

1 **Curation and Expansion of Human Phenotype Ontology for Defined Groups of Inborn**
2 **Errors of Immunity**

3

4

5 Matthias Haimel PhD^{1-3*}, Julia Pazmandi MSc^{1-3*}, Raúl Jiménez Heredia MSc¹⁻³, Jasmin
6 Dmytrus MSc¹⁻³, Sevgi Köstel Bal M.D.,PhD¹⁻³, Samaneh Zoghi PhD¹⁻³, Paul van Daele
7 M.D.⁴, Tracy A. Briggs PhD^{5,6}, Carine Wouters M.D.^{7,8}, Brigitte Bader-Meunier M.D.^{9,10},
8 Florence A. Aeschlimann M.D.^{9,10}, Roberta Caorsi M.D.¹¹, Despina Eleftheriou M.D.^{12,13},
9 Esther Hoppenreijns M.D.¹⁴, Elisabeth Salzer M.D.,PhD¹⁻³, Shahrzad Bakhtiar M.D.¹⁵, Beata
10 Derfalvi M.D.¹⁶, Francesco Saettini M.D.¹⁷, Maaïke A. A. Kusters M.D.,PhD^{12,13}, Reem Elfeky
11 M.D.^{12,13}, Johannes Trück M.D.,Phil¹⁸, Jacques G. Rivière M.D.^{19,20}, Mirjam van der Burg
12 PhD^{21,22}, Marco Gattorno M.D.¹¹, Markus G. Seidel M.D.²³, Siobhan Burns M.D.²⁴, Klaus
13 Warnatz M.D.^{25,26}, Fabian Hauck M.D.,PhD^{27,28}, Paul Brogan M.D.^{12,13}, Kimberly C. Gilmour
14 PhD¹³, Catharina Schuetz M.D.²⁹, Anna Simon M.D.,PhD³⁰, Christoph Bock PhD^{1,3,31}, Sophie
15 Hambleton PhD³², Esther de Vries M.D., PhD^{33,34}, Peter Robinson M.D.³⁵, Marielle van Gijn
16 PhD^{36,†#}, Kaan Boztug M.D.^{1-3,37,†#}

17

18

19 * and †, these authors contributed equally

20 # to whom correspondence should be addressed:

21

22 Kaan Boztug, Ludwig Boltzmann Institute for Rare and Undiagnosed Diseases (LBI-RUD)
23 and St. Anna Children's Cancer Research Institute (CCRI), Zimmermannplatz 10, A-1090
24 Vienna, kaan.boztug@rud.lbg.ac.at, Phone: +43 1-40470-4080, Fax: +43-1-40170-7280

25 Marielle Van Gijn, Department of Genetics, University Medical Center Groningen, Antonius
26 Deusinglaan 1, 9713AV Groningen, Netherlands, m.e.van.gijn@umcg.nl, +31-55256416

27

28 ¹Ludwig Boltzmann Institute for Rare and Undiagnosed Diseases, Vienna, Austria

29 ²St. Anna Children's Cancer Research Institute, Vienna, Austria

30 ³CeMM Research Center for Molecular Medicine of the Austrian Academy of Sciences,
31 Vienna, Austria

32 ⁴Department of Clinical Immunology, Erasmus University Medical Center, Rotterdam, The
33 Netherlands

34 ⁵NW Genomic Laboratory Hub, Manchester Centre for Genomic Medicine, St Mary's
35 Hospital, Manchester University NHS Foundation Trust, Manchester, United Kingdom

36 ⁶Division of Evolution and Genomic Sciences, School of Biological Sciences, University of
37 Manchester, United Kingdom.

38 ⁷Department of Microbiology and Immunology, Immunobiology, KU Leuven, Leuven,
39 Belgium

40 ⁸Department of Pediatrics, Division of Pediatric Rheumatology, University Hospitals Leuven,
41 Leuven, Belgium

42 ⁹Pediatric Immuno-Hematology and Rheumatology Unit, Necker Hospital for Sick Children -
43 AP-HP, Paris, France, EU.

44 ¹⁰Reference Center for Rheumatic, Autoimmune and Systemic Diseases in Children (RAISE),
45 Paris, France

46 ¹¹Center for Autoinflammatory diseases and Immunodeficiency, IRCCS Istituto Giannina
47 Gaslini, Genova, Italy

48 ¹²University College London Great Ormond Street Institute of Child Health, London, United
49 Kingdom

50 ¹³Department of immunology, Great Ormond Street (GOS) Hospital for Children NHS
51 Foundation Trust, London, United Kingdom

52 ¹⁴Department of Paediatric Rheumatology, Radboud University Medical Centre, Nijmegen,
53 The Netherlands

54 ¹⁵Department for Children and Adolescents, Division for Stem Cell Transplantation,
55 Immunology and Intensive Care Unit, Goethe University, Frankfurt, Germany

56 ¹⁶Department of Pediatrics, Division of Immunology, Dalhousie University/IWK Health
57 Centre Halifax, Nova Scotia, Canada

58 ¹⁷Pediatric Hematology Department, Fondazione MBBM, University of Milano Bicocca, via
59 Pergolesi 33, 20900, Monza, Italy

60 ¹⁸Division of Immunology, University Children's Hospital Zurich, Switzerland

61 ¹⁹Pediatric Infectious Diseases and Immunodeficiencies Unit, Vall d'Hebron Research
62 Institute, Hospital Universitari Vall d'Hebron, Universitat Autònoma de Barcelona, Barcelona,
63 Spain

64 ²⁰Jeffrey Model Foundation Excellence Center, Barcelona, Spain

65 ²¹Department of Immunology, University Medical Center Rotterdam, Rotterdam, The
66 Netherlands

67 ²²Laboratory for Pediatric Immunology, Department of Pediatrics, Leiden University Medical
68 Center, Leiden, The Netherlands

69 ²³Research Unit for Pediatric Hematology and Immunology, Division of Pediatric Hemato-
70 Oncology, Department of Pediatrics and Adolescent Medicine, Medical University Graz, Graz,
71 Austria

72 ²⁴Department Immunology, UCL Institute of Immunity & Transplantation, Department of
73 immunology, Royal Free Hospital NHS Foundation Trust, Pond Street, London, NW3 2QG,
74 UK

75 ²⁵Division of Immunodeficiency, Department of Rheumatology and Clinical Immunology,
76 Medical Center - University of Freiburg, Faculty of Medicine, University of Freiburg, Freiburg,
77 Germany

78 ²⁶Center for Chronic Immunodeficiency (CCI), Medical Center - University of Freiburg,
79 Faculty of Medicine, University of Freiburg, Freiburg, Germany

80 ²⁷Department of Pediatrics, Dr. von Hauner Children's Hospital, University Hospital, Ludwig-
81 Maximilians-Universität München, Munich, Germany

82 ²⁸Munich Centre for Rare Diseases (M-ZSE^{LMU}), University Hospital, Ludwig-Maximilians-
83 Universität München, Munich, Germany

84 ²⁹Department of Pediatrics, Medizinische Fakultät Carl Gustav Carus, Technische Universität
85 Dresden, Germany

86 ³⁰Radboudumc Expertise Centre for Immunodeficiency and Autoinflammation (REIA),
87 department of Internal Medicine, Radboud University Nijmegen Medical Centre, Nijmegen,
88 The Netherlands

89 ³¹Institute of Artificial Intelligence and Decision Support, Center for Medical Statistics,
90 Informatics, and Intelligent Systems, Medical University of Vienna, Vienna, Austria

91 ³²Immunity and Inflammation Theme, Translational and Clinical Research Institute, Newcastle
92 University, Newcastle upon Tyne, United Kingdom

93 ³³Tranzo, Tilburg University, Tilburg, The Netherlands

94 ³⁴Laboratory for Medical Microbiology and Immunology, Elisabeth-Tweesteden Hospital,
95 Tilburg, The Netherlands

96 ³⁵The Jackson Laboratory for Genomic Medicine, 10 Discovery Drive, Farmington, CT 06032,
97 USA

98 ³⁶Department of Genetics, University Medical Center Groningen, Groningen, The Netherlands

99 ³⁷Department of Pediatrics and Adolescent Medicine, Medical University of Vienna, Vienna,
100 Austria

101

102

103

104 **Competing interests:** The authors declare no conflict of interests.

105 **Funding:** The work was supported by the European Research Council (ERC Consolidator
106 Grant 820074 “iDysChart” to K.B. Additional financial support for the workshops was granted
107 by the Ludwig Boltzmann Institute for Rare and Undiagnosed Diseases (LBI-RUD), the
108 European Reference Network on Rare Primary Immunodeficiency, Autoinflammatory and
109 Autoimmune diseases (ERN-RITA), and the European Society for Immunodeficiencies
110 (ESID).

111 **Author contributions:** MH, JP, MVG, KB: study design and manuscript writing. MH, JP: data
112 acquisition, coordination of working groups. MH, JP analysis and interpretation of data. SH,
113 clinical cohort data extraction. KB, MVG: Study supervision. All co-authors participated in the
114 meetings and revision of terms. The manuscript was reviewed, edited and approved by all co-
115 authors.

116 **Abstract**

117 **BACKGROUND:** Accurate, detailed and standardized phenotypic descriptions are essential
118 to support diagnostic interpretation of genetic variants and to discover new diseases. The
119 Human Phenotype Ontology (HPO), extensively used in rare disease research, provides a rich
120 collection of vocabulary with standardized phenotypic descriptions in a hierarchical structure.
121 However, to date the use of HPO has not yet been widely implemented in the field of inborn
122 errors of immunity (IEIs), mainly due to a lack of comprehensive IEI-related terms.

123 **OBJECTIVES:** We sought to systematically review available terms in HPO for the depiction
124 of IEIs, to expand HPO yielding more comprehensive sets of terms, and to reannotate IEIs with
125 HPO terms to provide accurate, standardized phenotypic descriptions.

126 **METHODS:** We initiated a collaboration involving expert clinicians, geneticists, researchers
127 working on IEIs and bioinformaticians. Multiple branches of the HPO tree were restructured
128 and extended based on expert review. Our ontology-guided machine learning coupled with a
129 two-tier expert review was applied to reannotate defined subgroups of IEIs.

130 **RESULTS:** We revised and expanded four main branches of the HPO tree. Here, we
131 reannotated 73 diseases from four IUIS-defined IEI disease subgroups with HPO terms. We
132 achieved a 4.7-fold increase in number of phenotypic terms per disease. Given the new HPO
133 annotations, we demonstrated improved ability to computationally match selected IEI cases to
134 their known diagnosis, and improved phenotype-driven disease classification.

135 **CONCLUSION:** Our targeted expansion and reannotation presents enhanced precision of
136 disease annotation, will enable superior HPO-based IEI characterization and hence benefit both
137 IEI diagnostic and research activities.

138

139 **Key message**

140 HPO is a robust resource for supporting IEI diagnostics and genetics with adequate ontology

141 breadth and disease annotation depth.

142 **Capsule Summary**

143 Our newly formed expert consortium systematically reviewed and expanded existing HPO
144 terms of IEIs and reannotated IEIs with HPO terms. This will support diagnostic pipelines and
145 analysis of variants from next-generation sequencing.

146

147 **Key words**

148 HPO; ontology; phenotype; rare diseases; inborn errors of immunity; immune deficiencies;
149 disease classification; diagnostic support; patient matching; genetic analysis

150

151 **Abbreviations**

152 ALPS - Autoimmune Lymphoproliferative Syndrome

153 CVID – Common Variable Immunodeficiency Disorders

154 EBV – Epstein-Barr Virus

155 EHR - Electronic Health Record

156 ERN-RITA - European Reference Network on Rare Primary Immunodeficiency;

157 Autoinflammatory and Autoimmune diseases

158 ESID - European Society for Immunodeficiencies

159 HLH - Hemophagocytic Lymphohistiocytosis

160 HPO - Human Phenotype Ontology

161 IEI - Inborn Errors of Immunity

162 IUIS – International Union of Immunological Societies

163 ISSAID - International Society of Systemic Autoinflammatory Diseases

164 LBI-RUD - Ludwig Boltzmann Institute for Rare and Undiagnosed Diseases

- 165 OMIM – Online Mendelian Inheritance in Man
- 166 PAD – Primary Antibody Deficiencies
- 167 SCID – Severe Combined Immunodeficiency
- 168 TRAPS - Tumor necrosis factor receptor-associated periodic syndrome
- 169 UDNI - Undiagnosed Diseases Network International
- 170 UDP and UDN - Undiagnosed Disease Program and Network

171 Introduction

172

173 Rare and undiagnosed diseases pose challenges for affected patients, clinicians and researchers
174 working to improve diagnostic and therapeutic approaches. Because of the rarity, clinicians
175 often only see a few patients with specific rare phenotypes throughout their careers, leading to
176 considerable diagnostic delay (1). Genetic research on rare diseases often relies on single
177 pedigrees or a few patients, leaving many patients undiagnosed (1). Compiling a cohort of
178 patients - so-called patient matching - is often crucial to gain insight into the phenotypic
179 spectrum, natural/clinical history of the disease, and adequate monitoring and treatment
180 strategies. The rare disease community has recognized these challenges and established tools
181 enabling efficient data sharing across institutions and borders, including genetic data exchange
182 through the Matchmaker Exchange platform (2) to solve undiagnosed exomes and genomes
183 (3). These platforms however are highly dependent on accurately phenotyped and categorized
184 patients and standardized disease classifications.

185 To date, several nomenclatures and reference systems for diseases have been developed (4,5).

186 In parallel, ontologies were established to provide a more systematic, hierarchical classification
187 of diseases (6,7). However, these nomenclatures group patients by disease label and do not
188 describe the underlying phenotypic features. Consequently, clinical features, laboratory
189 measurements, anatomical and functional phenotypes of patients are often described with
190 variable quality and specificity, which hampers patient matching, diagnostic efficiency, genetic
191 variant prioritization in diagnostic pipelines and global data exchange.

192 Given these challenges and the need for accurate, standardized phenotyping, the Human
193 Phenotype Ontology (HPO) system was conceptualized and published with initial terminology
194 in 2008 (8,9). To date, HPO provides the most comprehensive deep phenotyping resource for
195 rare diseases for clinicians, researchers, bioinformaticians and electronic health record (EHR)

196 systems in the world. HPO is used in many projects including the 100,000 Genomes Project,
197 the NIH Undiagnosed Disease Program and Network (UDP and UDN), the Undiagnosed
198 Diseases Network International (UDNI), RD-CONNECT, and SOLVE-RD (*1,10-13*). HPO is
199 a community-based tool and is increasingly adapted as the standard to describe phenotypic
200 abnormalities for everyday use (*14*). Each term in HPO describes a distinct phenotypic feature
201 (e.g. lymphadenopathy, HP:0002716) and the HPO tree structure allows similarity measures
202 between patient phenotypes. HPO contains over 200,000 phenotypic annotations for hereditary
203 diseases, of which 2,120 are considered rare diseases. Inborn errors of immunity (IEIs) form a
204 subgroup of these rare diseases. Clinical experts in IEI agree that a major barrier to the adoption
205 of HPO terminology has not been used widely for IEIs partly due to the lack of disease specific
206 HPO terms for IEI patients (*15*). Adequate depiction of the complex clinical and
207 immunological phenotypes of IEI disease entities with HPO terms would allow discrimination
208 between heterogeneous groups of IEIs. Illustrating the lack of terms, in 2017 HPO contained
209 more than 11,000 terms, out of which 5,000 terms have been applied to the musculoskeletal
210 system, with only 1,000 terms related to IEIs (*9,15*). In addition, the phenotypic annotation of
211 IEIs often includes results of specific immunological assays, which pose a challenge to
212 accurately reflect in HPO terms (*15*). Because of the lack of specific HPO terms depicting
213 results of laboratory assays, often a non-specific broader term is used for the annotation of IEIs.
214 Therefore, HPOs are currently not specific enough to be used for genetic analysis and
215 diagnostic aid for IEIs. In a study addressing the clinical efficacy of genetic testing in IEI,
216 bioinformatics tools using existing HPO terms missed the disease causing gene in 37% of the
217 patients with known monogenic disorders (*16*). In this study, we set out to improve HPO
218 terminology for IEIs by applying established bioinformatic methodologies coupled with expert
219 review. The aims of this project were therefore to i) systematically review existing HPO terms
220 for IEIs, ii) revise ontology structures, to iii) add missing terms, as well as iv) reannotate

221 existing IEIs with HPO terms, to collectively enable systematic use of HPO by the IEI-
222 community.

223 Materials and Methods

224 Spearheaded by the European Reference Network on Rare Primary Immunodeficiency,
225 Autoinflammatory and Autoimmune diseases (ERN-RITA) and the European Society for
226 Immunodeficiencies (ESID), we set up working groups comprising members of the
227 participating immunodeficiency societies to revise and expand HPO terms for IEIs. Three
228 workshops, numerous teleconferences and joint task forces took place over the span of 2 years,
229 with over 30 participants including expert clinicians, geneticists, researchers working on IEIs
230 and bioinformaticians. All participating clinicians and geneticists identified through ERN-
231 RITA, ESID, and the International Society of Systemic Autoinflammatory Diseases (ISSAID)
232 are established experts in their fields from different European countries and North America.
233 Additional scientific support provided the indispensable bioinformatics expertise.

234

235 Establishment of working structure

236 A remote working structure (detailed in the Supplementary Methods) was launched to address
237 gaps in the HPO tree and in the annotation of IEI diseases.

238

239 Expansion and restructuring of disease-related branches of the HPO tree

240 Disease-specific HPO restructuring was discussed within four working groups. Each group
241 focused on a different HPO branch; the suggested changes were agreed on among all
242 participants. Differences between centers and countries in the use of terms and definitions were
243 highlighted during the face-to-face workshops. The results were summarized electronically in
244 Excel documents or pictures and flipchart drawings by the main coordinators before being
245 submitted to HPO. The full list of restructured tree elements is detailed in the Supplementary
246 Document 1. New submitted HPO terms can be found in Supplementary Document 2.
247 Additionally, missing terms describing pulmonary and gastro-intestinal complications of

248 primary antibody deficiency (PAD) were discussed during teleconferences and thereafter
249 submitted to update the HPO ontology.

250

251 *Standardized reannotation of rare, genetically diagnosed diseases*

252 A four-step process was developed for a standardized reannotation effort across working
253 groups and to consistently annotate IEIs (spanning over 300 different diseases in Online
254 Mendelian Inheritance in Man (OMIM)) with HPO terms (Fig 1). As IEIs represent a large and
255 heterogenous group of rare diseases, we here decided to selectively focus on defined subgroups
256 of IEI to test the feasibility and usefulness of such an endeavor. First, publications were
257 collected by experts for each disease within the subgroups (minimum of two articles per
258 disease), representing key phenotypic presentation(s) of the specific disease. In the second step,
259 HPO terms were extracted from the provided publications for each disease using machine
260 learning ((17), explained in detail in Supplementary Materials and Methods) and summarized
261 into Excel documents. Third, a two-tier expert review evaluated the text mined terms,
262 suggested additional terms if required and the responsible working group agreed (defined as at
263 least 80% agreement amongst group experts) on the final HPO annotations for each disease.
264 Fourth, the validated terms were submitted to HPO. Supplementary Document 2 contains the
265 reannotated diseases and the list of reannotated terms for each disease is available in
266 Supplementary Document 3.

267

268 *Standardized reannotation of genetically undiagnosed diseases*

269 The methods above were specifically designed for application in (very) rare diseases, where
270 the number of patients and therefore the described phenotypic spectrum and clinical
271 presentation is sparse. In case of diseases and disease groups where an adequate amount of
272 patient and phenotype data was available, in addition to a True/False annotation, the frequency

273 of each phenotypic item was assessed. The frequencies correspond to the following
274 representation in patients: common = Frequent (79-30%); sometimes = Occasional (29-5%);
275 rare = Very rare (<4-1%).

276

277 *Patient cohort*

278 We randomly selected 30 patients that harbored a genetic diagnosis in one of the reannotated
279 diseases from a large pediatric referral center research database. Clinical summaries of these
280 patients prior to genetic diagnosis were retrieved by an expert clinician. The clinical summaries
281 were parsed and HPO terms were extracted using machine learning as in the Supplementary
282 Methods.

283

284 *HPO information content measures, and disease patient similarity measures*

285 Information content of all HPO terms was assessed with the *R* package *ontologyIndex* v2.5
286 (18). The phenotypic similarity of diseases and patients before and after reannotation was
287 compared using the *R* package *ontologySimilarity* v2.3 (18). The Euclidean distances between
288 the diseases were computed based on similarity measures, clustered with hierarchical clustering
289 and visualized with *ggtree* using the *R* packages *ggtree* (19) and *ape* v5.2 (20).

290

291 A detailed description including the data processing pipeline and tools are available in the
292 Supplementary Materials and Methods.

293

294 *Supplementary Materials*

295 Supplementary Materials and Methods

296 Supplementary Document 1: HPO tree restructuring and list of new terms

297 Supplementary Document 2: Summary of diseases reannotated

- 298 Supplementary Document 3: List of all terms per disease after reannotation
- 299 Supplementary Document 4: List of cases used for phenotype to diagnosis matching

300 Results

301

302 Systematic evaluation and expansion of the HPO structure and terms relevant to IELs

303 Our approach has resulted in the restructuring of four main branches of the HPO tree, namely:

304 i) abnormality of the immune system (HP:0002715) ii) abnormality of metabolism/homeostasis

305 (HP:0001939) iii) abnormality of the integument (HP:0001574) and iv) abnormality of the

306 cardiovascular system. (Fig 2A, Supplementary Document 1). Together, this revision prompted

307 the replacement/restructuring of 67 terms, and the addition of 57 new terms to the HPO tree,

308 among them “recurrent fever”, “unusual infections”, “IgG levels in blood” (Fig 2B,

309 comprehensive list in Supplementary Documents 1 and 2).

310

311 Directed expansion of primary antibody deficiency (PAD) terms

312 Overall, the PAD working group focused on replacing broad and non-specific terms with terms

313 that describe phenotypes in more detail and accuracy (example: ‘partially absent total

314 IgG/IgA/IgM in blood’ and ‘(near) absent total IgG/IgA/IgM in blood’ instead of

315 ‘hypogammaglobulinemia’) Fig 2B. In addition, we proposed that the full detailed spectrum of

316 specific antibody as well as IgG-subclass deficiencies was described by separate HPO terms.

317 For example, we described individual terms related to ‘decreased specific antibody response to

318 vaccination in blood’ divided according to the response to different types of vaccination

319 (protein, protein-conjugated polysaccharide and unconjugated polysaccharide).

320

321 Standardized reannotation of rare, genetically diagnosed IELs

322 We started by a systematic review of four disease categories of the IUIS classification of IELs,

323 as proof of concept: diseases affecting cellular and humoral immunity (IUIS Table 1), diseases

324 of immune dysregulation (IUIS Table 4), autoinflammatory disorders (IUIS Table 7) and

325 genetically undiagnosed predominantly antibody deficiencies (IUIS Table 3), detailed in Table
326 1 and Supplementary Document 3. As a first step, we assessed the already available HPO
327 annotation for each disease in the 2019-06-03 HPO release. We found that 15% of diseases
328 considered (11 of 73 diseases in total) did not have any associated HPO terms (Fig 3A). Overall,
329 we found that on average 13.3 phenotype terms were available per disease (Fig 3B), later
330 referred to as “existing terms”.

331 The text mining and evaluation process was separated into four steps shown in Fig 3C. We
332 have first focused the reannotation of 72 genetically diagnosed IEs, and genetically
333 undiagnosed PADs. For genetically diagnosed IEs, text mining was based on 162 expert-
334 curated articles, on average 2.57 articles per disease (Fig 3D). This resulted in 4,517 extracted
335 phenotype terms, 66.42 terms per disease (Fig 3E). Of these terms, 3,242 - or 71% per disease
336 (47.67 out of 66.42) - were accepted as correctly attributed terms by the expert reviewers (Fig
337 3F). Expert suggestions added up to 529 additional HPO terms, in addition to the existing and
338 text mined terms.

339 After reannotation, a mean of 63.1 terms were available for each disease, resulting in a 4.7-fold
340 gain in the number of available annotations (Fig 3G). The mean information content as
341 measured by the overall frequency of terms in each disease’s annotations has increased from
342 6.17 to 8.3 (Fig 3H) after reannotation.

343

344 The new annotation of diseases consisted mainly of text mined terms (70.6%) (Fig 3I),
345 followed by already existing terms (9.3%) and additional suggestions by experts (9.3%, adding
346 a further 5.2 additional terms per disease) (Supplementary Document 3).

347

348 *Standardized reannotation of genetically undiagnosed primary antibody deficiencies (PADs)*
349 PADs form a heterogeneous group, and the majority of PADs do not (as yet) have a genetic
350 diagnosis. We collected articles describing the heterogeneous PADs related to common
351 variable immunodeficiency disorders (CVID), agammaglobulinemia, selective IgM deficiency,
352 selective IgA deficiency, IgG-subclass deficiency, specific antibody deficiency and
353 unclassified antibody deficiency subgroups. In total, 541 terms were text mined from these
354 articles, many of these in more than one PAD subgroup, and 245 of these terms (45.2%) were
355 annotated as correctly associated to the respective PAD subgroup by the expert reviewers (Fig
356 3J). Of these 245 terms, the experts annotated 16.3% as commonly found in PAD diseases,
357 48.97% as sometimes associated (albeit less commonly), and 34.7% as rarely associated with
358 PAD (Fig 3K).

359

360 *Patient-disease matching*

361 We set out to showcase the efficacy of our reannotation effort by highlighting the potential
362 diagnostic impact of optimized disease annotation. To do this, we have selected 30 clinical
363 cases from a large immunology referral center research database (Supplementary Document
364 4). HPO terms were matched to patient phenotypes by experts from the clinical synopsis and
365 the phenotypic similarity to all HPO-annotated diseases was calculated based on these selected
366 patient HPO terms (Fig 4A), as illustrated by one concrete clinical example of a patient with
367 Tumor Necrosis Factor Receptor Associated Periodic Syndrome (TRAPS, Fig 4B). Overall,
368 we show a significant 47% improvement in the specificity of patient phenotype matching to
369 correct diagnosis (from 0.49 to 0.72, p value = $1.8e-07$, Fig 4C), and a significantly better
370 ranking of the correct clinical diagnosis across all possible diseases after reannotation: in the
371 majority of cases, the correct diagnosis was in the top 10 of matched diseases (Fig 4D) after
372 reannotation, and the rank of the correct diagnosis for individual patients was highly

373 significantly improved, from a mean of 285 to 19 (14.9 fold improvement, p value = 9.1e-07,
374 Fig 4E).

375

376 *Phenotype-driven disease classification*

377 We tested the efficacy of our approach in selecting biologically and clinically meaningful
378 phenotypes by assessing the HPO-ontology based phenotypic similarity of diseases before and
379 after reannotation. In particular, we assessed whether the similarity was greater within or
380 between IUIS clinically defined groups. We found that the phenotype-driven disease
381 classification after reannotation has resulted in a clustering more in concordance with the IUIS-
382 based clinical classification (Fig 5A-B).

383 Discussion

384

385 Unified data standards, consistent classification and robustly verified clinical data are vital
386 pillars supporting diagnostic pipelines and data-driven research. Although databases and
387 vocabularies that aim to provide accurate phenotypic descriptions exist (5-9), there are still
388 major gaps in the depiction of IEs in these datasets. Here we used a cross-community
389 collaboration to review, expand and improve the depiction of IEs in HPO, and reannotate IEs
390 with HPO terms. We reviewed four separate branches of the HPO tree and submitted 57 new
391 and expanded HPO terms, the majority of which are now included in the official HPO dataset.
392 We introduced a semi-automated reannotation pipeline, that combines ontology-guided
393 machine learning and a two-tier expert review to reannotate four main categories of IEs. The
394 basis of the ontology-guided machine learning was the expert curated list of articles (162 in
395 total), that was submitted to the PanelApp (21) to serve as a public resource. The text mined
396 phenotypes were subjected to expert review to confer face validity or refute the putative new
397 HPO terms. IEs and their current HPO terms covered by the working groups were scrutinized
398 in-depth, resulting in high-quality annotations. Overall, we have achieved a 4.7-fold gain in
399 number of HPO terms annotating each disease. These annotations included unspecific
400 (frequently annotated) as well as specific (less frequently annotated) HPO terms holding less
401 and more information content respectively. Combined, the mean information content increased
402 from 6.17 to 8.3.

403 Each reannotated disease showed an increase in information content and a quantitative gain in
404 the number of available HPO terms. Through patient-disease matching and disease-similarity
405 examples we illustrated that these gains and increases translated to significant qualitative
406 improvement in patient-disease matching in an independent cohort of IEI patients (Figure 4),
407 and phenotype-driven classification of IEs that more closely resembles clinical consensus

408 (Figure 5). Although neither of these measures are systematic assessments of global patient-
409 disease matching and disease similarity comparisons, they highlight that there is considerable
410 benefit by the revision of specific subclasses of diseases. Once a near complete HPO phenotype
411 reannotation of almost all IEIs is available, it will be intriguing to assess how well patients with
412 genetic diagnoses match reannotated OMIM diseases in a clinical setting, how patient matching
413 to genetic diagnosis is transformed, and if these changes ultimately lead to an earlier diagnosis.
414 Finally, once a detailed and accurate phenotypic description is available for all IEIs,
415 identification phenotype-driven patient subgroups will be common practice, and a more
416 objective entirely phenotype-driven classification and ontology of IEIs can become a reality.

417

418 Accurate phenotypic description of patients holds promise for diagnostic utility and for the
419 discovery of novel diseases. Phenotype-driven genetic diagnostic tools now exist, but their full
420 clinical potential is hampered by the lack of complete phenotypic descriptions for most types
421 of IEIs. Phenotips (22) is a free and open source software for collecting and analyzing
422 phenotypic information of patients with genetic disorders that is widely used in the rare disease
423 community. Tools such as Exomiser use HPO terms to annotate and to prioritize potentially
424 casual variants (23). New integrative ‘omics approaches and the analysis of large-scale data
425 with artificial intelligence will allow us to go from a one-size-fits-all to a more personalized
426 medicine, including in IEIs. We see the potential to integrate the richer phenotyping of
427 previously undiagnosed groups of IEI patients with available sequencing data to accelerate
428 disease gene discovery and at the same time increase the diagnostic rate in new patients (24).

429 Novel disease-gene or phenotype associations depends on sufficient numbers of cases as well
430 as a control cohort of comparable quality. Cross-institute and cross-country collaborations for
431 cohorts of undiagnosed, but well-phenotyped patients could shed light on novel disease-
432 causing genes not only of the immune system. Trusted and accepted data and information

433 sharing platforms are already being developed (13, 22) to provide robust and sufficiently
434 granular HPO terms as a standardized way of phenotyping patients. Electronic health records
435 (EHR) (25) could facilitate the transfer of HPO terms by integrating with available sharing
436 platforms. Capturing HPO annotations of novel rare diseases or cases is an ongoing challenge
437 for a complete disease representation. Thus it is important that alongside of updating the official
438 IUIS classification, HPO descriptions of disorders are curated once every several years. We
439 suggest a community effort for such regular reviews of HPO regarding IEs, such as a team of
440 experts, part of big international groups of clinicians such as ESID or ERN RITA, the Clinical
441 Immunology Society (CIS) or other similar organizations. Publication standards that require
442 the submission of HPO annotations up-front would greatly improve this process.

443

444 Once phenotyped patients are available, robust and global approaches are accessible (2) to find
445 phenotypic similar cases. These comparisons are performed by advanced machine learning
446 algorithms. However, machine learning can also be a very powerful tool to automate the
447 identification of relevant phenotype information in publications or clinical notes. We applied
448 an ontology-guided machine learning tool to support the annotation of diseases and explored
449 the full spectrum of terms – from very relevant to not relevant at all. The same process can be
450 applied to unstructured clinical notes to accelerate in-depth annotation of patients. For patients
451 with EHR (25), abnormal clinical values can automatically be translated into HPO codes (26)
452 for a more precise diagnostic application and integrated with sharing platforms as mentioned
453 before. The foundation of these comparisons is an ontology with a comprehensive set of term,
454 which is widely used.

455

456 As there is currently no gold-standard on how to perform an expert-based review of ontologies.
457 guidance on annotating diseases with HPO phenotypes can vary between diseases, disease

458 classes and centers. IEIs are rare diseases, and often there are only a few patients described
459 (sometimes only one kindred in case of ultra-rare diseases). Therefore, the depth of currently
460 available published phenotypes is at times limited. The low number of patients and insufficient
461 depth of available phenotypes brings up a question as to which diseases to include in
462 phenotyping exercises of this nature. On the one hand, focusing on IEIs that are commonly
463 accepted, with multiple patients diagnosed and well described by multiple researchers can
464 increase the depth of phenotyping. However, this approach excludes at least 10% of IEIs (the
465 ultra-rare diseases). On the other hand, an all-inclusive approach including every disease
466 systematically means that we rely on sparsely phenotyped patients and perhaps insufficient
467 data for ultra-rare disorders. A warning of accuracy by indicating the frequency of each
468 phenotype for diseases could soon be possible, with the addition of phenotype frequency to the
469 HPO dataset, an expansion that is currently work in progress. This implies the need for a
470 responsive system, capable of assimilating new phenotypic information as the pool of
471 confidently diagnosed patients increases.

472

473 Our ongoing approach aims to address these gaps for IEIs and to provide an ontology that is
474 practical, useful and as complete as possible. However, the existence of a well-built ontology
475 and the awareness of clinicians and researchers itself does not guarantee a shift in the
476 community to fully adapt a standardized phenotyping approach. Our approach raised awareness
477 regarding the concept and importance of HPO amongst the IEI community. Moreover, the
478 process made the participating clinicians aware of the available terms and highlighted where
479 these were lacking. Moving forward, it is very important that official entities adopt HPO terms
480 as the unified means of patient phenotyping. We hypothesize that as soon as the widely used
481 registries such as the Undiagnosed Disease Network (11) or the IUIS (27) use HPO to refer to
482 phenotypic annotation, this will propel the IEI field towards adopting HPO as the main

483 nomenclature for phenotyping IEI patients. One promising move in this direction is the recent
484 expansion of the ESID registry working definitions for the clinical diagnosis of IEIs (28), which
485 derives HPO terms from OrphaNet using the ORDO Ontological Module (HOOM) platform
486 (29), prompted by our HPO initiative.

487

488 In summary, our work reviewed and expanded the phenotypic depiction of multiple subclasses
489 of IEIs, and to our knowledge, this initiative is the first endeavor of its kind with the aim of
490 standardizing IEI phenotypes. Our semi-automated annotation-based approach is scalable to
491 include all IEIs as illustrated herein. We propose our reannotation approach as a blueprint for
492 systematic HPO (re)annotation for additional immunological and non-immunological diseases.

493

494

495

496 **Fig 1: Pipeline for of standardized reannotation of IEI diseases.** First, scientific
497 publications were collected by experts for each disease within the subgroups. Second, HPO
498 terms were extracted from the provided publications for each disease using machine learning
499 and summarized into Excel documents. Third, a two-tier expert review evaluated the text mined
500 terms, suggested additional terms if required and the responsible working group agreed on the
501 final HPO annotations for each disease. Fourth, data were collated, and the agreed terms were
502 submitted to HPO.

503

504 **Fig 2: Revision and expansion of the HPO tree. A) Schematic representation of the**
505 **restructuring of the HPO tree.** Main branches of the HPO tree where restructuring was
506 performed are marked with light green. B) “Abnormality of temperature”, “Abnormality of
507 immunoglobulin level“ and “Unusual infections“ as examples of revised branches of the HPO
508 tree. New additions to the tree are marked with green, repositioned terms are marked with
509 yellow.

510

511 **Fig 3: Result of disease reannotation.** A) HPO annotation availability in the subset of 72
512 diseases. B) Distribution of number of available HPO terms per disease. C) Distribution of the
513 number of articles used per disease for the reannotation pipeline. D) Number of mined terms
514 per disease. Each dot represents a disease. E) All mined vs all accepted terms. F) Number of
515 available terms per disease before and after reannotation. Each dot represents a disease. G)
516 Mean information content available per disease before and after reannotation. H) The aggregate
517 mean annotation per disease after reannotation. I) All text mined terms from PAD publications
518 J) Frequency distribution of different PAD terms according to the experts. HPO: Human
519 Phenotype Ontology; PAD: Primary Antibody Deficiencies.

520

521 **Fig 4: Patient-disease matching.** A) Schematic overview of the different steps of patient-to-
522 disease matching. First, the phenotypes were identified in a patient's clinical history. Second,
523 these phenotypes were translated to HPO terms. Finally, patient phenotype to disease matching
524 was measured by Resnik similarity. B) Matching patient 1 to a diagnosis. C) Similarity of
525 patients in patient cohort to genetic diagnosis before and after reannotation. D) The rank of
526 correct clinical diagnosis more often is in the top 10 of matched diseases after reannotation. E)
527 Improvement of ranks of clinical diagnosis before and after reannotation. Significance was
528 assessed by Student t-test.

529

530 **Fig 5: Phenotypic similarity of diseases before and after reannotation.** Diseases are
531 annotated with the IUIS disease group (inner circle), sub-group (outer circle) and OMIM
532 identifier. A) Clustering of diseases based on phenotypic similarity before reannotation. B)
533 Clustering of diseases based on phenotypic similarity after reannotation. HPO: Human
534 Phenotype Ontology; IUIS: International Union of Immunological Societies, OMIM: Online
535 Mendelian Inheritance in Men; IEI: Inborn Errors of Immunity; EBV: Epstein-Barr Virus

536

537 **References**

538

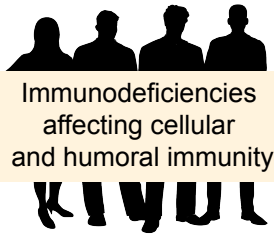
- 539 1. Gahl WA, Markello TC, Toro C, Fajardo KF, Sincan M, Gill F, et al. The National
540 Institutes of Health Undiagnosed Diseases Program: insights into rare diseases. *Genet*
541 *Med.* 2012 Jan;14(1):51–9.
- 542 2. Philippakis AA, Azzariti DR, Beltran S, Brookes AJ, Brownstein CA, Brudno M, et al.
543 The Matchmaker Exchange: a platform for rare disease gene discovery. *Hum Mutat.*
544 2015 Oct;36(10):915–21.
- 545 3. Sobreira N, Schiettecatte F, Valle D, Hamosh A. GeneMatcher: a matching tool for
546 connecting investigators with an interest in the same gene. *Hum Mutat.* 2015
547 Oct;36(10):928–30.
- 548 4. Hernandez-Ibarburu G, Perez-Rey D, Alonso-Oset E, Alonso-Calvo R, de Schepper K,
549 Meloni L, et al. ICD-10-CM extension with ICD-9 diagnosis codes to support
550 integrated access to clinical legacy data. *Int J Med Inform.* 2019;129:189–97.
- 551 5. Amberger J, Bocchini C, Hamosh A. A new face and new challenges for Online
552 Mendelian Inheritance in Man (OMIM®). *Hum Mutat.* 2011 May;32(5):564–7.
- 553 6. Pavan S, Rommel K, Mateo Marquina ME, Höhn S, Lanneau V, Rath A. Clinical
554 Practice Guidelines for Rare Diseases: The Orphanet Database. *PLoS ONE.*
555 2017;12(1):e0170365.
- 556 7. Schriml LM, Mitraka E, Munro J, Tauber B, Schor M, Nickle L, et al. Human Disease
557 Ontology 2018 update: classification, content and workflow expansion. *Nucleic Acids*
558 *Res.* 2019 08;47(D1):D955–62.
- 559 8. Robinson PN, Köhler S, Bauer S, Seelow D, Horn D, Mundlos S. The Human
560 Phenotype Ontology: a tool for annotating and analyzing human hereditary disease.
561 *Am J Hum Genet.* 2008 Nov;83(5):610–5.

- 562 9. Köhler S, Vasilevsky NA, Engelstad M, Foster E, McMurry J, Aymé S, et al. The
563 Human Phenotype Ontology in 2017. *Nucleic Acids Res.* 2017 04;45(D1):D865–76.
- 564 10. Ramoni RB, Mulvihill JJ, Adams DR, Allard P, Ashley EA, Bernstein JA, et al. The
565 Undiagnosed Diseases Network: Accelerating Discovery about Health and Disease.
566 *Am J Hum Genet.* 2017 02;100(2):185–92.
- 567 11. Taruscio D, Groft SC, Cederroth H, Melegh B, Lasko P, Kosaki K, et al. Undiagnosed
568 Diseases Network International (UDNI): White paper for global actions to meet patient
569 needs. *Mol Genet Metab.* 2015 Dec;116(4):223–5.
- 570 12. Gall T, Valkanas E, Bello C, Markello T, Adams C, Bone WP, et al. Defining Disease,
571 Diagnosis, and Translational Medicine within a Homeostatic Perturbation Paradigm:
572 The National Institutes of Health Undiagnosed Diseases Program Experience. *Front*
573 *Med (Lausanne).* 2017;4:62.
- 574 13. Thompson R, Johnston L, Taruscio D, Monaco L, Bérout C, Gut IG, et al. RD-
575 Connect: an integrated platform connecting databases, registries, biobanks and clinical
576 bioinformatics for rare disease research. *J Gen Intern Med.* 2014 Aug;29 Suppl
577 3:S780-787.
- 578 14. Köhler S, Carmody L, Vasilevsky N, Jacobsen JOB, Danis D, Gourdine J-P, et al.
579 Expansion of the Human Phenotype Ontology (HPO) knowledge base and resources.
580 *Nucleic Acids Research.* 2019 Jan 8;47(D1):D1018–27.
- 581 15. Chinn IK, Chan AY, Chen K, Chou J, Dorsey MJ, Hajjar J, et al. Diagnostic
582 interpretation of genetic studies in patients with primary immunodeficiency diseases:
583 A working group report of the Primary Immunodeficiency Diseases Committee of the
584 American Academy of Allergy, Asthma & Immunology. *Journal of Allergy and*
585 *Clinical Immunology.* 2020 Jan;145(1):46–69.

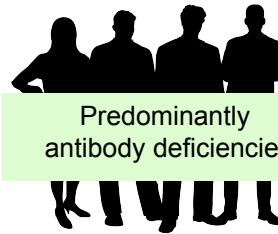
- 586 16. Rae W, Ward D, Mattocks C, Pengelly RJ, Eren E, Patel SV, et al. Clinical efficacy of
587 a next-generation sequencing gene panel for primary immunodeficiency diagnostics.
588 *Clin Genet.* 2018 Mar;93(3):647–55.
- 589 17. Arbabi A, Adams DR, Fidler S, Brudno M. Identifying Clinical Terms in Medical Text
590 Using Ontology-Guided Machine Learning. *JMIR Med Inform.* 2019 May
591 10;7(2):e12596.
- 592 18. Greene D, Richardson S, Turro E. ontologyX: a suite of R packages for working with
593 ontological data. *Bioinformatics.* 2017 01;33(7):1104–6.
- 594 19. Yu G, Lam TT-Y, Zhu H, Guan Y. Two Methods for Mapping and Visualizing
595 Associated Data on Phylogeny Using Ggtree. *Mol Biol Evol.* 2018 01;35(12):3041–3.
- 596 20. Paradis E, Schliep K. ape 5.0: an environment for modern phylogenetics and
597 evolutionary analyses in R. *Bioinformatics.* 2019 Feb 1;35(3):526–8.
- 598 21. Martin AR, Williams E, Foulger RE, Leigh S, Daugherty LC, Niblock O, et al.
599 PanelApp crowdsources expert knowledge to establish consensus diagnostic gene
600 panels. *Nat Genet.* 2019;51(11):1560–5.
- 601 22. Girdea M, Dumitriu S, Fiume M, Bowdin S, Boycott KM, Chénier S, et al. PhenoTips:
602 patient phenotyping software for clinical and research use. *Hum Mutat.* 2013
603 Aug;34(8):1057–65.
- 604 23. Smedley D, Jacobsen JOB, Jäger M, Köhler S, Holtgrewe M, Schubach M, et al. Next-
605 generation diagnostics and disease-gene discovery with the Exomiser. *Nat Protoc.*
606 2015 Dec;10(12):2004–15.
- 607 24. Westbury SK, Turro E, Greene D, Lentaigne C, Kelly AM, Bariana TK, et al. Human
608 phenotype ontology annotation and cluster analysis to unravel genetic defects in 707
609 cases with unexplained bleeding and platelet disorders. *Genome Medicine.* 2015 Apr
610 9;7(1):36.

- 611 25. Lehne M, Luijten S, Vom Felde Genannt Imbusch P, Thun S. The Use of FHIR in
612 Digital Health - A Review of the Scientific Literature. *Stud Health Technol Inform.*
613 2019 Sep 3;267:52–8.
- 614 26. Shefchek KA, Harris NL, Gargano M, Matentzoglou N, Unni D, Brush M, et al. The
615 Monarch Initiative in 2019: an integrative data and analytic platform connecting
616 phenotypes to genotypes across species. *Nucleic Acids Res.* 2020 08;48(D1):D704–15.
- 617 27. Bousfiha A, Jeddane L, Picard C, Ailal F, Bobby Gaspar H, Al-Herz W, et al. The
618 2017 IUIS Phenotypic Classification for Primary Immunodeficiencies. *J Clin*
619 *Immunol.* 2018;38(1):129–43.
- 620 28. Seidel MG, Kindle G, Gathmann B, Quinti I, Buckland M, van Montfrans J, et al. The
621 European Society for Immunodeficiencies (ESID) Registry Working Definitions for
622 the Clinical Diagnosis of Inborn Errors of Immunity. *J Allergy Clin Immunol Pract.*
623 2019 Aug;7(6):1763–70.
- 624 29. Gasteiger LM, Robinson PN, Pazmandi J, Boztug K, Seppänen MRJ, Seidel MG.
625 Supplementation of the ESID registry working definitions for the clinical diagnosis of
626 inborn errors of immunity with encoded human phenotype ontology (HPO) terms. *The*
627 *Journal of Allergy and Clinical Immunology: In Practice.* 2020 May 1;8(5):1778.
- 628 30. Groza T, Köhler S, Doelken S, Collier N, Oellrich A, Smedley D, et al. Automatic
629 concept recognition using the human phenotype ontology reference and test suite
630 corpora. *Database (Oxford).* 2015;2015.

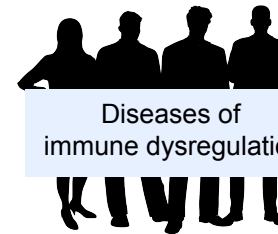
IUIS
Table 1



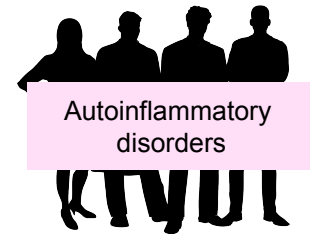
IUIS
Table 3



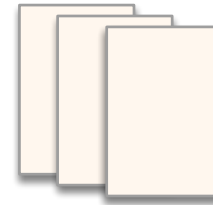
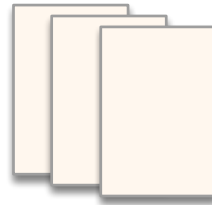
IUIS
Table 4



IUIS
Table 7



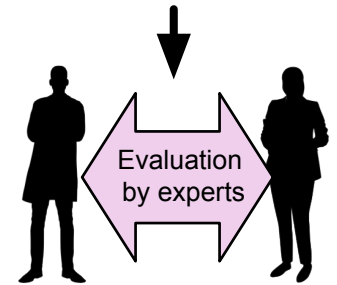
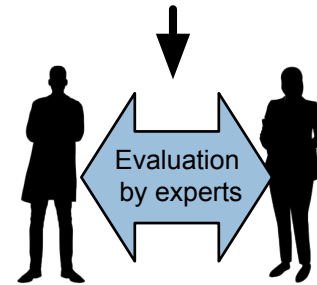
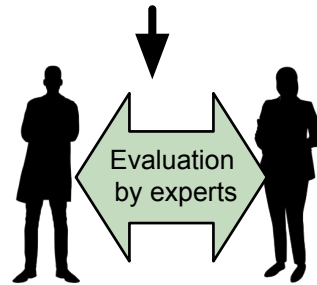
1. Collection of reviews/case reports



2. Text mining, data processing

Extraction of HPO terms with text miner
Preparation of summaries per disease

3. Two-tier expert review



4. Data processing, submission

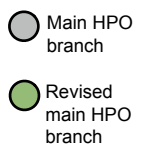
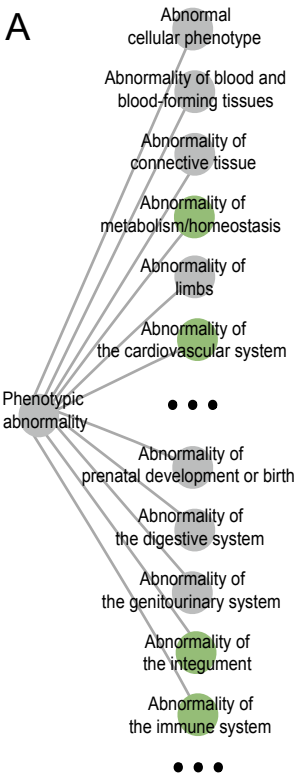
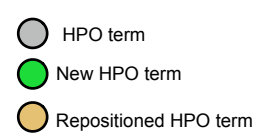
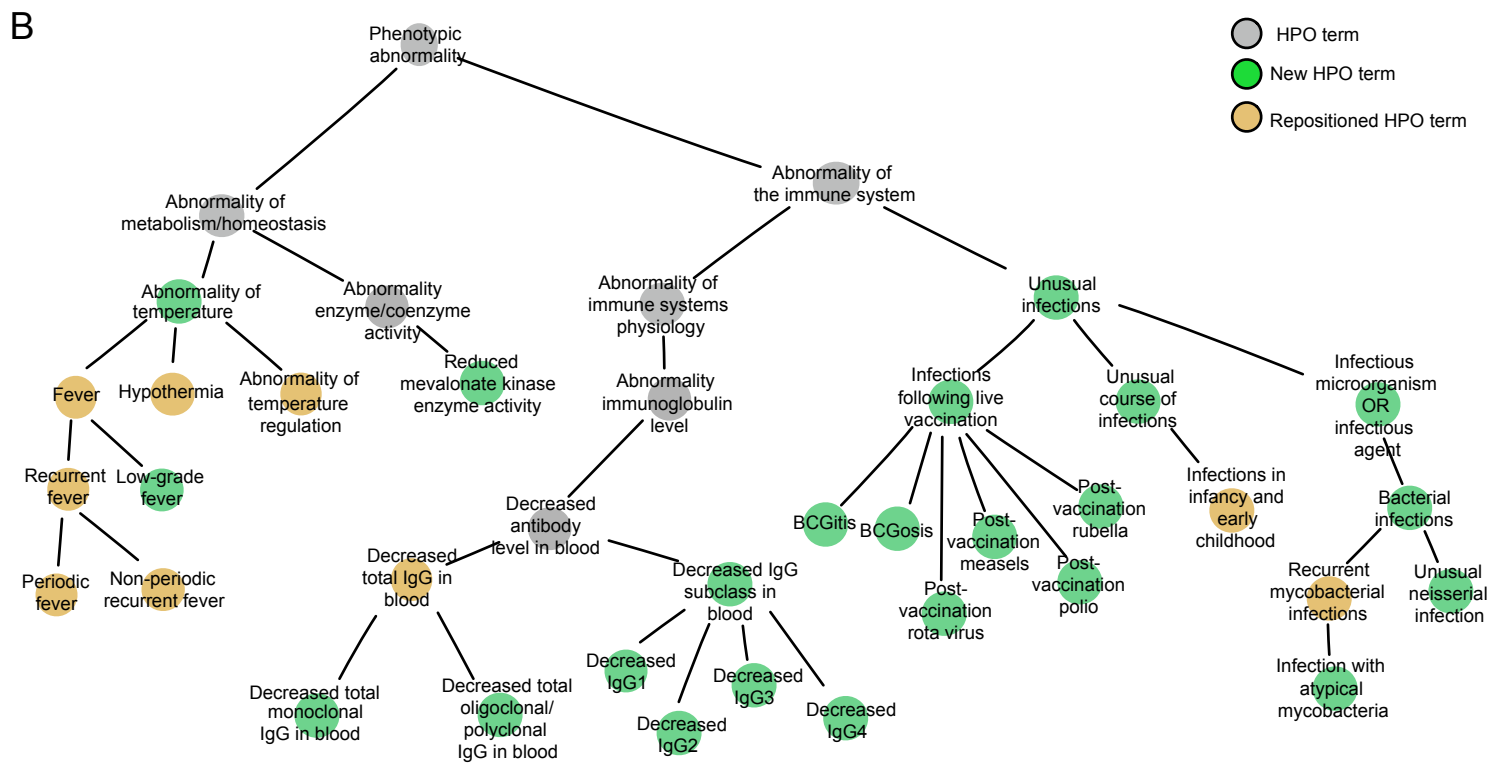
Data filtering and summary
Main coordinators

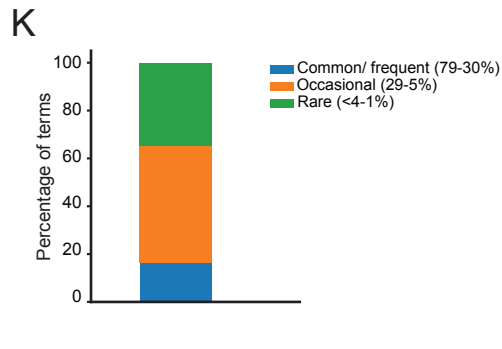
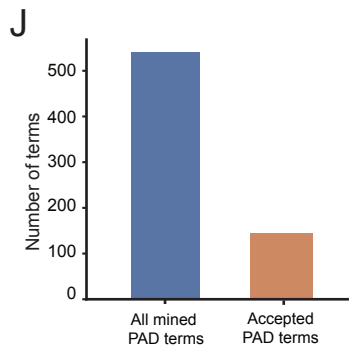
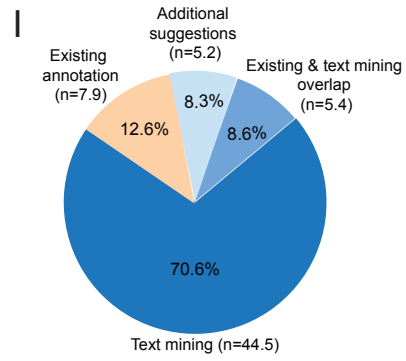
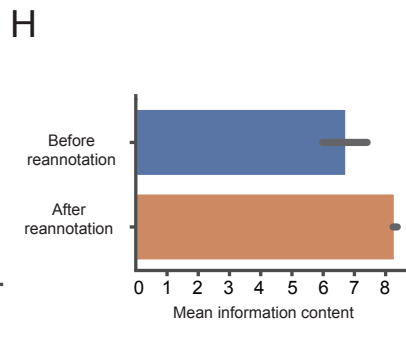
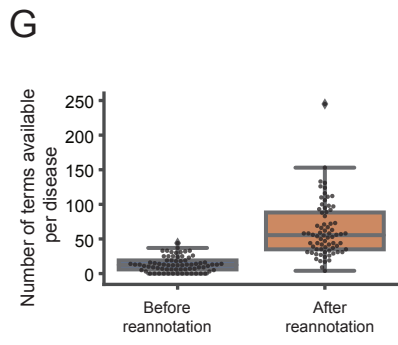
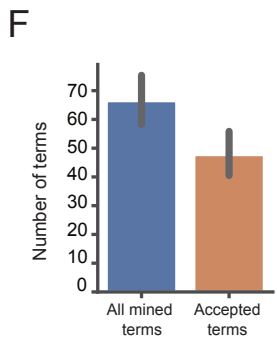
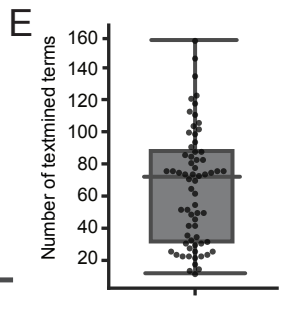
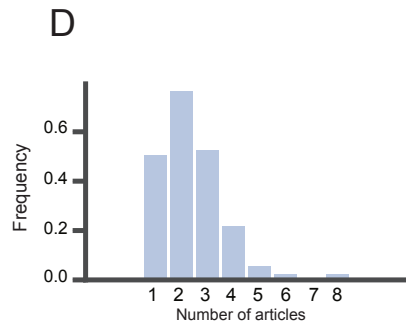
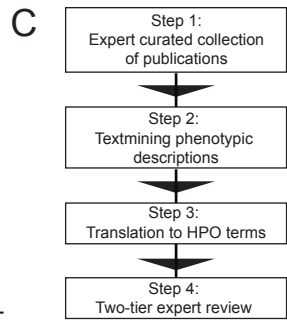
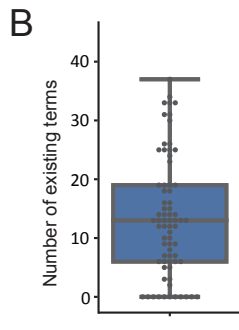
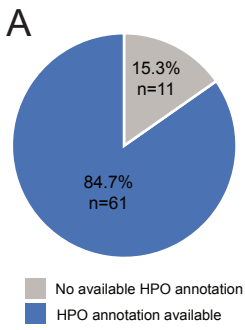


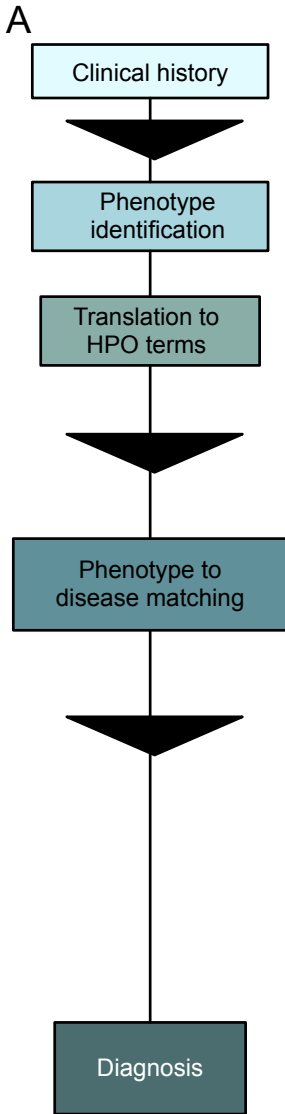
Updates to HPO

Updates to OMIM



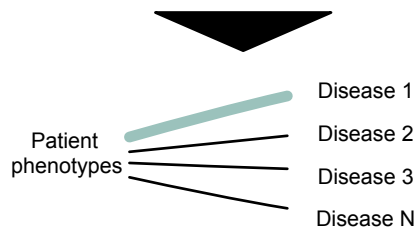
A**B**



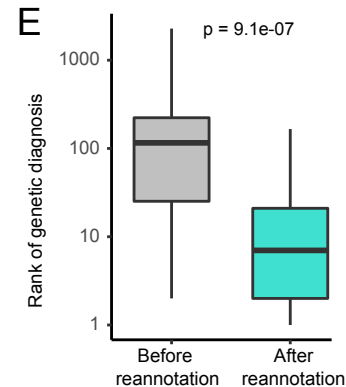
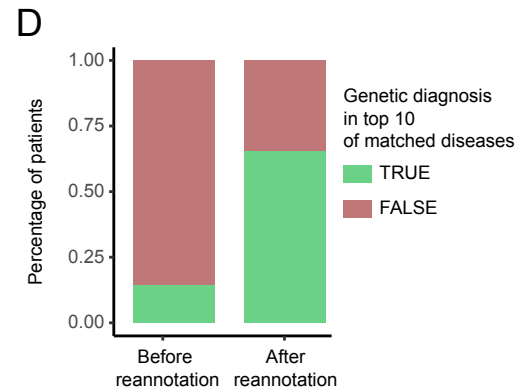
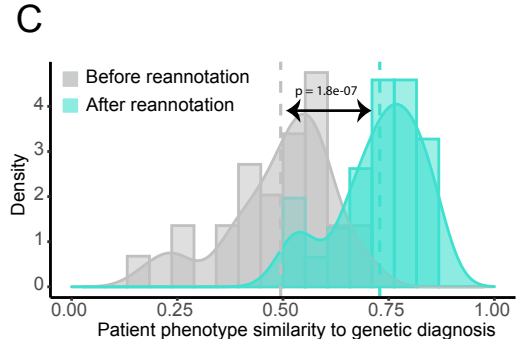
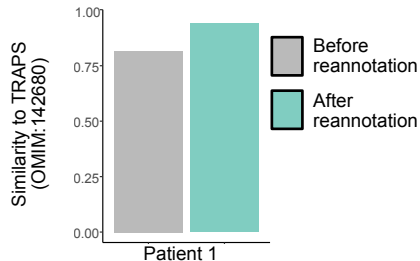


B Patient 1

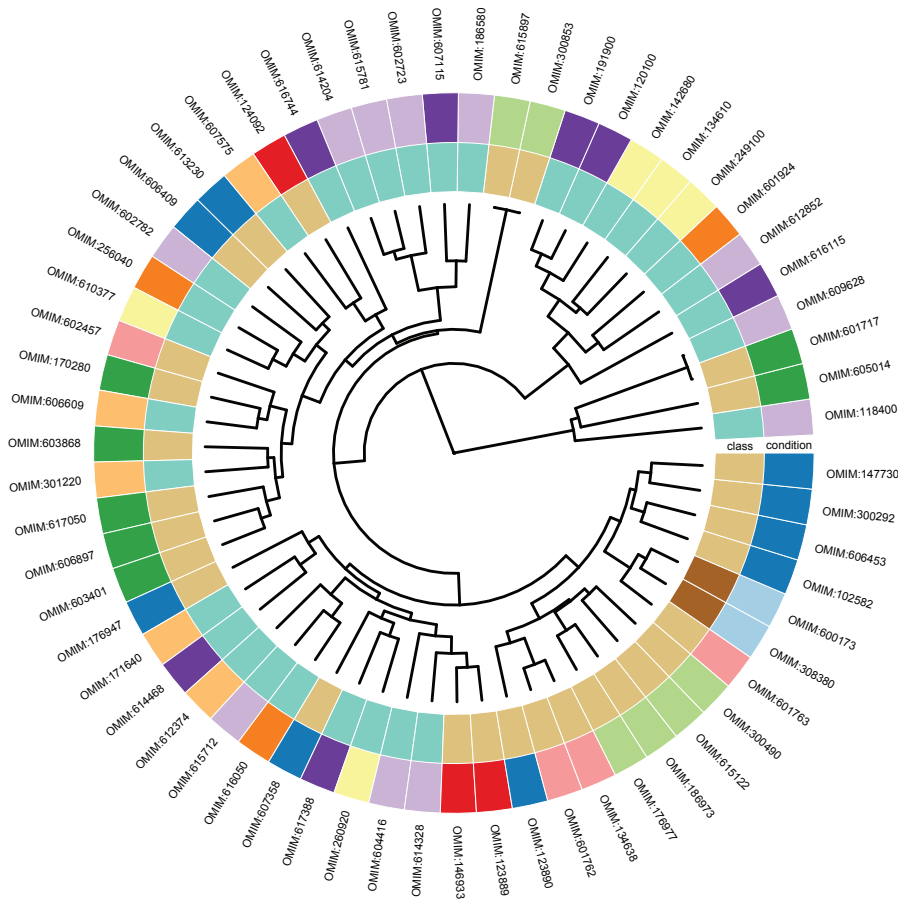
In patient 1, a boy of 4 years of age was seen by the paediatrician with prolonged fever episodes (HP:0001945) that lasted longer than ten weeks (HP:0001954), accompanied by episodes of rigor (HP:0002027) and returning rash (HP:0025145). In addition, the boy regularly experienced abdominal pain (HP:0000988).



Tumor necrosis factor receptor-associated periodic syndrome (TRAPS)
OMIM:142680



A



IUIS classification

- Autoinflammatory disorders
- Diseases of immune dysregulation
- Immunodeficiencies affecting cellular and humoral immunity
- Primary antibody deficiencies

Disease subgroup

- Autoimmune Lymphoproliferative Syndrome (ALPS)
- Hemophagocytic Lymphohistiocytosis (HLH)
- Immune dysregulation with colitis
- Others
- Primary antibody deficiency
- Recurrent inflammation
- SCID T-B+
- Sterile inflammation (skin / bone / joints)
- Susceptibility to EBV
- Syndromes with autoimmunity
- Systemic inflammation with urticaria rash
- Type 1 Interferonopathies

B

