1 **Whole Genome Sequencing of Primary Immunodeficiency reveals a role for common and rare**

2 **variants in coding and non-coding sequences**

3

4 James E. D. Thaventhiran[1,2,29], Hana Lango Allen[3,4,5,29], Oliver S. Burren[1,29], James H. R. Farmery[6,30], Emily

5 Staples[1,30], Zinan Zhang[1,30], William Rae[1], Daniel Greene[3,5,7], Ilenia Simeoni[3,5], Jesmeen Maimaris[8,9], Chris

6 Penkett[3,4,5], Jonathan Stephens[3,4,5], Sri V.V. Deevi[3,4,5], Alba Sanchis-Juan[3,4,5], Nicholas S Gleadall[3,4,5,] Moira

7 J. Thomas[10], Ravishankar B. Sargur[11,12], Pavels Gordins[13], Helen E. Baxendale[1,14], Matthew Brown[3,4,5], Paul

8 Tuijnenburg[15], Austen Worth[8,9], Steven Hanson[16,17], Rachel Linger[3,4,5], Matthew S. Buckland[16,17], Paula J.

9 Rayner-Matthews[3,4,5], Kimberly C. Gilmour[8,9], Crina Samarghitean[3,4,5], Suranjith L. Seneviratne[16,17], Paul

10 A. Lyons[1], David M. Sansom[16,17], Andy G. Lynch[6,18], Karyn Megy[3,4,5], Eva Ellinghaus[19], David Ellinghaus[20,21],

11 Silje F. Jorgensen[22,23] , Tom H Karlsen[19], Kathleen E. Stirrups[3,4,5], Antony J. Cutler[24], Dinakantha S.

12 Kumararatne[25], NBR-RD PID Consortium, NIHR BioResource[5], Sinisa Savic[26,27], Siobhan O. Burns[16,17], Taco

13 W. Kuijpers[15,31], Ernest Turro[3,4,5,7,31], Willem H. Ouwehand[3,4,5,28,31], Adrian J. Thrasher[8,9,31], Kenneth G. C.

14 Smith[1]

15

16 1.   Department of Medicine, School of Clinical Medicine, University of Cambridge, Cambridge Biomedical Campus, Cambridge,
17      UK.
18 2.   Medical Research Council, Toxicology Unit, School of Biological Sciences, University of Cambridge, Cambridge, UK
19 3.   Department of Haematology, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK.
20 4.   NHS Blood and Transplant, Cambridge Biomedical Campus, Cambridge, UK.
21 5.   NIHR BioResource, Cambridge University Hospitals, Cambridge Biomedical Campus, Cambridge, UK.
22 6.   Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Robinson Way, Cambridge, UK.
23 7.   Medical Research Council Biostatistics Unit, Cambridge Biomedical Campus, Cambridge, UK.
24 8.   UCL Great Ormond Street Institute of Child Health, London, UK.
25 9.   Great Ormond Street Hospital for Children NHS Foundation Trust, London, UK.
26 10.  Department of Immunology, Queen Elizabeth University Hospital, Glasgow, UK.
27 11.  Sheffield Teaching Hospitals NHS Foundation Trust, Sheffield, UK.
28 12.  Department of Infection Immunity and Cardiovascular Disease, University of Sheffield, Sheffield, UK.
29 13.  East Yorkshire Regional Adult Immunology and Allergy Unit, Hull Royal Infirmary, Hull and East Yorkshire Hospitals NHS
30      Trust, Hull, UK
31 14.  Royal Papworth Hospital NHS Foundation Trust, Cambridge, UK.
32 15.  Department of Pediatric Immunology, Rheumatology and Infectious Diseases, Emma Children's Hospital & The Department
33      of Experimental Immunology, Amsterdam University Medical Center (AMC), University of Amsterdam, Amsterdam, The
34      Netherlands.
35 16.  Institute of Immunity and Transplantation, University College London, London, UK.
36 17.  Department of Immunology, Royal Free London NHS Foundation Trust, London, UK.
37 18.  School of Mathematics and Statistics/School of Medicine, University of St Andrews, St Andrews, UK.
38 19.  K.G. Jebsen Inflammation Research Centre, Institute of Clinical Medicine, University of Oslo, Oslo University Hospital,
39      Rikshospitalet, Oslo, Norway.
40 20.  Department of Transplantation, Institute of Clinical Medicine, University of Oslo, Oslo University Hospital, Rikshospitalet,
41      Oslo, Norway.
42 21.  Institute of Clinical Molecular Biology, Christian Albrechts University of Kiel, Kiel, Germany.
43 22.  Section of Clinical Immunology and Infectious Diseases, Department of Rheumatology, Dermatology and Infectious
44      Diseases, Oslo University Hospital, Rikshospitalet, Norway.
45 23.  Research Institute of Internal Medicine, Division of Surgery, Inflammatory Diseases and Transplantation, Oslo University
46      Hospital, Rikshospitalet, Norway.
47 24.  JDRF/Wellcome Diabetes and Inflammation Laboratory, Wellcome Centre for Human Genetics, Nuffield Department of
48      Medicine, NIHR Oxford Biomedical Research Centre, University of Oxford, Oxford, UK.
49 25.  Department of Clinical Biochemistry and Immunology, Cambridge University Hospitals, Cambridge Biomedical Campus,
50      Cambridge, UK.
51 26.  The Department of Clinical Immunology and Allergy, St James's University Hospital, Leeds, UK.
52 27.  The NIHR Leeds Biomedical Research Centre and Leeds Institute of Rheumatic and Musculoskeletal Medicine, Leeds, UK.
53 28.  Department of Human Genetics, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK.

54

55 29.  These authors led the analysis: James E. D. Thaventhiran, Hana Lango Allen, Oliver S. Burren.
56 30.  These authors contributed equally: J. Henry R. Farmery, Emily Staples, Zinan Zhang
57 31.  These authors supervised this work: Taco W. Kuijpers, Ernest Turro, Willem H. Ouwehand, Adrian J. Thrasher

58

**Abstract**

Primary immunodeficiency (PID) is characterised by recurrent and often life-threatening infections, autoimmunity and cancer, and it presents major diagnostic and therapeutic challenges. Although the most severe forms present in early childhood, the majority of patients present in adulthood, typically with no apparent family history and a variable clinical phenotype of widespread immune dysregulation: about 25% of patients have autoimmune disease, allergy is prevalent, and up to 10% develop lymphoid malignancies. Consequently, in sporadic PID genetic diagnosis is difficult and the role of genetics is not well defined. We addressed these challenges by performing whole genome sequencing (WGS) of a large PID cohort of 1,318 subjects. Analysis of coding regions of 886 index cases found disease-causing mutations in known monogenic PID genes in 8.2%, while a Bayesian approach (BeviMed[1]) identified multiple potential new disease-associated genes. Exploration of the non-coding space revealed deletions in regulatory regions which contribute to disease causation. Finally, a genome-wide association study (GWAS) identified novel PID-associated loci and uncovered evidence for co-localisation of, and interplay between, novel high penetrance monogenic variants and common variants (at the *PTPN2* and *SOCS1* loci). This begins to explain the contribution of common variants to variable penetrance and phenotypic complexity in PID. Thus, a cohort-based WGS approach to PID diagnosis can increase diagnostic yield while deepening our understanding of the key pathways determining variation in human immune responsiveness.

79   The phenotypic heterogeneity of PID leads to diagnostic difficulty, and almost certainly to an
80   underestimation of its true incidence. Our cohort reflects this heterogeneity, though it is dominated by
81   adult onset, sporadic antibody deficiency associated PID (AD-PID: comprising Common Variable
82   Immunodeficiency (CVID), Combined Immunodeficiency (CID) and isolated antibody deficiency).
83   Identifying a specific genetic cause of PID can facilitate definitive treatment including haematopoietic
84   stem cell transplantation, genetic counselling, and the possibility of gene-specific therapy[2–4] while
85   contributing to our understanding of the human immune system[5]. Unfortunately, only 29% of patients
86   with PID receive a genetic diagnosis[6]. The lowest diagnosis rate is in patients who present as adults,
87   have no apparent family history, and in whom matching the clinical phenotype to a known genetic cause
88   is difficult, as the latter can be surprisingly variable even in patients with the same genetic defect (in the
89   UK PID cohort 78% of cases are adult and 76% sporadic[6]). Moreover, while over 300 monogenic causes
90   of PID have been described[7], the genotype-phenotype correlation in PID is complex. In CVID, for
91   example, pathogenic variants in *TACI* (*TNFRSF13B*) occur in 10% of patients but typically have low
92   disease penetration, appearing to act as disease modifiers[8]. Furthermore, a common variant analysis of
93   CVID identified two disease-associated loci, raising the possibility that common variants may impact
94   upon clinical presentation[9]. We therefore investigated whether applying WGS across a "real world" PID
95   cohort might illuminate the complex genetics of the range of conditions collectively termed PID.
96
97   **Patient cohort**
98
99   974 sporadic and familial PID patients, and 344 unaffected relatives, were recruited by collaborators as
100   part of the United Kingdom NIHR BioResource - Rare Diseases program (NBR-RD; **Supplementary Note**).
101   Of these, 886 were index cases who fell into one of the diagnostic categories of the European Society for
102   Immunodeficiencies (ESID) registry diagnostic criteria (**Fig. 1a; Supplementary Table 1**). This cohort
103   represents a third of CVID and half of CID patients registered in the UK[10]. Paediatric and familial cases
104   were less frequent, in part reflecting prior genetic testing of more severe cases (**Supplementary Fig. 1**).
105   Clinical phenotypes were dominated by adult-onset sporadic AD-PID: all had recurrent infections, 28%
106   had autoimmunity, and 8% had malignancy (**Fig. 1a-b, Supplementary Table 2**), mirroring the UK
107   national PID registry[6].
108
109   **Identification of Pathogenic Variants in Known Genes**
110
111   We analysed coding regions of genes previously causally associated with PID[11] (**Methods**). We identified
112   85 potentially causal variants in 73 index cases (8.2%) across 39 genes implicated in monogenic disease
113   (**Fig. 1c; Supplementary Table 3**). 60 patients (6.8%) had a previously reported pathogenic variant in the
114   disease modifier *TACI* (*TNFRSF13B*), increasing the diagnostic yield to 15.0% (133 patients). Interestingly,
115   5 patients with a monogenic diagnosis (in *BTK*, *LRBA*, *MAGT1*, *RAG2*, *SMARCAL1*) also had a pathogenic
116   *TACI* variant. The diagnostic yield rose to 17.0% (151 patients) once novel causal variants in *NFKB1* and
117   *ARPC1B*, associated with PID only after our initial analysis, were included. Of the 85 monogenic variants
118   we reported, 51 (60%) had not been previously described (**Supplementary Table 3**), and 4 were
119   structural variants, including a single exon deletion, unlikely to have been detected by whole exome
120   sequencing[12].
121
122   We observed divergence from an expected clinical phenotype for causal variants in 14 genes: for
123   instance, only 4 of the 8 *STAT1* patients had the pathognomonic chronic mucocutaneous candidiasis[13,14].
124   A more remarkable example of phenotypic complexity was the case of a 40 year-old patient presenting
125   with specific antibody deficiency and a premature stop variant at Arg328 in X-linked *IL2RG*, a defect
126   expected to cause absent T and NK cells and death in infancy. We found that the mild phenotype could

127  be ascribed to several independent somatic changes that reversed the premature stop codon, restoring
128  both T and NK cell lineages (**Fig. 1d** and **Supplementary Fig. 2**).

129

130  Since many PID-associated genes were initially discovered in a small number of typically familial cases, it
131  is perhaps not surprising that the phenotypes described do not reflect true clinical diversity. Thus, a
132  cohort-based WGS approach to PID can provide a significant diagnostic yield even in a predominantly
133  pre-screened and sporadic cohort, allows diagnoses which are not constrained by pre-existing
134  assumptions about genotype-phenotype relationships, and suggests caution in the use of clinical
135  phenotype in targeted gene screening and when interpreting PID genetic data.

136

137  **An approach to identifying new PID-associated genes in a WGS cohort**

138

139  We next sought to determine whether the cohort-based WGS approach could identify new genetic
140  associations with PID. We developed a Bayesian inference procedure, named BeviMed[1], to determine
141  posterior probabilities of association (PPA) between each gene and case/control status of the 886 index
142  cases and 9,283 unrelated controls (**Methods**). For each gene, we analysed variants with gnomAD minor
143  allele frequency (MAF) <0.001 and Combined Annotation Dependent Depletion (CADD) score >=10.
144  Genes with PPA>=0.18 are shown in **Fig. 1e**. There was a strong enrichment for known PID genes
145  (Wilcoxon $P<1\times10^{-200}$), supporting this statistical approach. Two novel BeviMed-identified genes were
146  subsequently causally associated with PID. *NFKB1* had the strongest probability of disease association
147  (PPA=1.0), driven by truncating heterozygous variants in 13 patients. Subsequent assessment of co-
148  segregation, protein expression, and B cell phenotype in pedigrees established these as disease-causing
149  variants, and consequently loss of function variants in *NFKB1* as the most common monogenic cause of
150  CVID[15]. Evidence of association of *ARPC1B* with PID (PPA=0.18) was driven by 2 functionally validated
151  cases, one homozygous for a complex InDel[16] and the other described below.

152  The discovery of both known and subsequently validated new PID genes using BeviMed underlines its
153  effectiveness in cohorts of unrelated patients with sporadic disease. Many candidate genes identified by
154  BeviMed remain to be functionally validated and, as the PID cohort grows, even very rare causes of PID
155  (e.g. affecting 0.2% of cases) will be detectable with a high positive predictive value (**Supplementary Fig.
156  3**).

157

158  **Identification of regulatory elements contributing to PID**

159

160  Sequence variation within non-coding regions of the genome can have profound effects on spatial and
161  temporal gene expression[17] and would be expected to contribute to PID susceptibility. We combined
162  rare variant and deletion events with a tissue-relevant catalogue of cis-regulatory elements (CREs)[18]
163  generated using promoter capture Hi-C (pcHi-C)[19] in matching tissues to prioritise putative causal PID
164  genes (**Fig. 2a**). Being underpowered to detect single nucleotide variants affecting CREs[20], we limited our
165  initial analysis to rare structural variants (SV) overlapping exon, promoter or 'super-enhancer' CREs of
166  known PID genes. No homozygous deletion events affecting CREs were identified, so we sought CRE SV
167  deletions that might cause disease through a candidate compound heterozygote (cHET) mechanism with
168  either a heterozygous rare coding variant or another SV in a pcHi-C linked gene (**Fig. 2a**). Out of 22,296
169  candidate cHET deletion events, after filtering by MAF, functional score and known PID gene status, we
170  obtained 10 events; the functional follow-up of three is described (**Fig. 2b**).

171

172  The *LRBA* and *DOCK8* cHET variants (**Supplementary Fig. 4**) were functionally validated; the former was
173  demonstrated to result in impaired surface CTLA-4 expression on Treg cells (**Supplementary Fig. 5**)

174  whilst the latter led to DOCK8 deficiency as confirmed by flow cytometry (data not shown). Although in
175  these two cases SV deletions encompassed both non-coding CREs and coding exons, the use of WGS PID
176  cohorts to detect a contribution of CREs confined to the non-coding space would represent a major
177  advance in PID pathogenesis and diagnosis. *ARPC1B* fulfilled this criterion, with its BeviMed association
178  partially driven by a patient cHET for a novel p.Leu247Glyfs*25 variant resulting in a premature stop,
179  and a 9Kb deletion spanning the promoter region including an untranslated first exon (**Fig. 2c**) that has
180  no coverage in the ExAC database (http://exac.broadinstitute.org). Two first-degree relatives were
181  heterozygous for the frameshift variant, and two for the promoter deletion (**Fig. 2d**). Western blotting
182  demonstrated complete absence of ARPC1B (**Fig. 2e**) and, consistent with previous reports[21], raised
183  ARPC1A in platelets. *ARPC1B* mRNA was almost absent from mononuclear cells in the cHET patient and
184  reduced in a clinically unaffected sister carrying the frameshift mutation (**Fig. 2f**). An allele specific
185  expression assay demonstrated that the promoter deletion essentially abolished mRNA expression (**Fig.
186  2g,h**).

187  These examples show the utility of WGS for detecting compound heterozygosity for a coding variant and
188  a non-coding CRE deletion, and demonstrate a further advantage of a WGS approach to PID diagnosis.
189  Improvements in analysis methodology, cohort size and better annotation of regulatory regions will be
190  required to explore the non-coding space more fully and discover new disease-causing genetic variants.

191
192  **WGS identifies PID-associated telomere shortening**
193
194  A striking example of WGS data providing more than just the linear genomic sequence is telomere
195  length estimation from mapped and unmapped reads[22]. We validated this method by showing
196  correlation with gender (**Fig. 3a**) and a particularly strong correlation with age (**Supplementary Fig. 6**) in
197  3,313 NBR-RD subjects (**Methods**). We demonstrated the effectiveness of this, the first large-scale
198  application of WGS-based telomere length estimation, by replicating an association with the telomerase
199  RNA component gene (*TERC*: **Supplementary Table 4**)[23] and identifying several PID cases with short
200  telomeres (**Fig. 3b**). Given that disruption of telomerase genes can cause PID[24], we looked for potentially
201  damaging coding variants in known telomere deficiency genes[25] in these PID cases, identifying 3 subjects
202  with novel variants potentially causative for telomerase deficiency (**Fig. 3b**). One had a homozygous
203  defect in telomerase reverse transcriptase (*TERT*), a subunit of the telomerase complex. Two male
204  siblings were found to have a hemizygous variant in dyskerin (*DKC1*), known to be associated with PID
205  and X-linked dyskeratosis congenita[26] (**Fig. 3c**). Therefore, WGS telomere length estimation can be used
206  as an effective approach to identify PID patients with novel variants causing telomere shortening.
207
208  **GWAS of the WGS cohort reveals novel PID-associated loci**
209
210  The diverse clinical phenotype and variable within-family disease penetrance of PID may be in part due
211  to stochastic events (e.g. unpredictable pathogen transmission) but may also have a genetic basis. We
212  therefore performed a GWAS of common SNPs (MAF>0.05), restricted to 733 AD-PID cases (**Fig. 1a**) to
213  reduce phenotypic heterogeneity, and 9,225 unrelated NBR-RD controls. We confirmed the known MHC
214  association and identified additional loci with suggestive association (**Fig. 4a, Supplementary Fig. 7**). A
215  GWAS of SNPs of intermediate frequency (0.005<MAF<0.05) identified a single locus incorporating
216  *TNFRSF13B* (**Fig. 4a, Supplementary Table 5, Extended Data Fig. 1**), for which the lead p.Cys104Arg
217  variant has been previously reported[27].

218  To increase power, we conducted a fixed effect meta-analysis of the AD-PID GWAS with summary
219  statistics data from an ImmunoChip study of 778 CVID cases and 10,999 controls[9] (**Fig. 4a,
220  Supplementary Table 5**). This amplified the MHC and 16p13.13 associations[9], found an additional locus

221  at 3p24.1 within the promoter region of *EOMES* (**Extended Data Fig. 2**), and a suggestive association at
222  18p11.21 proximal to *PTPN2* (**Extended Data Fig. 3**). Conditional analysis of the MHC locus revealed
223  independent signals at the Class I and Class II regions (**Supplementary Fig. 8**), driven by classical alleles
224  HLA-B*08:01 and HLA-DRB1*15:01 (**Methods**) with amino-acid changes known to impact upon peptide
225  binding (**Fig. 4b**).

226  We next sought to examine, genome-wide, the enrichment of non-MHC AD-PID associations in 9 other
227  diseases (**Extended Data Table 1**). We found significant enrichment for allergic (e.g. asthma) and
228  immune-mediated diseases (e.g. Crohn's disease), which was not evident in Type 2 diabetes or coronary
229  artery disease (**Fig. 4c**). This suggests that the common variant association between PID and other
230  immune-mediated diseases extends beyond the 4 genome-wide loci to multiple sub-genome-wide
231  associations, and that dysregulation of common pathways contributes to susceptibility to both.
232  Understanding the impact of these interrelationships will be a complex process. For example, while
233  variants in the *HLA-DRB1* and 16p13.13 loci increase the risk of both PID and autoimmunity, those at the
234  *EOMES* locus predispose to PID but protect from rheumatoid arthritis[28] (**Extended Data Fig. 2**).

236  Given this observed enrichment, we sought to investigate whether candidate genes identified through
237  large cohort association analysis of immune-mediated disease might have utility in prioritising novel
238  candidate genes harbouring rare coding variation causal for PID. We used the data-driven capture-HiC
239  omnibus gene score (COGS) approach[19] to prioritise putative causal genes across the 4 non-MHC AD-PID
240  loci identified by our meta-analysis, and assessed across 11 immune-mediated diseases (**Supplementary
241  Tables 5 and 6**). Hypothesising that causal PID genes would be intolerant to protein-truncating variation,
242  we computed an overall prioritisation score by taking the product of pLI (a measure of tolerance to loss
243  of gene function) and COGS gene scores for each disease. Six protein coding genes had an above
244  average prioritisation score in one or more diseases (**Fig. 4d**) which we examined for rare, potentially
245  causative variants within our cohort. We identified a single protein truncating variant in *ETS1, SOCS1*
246  and *PTPN2* genes, all occurring exclusively in PID patients in the NBR-RD cohort. None of the genes are
247  recognised causes of PID despite their involvement in immune processes (**Supplementary Discussion**).
248  The two cases with *SOCS1* and *PTPN2* variants were analysed further.

250  The patient with a heterozygous protein-truncating *SOCS1* variant (p.Met161Alafs*46) presented with
251  CVID complicated by lung and liver inflammation and B cell lymphopenia (**Supplementary Discussion,
252  Supplementary Fig. 9**). *SOCS1* limits phosphorylation of targets including STAT1, and is a key regulator of
253  IFN-γ signalling. SOCS1 haploinsufficiency in mice leads to B lymphopenia[29,30], immune-mediated liver
254  inflammation[31] and colitis[32]. In patient T cell blasts SOCS1 was deficient and IFN-γ induced STAT1
255  phosphorylation was abnormal (**Fig. 4e**), consistent with SOCS1 haploinsufficiency causing PID. The
256  patient also carries the *SOCS1* pcHiC-linked 16p13.13 risk-allele identified in the AD-PID GWAS
257  (**Extended Data Fig. 4**). Long read sequencing using Oxford Nanopore technology showed this to be in
258  *trans* with the novel *SOCS1*-truncating variant (**Methods**); such compound heterozygosity raises the
259  possibility that common and rare variants may combine to cause disease.

261  A more detailed example of an interplay between rare and common variants is provided by a family
262  containing a novel *PTPN2* premature stop-gain at p.Glu291 and a common autoimmunity-associated
263  variant (**Fig. 4f**). *PTPN2* encodes the non-receptor T-cell protein tyrosine phosphatase (TC-PTP) protein,
264  that negatively regulates immune responses by dephosphorylation of the proteins mediating cytokine
265  signalling. *PTPN2* deficient mice are B cell lymphopenic[33,34], while inducible haematopoietic deletion of
266  *PTPN2* leads to B and T cell proliferation and autoimmunity[35]. The novel truncating variant was
267  identified in a "sporadic" index case presenting with CVID at age 20; he had B lymphopenia
268  (**Supplementary Fig. 9**), low IgG, symmetrical rheumatoid-like polyarthropathy, severe recurrent

6

269 bacterial infections, splenomegaly and inflammatory lung disease. His mother, also heterozygous for the
270 *PTPN2* truncating variant, had systemic lupus erythematosus (SLE), insulin-dependent diabetes mellitus
271 diagnosed at 42, hypothyroidism and autoimmune neutropenia (**Supplementary Discussion**). Gain-of-
272 function variants in *STAT1* can present as CVID (**Supplementary Table 3**) and TC-PTP, like SOCS1,
273 reduces phosphorylated-STAT1 (**Fig. 4g**). Both mother and son demonstrated reduced TC-PTP expression
274 and STAT1 hyperphosphorylation in T cell blasts, similar to the SOCS1 haploinsufficient patient above
275 and to known STAT1 GOF patients; abnormalities that were more pronounced in the *PTPN2* index case
276 (**Fig. 4h**).

277

278 The index case, but not his mother, carried the G allele of variant rs2847297 at the *PTPN2* locus, an
279 expression quantitative trait locus (eQTL)[36] previously associated with rheumatoid arthritis[37]. His
280 brother, generally healthy apart from severe allergic nasal polyposis, was heterozygous at rs2847297
281 and did not inherit the rare variant (**Fig. 4f**). Allele-specific expression analysis demonstrated reduced
282 *PTPN2* transcription from the rs2847297-G allele, explaining the lower expression of TC-PTP and greater
283 persistence of pSTAT1 in the index case compared to his mother (**Fig. 4i**). This in turn could explain the
284 variable disease penetrance in this family, with *PTPN2* haploinsufficiency alone driving autoimmunity in
285 the mother, but with the additional impact of the common variant on the index case causing
286 immunodeficiency (and perhaps reducing the autoimmune phenotype). The family illustrates the power
287 of cohort-wide WGS approach to PID diagnosis, by revealing both a new monogenic cause of disease,
288 and how the interplay between common and rare genetic variants may contribute to the variable clinical
289 phenotypes of PID.

290

291 In summary, we show that cohort-based WGS in PID is a powerful approach to provide immediate
292 diagnosis of known genetic defects, and to discover new coding and non-coding variants associated with
293 disease. Intriguingly, even with a limited sample size, we could explore the interface between common
294 and rare variant genetics, explaining why PID encompasses such a complex range of clinical syndromes
295 of variable penetrance. Increasing cohort size will be crucial for powering the analyses needed to
296 identify both causal and disease-modifying variants, thus unlocking the potential of WGS for PID
297 diagnosis. Improved analysis methodology and better integration of parallel datasets, such as GWAS and
298 cell surface or metabolic immunophenotyping, will allow further exploration of the non-coding space
299 and enhance diagnostic yield. Such an approach promises to transform our understanding of genotype-
300 phenotype relationships in PID and related immune-mediated conditions, and could redefine the clinical
301 boundaries of immunodeficiency, add to our understanding of human immunology, and ultimately
302 improve patient outcomes.

303

304

305

317

Members of the NBR-RD PID Consortium: Zoe Adhya, Hana Alachkar, Ariharan Anantharachagan, Richard Antrobus, Gururaj Arumugakani, Chiara Bacchelli, Helen E Baxendale, Claire Bethune, Shahnaz Bibi, Barbara Boardman, Claire Booth, Matthew Brown, Michael J Browning, Mary Brownlie, Matthew S Buckland, Siobhan O Burns, Oliver S Burren, Anita Chandra, Hayley Clifford, Nichola Cooper, Godelieve J de Bree, E Graham Davies, Sarah Deacock, John Dempster, Lisa A Devlin, Elizabeth Drewe, J David M Edgar, William Egner, Tariq El-Shanawany, James H R Farmery, H Bobby Gaspar, Rohit Ghurye, Kimberly C Gilmour, Sarah Goddard, Pavels Gordins, Sofia Grigoriadou, Scott J Hackett, Rosie Hague, Lorraine Harper, Grant Hayman, Archana Herwadkar, Stephen Hughes, Aarnoud P Huissoon, Stephen Jolles, Julie Jones, Yousuf M Karim, Peter Kelleher, Sorena Kiani, Nigel Klein, Taco W Kuijpers, Dinakantha S Kumararatne, James Laffan, Hana Lango Allen, Sara E Lear, Hilary Longhurst, Lorena E Lorenzo, Paul A Lyons, Jesmeen Maimaris, Ania Manson, Elizabeth M McDermott, Hazel Millar, Anoop Mistry, Valerie Morrisson, Sai H K Murng, Iman Nasir, Sergey Nejentsev, Sadia Noorani, Eric Oksenhendler, Mark J Ponsford, Waseem Qasim, Ellen Quinn, Isabella Quinti, Alex Richter, Crina Samarghitean, Ravishankar B Sargur, Sinisa Savic, Suranjith L Seneviratne, W A Carrock Sewell, Fiona Shackley, Ilenia Simeoni, Kenneth G C Smith, Emily Staples, Hans Stauss, Cathal L Steele, James E Thaventhiran, David C Thomas, Moira J Thomas, Adrian J Thrasher, John A Todd, Anton T J Tool, Rafal D Urniaz, Steven B Welch, Lisa Willcocks, Sarita Workman, Austen Worth, Nigel Yeatman, Patrick F K Yong

Correspondence and requests for materials should be addressed to J.E.D.T. (jedt2@cam.ac.uk) and K.G.C.S. (kgcs2@cam.ac.uk)

**References**

1. Greene, D., Richardson, S. & Turro, E. A Fast Association Test for Identifying Pathogenic Variants Involved in Rare Diseases. *Am. J. Hum. Genet.* **101,** 104–114 (2017).

2. Chaigne-Delalande, B. *et al.* Mg2+ Regulates Cytotoxic Functions of NK and CD8 T Cells in Chronic EBV Infection Through NKG2D. *Science (80-. ).* **341,** 186–191 (2013).

3. Lo, B. *et al.* Patients with LRBA deficiency show CTLA4 loss and immune dysregulation responsive to abatacept therapy. *Science* **349,** 436–40 (2015).

4. Rao, V. K. *et al.* Effective 'activated PI3Kδ syndrome'-targeted therapy with the PI3Kδ inhibitor leniolisib. *Blood* **130,** 2307–2316 (2017).

5. Casanova, J.-L. Human genetic basis of interindividual variability in the course of infection. *Proc. Natl. Acad. Sci. U. S. A.* **112,** E7118-27 (2015).

6. Edgar, J. D. M. *et al.* The United Kingdom Primary Immune Deficiency (UKPID) Registry: report of

367 the first 4 years' activity 2008-2012. *Clin. Exp. Immunol.* **175,** 68–78 (2014).

368 7. Bousfiha, A. *et al.* The 2017 IUIS Phenotypic Classification for Primary Immunodeficiencies. *J. Clin.*
369 *Immunol.* **38,** 129–143 (2018).

370 8. Pan-Hammarström, Q. *et al.* Reexamining the role of TACI coding variants in common variable
371 immunodeficiency and selective IgA deficiency. *Nat. Genet.* **39,** 429–430 (2007).

372 9. Li, J. *et al.* Association of CLEC16A with human common variable immunodeficiency disorder and
373 role in murine B cells. *Nat. Commun.* **6,** 6804 (2015).

374 10. Shillitoe, B. *et al.* The United Kingdom Primary Immune Deficiency (UKPID) registry 2012 to 2017.
375 *Clin. Exp. Immunol.* **192,** 284–291 (2018).

376 11. Bousfiha, A. *et al.* The 2015 IUIS Phenotypic Classification for Primary Immunodeficiencies. *J. Clin.*
377 *Immunol.* **35,** 727–38 (2015).

378 12. Fromer, M. *et al.* Discovery and statistical genotyping of copy-number variation from whole-
379 exome sequencing depth. *Am. J. Hum. Genet.* **91,** 597–607 (2012).

380 13. van de Veerdonk, F. L. *et al.* STAT1 Mutations in Autosomal Dominant Chronic Mucocutaneous
381 Candidiasis. *N. Engl. J. Med.* **365,** 54–61 (2011).

382 14. Liu, L. *et al.* Gain-of-function human STAT1 mutations impair IL-17 immunity and underlie chronic
383 mucocutaneous candidiasis. *J. Exp. Med.* **208,** 1635–48 (2011).

384 15. Tuijnenburg, P. *et al.* Loss-of-function nuclear factor κB subunit 1 (NFKB1) variants are the most
385 common monogenic cause of common variable immunodeficiency in Europeans. *J. Allergy Clin.*
386 *Immunol.* **142,** 1285–1296 (2018).

387 16. Kuijpers, T. W. *et al.* Combined immunodeficiency with severe inflammation and allergy caused
388 by ARPC1B deficiency. *J. Allergy Clin. Immunol.* **140,** 273–277.e10 (2017).

389 17. Lettice, L. A. *et al.* A long-range Shh enhancer regulates expression in the developing limb and fin
390 and is associated with preaxial polydactyly. *Hum. Mol. Genet.* **12,** 1725–35 (2003).

391 18. Hnisz, D. *et al.* Super-Enhancers in the Control of Cell Identity and Disease. *Cell* **155,** 934–947
392 (2013).

393 19. Javierre, B. M. *et al.* Lineage-Specific Genome Architecture Links Enhancers and Non-coding
394 Disease Variants to Target Gene Promoters. *Cell* **167,** 1369–1384.e19 (2016).

395 20. Short, P. J. *et al.* De novo mutations in regulatory elements in neurodevelopmental disorders.
396 *Nature* **555,** 611–616 (2018).

397 21. Kahr, W. H. A. *et al.* Loss of the Arp2/3 complex component ARPC1B causes platelet
398 abnormalities and predisposes to inflammatory disease. *Nat. Commun.* **8,** 14816 (2017).

399 22. Farmery, J. H. R., Smith, M. L. & Lynch, A. G. Telomerecat: A ploidy-agnostic method for
400 estimating telomere length from whole genome sequencing data. *Sci. Rep.* **8,** 1300 (2018).

401 23. Codd, V. *et al.* Identification of seven loci affecting mean telomere length and their association
402 with disease. *Nat. Genet.* **45,** 422–7, 427e1–2 (2013).

403 24. Jyonouchi, S., Forbes, L., Ruchelli, E. & Sullivan, K. E. Dyskeratosis congenita: a combined
404 immunodeficiency with broad clinical spectrum - a single-center pediatric experience. *Pediatr.*
405 *Allergy Immunol.* **22,** 313–9 (2011).

406 25. Tummala, H. *et al.* Poly(A)-specific ribonuclease deficiency impacts telomere biology and causes
407 dyskeratosis congenita. *J. Clin. Invest.* **125,** 2151–60 (2015).

408 26. Cossu, F. *et al.* A novel DKC1 mutation, severe combined immunodeficiency (T+B-NK- SCID) and
409 bone marrow transplantation in an infant with Hoyeraal-Hreidarsson syndrome. *Br. J. Haematol.*
410 **119,** 765–8 (2002).

411 27. Salzer, U. *et al.* Mutations in TNFRSF13B encoding TACI are associated with common variable
412 immunodeficiency in humans. *Nat. Genet.* **37,** 820–828 (2005).

413 28. Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery.
414 *Nature* **506,** 376–381 (2014).

415 29. Starr, R. *et al.* Liver degeneration and lymphoid deficiencies in mice lacking suppressor of
416 cytokine signaling-1. *Proc. Natl. Acad. Sci. U. S. A.* **95,** 14395–9 (1998).

417 30. Alexander, W. S. *et al.* SOCS1 is a critical inhibitor of interferon gamma signaling and prevents the
418 potentially fatal neonatal actions of this cytokine. *Cell* **98,** 597–608 (1999).

9

419   31.   Yoshida, T. *et al.* SOCS1 is a suppressor of liver fibrosis and hepatitis-induced carcinogenesis. *J.*
420          *Exp. Med.* **199,** 1701–7 (2004).

421   32.   Horino, J. *et al.* Suppressor of cytokine signaling-1 ameliorates dextran sulfate sodium-induced
422          colitis in mice. *Int. Immunol.* **20,** 753–62 (2008).

423   33.   Bourdeau, A. *et al.* TC-PTP-deficient bone marrow stromal cells fail to support normal B
424          lymphopoiesis due to abnormal secretion of interferon-{gamma}. *Blood* **109,** 4220–8 (2007).

425   34.   You-Ten, K. E. *et al.* Impaired bone marrow microenvironment and immune function in T cell
426          protein tyrosine phosphatase-deficient mice. *J. Exp. Med.* **186,** 683–93 (1997).

427   35.   Wiede, F., Sacirbegovic, F., Leong, Y. A., Yu, D. & Tiganis, T. PTPN2-deficiency exacerbates T
428          follicular helper cell and B cell responses and promotes the development of autoimmunity. *J.*
429          *Autoimmun.* **76,** 85–100 (2017).

430   36.   Kilpinen, H. *et al.* Common genetic variation drives molecular heterogeneity in human iPSCs.
431          *Nature* **546,** 370–375 (2017).

432   37.   Okada, Y. *et al.* Meta-analysis identifies nine new loci associated with rheumatoid arthritis in the
433          Japanese population. *Nat. Genet.* **44,** 511–6 (2012).

434

**Figure Legends**

**Figure 1. Description of the immunodeficiency cohort and disease associations in coding regions. (a)** Number of index cases recruited under different phenotypic categories (red – adult cases, blue – paediatric cases). **(b)** Number of index cases with malignancy, autoimmunity and CD4+ lymphopenia. (black bar – total number of cases, blue bar - number of cases with AD-PID phenotype). **(c)** Number of patients with reported genetic findings subdivided by gene. Previously reported variants are those identified as immune disease-causing in the HGMD-Pro database. **(d)** Pie charts showing proportions of the germline p.Arg328* stop-gain variant and different somatic reversions in FACS-sorted blood cell populations from a male adult patient with an inherited *IL2RG* mutation that causes X-linked infantile fatality. **(e)** BeviMed assessment of enrichment for candidate disease-causing variants in individual genes, in the PID cohort relative to the rest of NBR-RD cohort. The top candidate genes (with BeviMed PPA>=0.18) are shown. Named genes are those in which the variants driving the association have been confirmed to be causal.

**Figure 2. Assessment of WGS data for regulatory region deletions that impact upon PID (a)** Schematic overview of configurations of large deletions and putative damaging variants that could lead to gene loss of function. **(b)** Flow-chart demonstrating filtering steps to prioritise patients with candidate compound heterozygous causal variants comprising of a rare (gnomAD v1 AF<0.001) damaging (CADD>20) coding variant within a known PID gene, and a structural deletion event (with internal MAF<0.03) over the gene's regulatory region. **(c)** Genomic configuration of the *ARPC1B* gene locus highlighting the compound heterozygous gene variants. ExAC shows that the non-coding deletion is outside of the exome-targeted regions. **(d)** Pedigree of patient in (c) and co-segregation of *ARPC1B* genotype (wt – wild-type, del – deletion, fs – frameshift). **(e)** ARPC1A and ARPC1B protein levels in neutrophils and platelets in the patient depicted in (c). **(f)** Histogram showing *ARPC1B* mRNA levels in patient depicted in (c), her sibling highlighted in (d), and healthy control. **(g)** Allele-specific expression assay showing the ratio of wt, del and fs alleles in genomic DNA (gDNA) from peripheral blood mononuclear cells of the patient and sibling. **(h)** Relative expression of *ARPC1B* mRNA from each allele in the patient and sibling. Allele-specific expression assessed in complementary DNA (cDNA; synthesized from pre-mRNA).

**Figure 3. Telomere lengths calculated from whole-genome data can be used to identify causal rare and common genomic variants associated with telomere variation. (a)** Telomerecat calculated telomere lengths (TLs) against age and sex in 3,313 NBR-RD recruited subjects. The Boxplot summarises the distribution of TLs within an age and gender bin; the lower, mid and upper box bounds represent the first, second (median) and third quartile respectively. Lines extend to 1.5 times the interquartile range, and outliers are marked as individual points. **(b)** Centiles of telomere lengths against age in PID cases. Symbols represent subjects with rare genomic homozygous/hemizygous single nucleotide variants (SNV) in *TERT* and *DKC1*. **(c)** Top: Pedigree of individuals with *DKC1* variants showing co-segregation with disease phenotypes. The four individuals assayed by Flow-FISH are marked by dotted line. Bottom: Flow-FISH assessment of telomere length in *DKC1* variant carrying siblings and their spouses in granulocytes and lymphocytes.

**Figure 4. Antibody deficiency (AD-PID) GWAS identifies common variants that mediate disease risk and suggests novel monogenic candidate genes. (a)** A composite Manhattan plot for the AD-PID GWAS. Blue – common variants (MAF>0.05) analysed in this study (NBR-RD) only, red – meta-analysed with

11

480   data from Li *et al.*; and purple – genome-wide significant low frequency (0.005<MAF<0.05) variants in
481   *TNFRSF13B* locus. Loci of interest are labelled with putative causal protein coding gene names. **(b)**
482   Protein modelling of two independent MHC locus signals: residue E71 on HLA-DRB1*1501 and residue
483   N114 on HLA-B*0801 using PDB 1BX2 and PDB 4QRQ respectively. Protein is depicted in white,
484   highlighted residue in red, and peptide is in green. **(c)** Immune mediated trait enrichment of AD-PID
485   association signals.  CAD – coronary artery disease, CRO – Crohn's disease, RA – rheumatoid arthritis,
486   SLE – systemic lupus erythematosus, T1D – type 1 diabetes, T2D – type 2 diabetes and UC – ulcerative
487   colitis (**See Extended Data Table 1**). **(d)** COGS prioritisation scores of candidate monogenic causes of
488   PID using previous autoimmune targeted genotyping studies (See Supplementary Table 6) across
489   suggestive AD-PID loci (n=4). For clarity, only diseases prioritising one or more genes are shown. CEL –
490   coeliac disease, CRO- Crohn's disease, UC – ulcerative colitis, MS – multiple sclerosis, PBC – primary
491   biliary cirrhosis and T1D – type 1 diabetes **(e)** T cells from the *SOCS1* mutation patient and healthy
492   control were cultured following TCR/CD28 stimulation in the presence of anti-IFN-γ and anti-IFN-γR
493   antibodies. At day 4 post-stimulation cells were washed and re-cultured without IFN-γ blockade. At day
494   6 cells were stimulated for 2 hours with IFN-γ and protein-lysates assessed for the indicated protein
495   expression. (Left) Representative western blot. (Right) The pSTAT1 and SOCS1 levels calculated from
496   image quantification of the western blots in 4 replicate samples. Error bars represent standard error of
497   mean. **(f)** The pedigree of the CVID patient identified with a premature stop mutation in *PTPN2*. Carriers
498   of the rs2847297-G risk allele are indicated. **(g)** Simplified model of how SOCS1 and TC-PTP limit the
499   phosphorylated-STAT1 triggered by interferon signalling. **(h)** T cells from the indicated members of the
500   *PTPN2* pedigree, 3 healthy controls, the *SOCS1* mutation patient and a *STAT1* gain of function (GOF)
501   patient were cultured for 4 days and treated +/- IFN-γ for 2 hours and protein-lysates assessed for
502   protein levels. (Left) PTPN2 protein levels normalised to Tublin level (loading control). (Right) pSTAT1
503   protein levels normalised to total STAT1 level. **(i)** Relative expression from each allele of the *PTPN2*
504   rs2847297 locus in the sibling II.3 of the CVID patient II.1 in (f). Shown are the proportion of directly
505   genotyped individual bacterial colonies, transformed with the PCR product containing the rheumatoid
506   arthritis risk allele rs2847297-G generated from either gDNA or cDNA.

507

508 **Methods**

509

510 <u>PID cohort</u>

511 The PID patients and their family members were recruited by specialists in clinical immunology across 26
512 hospitals in the UK, and one each from the Netherlands, France and Germany. The recruitment criteria
513 were intentionally broad, and included the following: clinical diagnosis of common variable
514 immunodeficiency disorder (CVID) according to internationally established criteria (**Supplementary**
515 **Table 1**); extreme autoimmunity; or recurrent and/or unusual severe infections suggestive of defective
516 innate or cell-mediated immunity. Patients with known secondary immunodeficiencies caused by cancer
517 or HIV infection were excluded. Although screening for more common and obvious genetic causes of PID
518 prior to enrolment into this WGS study was encouraged, it was not a requirement. Consequently, a
519 minority of patients (16%) had some prior genetic testing, from single gene Sanger sequencing or MLPA
520 to a gene panel screen.

521 To expedite recruitment a minimal clinical dataset was required for enrolment, though more detail was
522 often provided. There was a large variety in patients' phenotypes, from simple "chest infections" to
523 complex syndromic features, and the collected phenotypic data of the sequenced individuals ranged
524 from assigned disease category only to detailed clinical synopsis and immunophenotyping data. The
525 clinical subsets used to subdivide PID patients were based on ESID definitions, as shown in
526 **Supplementary Table 1**.

527 To facilitate analysis by grouping patients with a degree of phenotypic coherence while excluding some
528 distinct and very rare clinical subtypes of PID that may have different aetiologies, a group of patients
529 was determined to have antibody deficiency-associated PID (AD-PID). This group comprised 733 of the
530 886 unrelated index cases, and included all patients with CID, CVID or Antibody Defect ticked on the
531 recruitment form, together with patients requiring IgG replacement therapy and those with specified
532 low levels of IgG/A/M. SCID patients satisfying these criteria were not assigned to the AD-PID cohort.

533

534 <u>WGS data processing</u>

535 Details of DNA sample processing, whole genome sequencing, data processing pipeline, quality checks,
536 alignment and variant calling, ancestry and relatedness estimation, variant normalisation and
537 annotation, large deletion calling and filtering, and allele frequency calculations, are fully described in
538 [NIHR BioResource, in preparation; see Cover Letter]. Briefly, DNA or whole blood EDTA samples were
539 processed and quality checked according to standard laboratory practices and shipped on dry ice to the
540 sequencing provider (Illumina Inc, Great Chesterford, UK). Illumina Inc performed further QC array
541 genotyping, before fragmenting the samples to 450bp fragments and processing with the Illumina
542 TruSeq DNA PCR-Free Sample Preparation kit (Illumina Inc., San Diego, CA, USA). Over the three-year
543 duration of the sequencing phase of the project, different instruments and read lengths were used: for
544 each sample, either 100bp reads on three HiSeq2500 lanes; or 125bp reads on two HiSeq2500 lanes; or
545 150bp reads on a single HiSeq X lane. Each delivered genome had a minimum 15X coverage over at least
546 95% of the reference autosomes. Illumina performed the alignment to GRCh37 genome build and
547 SNV/InDel calling using their Isaac software, while large deletions were called with their Manta and
548 Canvas algorithms. The WGS data files were received at the University of Cambridge High Performance
549 Computing Service (HPC) for further QC and processing by our Pipeline team.

550 For each sample, we estimated the sex karyotype and computed pair-wise kinship coefficients using
551 PLINK, which allowed us to identify sample swaps and unintended duplicates, assign ethnicities,
552 generate networks of closely related individuals (sometimes undeclared relatives from across different
553 disease domains) and a maximal unrelated sample set (for the purposes of allele frequency estimation
554 and control dataset in case-control analyses). Variants in the gVCF files were normalised and loaded into

13

555 an HBase database, where Overall Pass Rate (OPR) was computed within each of the three read length
556 batches, and the lowest of these OPR values (minOPR) assigned to each variant.

557 Large deletions were merged and analysed collectively, as described in [NIHR BioResource, in
558 preparation]. The analyses presented here are based on SNVs/InDels with OPR>0.98, and a set of
559 deletions found through the SVH method to have high specificity after extensive manual inspection of
560 individual deletion calls. Variants were annotated with Sequence Ontology terms according to their
561 predicted consequences, their frequencies in other genomic databases (gnomAD, UK10K, 1000
562 Genomes), if they have been associated with a disease according to the HGMD Pro database, and
563 internal metrics (AN, AC, AF, OPR).

564

### Diagnostic reporting

566 We screened all genes in the IUIS 2015 classification for potentially causal variants. SNVs and small
567 InDels were filtered based on the following criteria: OPR>0.95; having a protein-truncating consequence,
568 gnomAD AF<0.001 and internal AF<0.01; or present in the HGMD Pro database as DM variant. Large
569 deletions called by both Canvas and Manta algorithms, passing standard Illumina quality filters,
570 overlapping at least one exon, and classified as rare by the SVH method were included in the analysis. In
571 order to aid variant interpretation and consistency in reporting, phenotypes were translated into Human
572 Phenotype Ontology (HPO) terms as much as possible. Multi-Disciplinary Team (MDT) then reviewed
573 each variant for evidence of pathogenicity and contribution to the phenotype, and classified them
574 according to the American College of Medical Genetics (ACMG) guidelines[38]. Only variants classified as
575 Pathogenic or Likely Pathogenic were systematically reported, but individual rare (gnomAD AF<0.001) or
576 novel missense variants that BeviMed analysis (see below) highlighted as having a posterior probability
577 of pathogenicity >0.2 were additionally considered as Variants of Unknown Significance (VUS). If the
578 MDT decided that they were likely to be pathogenic and contribute to the phenotype, they were also
579 reported and counted towards the overall diagnostic yield. All variants and breakpoints of large
580 deletions reported in this study were confirmed by Sanger sequencing using standard protocols.

581

### BeviMed

583 We used BeviMed[1] to evaluate the evidence for association between case/control status and rare
584 variant allele counts in each gene. We inferred a posterior probability of association (PPA) under
585 Mendelian inheritance models (dominant and recessive), and different variant selection criteria
586 ("moderate" and "high" impact variants based on functional consequences predicted by the Variant
587 Effect Predictor[39]). All genes were assigned the same prior probability of association with the disease of
588 0.01, regardless of their previously published associations with an immune deficiency phenotype. Genes
589 for which BeviMed inferred a PPA to be >=0.18 when summed over all four combinations of inheritance
590 model and variant selection criteria (each configuration being given a prior probability of association of
591 0.0025) are shown in **Fig. 1f**. Given each of the association models, the posterior probability that each
592 variant is pathogenic is also computed. We used a variant-level posterior probability of pathogenicity
593 >0.2 to select potentially pathogenic missense variants in known PID genes to report back.

594

### Telomerecat

596 Average telomere length was calculated from whole-genome sequence data using Telomerecat, as
597 reported previously[22]. Batch differences caused by changes in sequencing platform differences were
598 normalised by using a linear model. The linear model was defined as:

$$length = \beta_0 + \beta_1 batch_2 + \beta_2 batch_3 \ldots \beta_N batch_N + \epsilon$$

599

600 where β are regression coefficients, and batch represents a dummy variable denoting the plate a sample
601 was sequenced on. For each plate the relevant coefficient was subtracted from all of the observed
602 telomere lengths within each plate.

603 After adjusting for batch effects, telomere length was compared to age in 3,313 NBR-RD subjects. We
604 obtained a strong negative correlation with age (r = −0.56, Pearson's correlation), thus validating
605 Telomerecat as a reliable method for estimating telomere lengths. We found that each year of
606 additional age was equivalent to a 33bp deterioration in telomere length (**Supplementary Fig. 6**).
607 Although this observed negative correlation is well established within the literature, we obtain a
608 particularly high correlation with our method, which could be partly driven by the wide age range of our
609 sample set.

610 To normalise telomere lengths for comparison of samples from disparate age and gender, the following
611 linear model was fitted to the data using age as a continuous variable and gender as a dummy variable:

$$length = \beta_0 + \beta_1 age + \beta_2 age^2 + \beta_3 age^3 + \beta_4 gender + \epsilon$$

612

613 The relevant residuals produced by the cubic model were subtracted from the mean telomere length of
614 the cohort. These adjusted telomere lengths were used in the GWAS analysis.

615 To assess for monogenic causes of telomere shortening, subjects were identified within the PID cohort
616 that had telomere lengths below the 10th centile of age adjusted values and had hemizygous or
617 homozygous SNVs that occurred gnomAD AF<0.001 in *TERC*, *TERT*, *NHP2*, *TINF2*, *NOP10*, *PARN*, *ACD*,
618 *WRAP53*, *CTC1*, *RTEL1* or *DKC1* genes.

619

620 AD-PID GWAS

621 GWAS was performed both on the whole PID cohort (N cases = 886) and on a subset of AD-PID cases (N
622 cases = 733); here we present the results of the latter analysis, which was cleaner and less noisy despite
623 a reduced sample size. We used 9225 unrelated samples from non-PID NBR-RD cohorts as controls.

624 Variants were selected from a merged VCF file were filtered to include bi-allelic SNPs with overall
625 MAF>=0.05 and minOPR=1 (100% pass rate). We ran PLINK logistic association test under an additive
626 model using the read length, sex, and first 10 principal components from the ethnicity analysis as
627 covariates. After filtering out SNPs with HWE p<10$^{-6}$, we were left with the total of 4,993,945 analysed
628 SNPs. There was minimal genomic inflation of the test statistic (lambda = 1.027), suggesting population
629 substructure and sample relatedness had been appropriately accounted for. The only genome-wide
630 significant (p<5x10$^{-8}$) signal was at the MHC locus, with several suggestive (p<1x10$^{-5}$) signals
631 (**Supplementary Fig. 7**). We repeated the analysis with more relaxed SNP filtering criteria using
632 MAF>=0.005 and minOPR>0.95. The only additional signal identified were the three *TNFRSF13B* variants
633 shown in **Extended Data Fig. 1**.

634 We obtained summary statistics data from the Li et al. CVID Immunochip case-control study[9] and
635 performed a fixed effects meta-analysis on 95,417 variants shared with our AD-PID GWAS. For each of
636 the genome-wide and suggestive loci after meta-analysis, we conditioned on the lead SNP by including it
637 as an additional covariate in the logistic regression model, to determine if the signal is driven by the
638 single or multiple hits at those loci. Only the MHC locus showed evidence of multiple independent
639 signals (**Supplementary Fig. 8**).

640

641 MHC locus imputation

642 We imputed classical HLA alleles using the method implemented in the SNP2HLA v1.0.3 package[40],
643 which uses Beagle v3.0.4 for imputation and the HapMap CEU reference panel. We imputed allele

15

644 dosages and best-guess genotypes of 2-digit and 4-digit classical HLA alleles, as well as amino acids of
645 the MHC locus genes *HLA-A*, *HLA-B*, *HLA-C*, *HLA-DRB1*, *HLA-DQA1* and *HLA-DQB1*. We tested the
646 association of both allele dosages and genotypes using the logistic regression implemented in PLINK,
647 and obtained similar results. We then used the best-guess genotypes to perform the conditional analysis
648 in PLINK, since conditioning is not implemented in a model with allele dosages.

649

650 <u>Allele Specific Expression</u>

651 RNA and gDNA were extracted from PBMCs using the AllPrep kit (Qiagen) as per the manufacturer's
652 instructions. RNA was reverse transcribed to make cDNA using the SuperScript™ VILO™ cDNA synthesis
653 kit with appropriate minus reverse transcriptase controls, as per the manufacturer's instructions. The
654 region of interest in the gDNA and 1:10 diluted cDNA was amplified using Phusion (Thermo Fisher) and
655 the following primers on a G-Storm thermal cycler with 30 seconds at 98$^o$C then 35 cycles of 98$^o$C 10
656 seconds, 60$^o$C 30 seconds, 72$^o$C 15 seconds.

657 ***ARPC1B***

658 The region of interest spanning the frameshift variant was amplified using the following primers:
659 Forward: GGGTACATGGCGTCTGTTTC / Reverse: CACCAGGCTGTTGTCTGTGA

660 PCR products were run on a 3.5% agarose gel. Bands were cut out and product extracted using the QIA
661 Quick Gel Extraction Kit (Qiagen), as per protocol. Expected products were confirmed by Sanger
662 sequencing. 4ul fresh PCR product was used in a TOPO®cloning reaction (Invitrogen) and used to
663 transform One Shot™ TOP10 chemically competent E. coli. These were cultured overnight then spread
664 on LB agar plates. Individual colonies were picked and genotyped. ARPC1B mRNA expression was
665 assessed using a Taqman gene expression assay with 18S and EEF1A1 as control genes. Each sample was
666 run in triplicate for each gene with a no template control. PCR was run on a LightCycler® (Roche) with 2
667 mins 50$^o$C, 20 seconds 95$^o$C then 45 cycles of 95$^o$C 3 seconds, 60$^o$C 30 seconds.

668 ***PTPN2***

669 PTPN2 ASE protocol is modified from above. RNA and genomic DNA were extracted from PBMCs using
670 the AllPrep Kit (Qiagen). RNA was treated with Turbo DNAse (Thermo) and reverse transcribed to
671 generate cDNA using the SuperScript IV VILO master mix (Thermo). The intronic region of interest in
672 gDNA and cDNA was amplified by two nested PCR reactions using Phusion enzyme (Thermo). The
673 primers (F1/R1) and nested primers (F2/R2) used were:
674 Forward_1: aaagtctggagcaggcagag / Reverse_1: tgggggaactggttatgctttc
675 Forward_2: ggagctatgatcacgccacatg / Reverse_2: atgctttctggttgggctgac
676

677 PCR products were run on a 1% agarose gel. Bands were cut out and product extracted using the QIA
678 Quick Gel Extraction Kit (Qiagen), as per protocol. Expected products were confirmed by Sanger
679 sequencing. 5ng fresh PCR product was used in a TOPO®cloning reaction (Invitrogen) and used to
680 transform One Shot™ TOP10 chemically competent E. coli. These were cultured overnight then spread
681 on LB agar plates. Individual colonies were picked and genotyped. PTPN2 mRNA expression was
682 assessed using a Taqman SNP genotyping assay and on a LightCycler (Roche).
683

684 <u>PAGE and Western Blot analysis</u>

685 Samples were separated by SDS polyacrylamide gel electrophoresis and transferred onto a nitrocellulose
686 membrane. Individual proteins were detected with antibodies against ARPC1b (goat polyclonal
687 antibodies, ThermoScientific, Rockford, IL, USA), against ARPC1a (rabbit polyclonal antibodies, Sigma, St
688 Louis, USA) and against actin (mouse monoclonal antibody, Sigma). Secondary antibodies were either
689 donkey-anti-goat-IgG IRDye 800CW, Goat-anti-mouse-IgG IRDye 800CW or Donkey-anti-rabbit-IgG IRDye

690    680CW (LI-COR Biosciences, Lincoln, NE, USA). Quantification of bound antibodies was performed on an
691    Odyssey Infrared Imaging system (LI-COR Biosciences, Lincoln, NE, USA).

692

693    Phasing of *SOCS1* variants

694    To phase common rs2286974 variant with the novel stop-gain *SOCS1* variant (chr16:11348854
695    T>TGCGGC) identified in the same patient, we performed long-read WGS with Oxford Nanopore
696    Technologies (ONT). The sample was prepared using the 1D ligation library prep kit (SQK-LSK108), and
697    genomic libraries were sequenced on R9.4 flowcells. Sequencing was carried out on GridION system,
698    read sequences were extracted from base-called FAST5 files by Guppy (v0.5.1) to generate FASTQ files,
699    which were then aligned against the GRCh37/hg19 human reference genome using minimap2 (v2.2).
700    Four runs were performed in order to reach an average coverage of 14x, with a median read length of
701    5006 ± 3981. Haplotyping and genotyping was performed with MarginPhase.

702

703    Structural deletion analysis

704    Structural (length >50bp) deletions (MAF>0.03) were called as previously described[41]. For all
705    downstream analysis we used gencode  v26 annotations downloaded from
706    [ftp://ftp.sanger.ac.uk/pub/gencode/Gencode_human/release_26/GRCh37_mapping/gencode.v26lift37
707    .annotation.gtf.gz]. We defined promoters as a window +/- 500bp of any protein coding gene
708    transcriptional start site (TSS). In order to associate cis regulatory elements (cRE) with putative target
709    genes we combined by physical location overlap, super enhancer cRE annotations from [18], with
710    promoter capture Hi-C (pcHi-C) from [19], matching by tissue. We next computed the overlap of structural
711    variants occurring in the PID cohort with cREs for which putative target genes were available. We
712    classified overlaps between deletions and functional annotations into three non-mutually exclusive
713    categories; `prom' - overlaps focal gene promoter, `exon' - overlaps focal gene exon, `pse' - overlaps
714    Hnisz *et al.*[18] SE annotation linked to focal gene by pcHi-C. We compiled a catalogue of compound
715    heterozygous deletions where there was evidence in the same individual for a damaging (CADD>20) rare
716    (gnomAD AF<0.001) variant within the same gene.

717

718    AD-PID GWAS Enrichment

719    Due to the size of the AD-PID cohort, we were unable to use LD-score regression[42] to assess genetic
720    correlation between distinct and related traits. We therefore adapted the previous enrichment method
721    `blockshifter`[43] in order to assess evidence for the enrichment of AD-PID association signals in a
722    compendium of 9 GWAS European Ancestry summary statistics was assembled from publicly available
723    data. We removed the MHC region from all downstream analysis [GRCh37 chr6:25-45Mb]. To adjust for
724    linkage disequilibrium (LD), we split the genome into 1cM recombination blocks based on HapMap
725    recombination frequencies [44]. For a given GWAS trait, for $n$ variants within LD block $b$ we used
726    Wakefield's synthesis of asymptotic Bayes factors (aBF)[45] to compute the posterior probability that the
727    $i^{th}$ variant is causal ($PPCV_i$) under single causal variant assumptions[46] :

728

$$PPCV_i = \frac{aBF_i \pi_i}{\sum_{j=1}^{n}(aBF_j \pi_j) + 1}$$

730    Here $\pi_i = \pi_j$ are flat prior probabilities for a randomly selected variant from the genome to be causal
731    and we use the value 1x10^{-4} [47]. We sum over these PPCV within an LD block, $b$ to obtain the posterior
732    probability that $b$ contains a single causal variant (PPCB).

733    To compute enrichment for trait $t$, we convert PPCBs into a binary label by applying a threshold such
734    that $PPCB_t > 0.95$. We apply these block labels for trait $t$, to PPCBs (computed as described above) for
735    our AD-PID cohort GWAS, using them to compute a non-parametric Wilcoxon rank sum statistic, W
736    representing the enrichment. Whilst the aBF approach naturally adjusts for LD within a block, residual
737    LD between blocks may exist. In order to adjust for this and other confounders (e.g. block size) we use a
738    circularised permutation technique[48] to compute $W_{null}$. To do this, for a given chromosome, we select
739    recombination blocks, and circularise such that beginning of the first block adjoins the end of the last.
740    Permutation proceeds by rotating the block labels, but maintaining AD-PID PPCB assignment. In this way
741    many permutations of $W_{null}$ can be computed whilst conserving the overall block structure.

742    For each trait we used $10^4$ permutations to compute adjusted Wilcoxon rank sum scores using *wgsea*
743    [https://github.com/chr1swallace/wgsea] R package.

744

745    <u>PID monogenic candidate gene prioritisation</u>

746    We hypothesised, given the genetic overlap with antibody associated PID, that common regulatory
747    variation, elucidated through association studies of immune-mediated disease, might prioritise genes
748    harbouring damaging LOF variants underlying PID.  Firstly, using summary statistics from our combined
749    fixed effect meta-analysis of AD-PID, we compiled a list of densely genotyped ImmunoChip regions
750    containing one or more variant where $P<1x10^{-5}$. Next, we downloaded ImmunoChip (IC) summary
751    statistics from ImmunoBase (accessed 30/07/2018) for all 11 available studies.  For each study we
752    intersected PID suggestive regions, and used COGS (https://github.com/ollyburren/rCOGS) in
753    conjunction with promoter-capture Hi-C datasets for 17 primary cell lines[19,43] in order to prioritise genes.
754    We filtered by COGS score to select genes with a COGS score >0.5 [19,43] to obtain a list of 11 protein
755    coding genes.

756    We further hypothesised that genes harbouring rare LOF variation causal for PID would be intolerant to
757    variation. We thus downloaded pLI scores[49] and took the product between these and the COGS scores
758    to compute an `overall' prioritisation score across each trait and gene combination. We applied a final
759    filter taking forward only those genes having an above average `overall' score to obtain a final list of 6
760    candidate genes (Fig. 4d).  Finally, we filtered the cohort for damaging rare (gnomAD AF<0.001) protein-
761    truncating variants (frameshift, splice-site, nonsense) within these genes in order to identify individuals
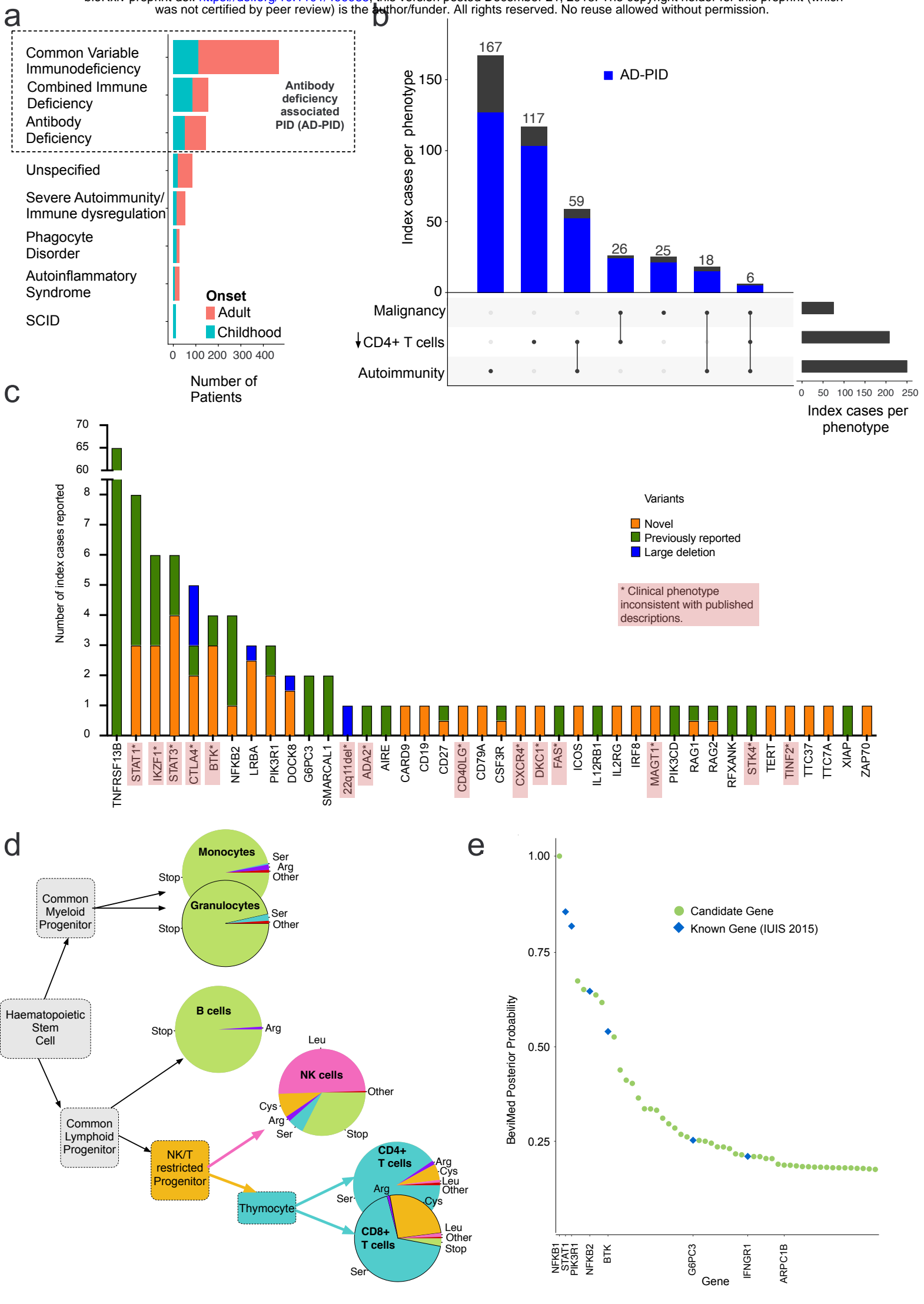762    for functional follow up.

763

764    <u>Statistical analysis</u>

765    Statistical analysis was carried out using R (3.3.3 – "Another Canoe") and Graphpad Prism (version 7)
766    unless otherwise stated. R code for running major analyses are available at
767    https://github.com/ollyburren/pid_thaventhiran_et_al.
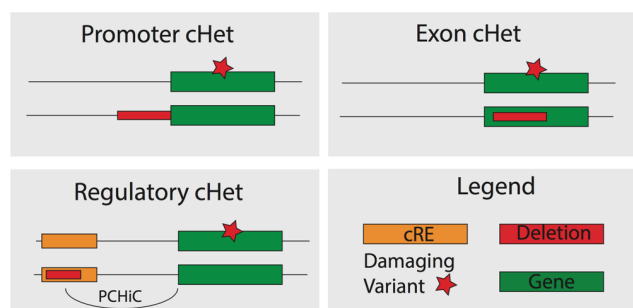
768

769    **Methods References**
770

771    38.    Richards, S. *et al.* Standards and guidelines for the interpretation of sequence variants: a joint
772           consensus recommendation of the American College of Medical Genetics and Genomics and the
773           Association for Molecular Pathology. *Genet. Med.* **17,** 405–423 (2015).
774    39.    McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17,** 122 (2016).
775    40.    Jia, X. *et al.* Imputing amino acid polymorphisms in human leukocyte antigens. *PLoS One* **8,**
776           e64683 (2013).
777    41.    Carss, K. J. *et al.* Comprehensive Rare Variant Analysis via Whole-Genome Sequencing to
778           Determine the Molecular Pathology of Inherited Retinal Disease. *Am. J. Hum. Genet.* **100,** 75–90
779           (2017).
780    42.    Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in
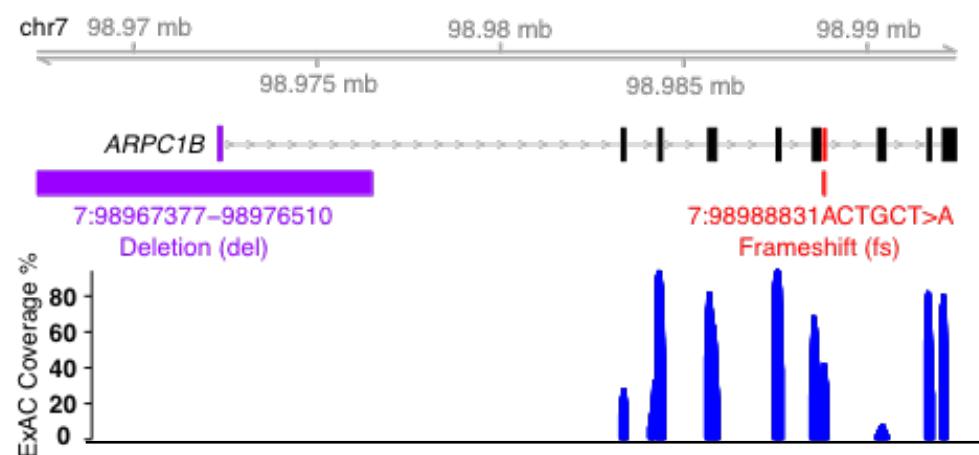
781           genome-wide association studies. *Nat. Genet.* **47,** 291–5 (2015).

782    43.   Burren, O. S. *et al.* Chromosome contacts in activated T cells identify autoimmune disease
783           candidate genes. *Genome Biol.* **18,** 165 (2017).

784    44.   International HapMap Consortium, T. I. H. *et al.* A second generation human haplotype map of
785           over 3.1 million SNPs. *Nature* **449,** 851–61 (2007).

786    45.   Wakefield, J. Bayes factors for genome-wide association studies: comparison with *P* -values.
787           *Genet. Epidemiol.* **33,** 79–86 (2009).

788    46.   Wellcome Trust Case Control Consortium, J. B. *et al.* Bayesian refinement of association signals
789           for 14 loci in 3 common diseases. *Nat. Genet.* **44,** 1294–301 (2012).

790    47.   Huang, H. *et al.* Fine-mapping inflammatory bowel disease loci to single-variant resolution.
791           *Nature* **547,** 173–178 (2017).

792    48.   Trynka, G. *et al.* Disentangling the Effects of Colocalizing Genomic Annotations to Functionally
793           Prioritize Non-coding Variants within Complex-Trait Loci. *Am. J. Hum. Genet.* **97,** 139–52 (2015).

794    49.   Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536,** 285–91
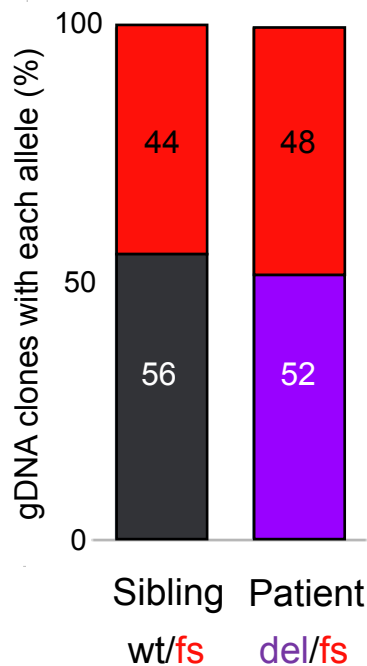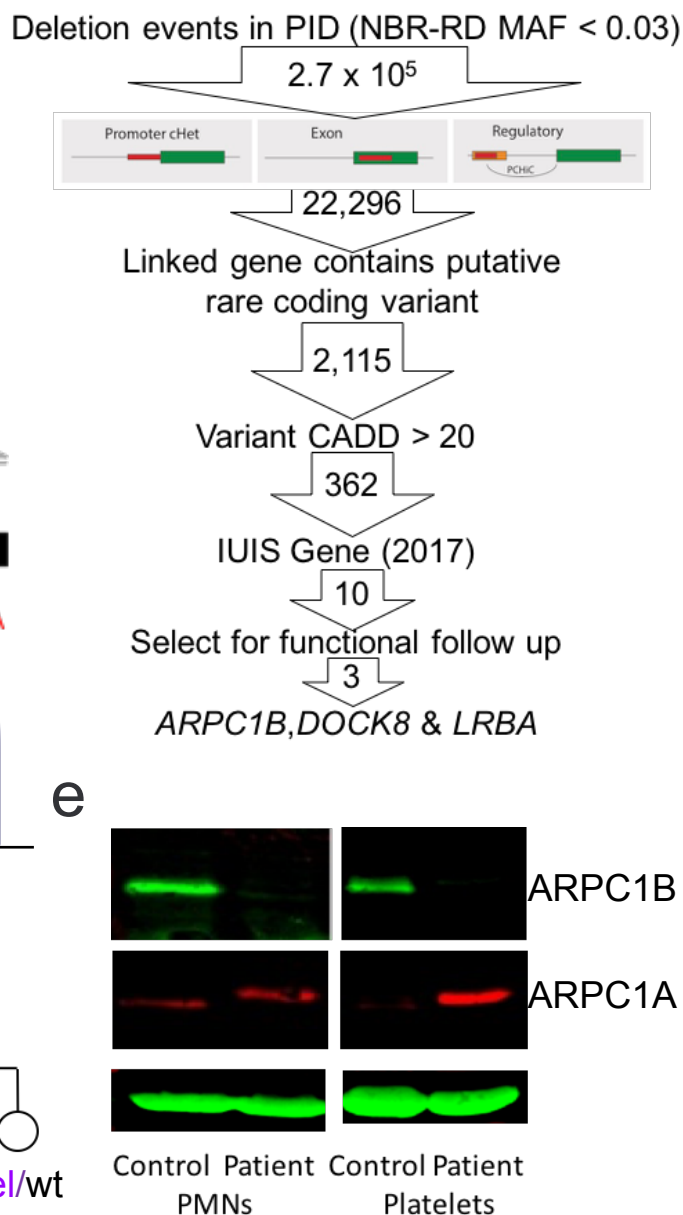795           (2016).

**a**

Promoter cHet    Exon cHet

Regulatory cHet    Legend

cRE    Deletion
Damaging
Variant    Gene

**b**

Deletion events in PID (NBR-RD MAF < 0.03)

$2.7 \times 10^5$

Promoter cHet    Exon    Regulatory

PCHiC

22,296

Linked gene contains putative
rare coding variant

2,115

Variant CADD > 20

362

IUIS Gene (2017)

10

Select for functional follow up

3

*ARPC1B, DOCK8 & LRBA*

**c**

chr7    98.97 mb    98.98 mb    98.99 mb

98.975 mb    98.985 mb

*ARPC1B*

7:98967377–98976510
Deletion (del)

7:98988831ACTGCT>A
Frameshift (fs)

ExAC Coverage %

**d**

wt/fs    del/wt

Patient    Sibling

del/fs    wt/wt    wt/wt    wt/fs    wt/wt    wt/wt    wt/fs    del/wt

**e**

ARPC1B

ARPC1A

Control Patient    Control Patient
PMNs    Platelets

**f**

ARPC1B mRNA
$2^{-\Delta\Delta Ct}$

Control    Sibling    Patient
wt/wt    wt/fs    del/fs

**g**

gDNA clones with each allele (%)

44    48

56    52

Sibling    Patient
wt/fs    del/fs

**h**

ARPC1B mRNA from each allele
$2^{-\Delta\Delta Ct}$

10

90    71

29

Sibling    Patient
wt/fs    del/fs

**a**

**b**

Homozygous
TERT SNV

Hemizygous
DKC1 SNVs

Centile
99th
90th
50th
10th
1st

**c**

Phenotype: Healthy    PCP    CMV retinitis    Healthy

X chromosome
DKC1:        wt/wt    R449Q    R449Q    wt/wt

Granulocytes

Lymphocytes
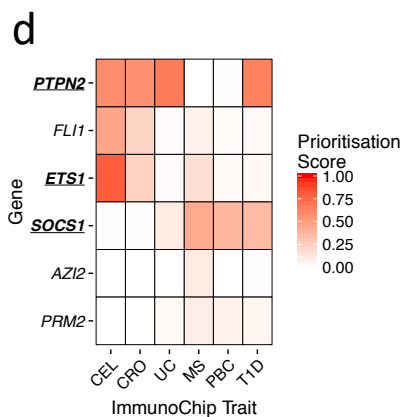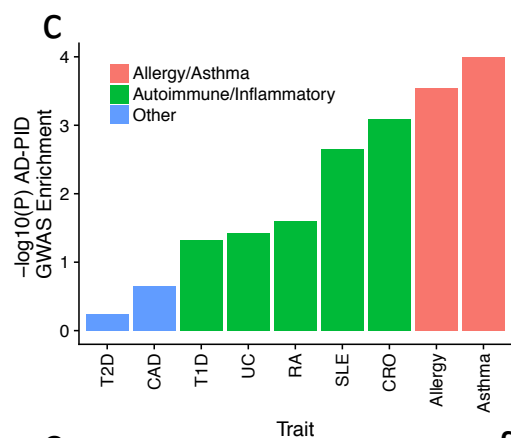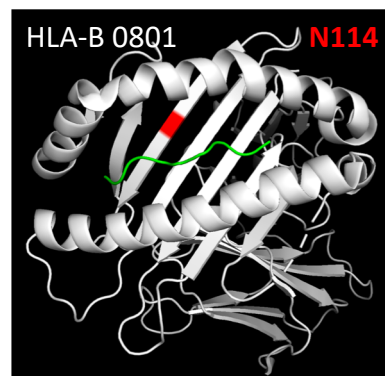
Centile
99th
90th
50th
10th
1st

a



b

c

d

e

f

g
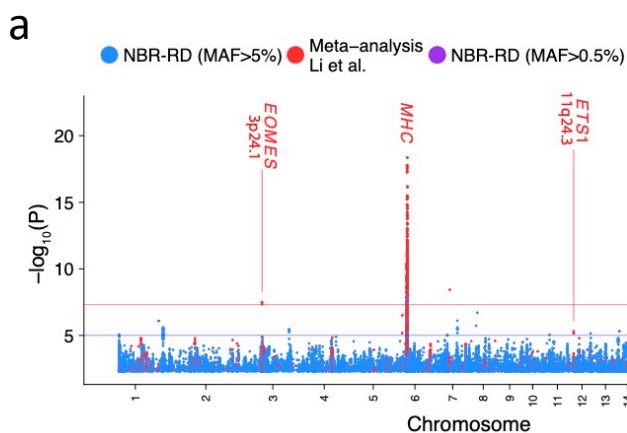
h

i

## Extended Data Items Legends

**Extended Data Table 1.** GWAS studies used in enrichment analysis

**Extended Data Fig. 1.** Regional AD-PID association plot of 17p11.2 (*TNFRSF13B/TACI)* region. Tracks are as follows: **AD-PID -log10(P)** dot plot of AD-PID association, index SNP ( is purple, others are coloured based on LD information from UK10K project with red indicating high LD (r^2>0.9), blue low (r^2<0.2) and grey where no information available. **Gene** - Cannonical gene annotation (Ensembl V75), **GWAS** - location of index variants from other immune-mediated disease, **CD4**, **B**, **Mon** putative regulatory regions in CD+T-cells, Total B cells and Monocytes computed from the union of ATAC-Seq and H3K27ac ChIP-Seq data, **pcHi-C** - Promoter Capture Hi-C interactions, in above primary cell types.

**Extended Data Fig. 2.** Regional AD-PID meta-analysis association plot of 3p24.1 (*EOMES*) region. Tracks as described for Extended Data Fig.1 with the exception of **Meta -log(P)** which shows dot plot of AD-PID association meta-analysis with Li *et al.* Detail shows location of RA index SNP that overlaps AD-PID index variant and its promoter proximity.

**Extended Data Fig. 3.** Comparison of GWAS association signals at 18p11.21 for **Li *et al.***, **NBR-RD AD-PID** (this study), and **Meta** (Meta-analysis). Y-axis is -log10(P) of univariate association statistic. Top SNP in each study is marked in cyan.

**Extended Data Fig. 4.** Locus plot of Regional AD-PID meta-analysis association plot of 16p13.13 (*CLEC16A/SOCS1)*. Tracks as described in Extended Data Fig. 2.