

2019/11/11

HwaZhong Argicultural University

PubCaseFinder : a case-report-based, phenotype-driven differential-diagnosis system for rare diseases

Jin-Dong Kim

Research Organization of Information and systems (ROIS)

Database Center for Life Science (DBCLS)

ARTICLE | [VOLUME 103, ISSUE 3, P389-399, SEPTEMBER 06, 2018](#)

PubCaseFinder: A Case-Report-Based, Phenotype-Driven Differential-Diagnosis System for Rare Diseases

[Toyofumi Fujiwara](#)   • [Yasunori Yamamoto](#) • [Jin-Dong Kim](#) • [Orion Buske](#) • [Toshihisa Takagi](#)

[Open Archive](#) • Published: August 30, 2018 • DOI: <https://doi.org/10.1016/j.ajhg.2018.08.003> •



What is a RARE DISEASE ?

Any disease, disorder, illness or condition affecting **fewer than 200,000** people in the United States is considered RARE.¹

1 in 10

Americans has a **RARE DISEASE** } 30 million people have a serious, lifelong condition.



Holding hands, they would circle the globe about **1.5 times**



More than half are children¹



7,000

RARE DISEASES exist, with less than 500 FDA-approved treatments²

ONLY 5% of RARE DISEASES have treatments.²

Patients with RARE DISEASES are frequently **misdiagnosed or undiagnosed.**

80% of RARE DISEASES ARE GENETICALLY BASED.²



Many RARE DISEASES result in premature death of infants & young children or are fatal in early adulthood.²

Families & private foundations provide about **3%** of ALL medical research funding in the U.S.⁶



Rare diseases

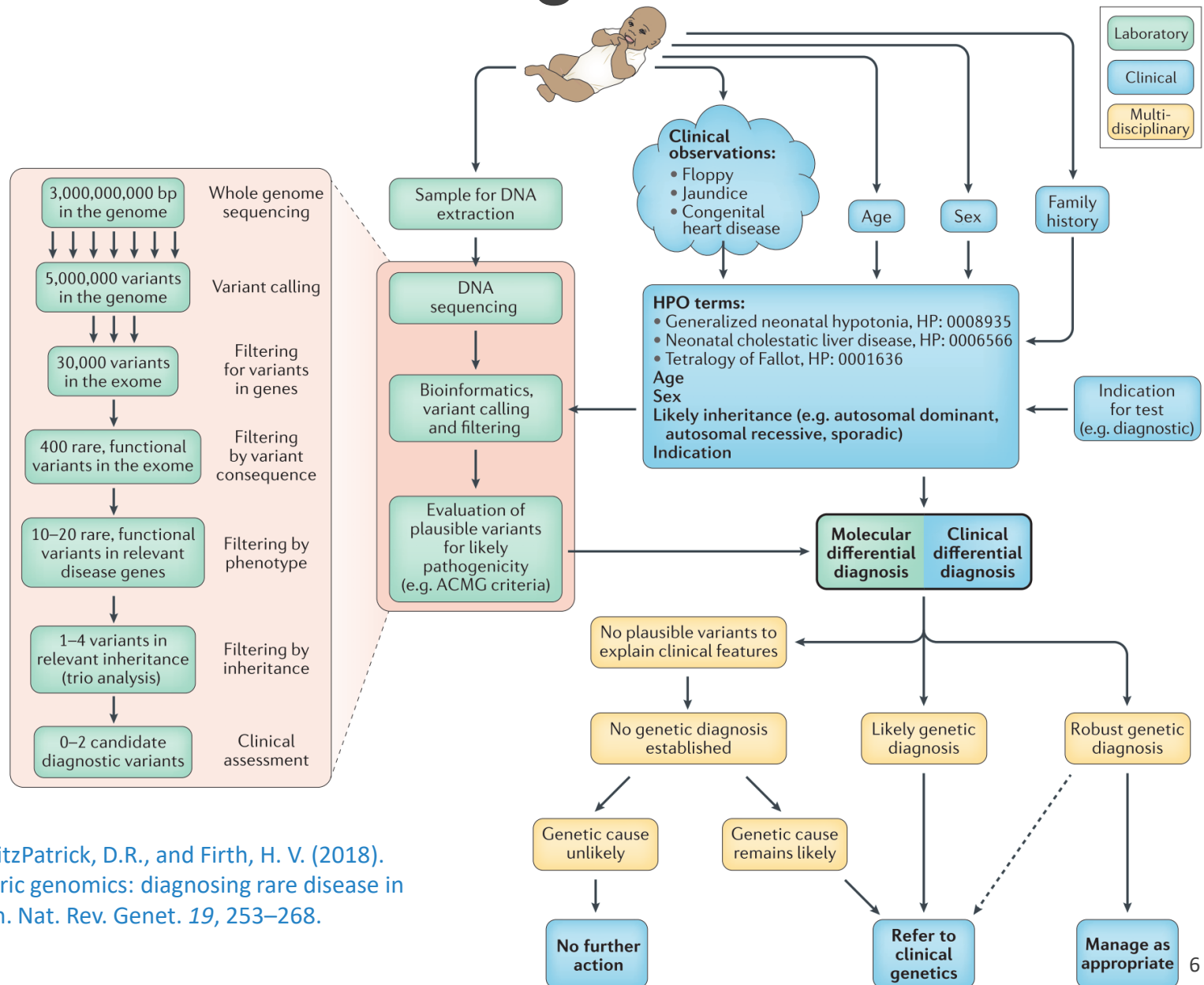
- ❑ For patients of rare diseases, getting diagnosis at an early stage is important
- ❑ E.g., Phenylketonuria (PKU)
 - ✓ An inborn error of metabolism
 - ✓ Untreated, it can lead to intellectual disability, seizures, behavioral problems, and mental disorders.
 - ✓ If it is diagnosed early enough, an affected newborn can grow up with normal development through diet, or a combination of diet and medication.

NGS-based diagnosis

□ Sawyer et al., 2016

- ✓ “Recent advances in sequencing, in particular whole-exome sequencing (WES), are identifying the genetic basis of disease for 25–40% of patients.”

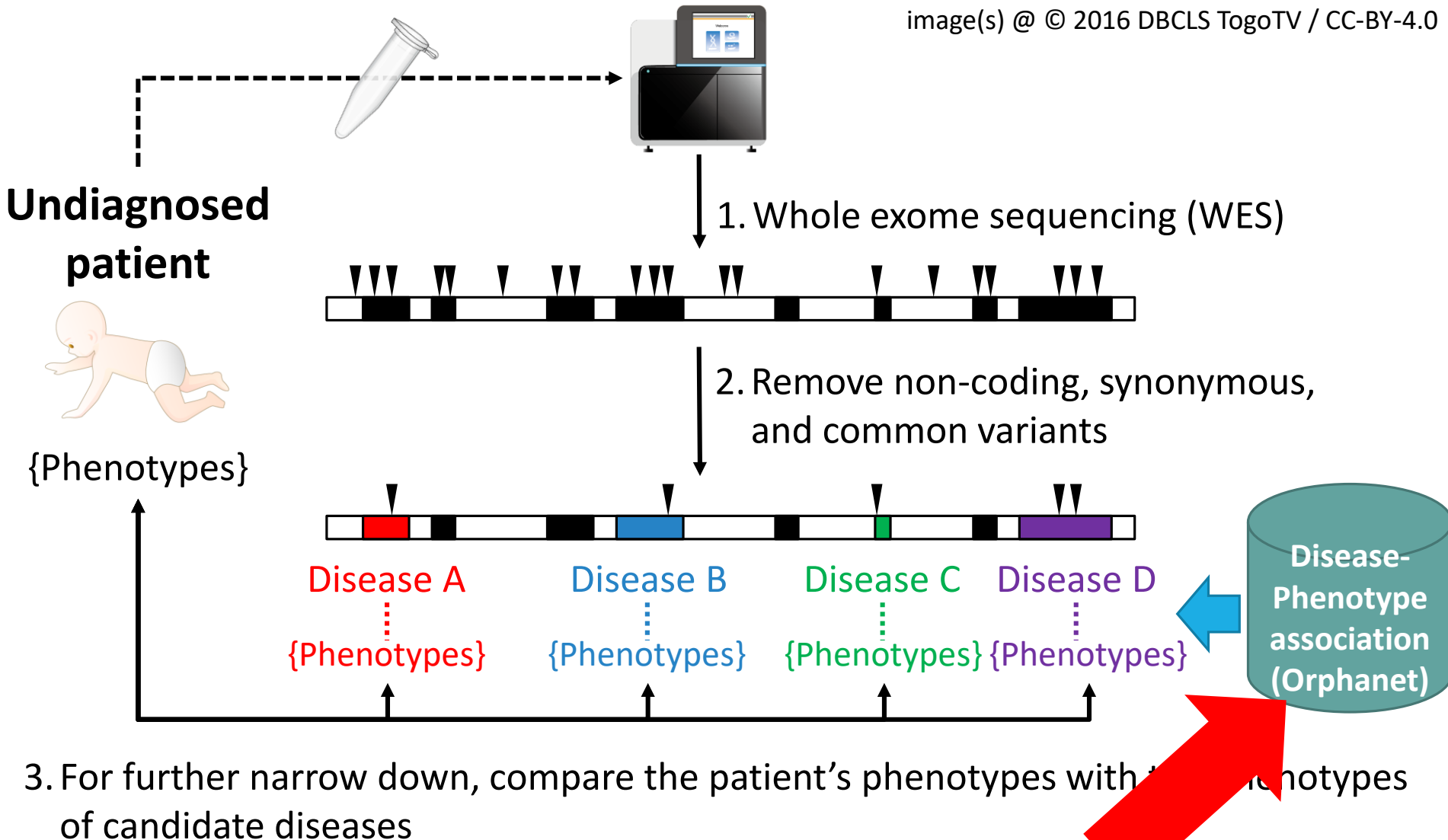
NGS-based diagnosis



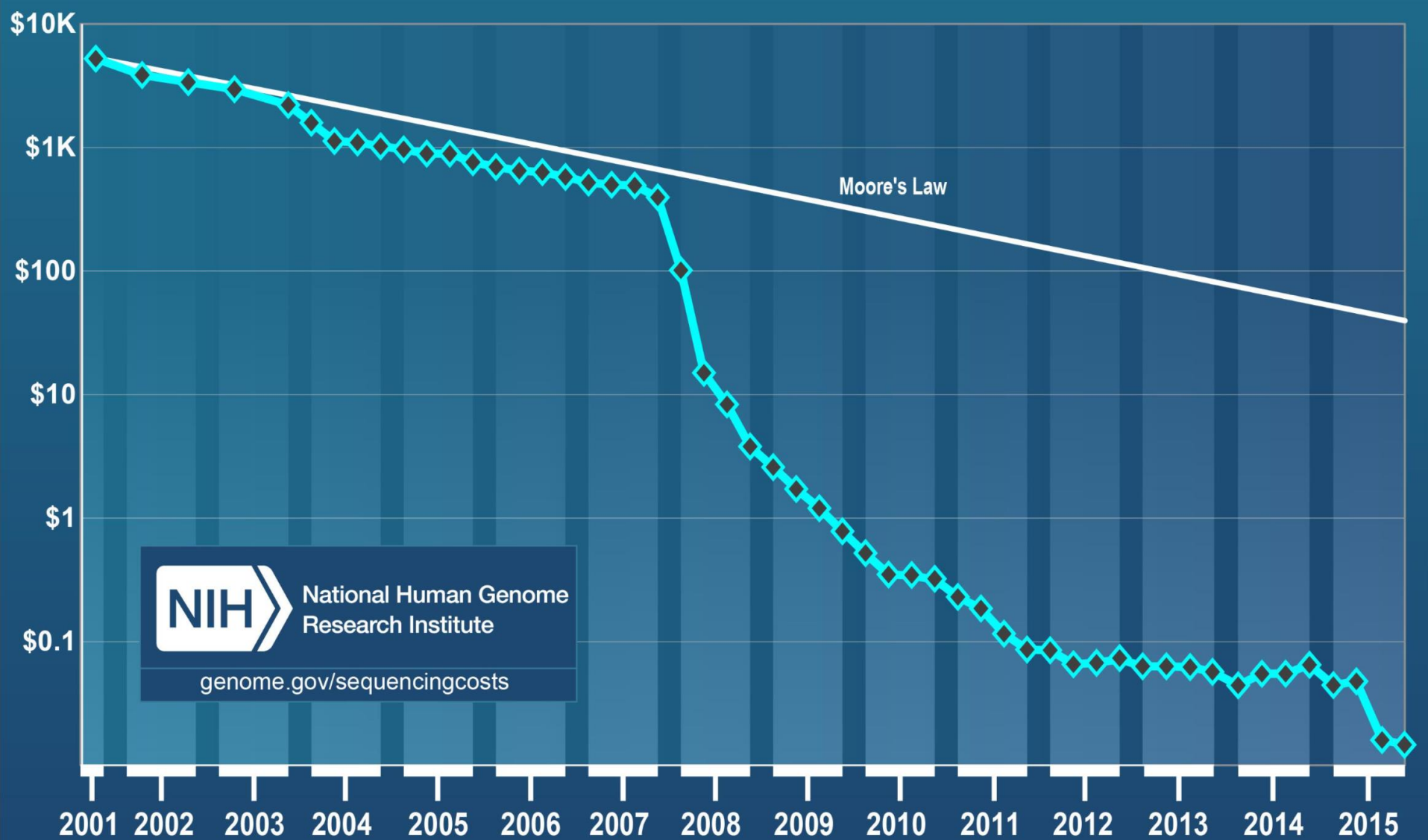
Wright, C.F., FitzPatrick, D.R., and Firth, H. V. (2018). Paediatric genomics: diagnosing rare disease in children. *Nat. Rev. Genet.* 19, 253–268.

For undiagnosed patients

image(s) @ © 2016 DBCLS TogoTV / CC-BY-4.0



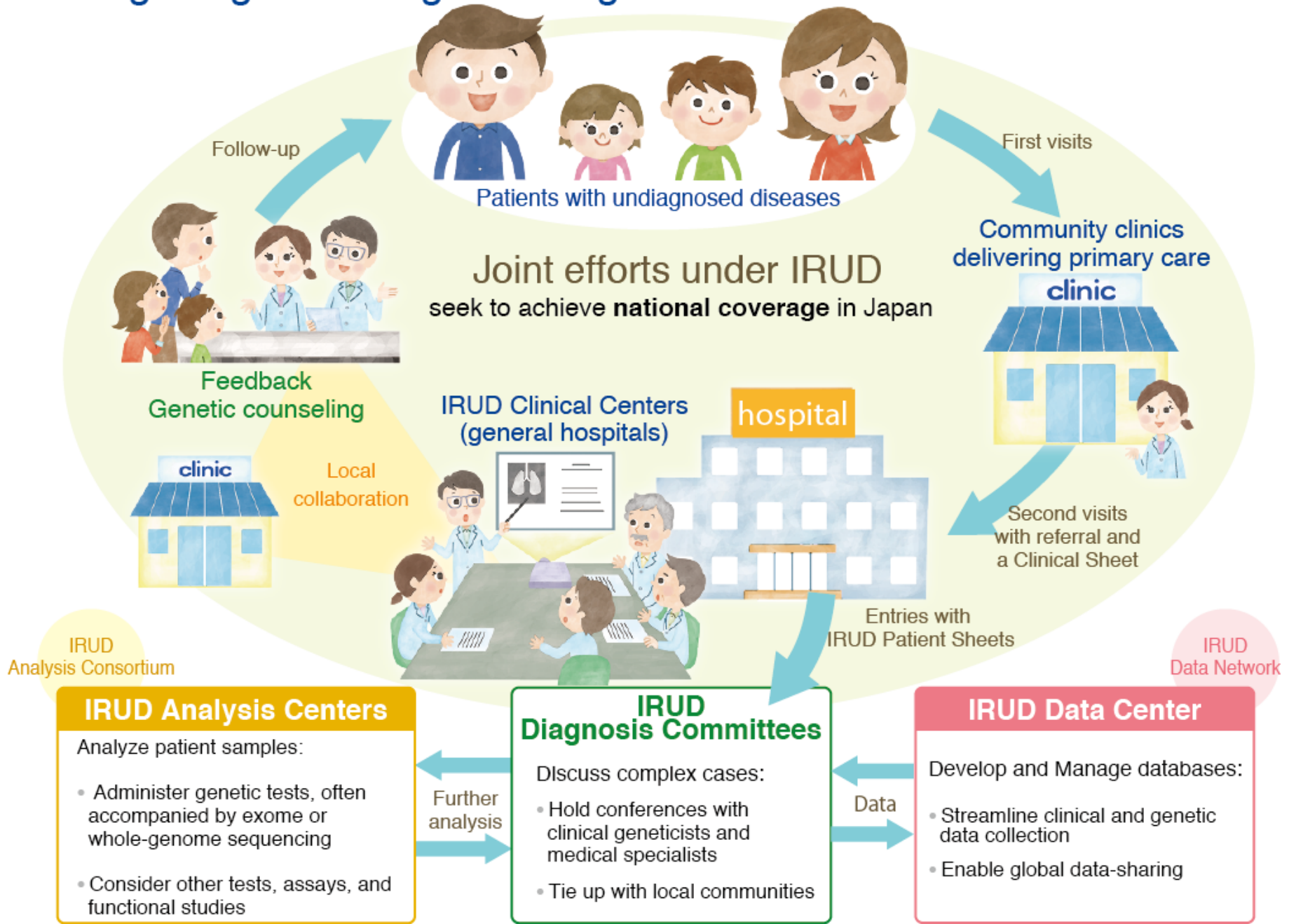
Cost per Raw Megabase of DNA Sequence



IRUD

- ❑ Initiative on Rare and Undiagnosed Diseases
- ❑ A national system to find diagnoses for undiagnosed patients across Japan, by applying NGS-based diagnostic tests
- ❑ Applied to >6,000 cases, so far
- ❑ Achieved diagnosis rate of 33%

Initiative on Rare and Undiagnosed Diseases (IRUD) : Integrating Knowledge for Diagnoses



IRUD Clinical Centers

第2期IRUD拠点体制地図

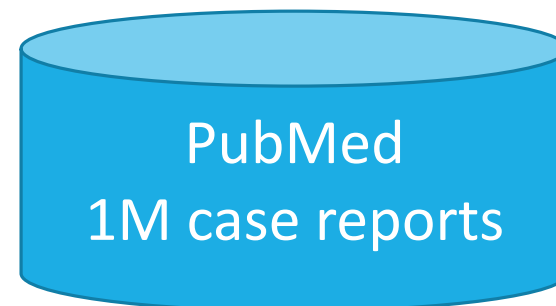
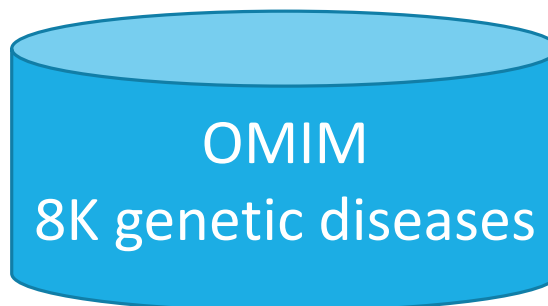
2018年4月 全国37拠点

- ⑱山梨大学
- ⑲浜松医科大学
- ⑳名古屋市立大学
- ㉑名古屋市立大学
- ㉒藤田保健衛生大学
- ㉓大阪市立大学
- ㉔大阪大学
- ㉕国立循環器病研究センター
- ㉖大阪母子医療センター
- ㉗京都大学
- ㉘神戸大学
- ㉙鳥取大学
- ㉚川崎医療福祉大学
- ㉛広島大学
- ㉜徳島大学
- ㉝愛媛大学
- ㉞長崎大学
- ㉟熊本大学
- ㊱琉球大学

- ①札幌医科大学
- ②北海道大学
- ③旭川医科大学
- ④秋田大学
- ⑤東北大学
- ⑥千葉大学
- ⑦東京医科歯科大学
- ⑧東京大学
- ⑨国立成育医療研究センター
- ⑩慶応義塾大学
- ⑪東京女子医科大学
- ⑫東京都立小児総合医療センター
- ⑬神奈川県立こども医療センター
- ⑭横浜市立大学
- ⑮新潟大学
- ⑯金沢大学
- ⑰信州大学

IRUDコーディネーティングセンター
国立精神・神経医療研究センター

Needs for access to relevant resources



患者の徴候

HP:0001878 溶血性貧血

疾患を絞り込む

結果の要約をダウンロード

希少疾患 (Orphanet)

合計: 4,066 件

順位 (類似度)

1 (100.0%)

1 (100.0%)

3 (97.9%)

患者の徴候

HP:0001878 溶血性貧血

疾患を絞り込む

結果の要約をダウンロード

希少疾患 (Orphanet)

合計: 6969 件

順位 (類似度)

1 (100.0%)

2 (91.9%)

2 (91.9%)

4 (90.1%)

患者の徴候

HP:0001878 溶血性貧血

疾患を絞り込む

ENT:286 ANK1 (SPH)
ENT:6521 SLC4A1 (C)
ENT:57674 RNF213
ENT:5906 RAP1A (K)

結果の要約をダウンロード

希少疾患 (Orphanet)

合計: 9 件

順位 (類似度)

1 (100.0%)

2 (81.3%)

3 (80.1%)

症例報告 徴候・症状 疾患原因遺伝子

患者の徴候・症状を入力

HP:0001878 溶血性貧血 HP:0000952 Jaundice HP:0001744 脾腫 HP:0004444 球状赤血球症 HP:0001081 胆石症

症例報告を絞り込む

症例報告を検索

クリア

合計: 205 (症例報告)

1 2 3 ... 21 >>

10 (表示件数)

対応する徴候・症状 遺伝子 変異 キーワード (MeSH)

順位 (類似度) PMID (PMCID)

1 (91.9%) 19763011 Hereditary spherocytosis in a 27-year-old woman: case report. Hassan A, Babadoko AA, Isa AH, Abunimye P. Ann Afr Med. 2009;8(1):61-3. [抄録を表示]

黄疸 脾腫 溶血性貧血 球状赤血球症

ヒト 女 成人 栄養補助食品 脾腫

2 (90.1%) 27423290 Hereditary Spherocytosis with Splenomegaly and Cholelithiasis in a Young Male of Western Region of Nepal - A Case Report. Ghimire P, Gurung NV, Shrestha S, Poudel SR, Chapagain A. Kathmandu Univ Med J (KUMJ). 2015;13(52):366-8. [抄録を表示]

黄疸 胆石症 脾腫 溶血性貧血

3 (89.3%) 8717295 [Intrathoracic extramedullary hematopoiesis in a case of hereditary spherocytosis]. Takahashi R, Igarashi T, Nakagawa A, Ohuchi H, Nishino M, Murakami S, Yoshida Y, Abe S. Nihon Kyobu Shikkan Gakkai Zassh. 1996;34(1):71-5. [抄録を表示]

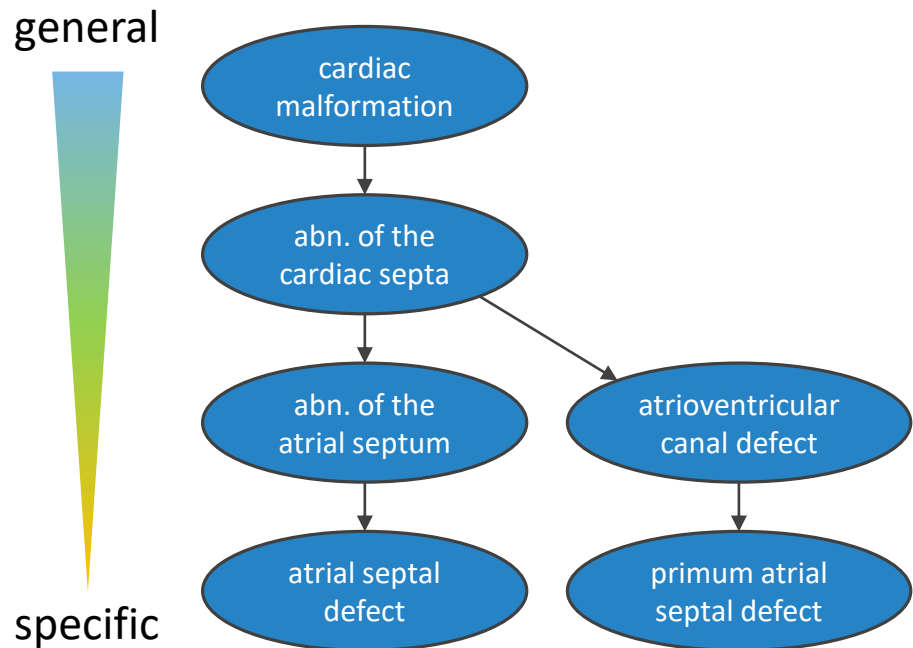
胆嚢炎 脾腫 貧血 球状赤血球症

ヒト 中年 男 縦隔腫瘍 胸部

4 (86.9%) 26073240 Hereditary Spherocytosis in a Middle-aged Man Complicated with Common Bile Duct Stones.

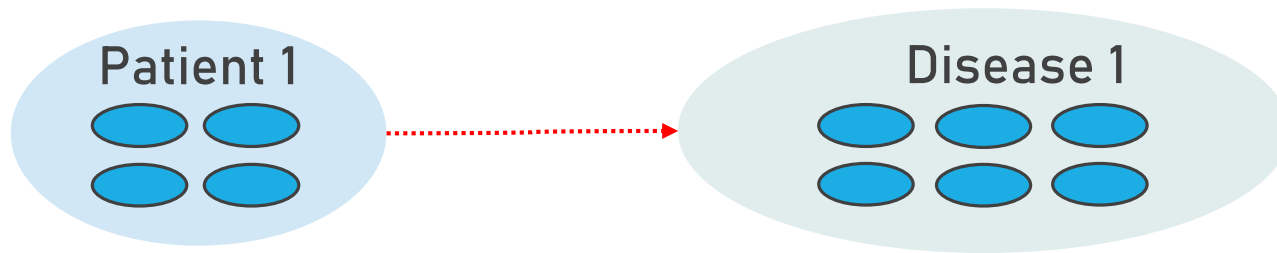
Human Phenotype Ontology (HPO)

- The Human Phenotype Ontology in 2017, Köhler et al., NAR, 2017
 - ✓ HPO is used by many data sets
 - Orphanet, OMIM, ClinVar, MedGen, GARD, ...



Similarity Measures

- Similarity between two sets of phenotypes can be computed in many ways



Measure	Equation	Variations	Reference
Resnik(a,b)	$\max_{t \in g^a \cap g^b} IC(t)$	Avg, Max	Rensik, 1995
Lin(a,b)	$\frac{2 * Resnik(a,b)}{IC(a) + IC(b)}$	Avg, Max	Lin, 1998
Jiang-Conrath(a,b)	$\frac{1}{IC(a) + IC(b) - 2 * Resnik(a,b) + 1}$	Avg, Max	Jiang, 1997
simGIC(P,Q)	$\frac{\sum_{t \in g^P \cap g^Q} IC(t)}{\sum_{t \in g^P \cup g^Q} IC(t)}$		Pesquita, 2007
GeneYenta(P,Q)			Gottlieb, 2015

Gene Yenta

□ Information Content (IC)

$$P(t) = \frac{|annot_t|}{|annot_{all}|}$$

$$IC_t = -\log P(t)$$

□ Similarity between two terms

$$sim_{terms}(t, t') = \max_{a_t \in A_t \cap A_{t'}} IC_{a_t}$$

□ Similarity between a case and a disease

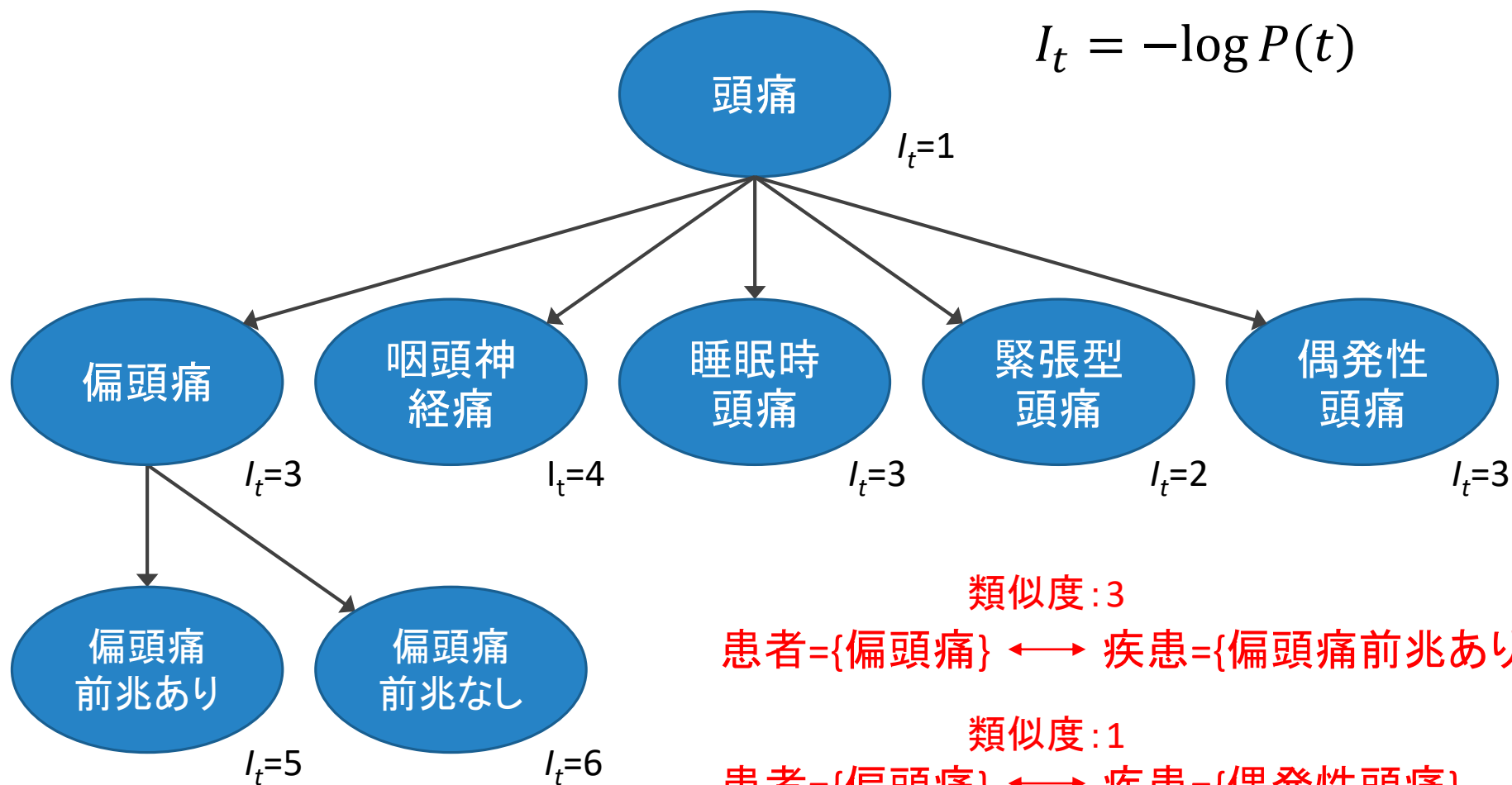
$$sim_{case_disease}(c, d) = \frac{\sum_{t \in T_c} R_t \times \max_{t' \in T_d} sim_{terms}(t, t')}{\sum_{t \in T_c} R_t \times IC_t}$$

症状セットの類似度計算手法

□ 症状間の類似度

$$P(t) = \frac{|annot_t|}{|annot_{all}|}$$

$$I_t = -\log P(t)$$



類似度:3

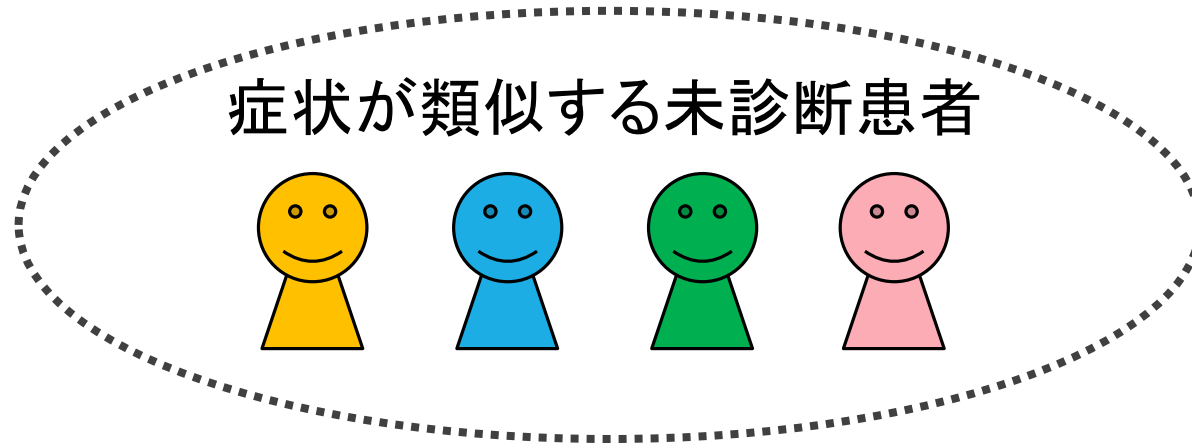
患者={偏頭痛} ↔ 疾患={偏頭痛前兆あり}

類似度:1

患者={偏頭痛} ↔ 疾患={偶発性頭痛}

HPOと類似度計算手法の活用事例

- 希少疾患分野では、HPOと類似度計算手法を用いて、症例データの共有が盛んに行われている
 - 症状が類似する複数の未診断患者が集まることで
 - **新規疾患の定義**
 - **疾患原因遺伝子の同定**



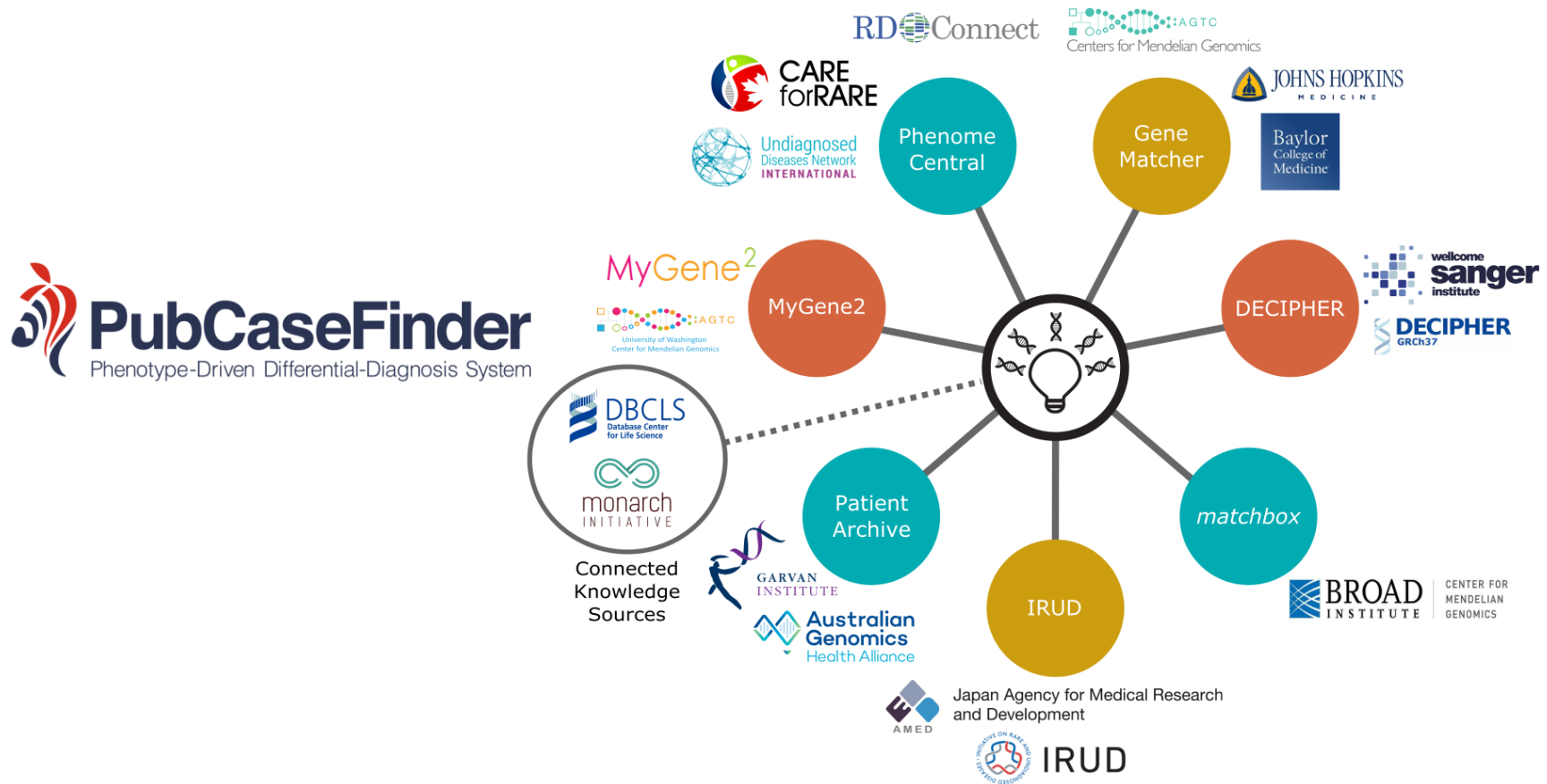
各未診断患者の症状はHPOで表現

Databases which use HPO

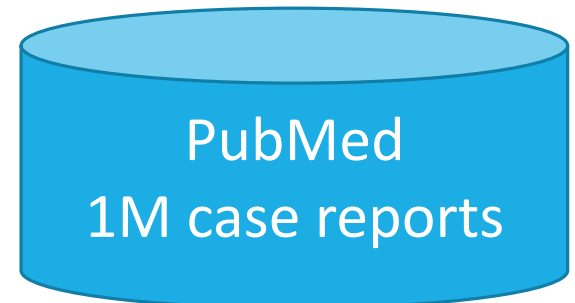
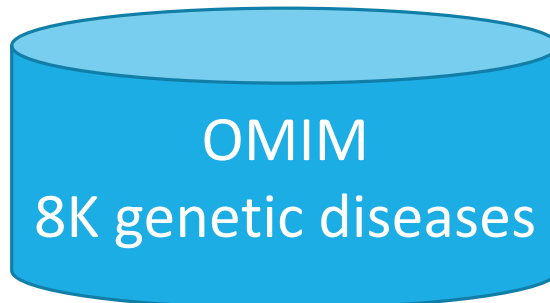
Name	URL
PhenomeCentral	phenomecentral.org
DDD (Deciphering Developmental Disorders)	www.ddduk.org
DECIPHER (DatabasE of genomiC varIation and Phenotype in Humans using Ensembl Resources)	decipher.sanger.ac.uk
ECARUCA (European Cytogeneticists Association Register of Unbalanced Chromosome Aberrations)	http://umcecaruca01.extern.umcn.nl:8080/ecaruca/ecaruca.jsp
The 100 000 Genomes Project	https://www.genomicsengland.co.uk/
Geno2MP (Exome sequencing data linked to phenotypic information from a wide variety of Mendelian gene discovery projects)	http://geno2mp.gs.washington.edu
NIH UDP (Undiagnosed Diseases Program)	available via phenomecentral.org
NIH UDN (Undiagnosed Diseases Network)	available via phenomecentral.org
HDG (Human Disease Gene Website series)	www.humandiseasegenes.com
Phenopolis (An open platform for harmonization and analysis of sequencing and phenotype data)	https://phenopolis.github.io
GenomeConnect (Patient portal developed by ClinGen (67))	www.genomeconnect.org
FORGE Canada & Care4Rare Consortium	available via phenomecentral.org
RD-Connect	platform.rd-connect.eu
Genesis	thegenesisprojectfoundation.org

MatchMaker Exchange Program

□ An official program of GA4GH



Needs for access to relevant resources



Limited covered of DBs

- ❑ Manually curated databases like Orphanet inherently have a limited coverage
 - ✓ For example, only 1/3 of rare diseases in Orphanet are associated with phenotypes

6K Diseases (Orphanet)

Moyamoya syndrome

Tyrosinemia type 2

Truncus arteriosus

...

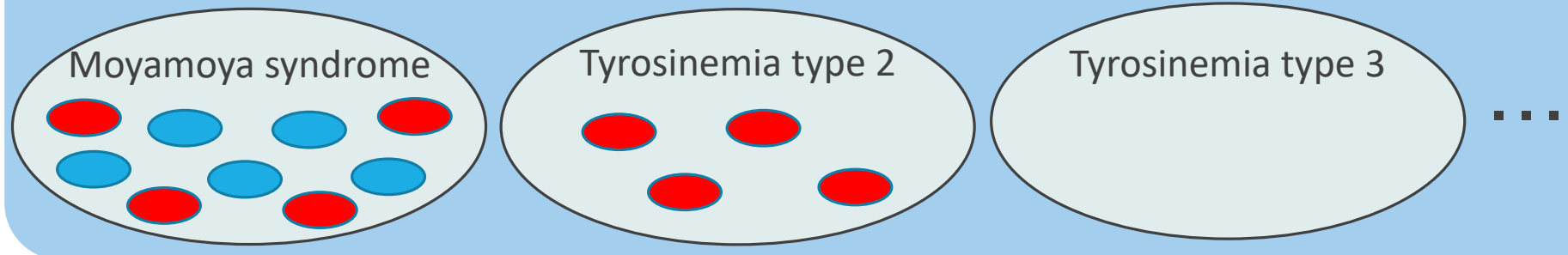
Text Mining

- Text mining-based approach is effective to improve the coverage of disease-phenotype associations
- To automatically extract the associations from more than one million **case reports** in PubMed



As a result, **2/3** of rare diseases in Orphanet could be associated with phenotypes

6,000 Diseases (Orphanet)



PubCaseFinder

- A diagnosis support system for rare diseases

PubCaseFinder Home About API

Query phenotype(s) + Upload File (HPO List):

HP:0004444 Spherocytosis × HP:0001744 Splenomegaly × HP:0001903 Anemia × HP:0000952 Jaundice × HP:0001297 Stroke × HP:0002721 Immunodeficiency ×

Narrow down the diseases + Upload File (Entrez Gene ID List):

ENT:286 ANK1 (SPH1) × ENT:2038 EPB42 (MGC116735 | MGC116737 | PA) × ENT:6521 SLC4A1 (CD233 | FR | RTA1A | SW | WR) × ENT:6708 SPTA1 (EL2) ×

Re-search Clear

Total: 5 20 (per page)

Specify causative disease genes to narrow down candidate diseases

A ranked list of rare diseases based on phenotypic similarity

Similarity	Disease Name	phenotype	Causative Gene
100.0%	Hereditary spherocytosis (ORDO:822)	Hemolytic anemia Immunodeficiency Jaundice Spherocytosis Splenomegaly Stroke	ANK1 EPB42 SLC4A1 SPTA1 SPTB
		ICD-10 (D58.0) OMIM (182900 , 270970 , 612653 , 612690 , 616649)	
81.71%	8p11.2 deletion syndrome (ORDO:251066)	Hemolytic anemia Microcephaly Sacral dimple Spherocytosis Splenomegaly	ANK1

PubCases: A diagnosis assistant tool for rare diseases based on disease-phenotype associations extracted from published case reports.,

T. Fujiwara; Y. Yamamoto; J.D. Kim; T. Takagi, ASHG2017

<https://ep70.eventpilotadmin.com/web/page.php?page=IntHtml&project=ASHG17&id=170121395>

文献から疾患に関連する症状を抽出

□ 課題: 「疾患-症状」関連データの不足が大きな課題

➤ 約100万件の症例報告から、テキストマイニング技術で自動取得



例) 遺伝性球状赤血球症 (Orphanet:822)

- 溶血性貧血 (HP:0001878)
- 黄疸 (HP:0000952)
- 脾腫 (HP:0001744)



症例報告 PMID: 12355853

青: 疾患名 赤: 症状

Hereditary spherocytosis is a genetic, frequently familial hemolytic blood disease characterized by varying degrees of hemolytic anemia, splenomegaly, and jaundice. . . .



Publication of Case Reports



文献から疾患に関連する症状を抽出

□ オントロジーを用いたアノテーションツール

- ConceptMapper (Tanenblatt, 2010)、MetaMap (Aronson, 2010)、NCBO Annotator (Jonquet, 2009)

□ CRAFT Corpus を利用した、パフォーマンス比較 (Christopher, 2014)

- 8つのオントロジーにおいてF-measureを比較結果、7つのオントロジーでConceptMapperのF-measureが最も高かった

□ HPO gold standard (Tudor, 2015)を用いたツール評価 (藤原, JSAI2017)

System	F-measure	Precision	Recall
NCBO Annotator	0.51	0.54	0.47
MetaMap	0.56	0.51	0.61
ConceptMapper	0.52	0.52	0.51

System	Processing time (sec)
NCBO Annotator	206.0
MetaMap	351.0
ConceptMapper	4.3

100万件の症例報告
の処理に要する時間



17.7 (day)



5.2 (hour)

Result of IE from case reports

【Orphanet】

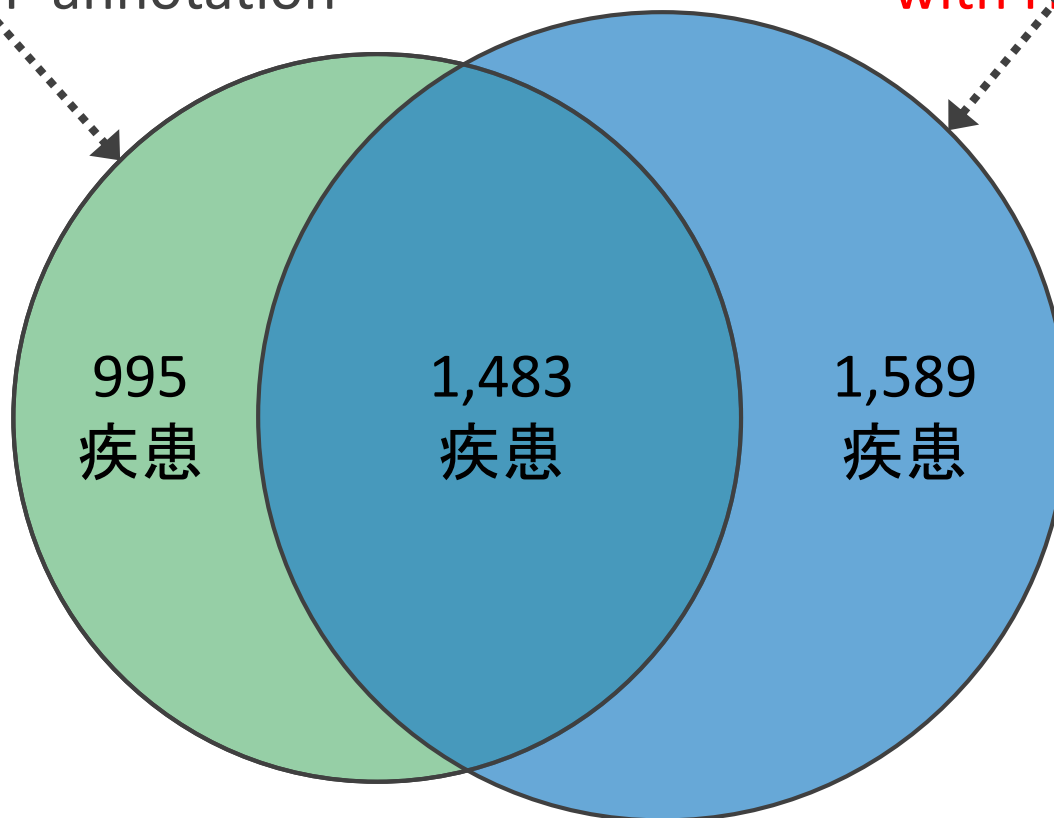
2,478 diseases

with HP annotation

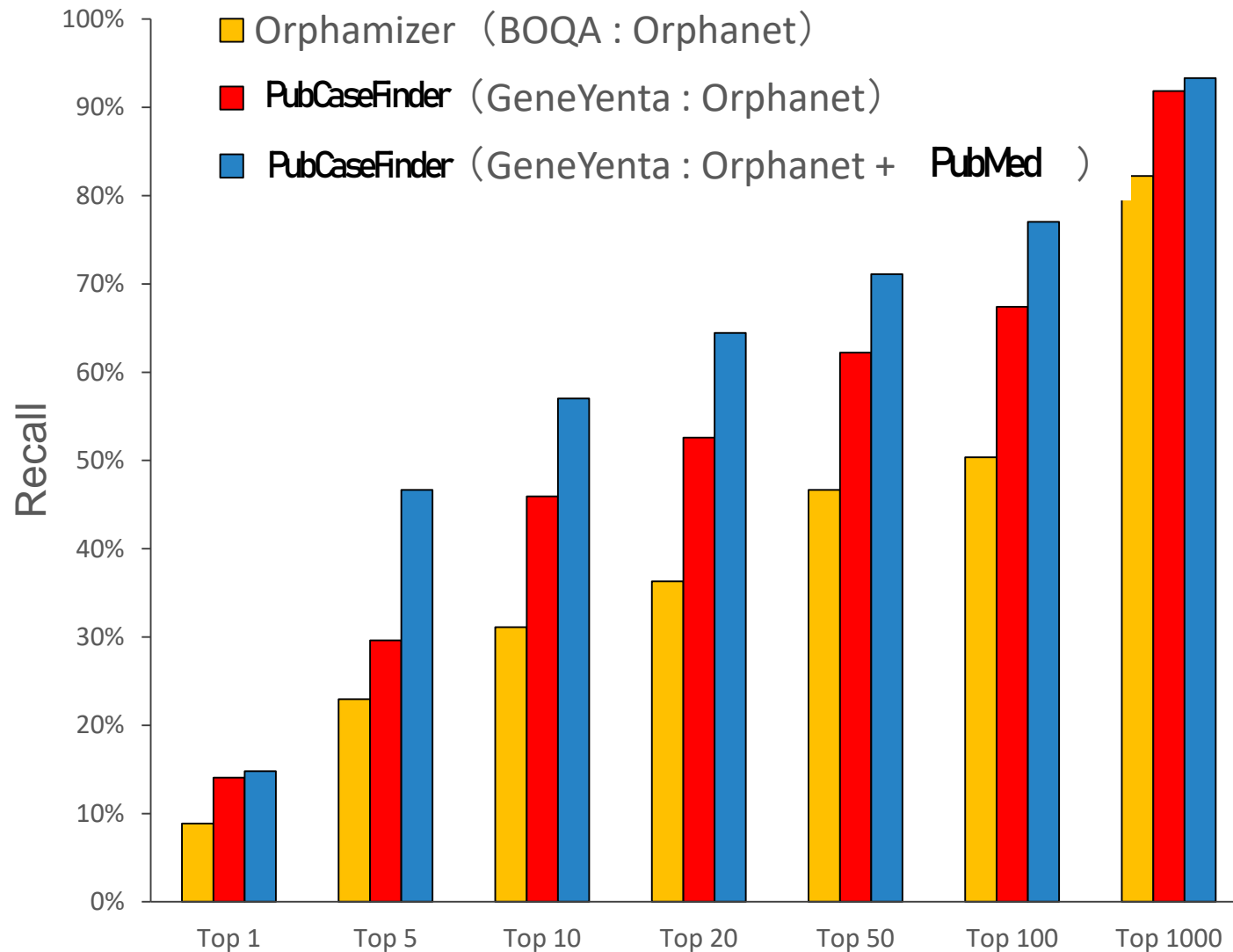
【Case Reports】

3,072 diseases

with HP annotation



Performance Evaluation



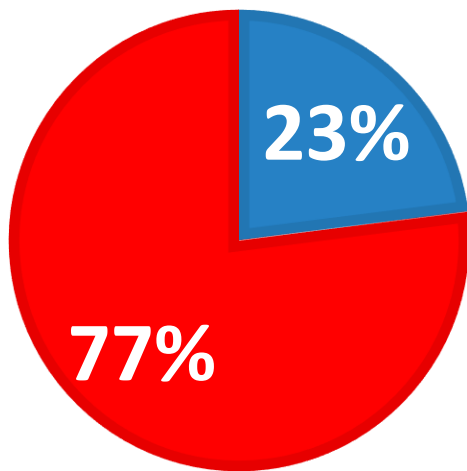
Performance Evaluation

□ For 135 test cases

- ✓ collected from PhenomeCentral
 - maintained by Care4Rare Canada Consortium
- ✓ Recall @ top5 was calculated

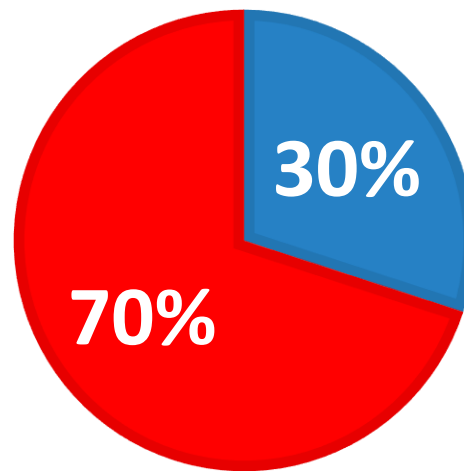
Orphamizer
(Orphanet)

■ 正解症例 ■ 不正解症例



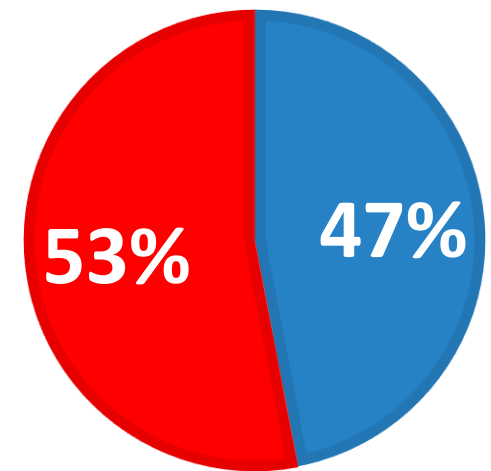
PubCaseFidner
(Orphanet)

■ 正解症例 ■ 不正解症例



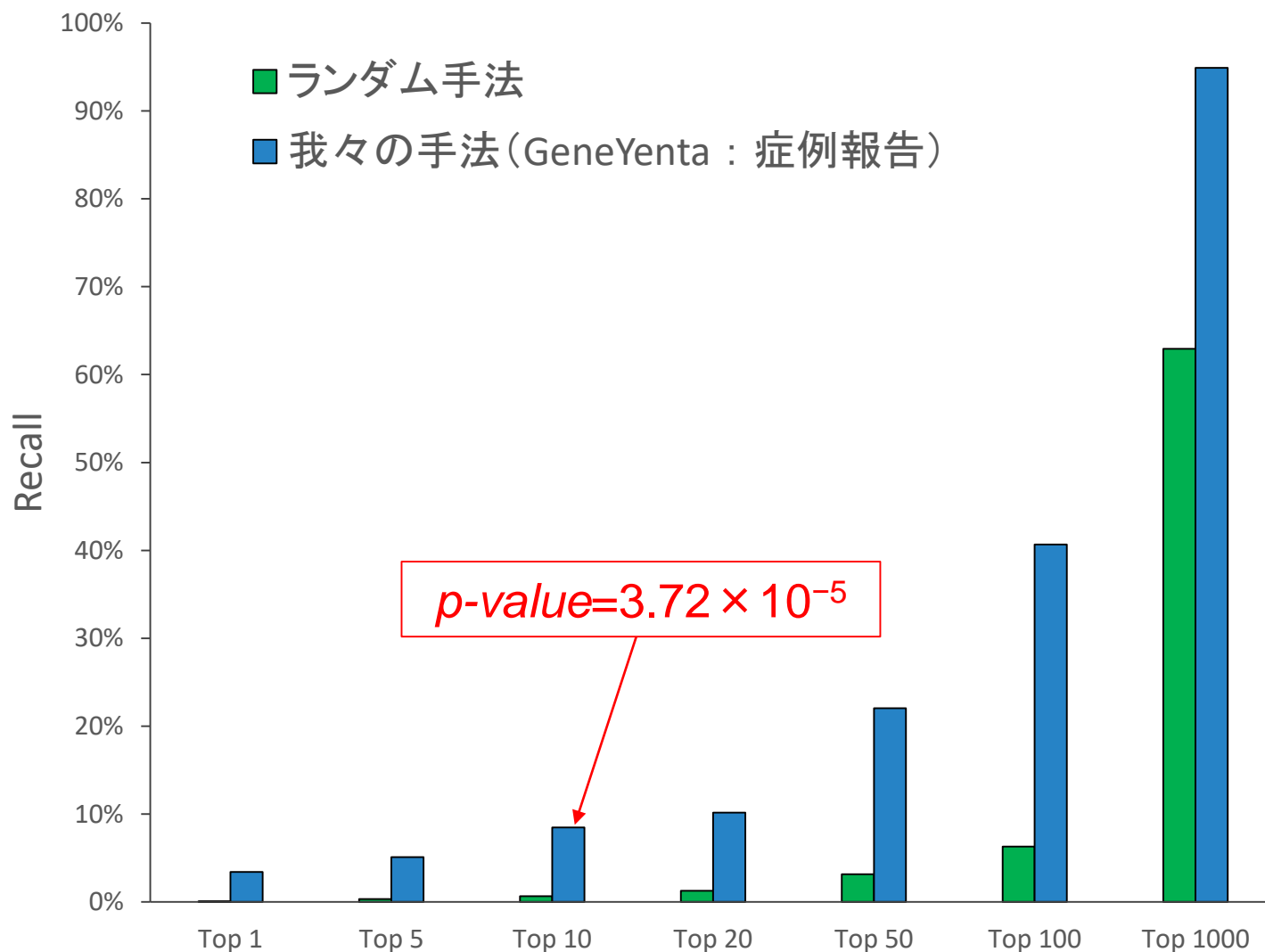
PubCaseFidner
(Orphanet + PubMed)

■ 正解症例 ■ 不正解症例



疾患ランキング精度の評価

- 評価用の59症例を利用して、Recall1,5,10,20,50,100,1000を比較



Conclusion

医療関係者

複数のデータベースを一度に検索

日本語での検索可



高精度な疾患
ランキング機能
を提供

特定の遺伝子
に関連した疾患
を対象に検索可



PubCaseFinder

Phenotype-Driven Differential-Diagnosis System



GA4GHへの
貢献



Matchmaker
Exchange

API

約2万件
の検索
(月間)