 nmol l⁻¹). Primers for human *ACTB* (encoding β-actin), *AQP1*, *ATP13A3*,
587 *B2M*, *HPRT* and *SOX17* were designed using PrimerBLAST
588 (<https://www.ncbi.nlm.nih.gov/tools/primer-blast/>) (Supplementary Table 9). Reactions were
589 amplified on a Quantstudio 6 Real-Time PCR system (Applied Biosystems). The relative
590 abundance of each target gene in different cell lines was compared using the equation $2^{-\text{CtGOI}-\text{Ct3HK}}$,
591 where Ct3HK corresponded to the arithmetic mean of the Cts for *ACTB*, *B2M* and
592 *HPRT* for each sample. For expression analysis of siRNA knockdown, the $2^{-(\Delta\Delta\text{Ct})}$ method
593 was used and fold expression determined relative to the DH1 control.

594 siRNA transfection

595 Prior to transfection, cells were preincubated in Opti-MEM-I reduced serum media
596 (Invitrogen) for 2h before transfection with 10nM siRNA that had been lipoplexed for 20 min
597 at RT with DharmaFECT1 (GE Dharmacon, Lafayette, CO). Cells were then incubated with
598 the siRNA/DharmaFECT1 complexes for 4h at 37°C before replaced by full growth media.
599 Cells were kept in growth media for 24h before further treatment. Knockdown efficiency was
600 confirmed by mRNA expression or immunoblotting. For proliferation assays, parallel RNA
601 samples were collected both on day0 and day6, confirming that *ATP13A3* expression was
602 reduced by >90% on Day 0 and still reduced by >70% at Day 6. For all other assays, parallel
603 RNA samples were collected on the day of the experiment to confirm knockdown, which was
604 >90%. The siRNAs used were oligos targeting *ATP13A3* (SASI_Hs02_00356805) from
605 Sigma-Aldrich and ON-TARGET plus non-targeting Pool (siCP; GE Dharmacon).

606 Flow cytometric apoptosis assay

607 BOECs were plated 150,000/well into 6-well plates and transfected with si*ATP13A3* or siCP
608 lipoplexed with DharmaFECT1. Cells were then serum-starved in EBM-2 (Lonza) containing
609 0.1% FBS and A/A for 8 hours before treating with EBM-2 and A/A containing either
610 0.1%FBS or 5%FBS for another 24 hours. Cells were then trypsinized and after washing
611 with PBS, stained using the FITC Annexin V Apoptosis Detection Kit I (BD Biosciences).
612 For each condition, dual-staining of 5µl FITC conjugated Annexin V and 5µl propidium iodide
613 (PI) were added and incubated at room temperature for 15 minutes. For the single staining
614 controls for compensation, either 5µl FITC Annexin V or 5µl PI was added into non-
615 transfected cells. All samples were analysed on BD Accuri™ C6 Plus platform (BD
616 Biosciences). Data were collected and analysed using FlowJo software, with AnnexinV⁺/PI
617 cells defined as early apoptotic (Treestar).
618

619 Caspase-Glo 3/7 assay

620 BOECs were seeded at a density of 150,000/well into 6-well plates and transfected with
621 si*ATP13A3* or siCP lipoplexed with DharmaFECT1. For each condition, cells were
622 trypsinized from 6-well plates and reseeded in triplicates into a 96-well plate at a density of
623 15,000-20,000/well and left to adhere overnight. Cells were quiesced in EBM-2 containing
624 0.1%FBS for 24h before treating with or without EBM-2 and A/A containing either 0.1%FBS
625 or 5%FBS for 16 hours. For measuring caspase activities, 100ul Caspase-Glo® 3/7 Reagent
626 (G8091 Promega) was added into each well, incubated and mixed on a plate shaker in the
627 dark for 30 minutes at room temperature. The lysates were transferred to a white-walled 96-
628 well plate and luminescence was read in a GloMax® luminometer (Promega).

629 Data availability

630 WGS data of PAH cases included in this manuscript and eligible for public release according
631 to the UK Research Ethics rules have been deposited in the European Genome-phenome
632 Archive (EGA) at the EMBL - European Bioinformatics Institute under accession number
633 EGAD00001003423 (<https://www.ebi.ac.uk/ega/studies/EGAD00001003423>).

634 Acknowledgements

635 The UK National Institute for Health Research BioResource (NIHRBR) and the BHF/MRC
636 UK National Cohort of Idiopathic and Heritable PAH made this study possible. We gratefully
637 acknowledge the participation of patients recruited to the NIHRBR. We thank the NIHR BR-
638 RD staff and co-ordination teams at the University of Cambridge, and the research nurses
639 and coordinators at the specialist pulmonary hypertension centres involved in this study. The
640 UK National Cohort of Idiopathic and Heritable PAH is supported by the NIHR BR-RD, the
641 British Heart Foundation (BHF) (SP/12/12/29836), the BHF Cambridge Centre of
642 Cardiovascular Research Excellence, the UK Medical Research Council (MR/K020919/1),
643 the Dinosaur Trust, BHF Programme grants to RCT (RG/08/006/25302) and NWM
644 (RG/13/4/30107), and the UK NIHR Cambridge Biomedical Research Centre. Funding for
645 whole-exome sequencing was provided through a Bart's Charity award (MGU0205) to RCT
646 and DvH. NWM is a BHF Professor and NIHR Senior Investigator. CH is a NIHR Rare
647 Disease Translational Research Collaboration Clinical PhD Fellow. LS is supported by the
648 Wellcome Trust Institutional Strategic Support Fund (204809/Z/16/Z) awarded to St.
649 George's, University of London. CJR is supported by a BHF Intermediate Basic Science
650 Research Fellowship (FS/15/59/31839). AL is supported by a BHF Senior Basic Science
651 Research Fellowship (FS/13/48/30453). We acknowledge the support of the Imperial NIHR
652 Clinical Research Facility, the Netherlands CardioVascular Research Initiative, the Dutch
653 Heart Foundation, Dutch Federation of University Medical Centres, the Netherlands
654 Organisation for Health Research and Development and the Royal Netherlands Academy of
655 Sciences. We also gratefully acknowledge Dr Claudia Cabrera in the NIHR Barts
656 Cardiovascular Biomedical Research Centre for bioinformatics support. We thank all the
657 patients and their families who contributed to this research and the Pulmonary Hypertension
658 Association (UK) for their support.

659 Author contributions

660 S.G., N.W.M. and W.H.O. conceived and designed the research. S.G., M.H, M.B. and C.H.
661 processed the data and performed the statistical analysis. S.G., M.H, M.B., C.H. and N.W.M.
662 drafted the manuscript. L.S., R.D.M. and R.C.T. conducted the *SOX17* familial segregation
663 analyses. W.L. performed the structural analysis of the rare variants. R.S. generated the
664 mutant cells. J.H., R.M.S, B.L. and P.D.U. conducted the functional experiments on the
665 novel disease genes. M.S. performed the immunohistochemistry for novel gene products.
666 L.C.D. helped with the assessment of pertinent findings. O.S. was involved with data
667 analysis. D.W. participated in DNA extraction, sample QC and plating. L.S, R.D.M, S.H,
668 M.A., C.J.R., W.H.O., N.S., A.L., R.C.T. and M.R.W. helped with data analysis and
669 interpretation and made critical revision of the manuscript for important intellectual content.
670 J.M.M., C.M.T. and K.Y. coordinated data collection. N.W.M. and W.H.O. handled the
671 funding for the study. All other authors were responsible for data acquisition and recruitment
672 of subjects to the study and helped to draft the final version of the manuscript.

673 Competing interests

674 The authors declare no competing financial and non-financial interests.

675 References

- 676 1. Wagenvoort CA. The pathology of primary pulmonary hypertension. *J Pathol.* 1970
677 Aug;101(4):Pi
- 678 2. McGoon MD, Benza RL, Escribano-Subias P, Jiang X, Miller DP, Peacock AJ, *et al.*
679 Pulmonary arterial hypertension: epidemiology and registries. *J Am Coll Cardiol.* 2013 Dec
680 24;62(25 Suppl):D51-9
- 681 3. McLaughlin VV, Shillington A, Rich S. Survival in primary pulmonary hypertension: the
682 impact of epoprostenol therapy. *Circulation.* 2002 Sep 17;106(12):1477-82
- 683 4. Lane KB, Machado RD, Pauciuolo MW, Thomson JR, Phillips JA 3rd, Loyd JE, *et al.*
684 Heterozygous germline mutations in *BMPR2*, encoding a TGF-beta receptor, cause familial
685 primary pulmonary hypertension. *Nat Genet.* 2000 Sep;26(1):81-4
- 686 5. Deng Z, Morse JH, Slager SL, Cuervo N, Moore KJ, Venetos G, *et al.* Familial primary
687 pulmonary hypertension (gene *PPH1*) is caused by mutations in the bone morphogenetic
688 protein receptor-II gene. *Am J Hum Genet.* 2000 Sep;67(3):737-44
- 689 6. Evans JD, Girerd B, Montani D, Wang XJ, Galie N, Austin ED, *et al.* *BMPR2* mutations
690 and survival in pulmonary arterial hypertension: an individual participant data meta-analysis.
691 *Lancet Respir Med.* 2016 Feb;4(2):129-37
- 692 7. Trembath RC, Thomson JR, Machado RD, Morgan NV, Atkinson C, Winship I, *et al.*
693 Clinical and molecular genetic features of pulmonary hypertension in patients with hereditary
694 hemorrhagic telangiectasia. *N Engl J Med.* 2001 Aug 2;345(5):325-34
- 695 8. Harrison RE, Berger R, Haworth SG, Tulloh R, Mache CJ, Morrell NW, *et al.* Transforming
696 growth factor-beta receptor mutations and pulmonary arterial hypertension in childhood.
697 *Circulation.* 2005 Feb 1;111(4):435-41
- 698 9. Nasim MT, Ogo T, Ahmed M, Randall R, Chowdhury HM, Snape KM, *et al.* Molecular
699 genetic characterization of SMAD signaling molecules in pulmonary arterial hypertension.
700 *Hum Mutat.* 2011 Dec;32(12):1385-9
- 701 10. Shintani M, Yagi H, Nakayama T, Saji T, Matsuoka R. A new nonsense mutation of
702 *SMAD8* associated with pulmonary arterial hypertension. *J Med Genet.* 2009 May;46(5):331-
703 7
- 704 11. Austin ED, Ma L, LeDuc C, Berman Rosenzweig E, Borczuk A, Phillips JA 3rd, *et al.*
705 Whole exome sequencing to identify a novel gene (*caveolin-1*) associated with human
706 pulmonary arterial hypertension. *Circ Cardiovasc Genet.* 2012 Jun;5(3):336-43
- 707 12. Ma L, Roman-Campos D, Austin ED, Eyries M, Sampson KS, Soubrier F, *et al.* A novel
708 channelopathy in pulmonary arterial hypertension. *N Engl J Med.* 2013 Jul 25;369(4):351-
709 361
- 710 13. Kerstjens-Frederikse WS, Bongers EM, Roofthoof MT, Leter EM, Douwes JM, Van Dijk
711 A, *et al.* *TBX4* mutations (small patella syndrome) are associated with childhood-onset
712 pulmonary arterial hypertension. *J Med Genet.* 2013 Aug;50(8):500-6
- 713 14. Eyries M, Montani D, Girerd B, Perret C, Leroy A, Lonjou C, *et al.* *EIF2AK4* mutations
714 cause pulmonary veno-occlusive disease, a recessive form of pulmonary hypertension. *Nat*
715 *Genet.* 2014 Jan;46(1):65-9
- 716 15. Best DH, Sumner KL, Austin ED, Chung WK, Brown LM, Borczuk AC, *et al.* *EIF2AK4*
717 mutations in pulmonary capillary hemangiomatosis. *Chest.* 2014 Feb;145(2):231-236
- 718 16. Abenhaim L, Moride Y, Brenot F, Rich S, Benichou J, Kurz X, *et al.* Appetite-suppressant
719 drugs and the risk of primary pulmonary hypertension. International Primary Pulmonary
720 Hypertension Study Group. *N Engl J Med.* 1996 Aug 29;335(9):609-16

721 17. Machado RD, Southgate L, Eichstaedt CA, Aldred MA, Austin ED, Best DH, *et al.*
722 Pulmonary Arterial Hypertension: A Current Perspective on Established and Emerging
723 Molecular Genetic Defects. *Hum Mutat.* 2015 Dec;36(12):1113-27
724 18. Thorvaldsdottir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV): high-
725 performance genomics data visualization and exploration. *Brief Bioinform.* 2013
726 Mar;14(2):178-92
727 19. Hadinnapola C, Bleda M, Haimel M, Screatton N, Swift A, Dorfmuller P, *et al.* Phenotypic
728 Characterization of EIF2AK4 Mutation Carriers in a Large Cohort of Patients Diagnosed
729 Clinically With Pulmonary Arterial Hypertension. *Circulation.* 2017 Nov 21;136(21):2022-
730 2033
731 20. Hamid R, Cogan JD, Hedges LK, Austin E, Phillips JA 3rd, Newman JH, *et al.*
732 Penetrance of pulmonary arterial hypertension is modulated by the expression of normal
733 BMPR2 allele. *Hum Mutat.* 2009 Apr;30(4):649-54
734 21. David L, Mallet C, Mazerbourg S, Feige JJ, Bailly S. Identification of BMP9 and BMP10
735 as functional activators of the orphan activin receptor-like kinase 1 (ALK1) in endothelial
736 cells. *Blood.* 2007 Mar 1;109(5):1953-61
737 22. Mi LZ, Brown CT, Gao Y, Tian Y, Le VQ, Walz T, *et al.* Structure of bone morphogenetic
738 protein 9 procomplex. *Proc Natl Acad Sci U S A.* 2015 Mar 24;112(12):3710-5
739 23. David L, Mallet C, Keramidas M, Lamande N, Gasc JM, Dupuis-Girod S, *et al.* Bone
740 morphogenetic protein-9 is a circulating vascular quiescence factor. *Circ Res.* 2008 Apr
741 25;102(8):914-22
742 24. Thever MD, Saier MH Jr. Bioinformatic characterization of p-type ATPases encoded
743 within the fully sequenced genomes of 26 eukaryotes. *J Membr Biol.* 2009 Jun;229(3):115-
744 30
745 25. Kanai R, Ogawa H, Vilsen B, Cornelius F, Toyoshima C. Crystal structure of a Na⁺-
746 bound Na⁺,K⁺-ATPase preceding the E1P state. *Nature.* 2013 Oct 10;502(7470):201-6
747 26. Sui H, Han BG, Lee JK, Walian P, Jap BK. Structural basis of water-specific transport
748 through the AQP1 water channel. *Nature.* 2001 Dec 20-27;414(6866):872-8
749 27. Sinner D, Rankin S, Lee M, Zorn AM. Sox17 and beta-catenin cooperate to regulate the
750 transcription of endodermal genes. *Development.* 2004 Jul;131(13):3069-80
751 28. Remenyi A, Lins K, Nissen LJ, Reinbold R, Scholer HR, Wilmanns M. Crystal structure of
752 a POU/HMG/DNA ternary complex suggests differential assembly of Oct4 and Sox2 on two
753 enhancers. *Genes Dev.* 2003 Aug 15;17(16):2048-59
754 29. Palasingam P, Jauch R, Ng CK, Kolatkar PR. The structure of Sox17 bound to DNA
755 reveals a conserved bending topology but selective protein interaction platforms. *J Mol Biol.*
756 2009 May 8;388(3):619-30
757 30. Jauch R, Aksoy I, Hutchins AP, Ng CK, Tian XF, Chen J, *et al.* Conversion of Sox17 into
758 a pluripotency reprogramming factor by reengineering its association with Oct4 on DNA.
759 *Stem Cells.* 2011 Jun;29(6):940-51
760 31. Toshner M, Dunmore BJ, McKinney EF, Southwood M, Caruso P, Upton PD, *et al.*
761 Transcript analysis reveals a specific HOX signature associated with positional identity of
762 human endothelial cells. *PLoS One.* 2014;9(3):e91334
763 32. Schultheis PJ, Hagen TT, O'Toole KK, Tachibana A, Burke CR, McGill DL, *et al.*
764 Characterization of the P5 subfamily of P-type transport ATPases in mice. *Biochem Biophys*
765 *Res Commun.* 2004 Oct 22;323(3):731-8
766 33. Madan M, Patel A, Skruber K, Geerts D, Altomare DA, Iv OP. ATP13A3 and caveolin-1
767 as potential biomarkers for difluoromethylornithine-based therapies in pancreatic cancers.
768 *Am J Cancer Res.* 2016;6(6):1231-52

769 34. Taraseviciene-Stewart L, Kasahara Y, Alger L, Hirth P, Mc Mahon G, Waltenberger J, *et*
770 *al.* Inhibition of the VEGF receptor 2 combined with chronic hypoxia causes cell death-
771 dependent pulmonary endothelial cell proliferation and severe pulmonary hypertension.
772 *FASEB J.* 2001 Feb;15(2):427-38
773 35. Teichert-Kuliszewska K, Kutryk MJ, Kuliszewski MA, Karoubi G, Courtman DW, Zucco L,
774 *et al.* Bone morphogenetic protein receptor-2 signaling promotes pulmonary arterial
775 endothelial cell survival: implications for loss-of-function mutations in the pathogenesis of
776 pulmonary hypertension. *Circ Res.* 2006 Feb 3;98(2):209-17
777 36. Wang G, Fan R, Ji R, Zou W, Penny DJ, Varghese NP, *et al.* Novel homozygous BMP9
778 nonsense mutation causes pulmonary arterial hypertension: a case report. *BMC Pulm Med.*
779 2016 Jan 22;16:17
780 37. Long L, Ormiston ML, Yang X, Southwood M, Graf S, Machado RD, *et al.* Selective
781 enhancement of endothelial BMPR-II with BMP9 reverses pulmonary arterial hypertension.
782 *Nat Med.* 2015 Jul;21(7):777-85
783 38. Saadoun S, Papadopoulos MC, Hara-Chikuma M, Verkman AS. Impairment of
784 angiogenesis and cell migration by targeted aquaporin-1 gene disruption. *Nature.* 2005 Apr
785 7;434(7034):786-92
786 39. Schuoler C, Haider TJ, Leuenberger C, Vogel J, Ostergaard L, Kwapiszewska G, *et al.*
787 Aquaporin 1 controls the functional phenotype of pulmonary smooth muscle cells in hypoxia-
788 induced pulmonary hypertension. *Basic Res Cardiol.* 2017 May;112(3):30
789 40. Matsui T, Kanai-Azuma M, Hara K, Matoba S, Hiramatsu R, Kawakami H, *et al.*
790 Redundant roles of Sox17 and Sox18 in postnatal angiogenesis in mice. *J Cell Sci.* 2006
791 Sep 1;119(Pt 17):3513-26
792 41. Corada M, Orsenigo F, Morini MF, Pitulescu ME, Bhat G, Nyqvist D, *et al.* Sox17 is
793 indispensable for acquisition and maintenance of arterial identity. *Nat Commun.* 2013;4:2609
794 42. Lange AW, Haitchi HM, LeCras TD, Sridharan A, Xu Y, Wert SE, *et al.* Sox17 is required
795 for normal pulmonary vascular morphogenesis. *Dev Biol.* 2014 Mar 1;387(1):109-20
796 43. Sundin OH, Leppert GS, Silva ED, Yang JM, Dharmaraj S, Maumenee IH, *et al.* Extreme
797 hyperopia is the result of null mutations in MFRP, which encodes a Frizzled-related protein.
798 *Proc Natl Acad Sci U S A.* 2005 Jul 5;102(27):9553-8
799 44. The Genotype-Tissue Expression (GTEx) project. *Nat Genet.* 2013 Jun;45(6):580-5
800 45. Machado RD, Pauciulo MW, Thomson JR, Lane KB, Morgan NV, Wheeler L, *et al.*
801 BMPR2 haploinsufficiency as the inherited molecular mechanism for primary pulmonary
802 hypertension. *Am J Hum Genet.* 2001 Jan;68(1):92-102
803 46. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, *et al.* A framework
804 for variation discovery and genotyping using next-generation DNA sequencing data. *Nat*
805 *Genet.* 2011 May;43(5):491-8
806 47. Yang H, Wang K. Genomic variant annotation and prioritization with ANNOVAR and
807 wANNOVAR. *Nat Protoc.* 2015 Oct;10(10):1556-66
808 48. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, *et al.* Analysis of
809 protein-coding genetic variation in 60,706 humans. *Nature.* 2016 Aug 18;536(7616):285-91
810 49. Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, *et al.* A global
811 reference for human genetic variation. *Nature.* 2015 Oct 1;526(7571):68-74
812 50. Walter K, Min JL, Huang J, Crooks L, Memari Y, McCarthy S, *et al.* The UK10K project
813 identifies rare variants in health and disease. *Nature.* 2015 Oct 1;526(7571):82-90
814 51. Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, Shendure J. A general framework
815 for estimating the relative pathogenicity of human genetic variants. *Nat Genet.* 2014
816 Mar;46(3):310-5

817 52. Ng PC, Henikoff S. Predicting deleterious amino acid substitutions. *Genome Res.* 2001
818 May;11(5):863-74
819 53. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, *et al.* A
820 method and server for predicting damaging missense mutations. *Nat Methods.* 2010
821 Apr;7(4):248-9
822 54. Cooper GM, Stone EA, Asimenos G, Green ED, Batzoglou S, Sidow A. Distribution and
823 intensity of constraint in mammalian genomic sequence. *Genome Res.* 2005 Jul;15(7):901-
824 13
825 55. Conomos MP, Miller MB, Thornton TA. Robust inference of population structure for
826 ancestry prediction and correction of stratification in the presence of relatedness. *Genet*
827 *Epidemiol.* 2015 May;39(4):276-93
828 56. Conomos MP, Reiner AP, Weir BS, Thornton TA. Model-free Estimation of Recent
829 Genetic Relatedness. *Am J Hum Genet.* 2016 Jan 7;98(1):127-48
830 57. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, *et al.* PLINK: a
831 tool set for whole-genome association and population-based linkage analyses. *Am J Hum*
832 *Genet.* 2007 Sep;81(3):559-75
833 58. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, *et al.* The Sequence
834 Alignment/Map format and SAMtools. *Bioinformatics.* 2009 Aug 15;25(16):2078-9
835 59. Zhan X, Hu Y, Li B, Abecasis GR, Liu DJ. RVTESTS: an efficient and comprehensive
836 tool for rare variant association analysis using sequence data. *Bioinformatics.* 2016 May
837 1;32(9):1423-6
838 60. Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Kallberg M, *et al.* Manta:
839 rapid detection of structural variants and indels for germline and cancer sequencing
840 applications. *Bioinformatics.* 2016 Apr 15;32(8):1220-2
841 61. Roller E, Ivakhno S, Lee S, Royce T, Tanner S. Canvas: versatile and scalable detection
842 of copy number variants. *Bioinformatics.* 2016 Aug 1;32(15):2375-7
843 62. Zarrei M, MacDonald JR, Merico D, Scherer SW. A copy number variation map of the
844 human genome. *Nat Rev Genet.* 2015 Mar;16(3):172-83
845 63. Rainer J, Taliun D, D'Elia Y, Pattaro C, Domingues FS, Weichenberger CX. FamAgg: an
846 R package to evaluate familial aggregation of traits in large pedigrees. *Bioinformatics.* 2016
847 May 15;32(10):1583-5
848 64. Wei Z, Salmon RM, Upton PD, Morrell NW, Li W. Regulation of bone morphogenetic
849 protein 9 (BMP9) by redox-dependent proteolysis. *J Biol Chem.* 2014 Nov 7;289(45):31150-
850 9
851 65. Ormiston ML, Toshner MR, Kiskin FN, Huang CJ, Groves E, Morrell NW, *et al.*
852 Generation and Culture of Blood Outgrowth Endothelial Cells from Human Peripheral Blood.
853 *J Vis Exp.* 2015 Dec 23;(106):e53384
854

855 Figure legends

856 **Figure 1** Flow diagrams illustrating **(a)** the composition of the NIHR BioResource – Rare
857 Diseases PAH study and **(b)** the analysis strategy to identify novel PAH disease genes. **(a)**
858 The study comprised 1048 adult cases (aged 16 or over) attending specialist pulmonary
859 hypertension centres from the UK (n=731), and additional cases from France (n=142), The
860 Netherlands (n=45), Germany (n=82) and Italy (n=48). **(b)** A series of case-control
861 comparisons including and excluding cases with variants in previously reported disease
862 genes were undertaken using complementary filtering strategies.

863

864 **Figure 2** Analysis of copy number deletions. **(a)** Deletions affecting the *BMPR2* locus in 23
865 PAH cases. Genes are indicated in orange and labelled with their respective gene symbol.
866 Deletions are drawn as blue boxes above the genome axis (grey) showing the genomic
867 position on chromosome 2. The grey box highlights the location of *BMPR2*. **(b)** Locus zoom
868 on *BMPR2* highlighting the focal deletions affecting one or more exons. **(c)** WGS coverage
869 profiles of a selected set of smaller and larger deletions, visualised with the Integrative
870 Genomics Viewer (IGV)¹⁸, with deletions highlighted by red bars. **(d)** and **(e)** Manhattan plots
871 of the genome-wide case-control comparison of large deletions. In **(d)** all subject are
872 considered. In **(e)** subject with larger deletions affecting the *BMPR2* locus are excluded. The
873 adjusted *P* value threshold of 5×10^{-8} for genome-wide significance is indicated by the red
874 line.

875

876 **Figure 3** Manhattan plots of the rare variant analyses, having excluded cases carrying rare
877 variants in previously established PAH genes. Filtered variants were grouped per gene. We
878 tested for an excess of variants in PAH cases within genes using Fisher's exact test. The
879 negative decadic logarithm of unadjusted or adjusted *P*-values are plotted against the
880 chromosomal location of each gene. **(a)** Burden test of rare PTVs. **(b)** Burden test of rare
881 deleterious missense variants. **(c)** Burden test combining rare PTVs and likely deleterious
882 missense variants. **(d)** SKAT-O test of rare PTVs and missense variants.

883

884 **Figure 4** Pedigree structures and analysis of familial transmission of variants in *AQP1* and
885 *SOX17*. **(a)** Individual II.1 harbours a heterozygous *de novo* *SOX17* c.411C>G (p.Y137*)
886 PTV resulting in a premature termination codon, which has been transmitted to the affected
887 male (III.1). No unaffected family members carry the variant. No sample was available from
888 subject III.2. **(b)** Proband E011942 has inherited a heterozygous *AQP1* c.583C>T
889 (p.R195W) missense variant from her affected father. No sample was available from the
890 affected sister of the proband. The younger healthy uncle of the index case also carries the
891 *AQP1* variant. No samples or further clinical information was available for the grandparents,
892 who were not known to have cardiopulmonary disease. **(c)** Both the proband E012415 and
893 her father are affected and carry the rare *AQP1* c.527T>A (p.V176E) missense variant.
894 There was no further information available about the siblings of the father. **(d)** Subject
895 E010634 has inherited the heterozygous *AQP1* c.583C>T (p.R195W) missense variant from
896 her affected father. No rare variants in previously reported PAH genes were identified in any
897 of these families. Index cases are highlighted in red. yo: years old, mo: months old, d.:
898 death.

899

900 **Figure 5** Structural analysis of *GDF2* mutations. **(a)** Schematic diagram of *GDF2*
901 processing. The pre-pro-protein is processed into the mature growth factor domain (GFD)

902 bound to the prodomain upon secretion²³. **(b)** Plot of *GDF2* mutations found only in PAH
903 cases superimposed on the structure of prodomain bound *GDF2* (PDB: 4YCG
904 [<http://dx.doi.org/10.2210/pdb4YCG/pdb>])²². The *GDF2* growth factor domain is shown in
905 green and the prodomain in cyan. **(c)** Magnified view of the Arg110 and Glu143 mutations.
906 The wild type amino acids make double salt bridges to stabilise the prodomain conformation
907 at the interface between the growth factor domain and prodomain. The E143K and R110W
908 mutations both disrupt these interactions, destabilising the interaction between the growth
909 factor domain and prodomain. **(d)** *GDF2* levels secreted into supernatants of HEK293T cells
910 transfected with likely pathogenic variants found in PAH cases, compared with wild type
911 *GDF2* and cells transfected with an empty vector. *** $P < 0.001$ by ANOVA.
912

913 **Figure 6** Structural analysis of *ATP13A3* mutations. **(a)** Topology of *ATP13A3*, plotted
914 according to UniProtKB Q9H7F0 [<http://www.uniprot.org/uniprot/Q9H7F0>]. Frameshift and
915 stop-gained mutations identified in PAH cases are shown as khaki circles, and missense
916 mutations as red circles. Frameshift/stop-gained mutations are predicted to truncate the
917 protein prior to the catalytic domain and essential Mg binding sites, leading to loss of
918 ATPase activity. **(b)** Sequence alignment of *ATP13A3* with *ATP1A1* (P05024
919 [<http://www.uniprot.org/uniprot/P05024>]), of which the high resolution structure was used for
920 the structural analysis in **(c)**. The conserved regions of *ATP13A3* and *ATP1A1*, essential for
921 ATPase activity²⁴, show good alignment (data not shown). Only regions containing the
922 missense PAH mutations are shown, with positions of the four missense mutations
923 highlighted in yellow above the sequences. **(c)** Structural analysis of the 4 PAH missense
924 mutations plotted on the *ATP1A1* crystal structure based on the sequence alignment in **(b)**
925 (PDB: 3wgu [<http://dx.doi.org/10.2210/pdb3WGU/pdb>])²⁵. Green: α subunit (P05024
926 [<http://www.uniprot.org/uniprot/P05024>]), cyan: β subunit (P05027
927 [<http://www.uniprot.org/uniprot/P05027>]), grey: γ -subunit transcript variant a (Q58k79
928 [<http://www.uniprot.org/uniprot/Q58k79>]). Y535, Y677, R685 and I787 are the numbering in
929 *ATP1A1*. Positions of the four missense mutations found in PAH are labelled and highlighted
930 by red circles. **(d)** Magnified view of the cytoplasmic region of the ATPase, showing the
931 presence of ADP at the active site. The conserved regions essential for ATPase activity are
932 shown in light pink. The L675V and R858H mutations are located close to the ATP catalytic
933 region.
934

935 **Figure 7** Structural analysis of *AQP1* mutations. **(a)** Multiple sequence alignment of human
936 *AQP1* with seven other mammals. The bovine *AQP1* has the high resolution (2.2Å)
937 published structure. Mutations identified in PAH cases are highly conserved and highlighted
938 in yellow. **(b)** Crystal structure of bovine *AQP1* (PDB: 1j4n
939 [<http://dx.doi.org/10.2210/pdb1J4N/pdb>])²⁶. Left: side view; right: top view from the
940 extracellular direction. *AQP1* is shown as a semi-transparent cartoon and five water
941 molecules in the water channel are shown as red spheres. Key residues lining the water
942 channels are represented with stick structures. **(c)** Magnified view of the water channel, with
943 H-bonds connected to water molecules in the channel highlighted. Two asparagine-proline-
944 alanine (NPA) motifs, essential for the water transporting function of *AQP1*, are shown in
945 magenta. Conserved His180 that constricts the water channel is shown in yellow. Mutations
946 found in PAH cases, Arg195Trp and Val176Glu, are labelled and shown as orange stick
947 structures. Arg195 and His180 are highly conserved in the known water channels and are
948 strong indicators of water channel specificity. Arg195Trp and Val176Glu mutations are
949 predicted to disrupt the conformation of this conserved water channel.

950

951 **Figure 8** Structural analysis of *SOX17* mutations. **(a)** Schematic diagram of human *SOX17*
952 (Q9H6I2 [<http://www.uniprot.org/uniprot/Q9H6I2>]), based on UniProtKB annotation, and
953 published reports²⁷. Red arrows indicate PTVs and black arrows indicate missense
954 mutations identified in PAH patients. The blue bar illustrates the region that is covered in the
955 crystal structure (PDB: 3F27 [<http://dx.doi.org/10.2210/pdb3F27/pdb>])²⁹. The ability of
956 *SOX17* to activate transcription of target genes correlates with binding to β -catenin²⁷. As
957 illustrated, all PTVs lead to a loss of the β -catenin binding region. Two missense mutations
958 are located within and very close to the minimum β -catenin binding regions, and both are
959 highly conserved, indicating they are likely to be important for β -catenin binding. **(b)**
960 Structural analysis of HMG domain missense mutations found in PAH patients. Left,
961 Superposition of *SOX17*/DNA structure (*Sox17*: cyan, DNA: grey)²⁹ onto *SOX2*/DNA/*Oct1*
962 structure (PDB: 1GT0 [<http://dx.doi.org/10.2210/pdb1GT0/pdb>], *Sox2*: yellow, *Oct1*:
963 magenta, DNA: light blue)²⁸. Right: Magnified view of the interactions around Arg140 in the
964 *SOX2*/DNA/*Oct* structure. Arg140 in *SOX2* makes multiple H-bond interactions and mutating
965 this Arg in *SOX2* abolishes the interaction with transcription factors Pax6 and Oct4²⁸. *SOX2*
966 and *SOX17* both bind to Oct4³⁰ and *SOX17* K122E mutant can replace *SOX2* in maintaining
967 stem cell pluripotency³⁰, indicating this region in *SOX17* may interact with Oct4, similar to
968 *SOX2*. The three missense mutations in *SOX17* will likely disrupt interaction with Oct4. **(c)**
969 Supporting the analysis in **(b)**, sequence alignment shows that the HMG domain of *SOX2*
970 (P48431 [<http://www.uniprot.org/uniprot/P48431>]) and *SOX17* as well as *SOX8* (P57073
971 [<http://www.uniprot.org/uniprot/P57073>]) and *SOX18* (P35713
972 [<http://www.uniprot.org/uniprot/P35713>]) share high sequence identity and the three
973 mutations found in PAH (highlighted in yellow) are highly conserved emphasising their
974 functional importance. Similarly, the Gly and Thr that interact with Arg140 in *SOX2*
975 (highlighted in yellow) are also conserved between *Oct1* (PO2F1) and *Oct4* (PO5F1).

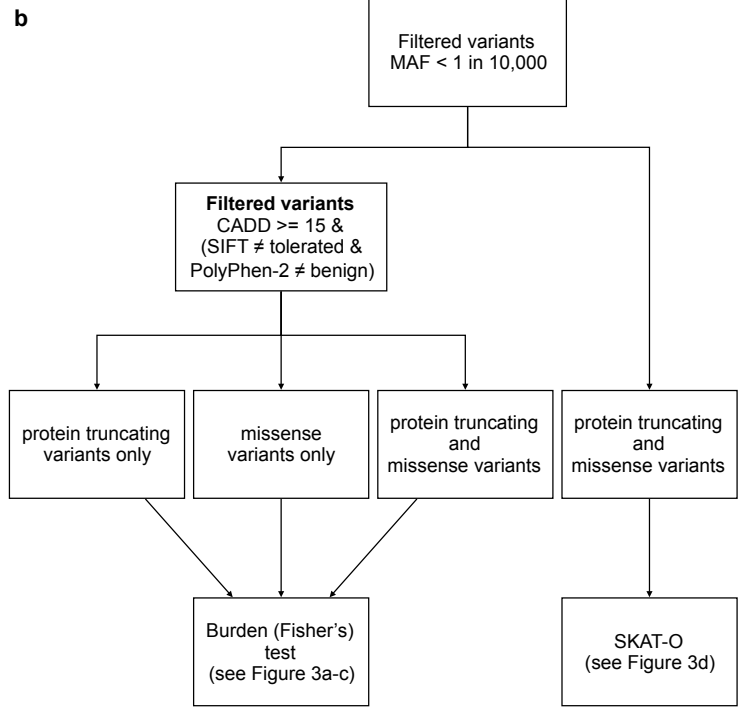
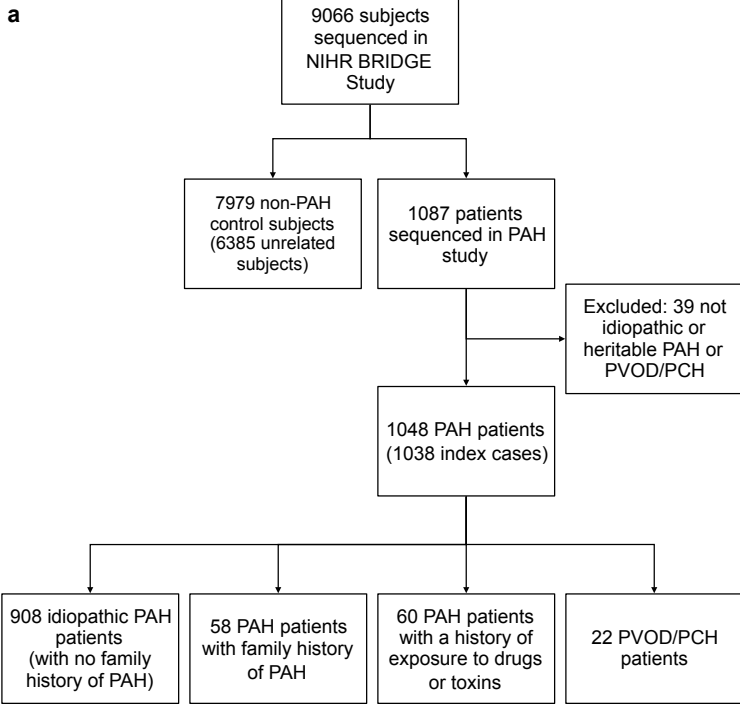
976

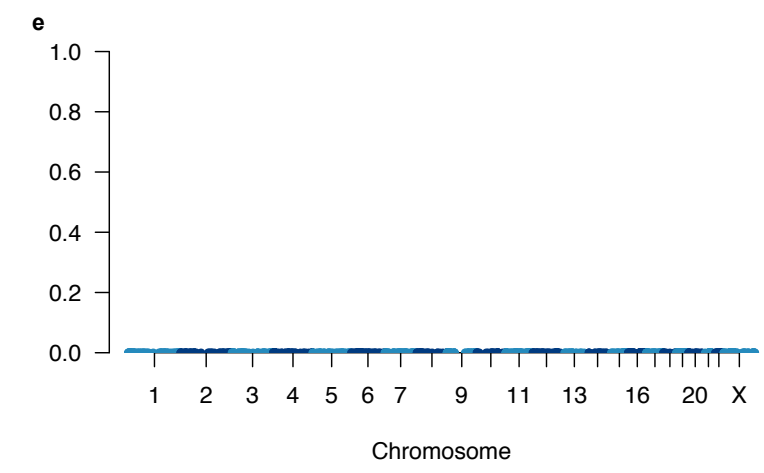
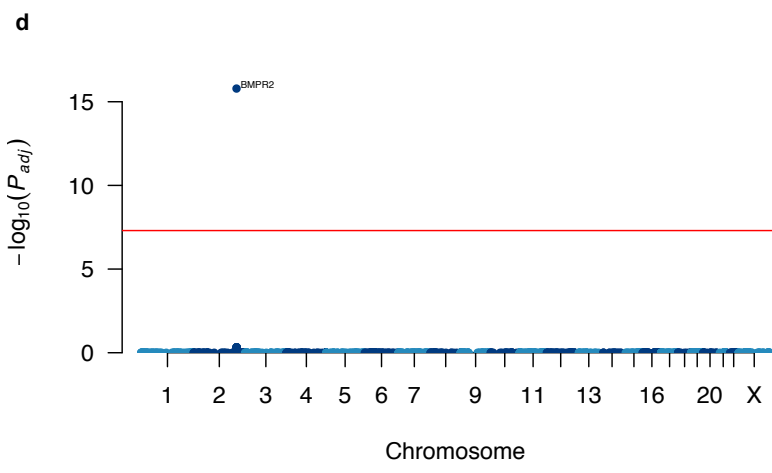
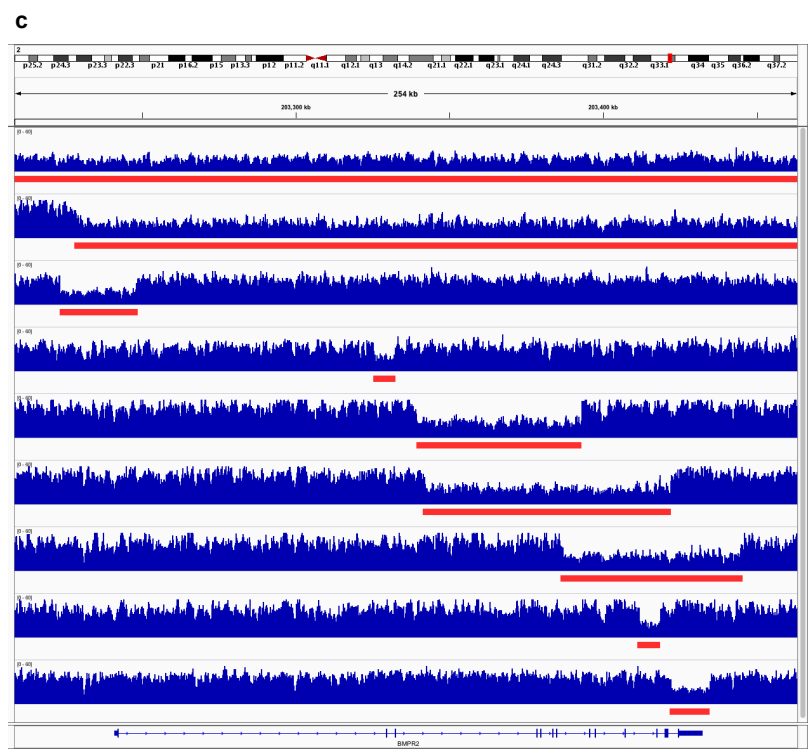
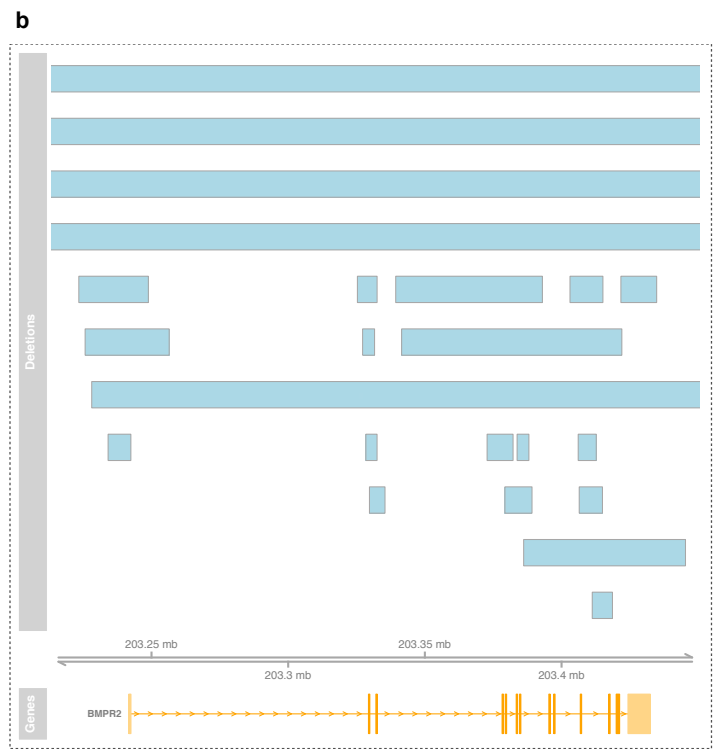
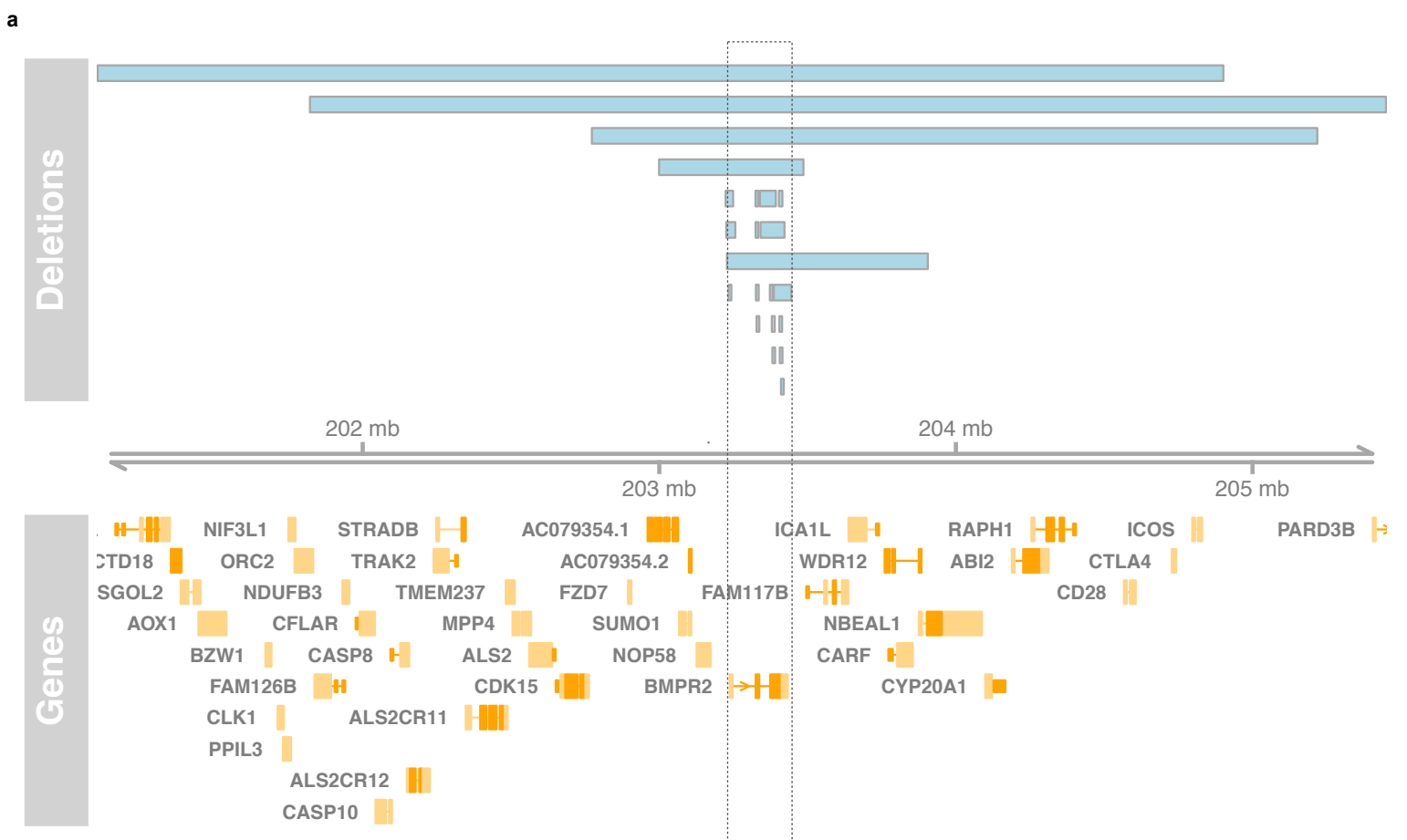
977 **Figure 9** Immunolocalisation of AQP1, ATP13A3 and *SOX17* in normal and PAH lung. The
978 typical histological findings (haematoxylin and eosin staining) of concentric vascular lesions
979 with associated plexiform lesions are shown **(a)**. Higher magnification images of plexiform
980 lesion **(b)**, with frequent endothelialised channels **(c)**; anti-CD31) surrounded by
981 myofibroblasts **(d)**; anti-SM α). Additional high magnification images demonstrating
982 endothelial expression of ATP13A3 **(e)**, AQP1 **(f)** and *SOX17* **(g)** in PAH lung. Controls lung
983 sections demonstrating predominantly endothelial expression of ATP13A3 **(h)**, AQP1 **(i)** and
984 *SOX17* **(j)**. (Scale bars = 50 μ m).

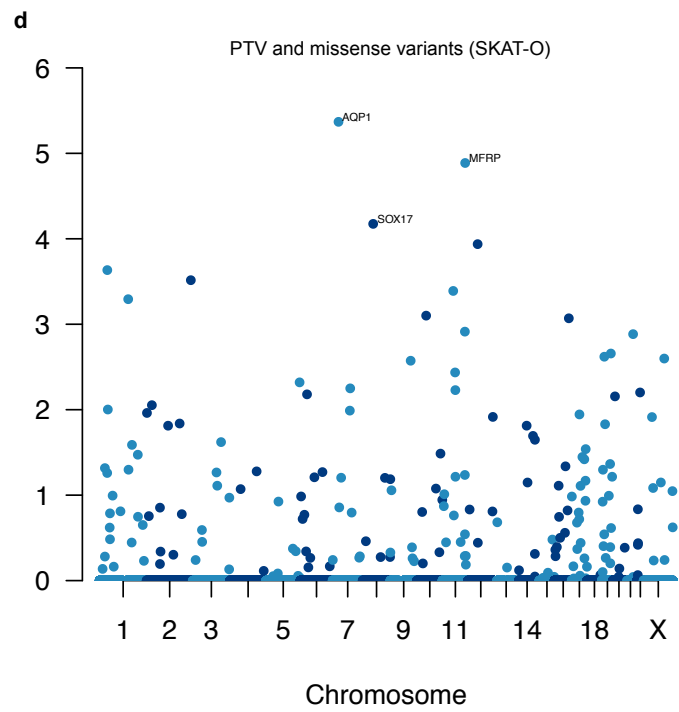
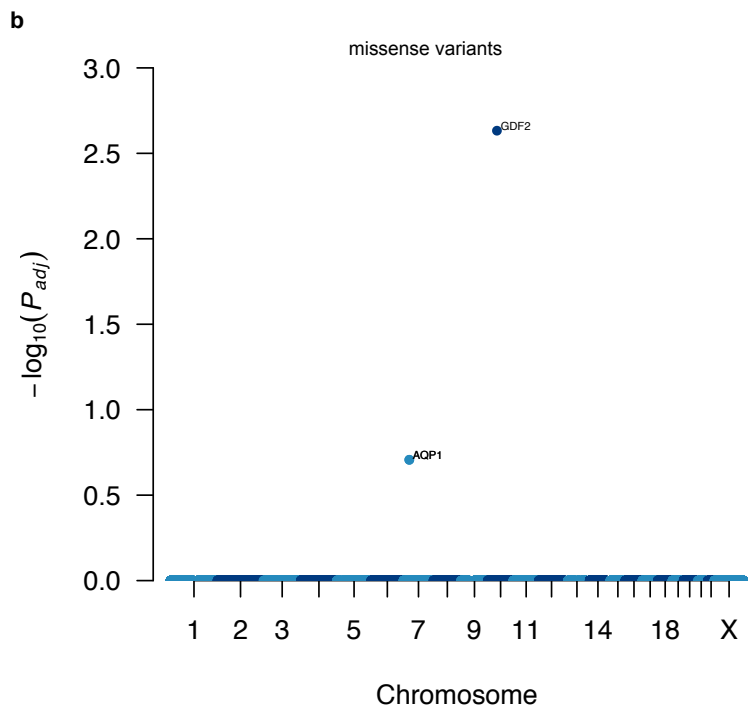
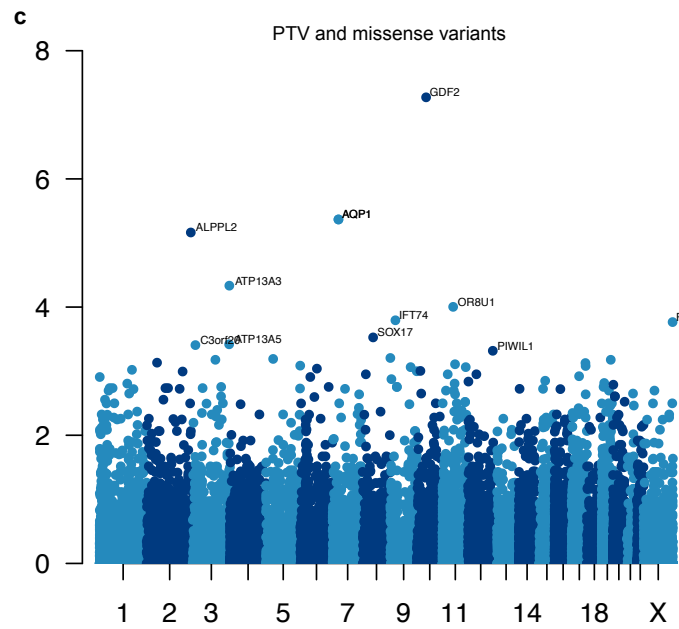
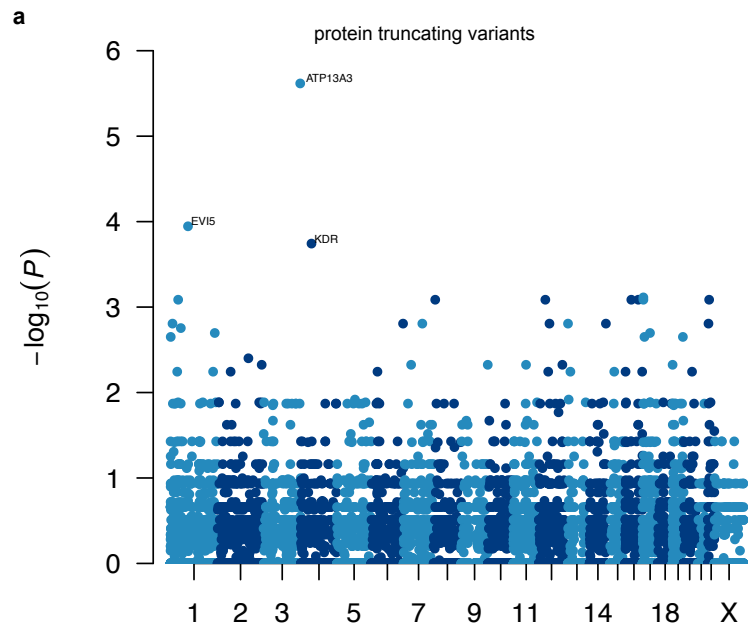
985

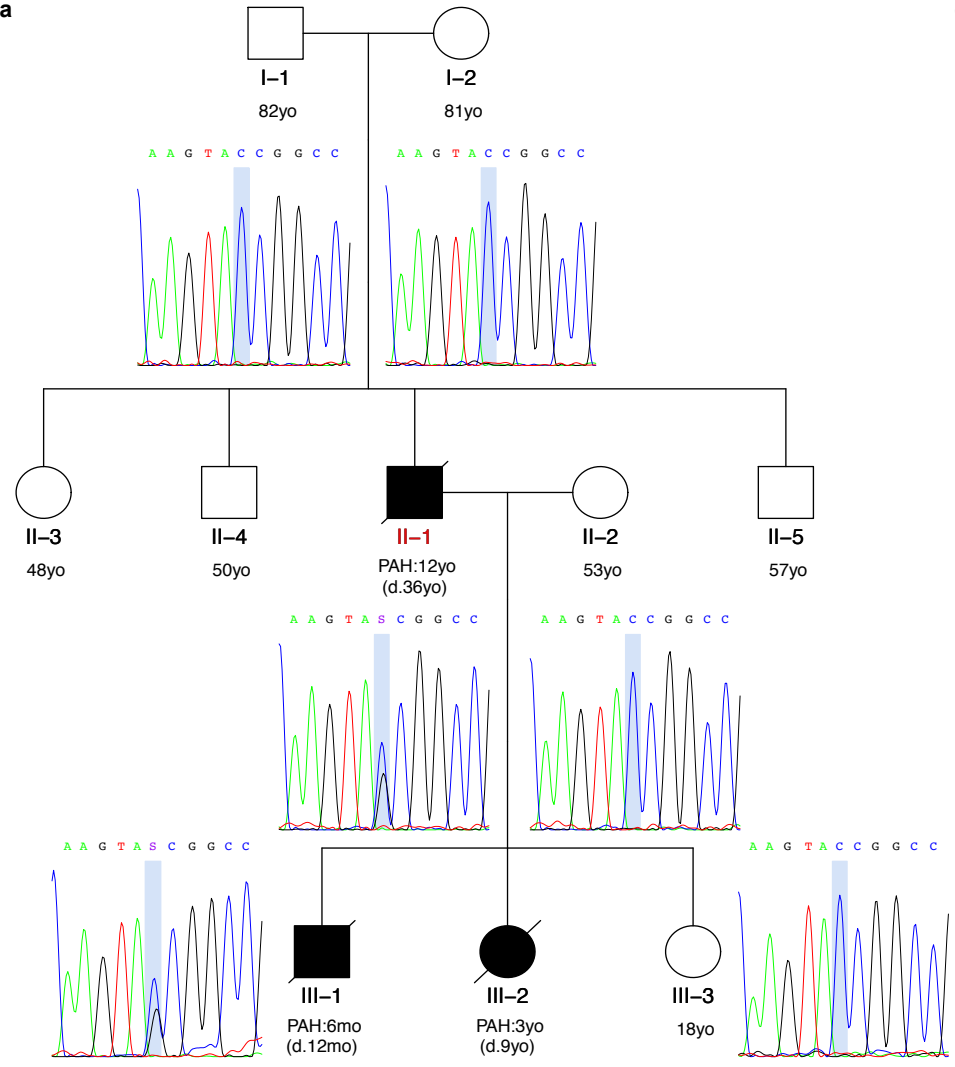
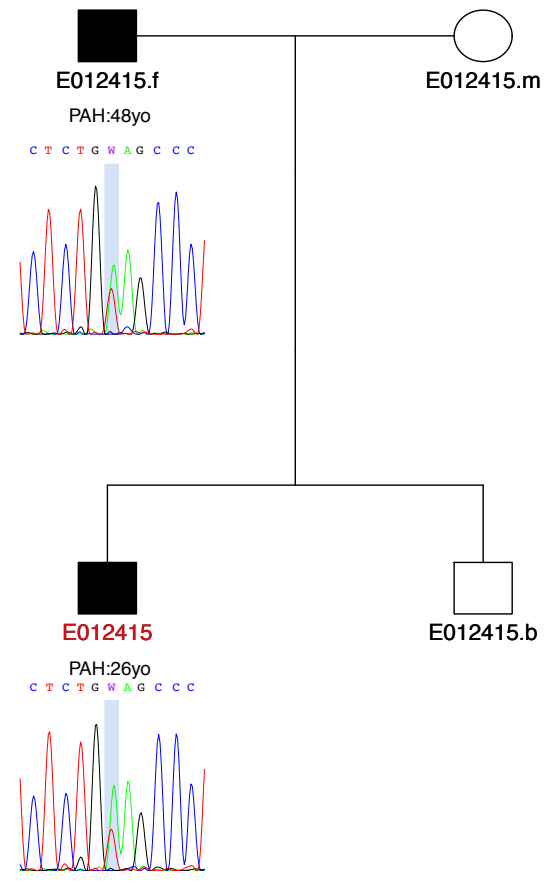
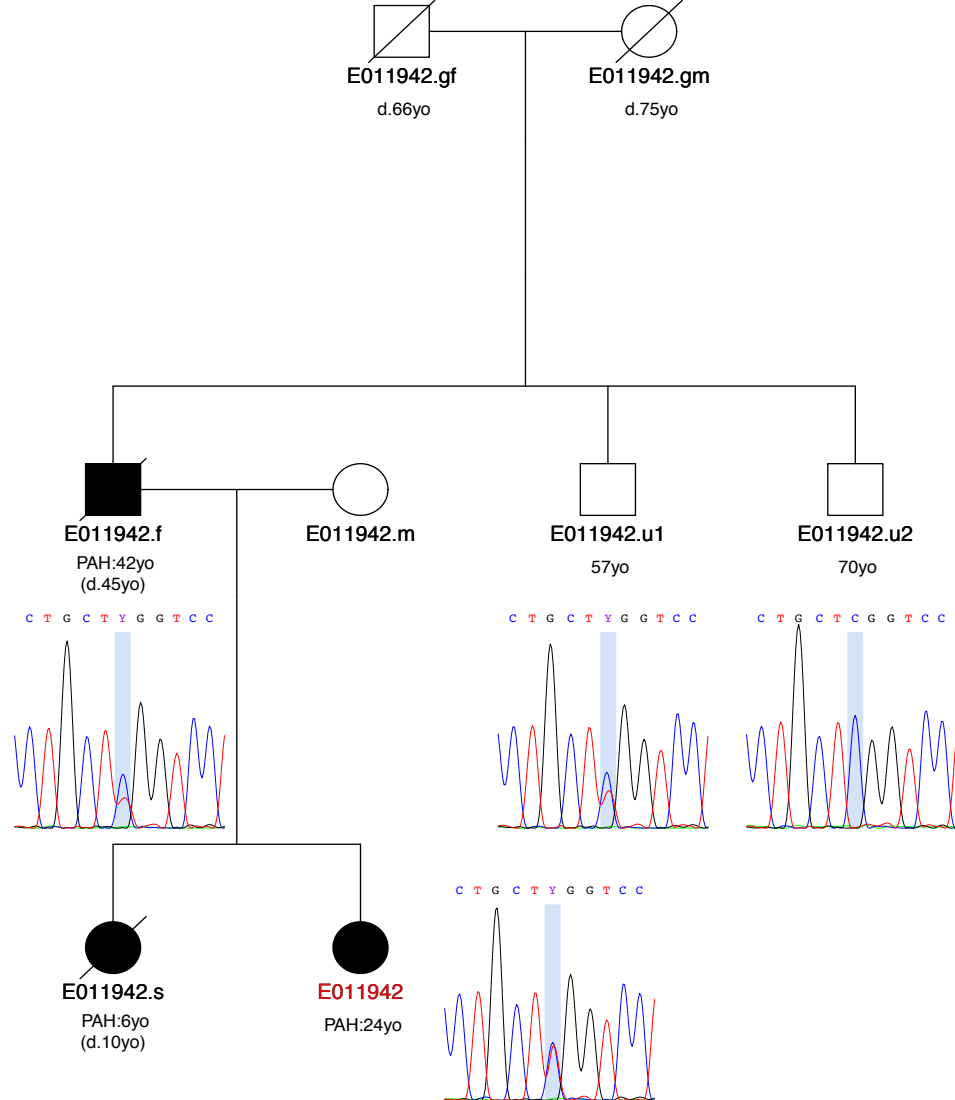
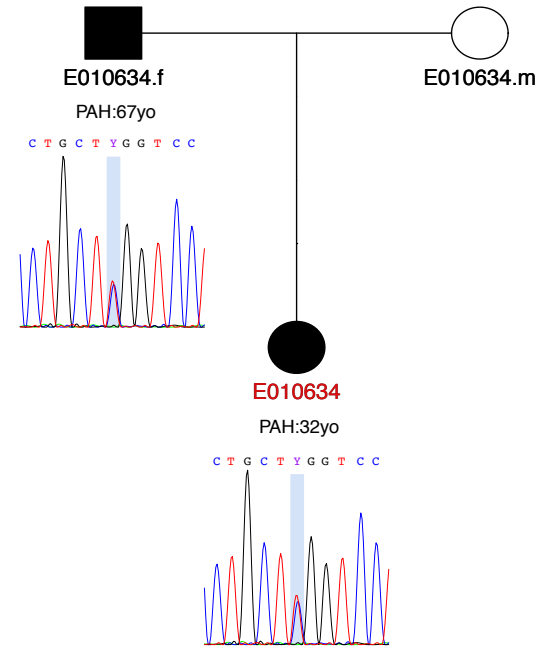
986 **Figure 10** Functional studies of novel genes. **(a-c)**. Expression of **(a)** AQP1, **(b)** ATP13A3
987 and **(c)** *SOX17* mRNA in human pulmonary artery smooth muscle cells, pulmonary artery
988 endothelial cells and blood outgrowth endothelial cells (BOECs) (n=4 biological replicates of
989 each). Relative expression of each transcript was normalised to three reference genes,
990 *ACTB*, *B2M* and *HPRT*. **(d)** Proliferation of BOECs in 5% FBS over 6 days. Cells were
991 transfected with DharmaFECT1 alone (DH1), siATP13A3 or non-targeting siRNA control
992 (siCP) **(e-f)** Quantification of apoptosis in BOECs, defined as Annexin V+/PI- cells, in
993 BOECs transfected with siATP13A3 or siCP in complex with DH1 followed by 24hr treatment
994 with 0.1% FBS or 5% FBS (n=4 biological repeats). **(f)** Measurement of apoptosis via
995 Caspase-Glo 3/7 activity measurements in BOECs transfected with siATP13A3 or siCP in
996 complex with DH1, followed by 16hr treatment in 0.1% FBS or 5% FBS. Data are from a
997 single experiment (n=4 wells) representative of 3 biological repeats. Data were analysed

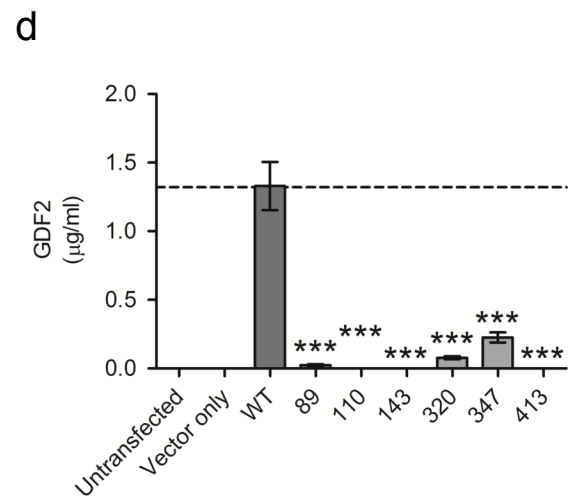
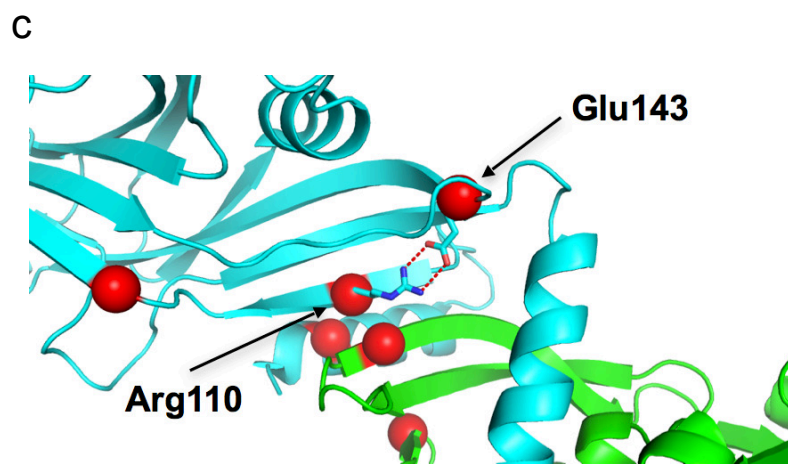
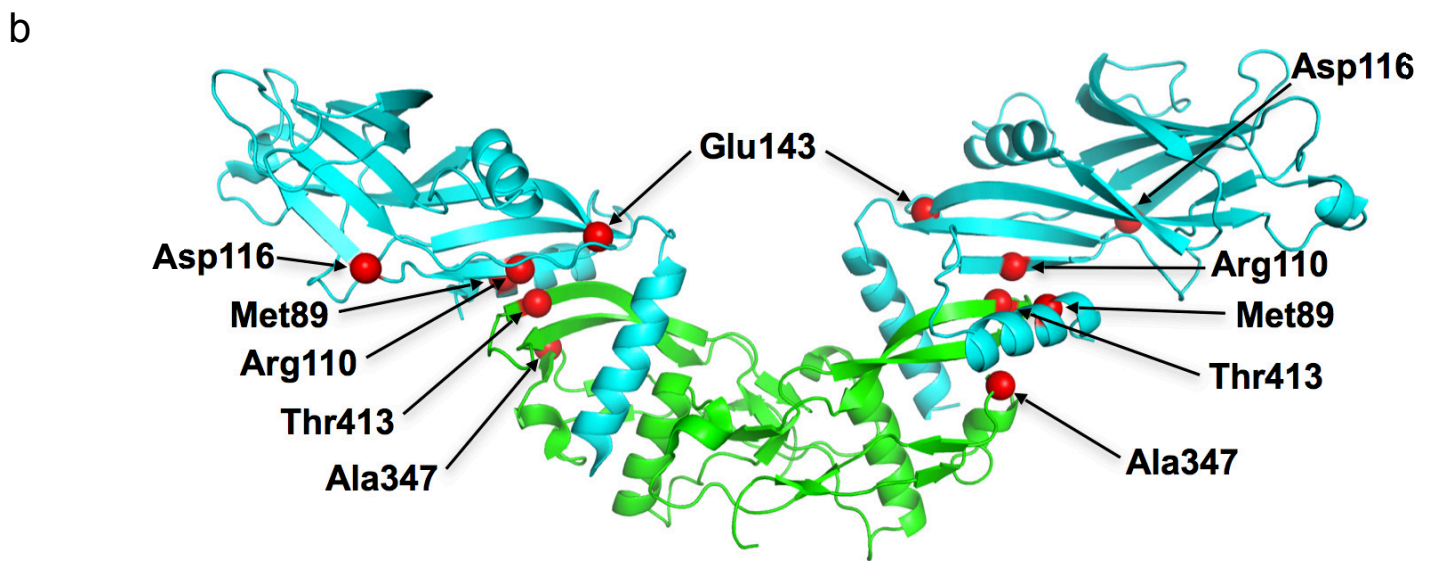
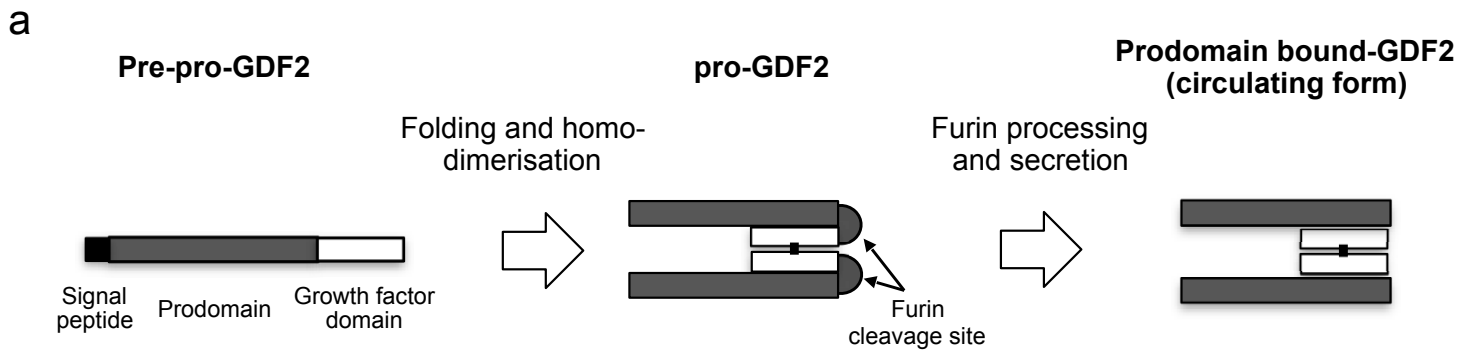
998 using a One-way analysis of variance with post hoc Tukey's test for multiple comparisons in
999 **d** and **f**. Data were analysed using a repeated measures One-way analysis of variance with
1000 post hoc Tukey's for multiple comparisons in **e**. *P<0.05, **P<0.01 within treatment groups.
1001 ###P<0.001 for effect of ligand against control for same transfection condition.



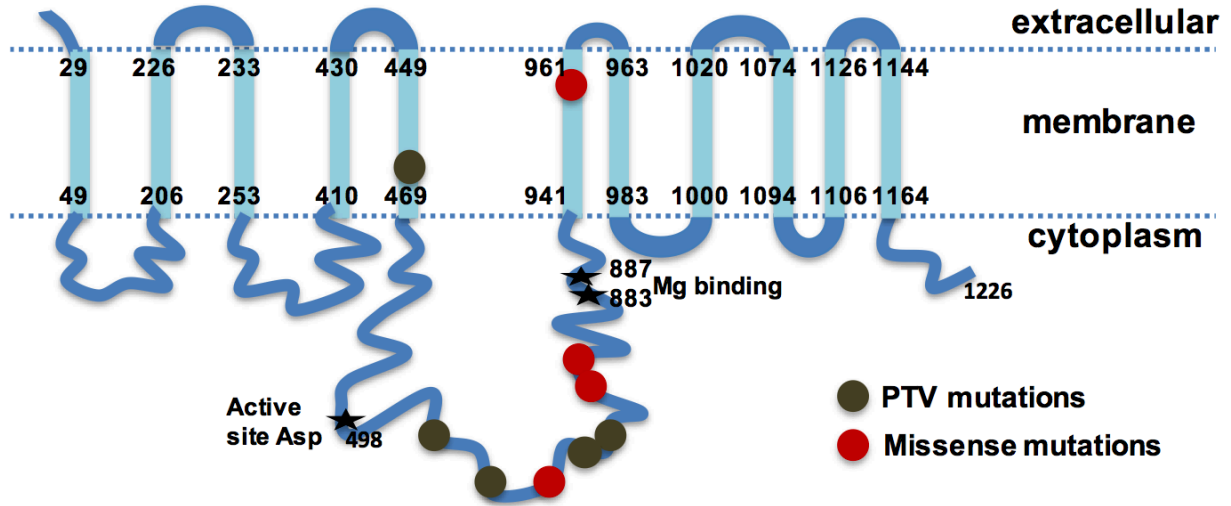




a**c****b****d**



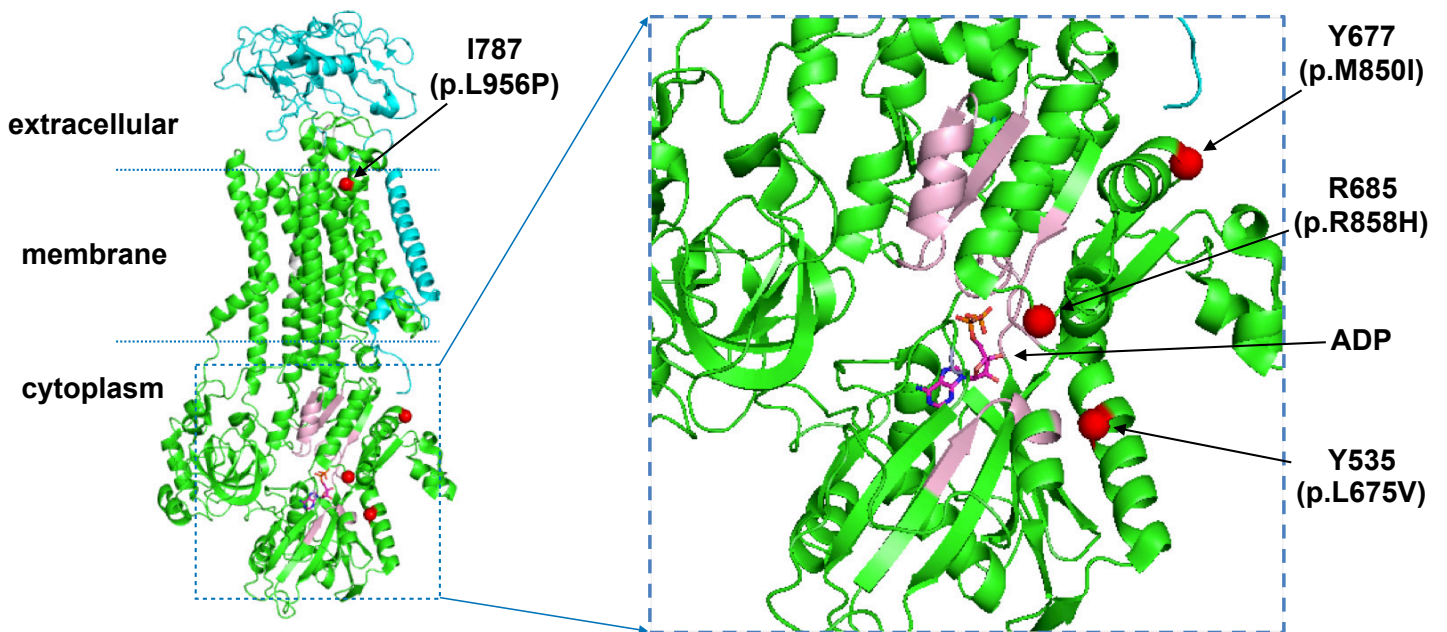
a

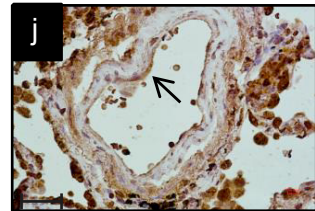
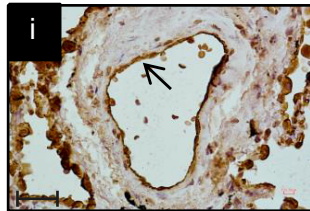
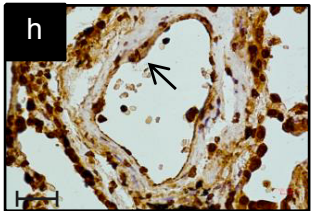
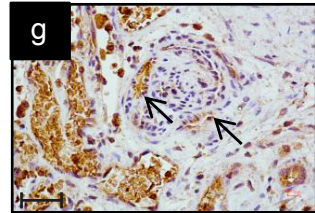
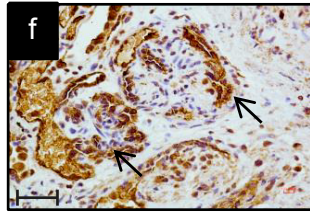
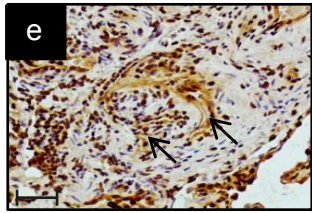
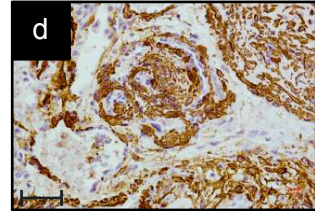
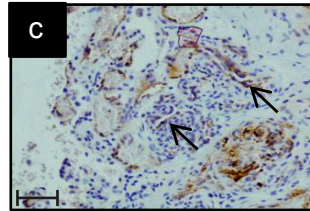
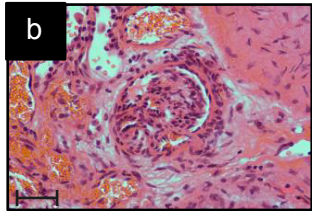
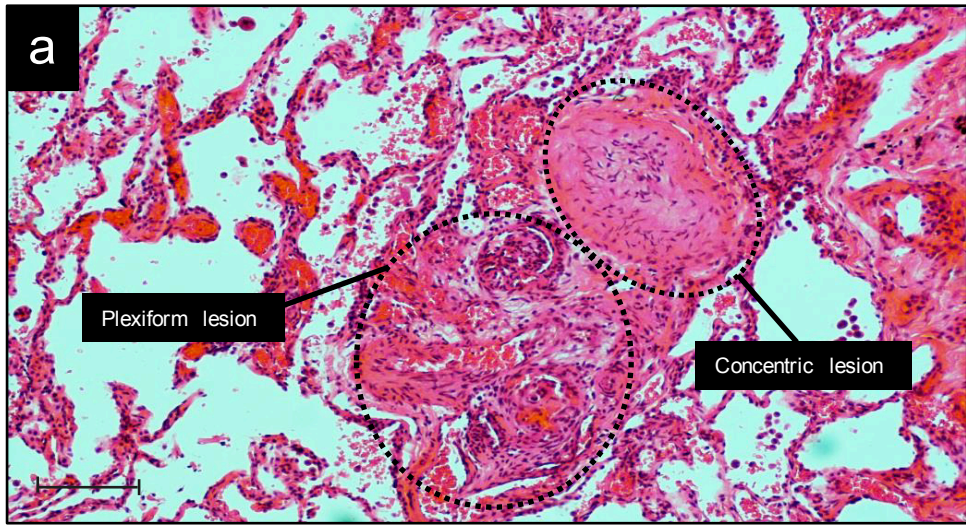


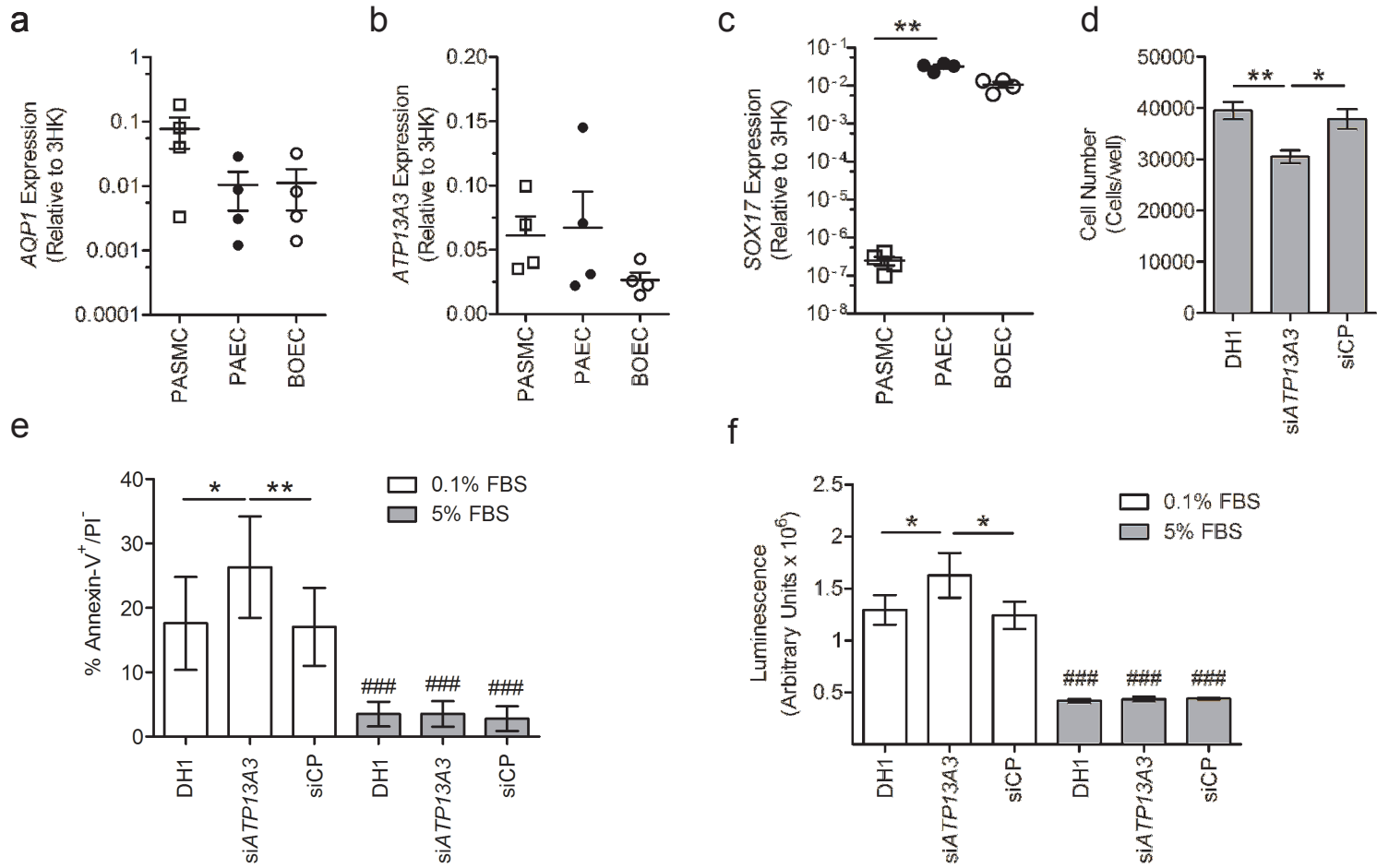
b

Q9H7F0 AT133_HUMAN	-----PVDFQNVLEDFTKQGF RVIALAHRKLESKLTWHK--VQNISRDAIENNM	714
P05023 AT1A1_HUMAN	HGKEQPLDEELKDAFQNAVLELGG LGERV LGFCHLFLPDEQFPEGFQFD TDDVNFPI DNL	583
	***. : : * * : : : * * : : : : * :	
Q9H7F0 AT133_HUMAN	SFSVILEHFQDLVPKMLLHGTVFAMAPDQKTQLIEALQNVDFVGMCGD GANDCGALKR	893
P05023 AT1A1_HUMAN	-----EQLDDIL--KYHTEIVFAN TSPQOKLIIVEGCQRQGAIVAVTGDGVNDSPALKK	727
	* : : * : : * * * * : : * * : : * . * . . : * : : * * * . * * * :	
Q9H7F0 AT133_HUMAN	QYFSVTL LYSILSNLGD FQFLFIDLAILVVVFTMSLNP AWKELVAQRPPSG----LISG	1004
P05023 AT1A1_HUMAN	TPFLIFIIANIP LPLGTVTILCIDLGTDMVPAISLAYEQAESDIMKRQPRNPKTDKLVNE	847
	* : : * . * * * . * * * * : : * . : : : : * : : : : * . * : : * :	

c







1 Pulmonary arterial hypertension (PAH) is a rare lung disorder characterised by narrowing and
2 obliteration of small pulmonary arteries ultimately leading to right heart failure. Here, the authors
3 sequence whole genomes of over 1000 PAH patients and identify likely causal variants in *GDF2*,
4 *ATP13A3*, *AQP1* and *SOX17*.

5