

# RDmap: A Map for Exploring Rare Diseases

Jian Yang, Bs<sup>#1,2</sup>; Cong Dong, Ms<sup>#1,2</sup>; Huilong Duan, PhD<sup>2</sup>; Qiang Shu, MD<sup>1</sup>; Haomin Li, PhD<sup>1\*</sup>

1. The Children's Hospital, Zhejiang University School of Medicine, National Clinical Research Center for Child Health, Zhejiang, China

2. The College of Biomedical Engineering and Instrument Science, Zhejiang University, Zhejiang, China

\*Address correspondence to: Haomin Li, The Children's Hospital, Zhejiang University School of Medicine, Binsheng Road 3333#, Hangzhou, China 310052, [hmli@zju.edu.cn], 086-13867445504;

# Jiang Yang and Cong Dong contributed equally to this study.

## Abstract

**Background:** The complexity of the phenotypic characteristics and molecular bases of many rare human genetic diseases make the diagnosis of such diseases a challenge for clinicians. A map for visualizing, locating and navigating rare diseases based on similarity will help clinicians and researchers understand and easily explore these diseases.

**Methods:** By defining the quantitative distance among phenotypes and pathogenic genes based on corresponding ontology systems, the distance matrix of rare diseases included in

23 Orphanet was calculated and mapped into Euclidean space. Enhanced by clustering classes  
24 and disease information, a rare disease map was developed based on ECharts.

25 **Results:** The rare disease map called RDmap was published at <http://rdmap.nbscn.org>. The  
26 phenotype-based map comprises 3,287 rare diseases and the gene-based map comprises 3,789  
27 rare genetic diseases and they were bridged by 1,718 overlapping diseases. RDmap works  
28 similar to the widely used Google map and supports zooming and panning. The phenotype  
29 similarity base disease location function performed better than traditional keyword search in  
30 an in-silico evaluation and 20 published cases of rare diseases also demonstrated that RDmap  
31 can be used by clinicians to improve diagnosis.

32 **Conclusion:** RDmap is the first user-interactive map-style rare disease knowledgebase. It  
33 will help clinicians and researchers explore the increasing complicated rare genetic diseases.

34 **Keywords:** rare disease; phenotype; pathogenetic gene; disease map; clinical decision  
35 support

## 36 **Background**

37 Rare diseases, most of which are caused by underlying genetic factors, often occur in infants  
38 or young children and affect the patients' whole life. Although rare, studies involving them

39 have revealed important insights about normal physiology that, in turn, have provided a better  
40 understanding of common disorders, universal mechanisms, critical pathways, and therapies  
41 that are useful to treat more than one disease. However, making the correct diagnosis for rare  
42 genetic diseases is extremely complicated and remains a challenge in both developed and  
43 developing countries. According to a survey from EURORDIS, it takes 5 to 30 years for a  
44 quarter of patients with rare genetic diseases from onset to diagnosis. During this period, the  
45 rate of first misdiagnosis is as high as 40%. These misdiagnoses would lead to a large number  
46 of invalid medical treatment or even unnecessary surgery, seriously endangering the health of  
47 the patients and wasting medical resources at the same time.

48 More than 7,000 known rare diseases have been identified, and more than 100 novel disease-  
49 gene associations have been identified per year since the introduction of next-generation  
50 sequencing technologies[1]. An important challenge of rare disease practice is to establish  
51 relationships among so many rare and diverse diseases from different levels. Accumulating  
52 studies have found that genetic diseases that are caused by similar molecules[2–4] can be  
53 diagnosed by similar phenotypic characteristics[5,6], and finally can be treated using similar  
54 drugs through corresponding targets[7–10]. Network-based medicine has emerged as a

55 complementary approach to identify disease-causing genes, genetic mediators, disruptions in  
56 the underlying cellular functions and for drug repositioning. Therefore, exploring the  
57 relationships of rare diseases can help to reveal the common attributes of similar rare genetic  
58 diseases. For example, the classification of rare diseases, phenotypic characteristics of  
59 diseases, and pathogenic genes of genetic diseases can improve the probability of discovering  
60 potential pathogenic mechanisms or drugs and, most importantly, can help with the clinical  
61 diagnosis of rare genetic disease and improve treatment plans.

62 In this study, we aimed to propose a method to construct two rare human disease maps based  
63 on the semantic similarities of both phenotypic characteristics and pathogenic genes of rare  
64 diseases. Using advanced visualization technologies, the disease map can be used to reveal  
65 the complex relationships among different rare human genetic diseases and support the  
66 clinical diagnosis process.

## 67 **Methods**

68 Methods to measure the distance between phenotypes

69 Human Phenotype Ontology (HPO)[11], proposed by Professor Robinson in 2008, provides a  
70 standardized vocabulary that covers all the common abnormal phenotypes in humans and has  
71 been recognized as a useful annotation of the phenotypic abnormalities of rare genetic  
72 diseases. As most of the modern ontology, HPO is structured as a directed acyclic graph  
73 (DAG), whereby the nodes of the DAG, also called HPO terms, represent abnormal  
74 phenotypic terms in humans. Additionally, these phenotypic terms are linked to their parents  
75 through subclass (“is a”) relationships. Therefore, subclass phenotypic terms have more  
76 accurate definitions than parent phenotypic terms, and each phenotypic term may have  
77 multiple parents, reflecting various semantic types.

78 In this study, we measured the distance between different phenotype terms based on the  
79 hierarchical structure of HPO. For any two HPO terms, the distance can be quantified by the  
80 shortest distance between the corresponding two nodes of the HPO DAG:

$$81 \quad \text{Dist}_p(p_1, p_2) = \frac{\min(d_1 + d_2)}{d_{max}} \quad (\text{Formula 1})$$

82 where  $d_1$  and  $d_2$  represent the distances between two child nodes and their common parent  
83 nodes of the HPO DAG, respectively. Additionally,  $d_{max}$  represents the maximum distance  
84 between nodes in the HPO DAG.

85 Method to measure the distance between genes

86 The Gene Ontology (GO) knowledgebase is the world's largest source of information on  
87 the functions of genes[12]. Similar to the above process, GO can be used to compute the  
88 distance between genes. The GO describes genes from three different aspects: *molecular*  
89 *function, biological process* and *cell component*. Thus, the distance between any two genes  
90 from GO can be defined as the mean value of the shortest distance between gene nodes of the  
91 GO DAG from these three aspects:

92 
$$Dist_g(g_1, g_2) = \frac{Dist_{cc} + Dist_{mf} + Dist_{bp}}{3} \quad (\text{Formula 2})$$

93 where  $Dist_{cc}$ ,  $Dist_{mf}$  and  $Dist_{bp}$  represent the distance between two genes calculated by  
94 Formula 1 from three different aspects.

95 Constructing the rare disease map based on Orphanet

96 Orphanet[13] was established in France in 1997 at the advent of the internet to gather  
97 scarce knowledge on rare diseases to improve the diagnosis, care and treatment of patients  
98 with rare diseases. Currently, Orphanet has become the reference source of information on  
99 rare diseases. In this study, 3,287 diseases with a known clinical phenotype and 3,789  
100 diseases with known pathogenic genes, including 1,718 overlapping diseases, were used to  
101 construct the rare disease map.

102 Because many rare diseases in Orphanet were annotated using HPO terms and frequency,  
103 each of these diseases can be represented by a set of phenotypes with weight. The phenotypic  
104 distance between disease  $d_1$  and disease  $d_2$  can be measured by Formula 3:

$$105 \quad Dist_{dp}(d_1, d_2) = \frac{1}{2} \left( \frac{\sum_{i=1}^m \min_{1 \leq j \leq n} (Dist_p(p_i, p_j)) * (w_i * w_j)}{m} + \frac{\sum_{i=1}^n \min_{1 \leq j \leq m} (Dist_p(p_i, p_j)) * (w_i * w_j)}{n} \right) \quad (\text{Formula 3})$$

106 where  $m$  and  $n$  represent the number of phenotypes contained in disease  $d_1$  and  $d_2$ ,  
107 respectively, and  $Dist(p_i, p_j)$  represents the distance between two phenotypes  $p_i$  and  $p_j$   
108 as shown in Formula 1, and  $w_i$  and  $w_j$  are the frequencies of two phenotypes  $p_i$  and  $p_j$   
109 in  $d_1$  and  $d_2$ , respectively.

110 Similarly, we extracted disease gene relationships from the Orphanet knowledgebase.

111 The genetic distance between diseases can then be transformed into the distance between

112 genes:

$$113 \quad Dist_{dg}(d_1, d_2) = \frac{1}{2} \left( \frac{\sum_{i=1}^m \min_{1 \leq j \leq n} (Dist_g(g_i, g_j))}{m} + \frac{\sum_{i=1}^n \min_{1 \leq j \leq m} (Dist_g(g_i, g_j))}{n} \right) \quad (\text{Formula 4})$$

114 where  $m$  and  $n$  represent the number of genes identified as pathogenic genes in disease  $d_1$

115 and  $d_2$ , respectively, and  $Dist_g(g_i, g_j)$  represents the distance between two genes  $g_i$  and

116  $g_j$  as shown in Formula 2.

117 By calculating these distances among all rare diseases from Orphanet, we generated two

118 distance matrices with the sizes of  $3287 \times 3287$  and  $3789 \times 3789$  for phenotype and gene,

119 respectively. We used multidimensional scaling[14] (*cmdscale* from the package *stats* in R)

120 to convert the distance matrix into 2D points which can be visualized as a map.

121 To further explore the internal relationship between phenotypes and genes of rare genetic

122 diseases, we divided the rare disease map into several disease clusters using the K-means

123 clustering method. To determine the optimal  $k$  for disease clustering, a bootstrap approach

124 implemented in the *clusterboot* function from the *fpc* package in R (version 3.4.0) was used.

125 Thus, we developed a web-based interactive rare disease map called RDMap  
126 (<http://RDMap.nbscn.org>) based on ECharts[15]. All the data precessions were under R, and  
127 the online web service was developed using Node.js and Python.

128 To evaluate the RDMap in clinical diagnosis, an in-silico test and a case report test were  
129 used. The targeted disease ranked in the recommend disease list was used to evaluate the  
130 performance of RDMap in clinical practice using only clinical information.

## 131 **Results**

132 In this study, 3,287 diseases in Orphanet with a clinical phenotype and 3,789 diseases  
133 with known pathogenic genes in Orphanet were plotted into the Euclidean space, as shown in  
134 Fig 1. In total, 17 phenotype-based disease clusters and 18 gene-based disease clusters were  
135 generated and highlighted by different colors. The detailed information of disease clustering  
136 is explained in the supplemental material.

137 We published RDMap online (<http://RDMap.nbscn.org>) to help the user to explore rare  
138 disease relationships interactively. The map supports zooming and panning to find and locate  
139 special diseases like the widely used google map (Fig 2). It also supports a feature-based  
140 exploration such as one or more phenotypes will locate the most likely rare diseases on the

141 map and filter by the similarity score (Fig 2A). The detailed information about the disease  
142 will be shown when the disease is selected on the RDmap (Fig 2B). When a disease was  
143 selected on RDmap, user can jump between phenotype map and gene map through a toolbar  
144 button. This will help user to explore the interested diseases from different levels. The user  
145 guide for RDmap is provided in the supplemental material.

146 In the traditional knowledge base, the entries were indexed by keywords, and users are  
147 required to use the exact term used in the knowledge base to query the knowledge. However,  
148 obtaining the exact phenotype features in a particular patient clinically and matching it with  
149 the standard phenotype terms used to annotated diseases in knowledgebases remain  
150 challenges[16]. Because thousands of genetic diseases are known, their clinical presentations  
151 often overlap in patients and are typically abridged with respect to classical descriptions. The  
152 incompleteness, heterogeneity, imprecision, and noise in phenotype description sometimes  
153 miss the right diagnosis and even to wrong diagnoses.

154 To compare the performance of RDMap and direct simple matching of phenotype vectors  
155 (Jaccard matching), we designed an in silico evaluation test in which 1,000 rare genetic  
156 diseases from the Orphanet database are taken as the target diseases, and each disease is

157 represented as a set of four characteristic phenotypes with the highest frequency of the disease.

158 Then, the adjacent node or parent node of the phenotype in the HPO DAG is defined as the

159 imprecise phenotype of the target phenotype. The rank of the target diseases in the query results

160 was used to evaluate the performance. The performance of the D-P method decreases

161 significantly as the number of imprecise phenotypes increases (Fig 3). This finding also

162 explains why it is very difficult to diagnose a rare genetic disease accurately in clinical practice

163 using only clinical phenotypes. The RDmap-proposed methods Similarity (one-way distance

164 calculation) and Similarity-Avg (average of two-way distance calculation) both have an

165 obvious advantage over the Jaccard matching method, particularly regarding imprecision

166 phenotypes. We also noticed that the one-way distance algorithm (Similarity) is more stable in

167 the disease recommendation than the Similarity-Avg in this scenario. This one-way distance

168 algorithm had been implemented in published RDmap.

169 To further evaluate the performance of RDmap in clinical practice, we collected 20 rare

170 disease cases reported by the Orphanet Journal of Rare Diseases as test cases. The case

171 presentations in publication were converted into HPO terms manually by one of the authors.

172 The targeted diseases ranked in the similarity search results on RDmap are shown in Table 1

173 (the detailed information of each test case is shown in the supplemental material). RDmap  
174 worked very well in most cases with clear clinical phenotype descriptions and the targeted  
175 disease rank average 1.8 and median top 1 in all diseases. If you check the detail test case in  
176 supplemental material, some other similar diseases for some test cases remain under  
177 consideration for the clinician. These results were all highlighted on RDmap, and a quick  
178 check of a typical phenotype and its frequency in these potential diagnoses on RDmap will  
179 support clinicians to make a decision.

## 180 **Discussion**

181 In this study, we constructed two maps of rare human genetic diseases based on phenotypic  
182 characteristics and genes and divided these genetic diseases into several disease clusters.  
183 Because diseases from the same cluster are related in phenotypic characteristics or gene  
184 functions, correlating clusters between two maps will be helpful to understand the  
185 physiological and pathological bases of related genetic diseases. Consistent with the results of  
186 Goh et al.[17], most of the diseases in the same phenotype-based cluster tend to have similar  
187 phenotypic characteristics. In total, 1,718 diseases overlapped in the two maps, and the  
188 relationship between 17 phenotype-based clusters and 18 gene-based clusters are shown in an

189 alluvial diagram in Fig 4 and supplemental material. The complicated branches among these  
190 clusters further confirmed the complicated relationships among the pathogenic genes and  
191 phenotypes of rare genetic diseases. Diseases with similar phenotypes may be divided into  
192 different gene-based disease clusters. However, diseases from the same gene-based clusters  
193 also present diverse phenotypes. However, we also noticed mainstreams among different  
194 clusters. RDmap also provide a button to jump from disease selected in phenotype-based map  
195 to same disease in gene-based map and vice versa. So, there are 1,718 bridges between two  
196 maps. These will inspire researchers to evaluate the inner relationships among pathogenic  
197 genes and phenotypes.

198 In recent years, to reveal the similar relationships between different human genetic diseases,  
199 many studies have used various ways to construct a human genetic disease network. For  
200 example, Goh et al. extracted known disease-gene associations from the OMIM database and  
201 constructed the human disease network[17]. The core idea of their method is that two  
202 diseases are related if they share at least one common gene. Lee et al. constructed a human  
203 disease network based on cell metabolism, and the core idea of this method is that two  
204 diseases are related if the related mutant enzyme catalyzes the adjacent metabolism

205 reaction[18]. Zhang et al. constructed a disease phenotype network using the similarity  
206 between phenotypes to obtain the gene function module[19]. However, RDmap shows a  
207 complicated disease relationship in a user-interactive map that we believe will be conducive  
208 to the discovery of potential relationships among pathogenic genes and phenotypic  
209 characteristics among many genetic diseases. The map-style visualization that reflects the  
210 distance of disease more intuitively will inspire investigators to understand the inner  
211 relationships among these diseases and their potential treatments and identify new pathogenic  
212 genes. Moreover, this tool can help the clinician or genetic counselor accurately diagnose rare  
213 genetic diseases effectively, especially when the clinical phenotypes are incomplete,  
214 imprecise or noise.

215       This study has some limitations. First, the two disease maps still not covering all rare  
216 genetic diseases. It based on a history version of Ophanet when this project started. Second,  
217 when a novel disease is enrolled in the map, all the disease maps and disease clustering need  
218 to be recalculated and updated. But we will annually update it based on the feedbacks from  
219 the community.

220 **Acknowledgements**

221 Not applicable

222 **Authors' contributions**

223 JY and CD developed the website and the initial draft of the manuscript. HD and QS  
224 provided data and developed the analysis protocol. HL conceived this project, analysis the  
225 data, design the website and revised the manuscript.

226 **Funding**

227 This study was supported by the National Natural Science Foundation of China (81871456)  
228 and National Key R&D Program of China (2016YFC0901905).

229 **Availability of data and materials**

230 All data generated or analyzed during this study are published online or included in this  
231 published article and its supplementary files.

232 **Ethics approval and consent to participate**

233 Not applicable

234 **Consent for publication**

235 Not applicable

236 **Competing interests**

237 The authors declare no conflict of interest.

238

239 **References**

240 1. Boycott KM, Rath A, Chong JX, Hartley T, Alkuraya FS, Baynam G, et al. International

241 Cooperation to Enable the Diagnosis of All Rare Genetic Diseases. *Am J Hum Genet.*

242 2017;100:695–705.

243 2. Aerts S, Lambrechts D, Maity S, Van Loo P, Coessens B, De Smet F, et al. *Gene*

244 prioritization through genomic data fusion. *Nat Biotechnol.* 2006;24:537–44.

245 3. Chavali S, Barrenas F, Kanduri K, Benson M. Network properties of human disease genes

246 with pleiotropic effects. *BMC Syst Biol.* 2010;4:78.

- 247 4. Franke L, Van Bakel H, Fokkens L, De Jong ED, Egmont-Petersen M, Wijmenga C.  
248 Reconstruction of a functional human gene network, with an application for prioritizing  
249 positional candidate genes. *Am J Hum Genet.* 2006;78:1011–25.
- 250 5. Robinson P, Mundlos S. The Human Phenotype Ontology. *Clin Genet* [Internet].  
251 2010;77:525–34. Available from: <http://doi.wiley.com/10.1111/j.1399-0004.2010.01436.x>
- 252 6. Robinson PN, Köhler S, Bauer S, Seelow D, Horn D, Mundlos S. The Human Phenotype  
253 Ontology: A Tool for Annotating and Analyzing Human Hereditary Disease. *Am J Hum*  
254 *Genet.* 2008;83:610–5.
- 255 7. Yu L, Ma X, Zhang L, Zhang J, Gao L. Prediction of new drug indications based on  
256 clinical data and network modularity. *Sci Rep.* 2016;6:32530.
- 257 8. Gottlieb A, Stein GY, Ruppin E, Sharan R. PREDICT: a method for inferring novel drug  
258 indications with application to personalized medicine. *Mol Syst Biol.* 2011;7:496.
- 259 9. Luo H, Wang J, Li M, Luo J, Peng X, Wu F-X, et al. Drug repositioning based on  
260 comprehensive similarity measures and Bi-Random walk algorithm. *Bioinformatics.*  
261 2016;32:2664–71.

- 262 10. Yu L, Wang B, Ma X, Gao L. The extraction of drug-disease correlations based on  
263 module distance in incomplete human interactome. *BMC Syst Biol.* 2016;10:111.
- 264 11. Köhler S, Doelken SC, Mungall CJ, Bauer S, Firth H V., Bailleul-Forestier I, et al. The  
265 Human Phenotype Ontology project: Linking molecular biology and disease through  
266 phenotype data. *Nucleic Acids Res.* 2014;
- 267 12. Carbon S, Douglass E, Dunn N, Good B, Harris NL, Lewis SE, et al. The Gene Ontology  
268 Resource: 20 years and still GOing strong. *Nucleic Acids Res.* 2019;47:D330–8.
- 269 13. Rath A, Olry A, Dhombres F, Brandt MM, Urbero B, Ayme S. Representation of rare  
270 diseases in health information systems: The orphanet approach to serve a wide range of end  
271 users. *Hum Mutat.* 2012;33:803–8.
- 272 14. Mead A. Review of the Development of Multidimensional Scaling Methods. *Stat.*  
273 1992;41:27.
- 274 15. Li D, Mei H, Shen Y, Su S, Zhang W, Wang J, et al. ECharts: A declarative framework  
275 for rapid construction of web-based visualization. *Vis Informatics.* 2018;2:136–46.

- 276 16. Eldomery MK, Coban-Akdemir Z, Harel T, Rosenfeld JA, Gambin T, Stray-Pedersen A,  
277 et al. Lessons learned from additional research analyses of unsolved clinical exome cases.  
278 *Genome Med.* 2017;9:26.
- 279 17. Goh K-I, Cusick ME, Valle D, Childs B, Vidal M, Barabasi A-L. The human disease  
280 network. *Proc Natl Acad Sci.* 2007;104:8685–90.
- 281 18. Lee DS, Park J, Kay KA, Christakis NA, Oltvai ZN, Barabasi A-L. The implications of  
282 human metabolic network topology for disease comorbidity. *Proc Natl Acad Sci.*  
283 2008;105:9880–5.
- 284 19. Zhang S-H, Wu C, Li X, Chen X, Jiang W, Gong B-S, et al. From phenotype to gene:  
285 Detecting disease-specific gene functional modules via a text-based human disease  
286 phenotype network construction. *FEBS Lett.* 2010;584:3635–43.
- 287 20. Al-Owain M, Mohamed S, Kaya N, Zagal A, Matthijs G, Jaeken J. A novel mutation and  
288 first report of dilated cardiomyopathy in ALG6-CDG (CDG-Ic): a case report. *Orphanet J*  
289 *Rare Dis.* 2010;5:7.

- 290 21. Böhm J, Yiş U, Ortaç R, Çakmakçı H, Kurul S, Dirik E, et al. Case report of intrafamilial  
291 variability in autosomal recessive centronuclear myopathy associated to a novel BIN1 stop  
292 mutation. *Orphanet J Rare Dis.* 2010;5:35.
- 293 22. Acién P, Galán F, Manchón I, Ruiz E, Acién M, Alcaraz LA. Hereditary renal adysplasia,  
294 pulmonary hypoplasia and Mayer-Rokitansky-Küster-Hauser (MRKH) syndrome: a case  
295 report. *Orphanet J Rare Dis.* 2010;5:6.
- 296 23. Mejia-Gaviria N, Gil-Pêa H, Coto E, Pérez-Menéndez TM, Santos F. Genetic and clinical  
297 peculiarities in a new family with hereditary hypophosphatemic rickets with hypercalciuria:  
298 A case report. *Orphanet J Rare Dis.* 2010;
- 299 24. Joy T, Cao H, Black G, Malik R, Charlton-Menys V, Hegele RA, et al. Alstrom  
300 syndrome (OMIM 203800): a case report and literature review. *Orphanet J Rare Dis.*  
301 2007;2:49.
- 302 25. Zhu Y, Zou Y, Yu Q, Sun H, Mou S, Xu S, et al. Combined surgical-orthodontic  
303 treatment of patients with cleidocranial dysplasia: case report and review of the literature.  
304 *Orphanet J Rare Dis.* 2018;13:217.

- 305 26. Zamel R, Khan R, Pollex RL, Hegele RA. Abetalipoproteinemia: two case reports and  
306 literature review. *Orphanet J Rare Dis.* 2008;3:19.
- 307 27. Vroegindeweyj LHP, Boon AJW, Wilson JHP, Langendonk JG. Effects of iron chelation  
308 therapy on the clinical course of aceruloplasminemia: an analysis of aggregated case reports.  
309 *Orphanet J Rare Dis.* 2020;15:105.
- 310 28. Zhou L, Ouyang R, Luo H, Ren S, Chen P, Peng Y, et al. Efficacy of sirolimus for the  
311 prevention of recurrent pneumothorax in patients with lymphangioleiomyomatosis: a case  
312 series. *Orphanet J Rare Dis.* 2018;13:168.
- 313 29. Dias RP, Buchanan CR, Thomas N, Lim S, Solanki G, Connor SEJ, et al. Os  
314 ontoideum in wolcott-rallison syndrome: A case series of 4 patients. *Orphanet J Rare Dis.*  
315 2016;
- 316 30. Valayannopoulos V, Nicely H, Harmatz P, Turbeville S. Mucopolysaccharidosis VI.  
317 *Orphanet J Rare Dis.* 2010;5:5.
- 318 31. Biesecker LG. The Greig cephalopolysyndactyly syndrome. *Orphanet J Rare Dis.* 2008;
- 319 32. Germain DP. Fabry disease. *Orphanet J Rare Dis.* 2010;5:30.

320 33. Drera B, Ritelli M, Zoppi N, Wischmeijer A, Gnoli M, Fattori R, et al. Loeys-Dietz  
321 syndrome type I and type II: Clinical findings and novel mutations in two Italian patients.  
322 Orphanet J Rare Dis. 2009;

323 34. Reibel A, Manière M-C, Clauss F, Droz D, Alembik Y, Mornet E, et al. Orofacial  
324 phenotype and genotype findings in all subtypes of hypophosphatasia. Orphanet J Rare Dis.  
325 2009;4:6.

326 35. Sarfati J, Bouvattier C, Bry-Gaillard H, Cartes A, Bouligand J, Young J. Kallmann  
327 syndrome with FGFR1 and KAL1 mutations detected during fetal life. Orphanet J Rare Dis.  
328 2015;10:71.

329 36. Weisfeld-Adams JD, Mehta L, Rucker JC, Dembitzer FR, Szporn A, Lublin FD, et al.  
330 Atypical Chédiak-Higashi syndrome with attenuated phenotype: three adult siblings  
331 homozygous for a novel LYST deletion and with neurodegenerative disease. Orphanet J Rare  
332 Dis. 2013;8:46.

333 37. Mowat DR, Wilson MJ, Goossens M. Mowat-Wilson syndrome. J. Med. Genet. 2003.

334 38. Chrzanowska KH, Gregorek H, Dembowska-Bagińska B, Kalina MA, Digweed M.  
335 Nijmegen breakage syndrome (NBS). Orphanet J Rare Dis. 2012;7:13.

336 39. Marshall BA, Paciorkowski AR, Hoekel J, Karzon R, Wasson J, Viehover A, et al.

337 Phenotypic characteristics of early Wolfram syndrome. *Orphanet J Rare Dis.* 2013;

338

339

340 **Figure Legends:**

341 **Fig. 1. Rare disease maps and clusters. The locations reflect the distance among diseases,**  
342 **and the size of the points reflect the prevalence of rare diseases** A. Rare disease map and  
343 clusters based on phenotype. B. Rare disease map and clusters based on gene.

344

345 **Fig. 2. Rare disease map zooming, panning, location, filtering and disease detail**  
346 **information.** A. The RDmap locates similar diseases based on phenotype search. The slider in  
347 the left bottom corner can control the similarity filtering threshold by the user. The prevalence  
348 options switch at the bottom right can filter the results based on prevalence. B. When a disease  
349 was selected on RDmap, its detail information will be shown like this.

350

351 **Fig. 3. In-silico test of RDMap.** Performance of RDMap under conditions with different  
352 numbers of imprecision phenotypes for the search

353

354 **Fig. 4. Alluvial diagram between 17 phenotype-based rare disease clusters and 18**  
355 **gene-based rare disease clusters.**

356

**Table 1 Evaluation of RDmap based on cases from publications**

Ref.	Disease (OMIM)	phenotypes	rank	Sim. score
Case1[20]	Congenital disorder of glycosylation (OMIM 603147)	HP:0025356 Psychomotor retardation/Psychomotor HP:0001252 Muscular hypotonia HP:0001644 Dilated cardiomyopathy HP:0001250 Seizures HP:0000486 Strabismus HP:0006610 Wide intermamillary distance	6	0.0625
Case2[21]	Centronuclear myopathies (OMIM 255200)	HP:0009073 Progressive proximal muscle weakness HP:0000297 Facial hypotonia HP:0000508 Ptosis HP:0000602 Ophthalmoplegia HP:0001315 Reduced tendon reflexes HP:0001256 Intellectual disability, mild	4	0.0486
Case3[22]	Mayer-Rokitansky-Küster-Hauser syndrome (OMIM 277000)	HP:0002089 Pulmonary hypoplasia HP:0000122 Unilateral renal agenesis HP:0000151 Aplasia of the uterus HP:0008726 Hypoplasia of the vagina	4	0.0937
Case4[23]	Hereditary hypophosphatemic rickets with hypercalciuria (OMIM 241530)	HP:0002148 Hypophosphatemia HP:0002150 Hypercalciuria	1	0
Case5[24]	Alström syndrome (OMIM 203800)	HP:0000662 Night blindness HP:0000618 Blindness HP:0012330 Pyelonephritis HP:0000822 Hypertension HP:0000819 Diabetes mellitus HP:0000510 Retinitis pigmentosa HP:0000518 Cataract	1	0.535
Case6[25]	Cleidocranial dysplasia (OMIM 119600)	HP:0000684 Delayed eruption of teeth HP:0000164 Abnormality of the teeth HP:0000316 Hypertelorism HP:0011069 Increased number of teeth	1	0
Case7[26]	Abetalipoproteinemia (OMIM 200100)	HP:0002630 Fat malabsorption HP:0001251 Ataxia HP:0001324 Muscle weakness HP:0001315 Reduced tendon reflexes	4	0.0416
Case8[27]	Aceruloplasminemia (OMIM 604290)	HP:0001935 Microcytic anemia HP:0001260 Dysarthria HP:0001288 Gait disturbance HP:0000819 Diabetes mellitus HP:0001903 Anemia HP:0001300 Parkinsonism	1	0.0416
Case9[28]	Lymphangi leiomyomatosis (OMIM 606690)	HP:0100749 Chest pain HP:0002094 Dyspnea HP:0002107 Pneumothorax	1	0
Case10[29]	Wolcott-Rallison syndrome (OMIM 226980)	HP:0006554 Acute hepatic failure HP:0001298 Encephalopathy HP:0000083 Renal insufficiency HP:0002654 Multiple epiphyseal dysplasia	1	0.0208
Case11[30]	Mucopolysaccharidosis type 6 (OMIM 253200)	HP:0000280 Coarse facial features HP:0000470 Short neck HP:0000158 Macroglossia HP:0002808 Kyphosis HP:0012471 Thick vermilion border	1	0.0083

Case12[31]	Greig cephalopolysyndactyly syndrome (OMIM 175700)	HP:0000256 Macrocephaly HP:0011304 Broad thumb HP:0001159 Syndactyly HP:0001162 Postaxial hand polydactyly HP:0005873 Polysyndactyly of hallux	1	0.016
Case13[32]	Fabry disease (OMIM 301500)	Angiokeratoma (HP:0001014)	1	0
Case14[33]	Loeys-Dietz syndrome (OMIM 609192)	Camptodactyly of finger (HP:0100490) Ulnar deviation of the hand or fingers of the hand (HP:0001193) Bilateral talipes equinovarus (HP:0001776) Blue sclerae (HP:0000592) Microretrognathia (HP:0000308) High palate (HP:0000218) Bifid uvula (HP:0000193)	1	0.0628
Case15[34]	Hypophosphatasia (OMIM 146300)	Recurrent fractures (HP:0002757) Craniosynostosis (HP:0001363) Premature loss of teeth (HP:0006480)	1	0.0138
Case16[35]	Kallmann syndrome (OMIM 308700)	Oligomenorrhea (HP:0000876) Breast hypoplasia (HP:0003187) Anosmia (HP:0000458) Hearing impairment (HP:0000365) Reduced number of teeth (HP:0009804)	1	0.0249
Case17[36]	Chédiak–Higashi syndrome (OMIM 214500)	Lower limb muscle weakness (HP:0007340) Dementia (HP:0000726) Ataxia (HP:0001251) Hypermetric saccades (HP:0007338) Bradykinesia (HP:0002067) Periodontitis (HP:0000704)	4	0.0972
Case18[37]	Mowat–Wilson syndrome (OMIM 235730)	Open mouth (HP:0000194) Abnormality of the eyebrow (HP:0000534) Frontal bossing (HP:0002007) Deeply set eye (HP:0000490) Wide nasal bridge (HP:0000431) Strabismus (HP:0000486)	1	0
Case19[38]	Nijmegen breakage syndrome (OMIM 251260)	Microcephaly (HP:0000252) Sloping forehead (HP:0000340) Retrognathia (HP:0000278) Macrotia (HP:0000400) Bulbous nose (HP:0000414)	1	0.0116
Case20[39]	Wolfram syndrome (OMIM 222300)	Diabetes mellitus (HP:0000819) Optic atrophy (HP:0000648) Diabetes insipidus (HP:0000873) Hearing impairment (HP:0000365) Gastroesophageal reflux (HP:0002020)	1	0.0333