

Five years of DisGeNET: Lessons learned and challenges ahead

Laura I. Furlong
IMIM-UPF

Big Data in Biomedicine
Barcelona
November 11, 2014



RESEARCH
PROGRAMME
ON BIOMEDICAL
INFORMATICS



Universitat
Pompeu Fabra
Barcelona



Institut Hospital del Mar
d'Investigacions Mèdiques

DisGeNET

- Knowledge platform on human diseases and their genes
- Aims to cover all disease therapeutic areas
- Developed by integration of data from expert-curated resources and from the literature by text mining
- Centered on gene-disease association (GDA) and its supporting evidence and provenance

Rett Syndrome
UMLS:C0035372

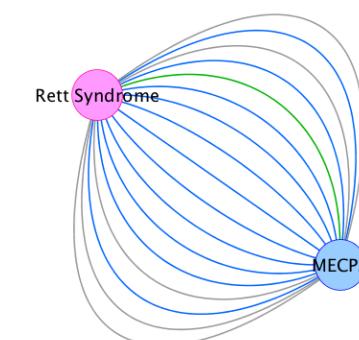


MCEP2
NCBI:4204

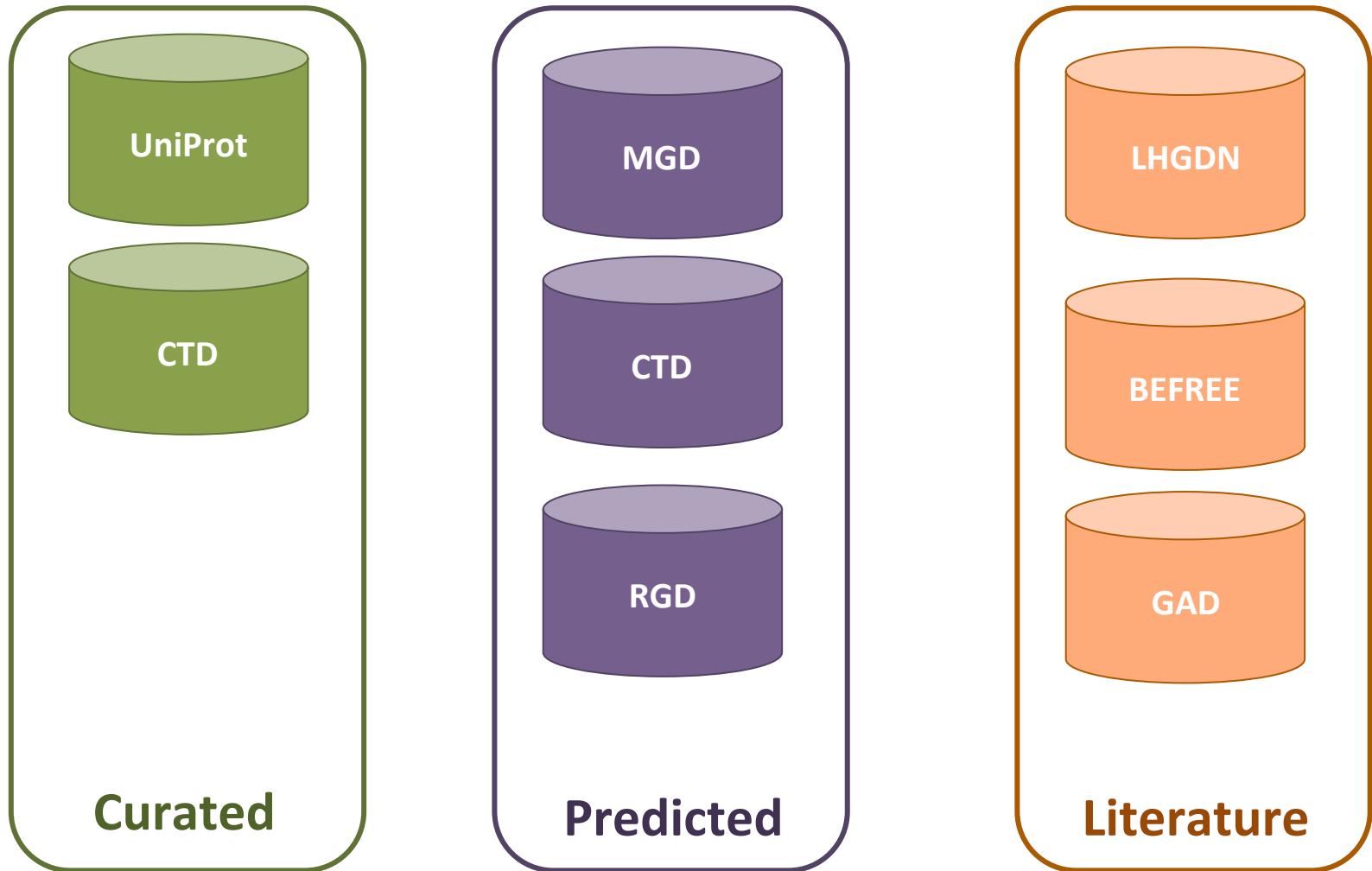
Rett Syndrome
UMLS:C0035372



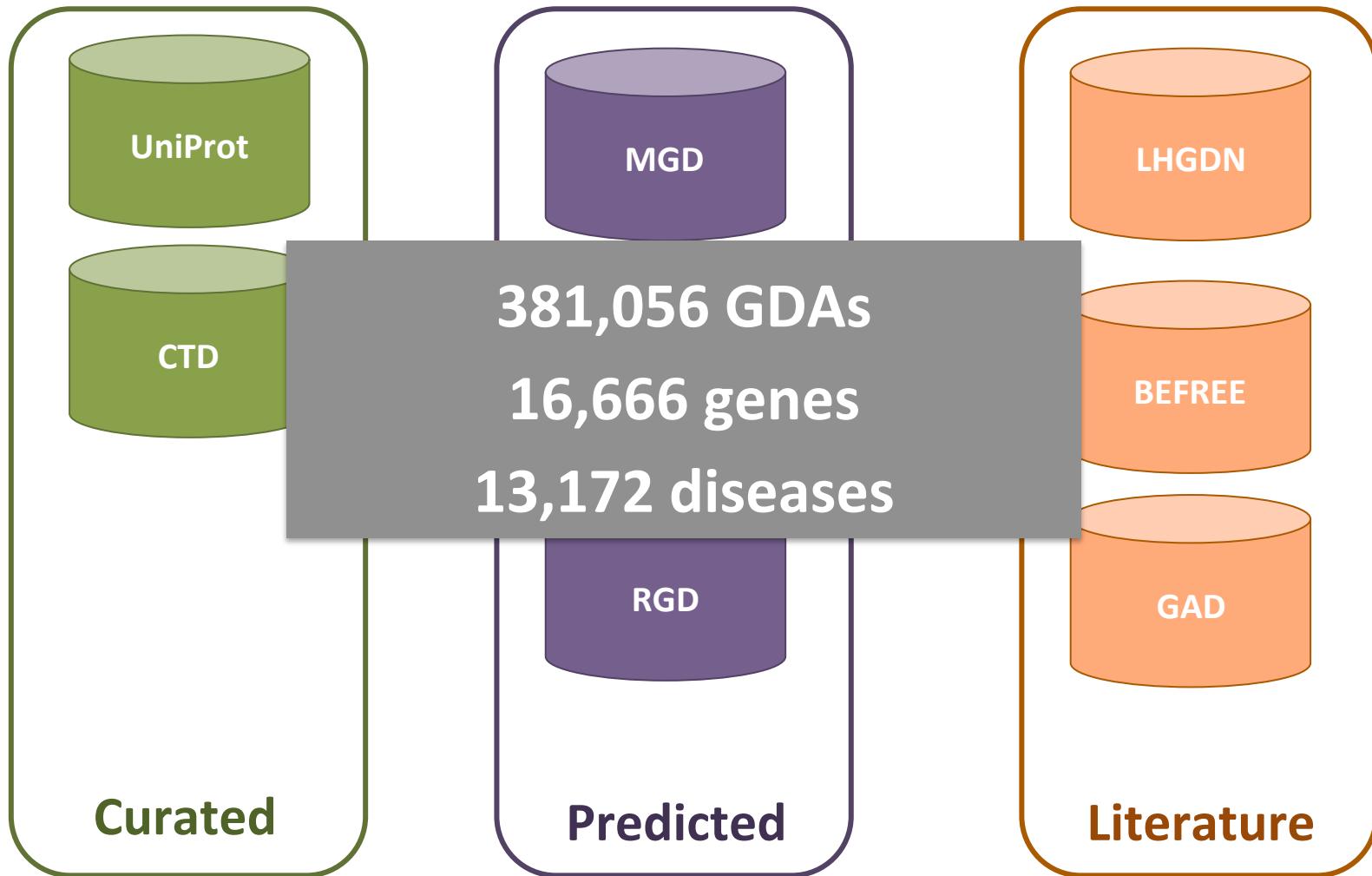
MCEP2
NCBI:4204



DisGeNET

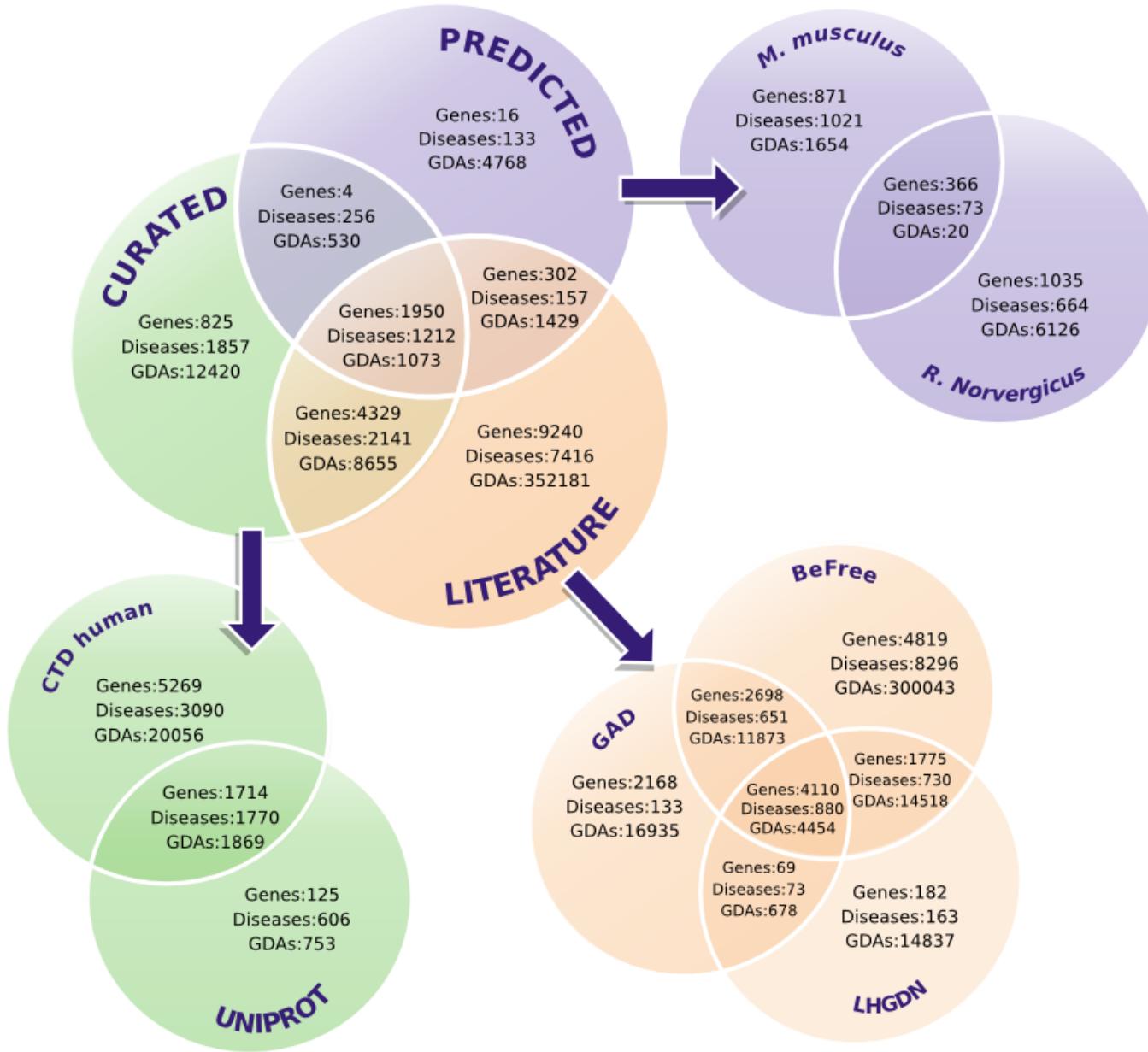


DisGeNET



Key points

- 1. Fragmentation of information as a barrier to knowledge on the mechanisms of human diseases**

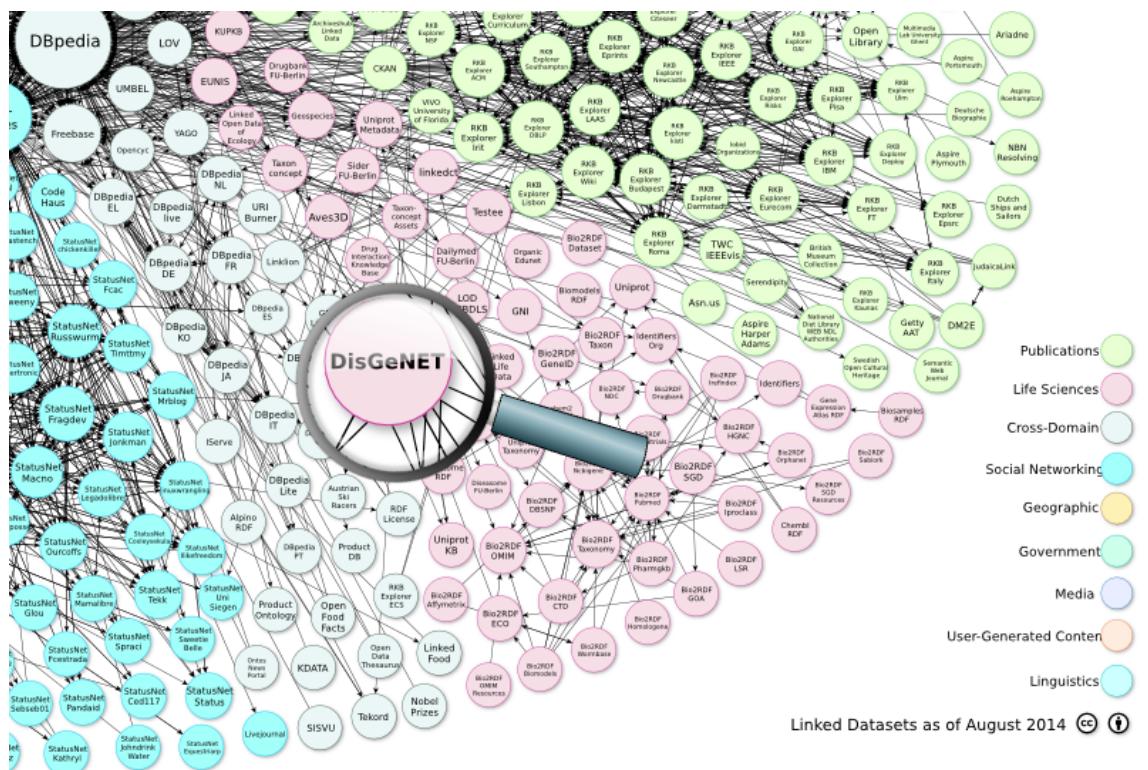


Standardization

- Controlled vocabularies and ontologies
- DisGeNET association type ontology

Open access

- RDF and nanopublications
- Data distributed under the Open database commons license



<http://semanticscience.org>

<http://opendatacommons.org/>

<http://lod-cloud.net/>

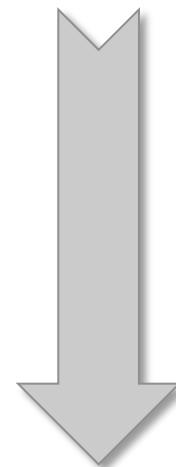
Key points

- 2. High rate of data generation on GDAs poses challenges to biocuration pipelines**

Publications on diseases and genes from 1980
(737,712 publications)



Text Mining



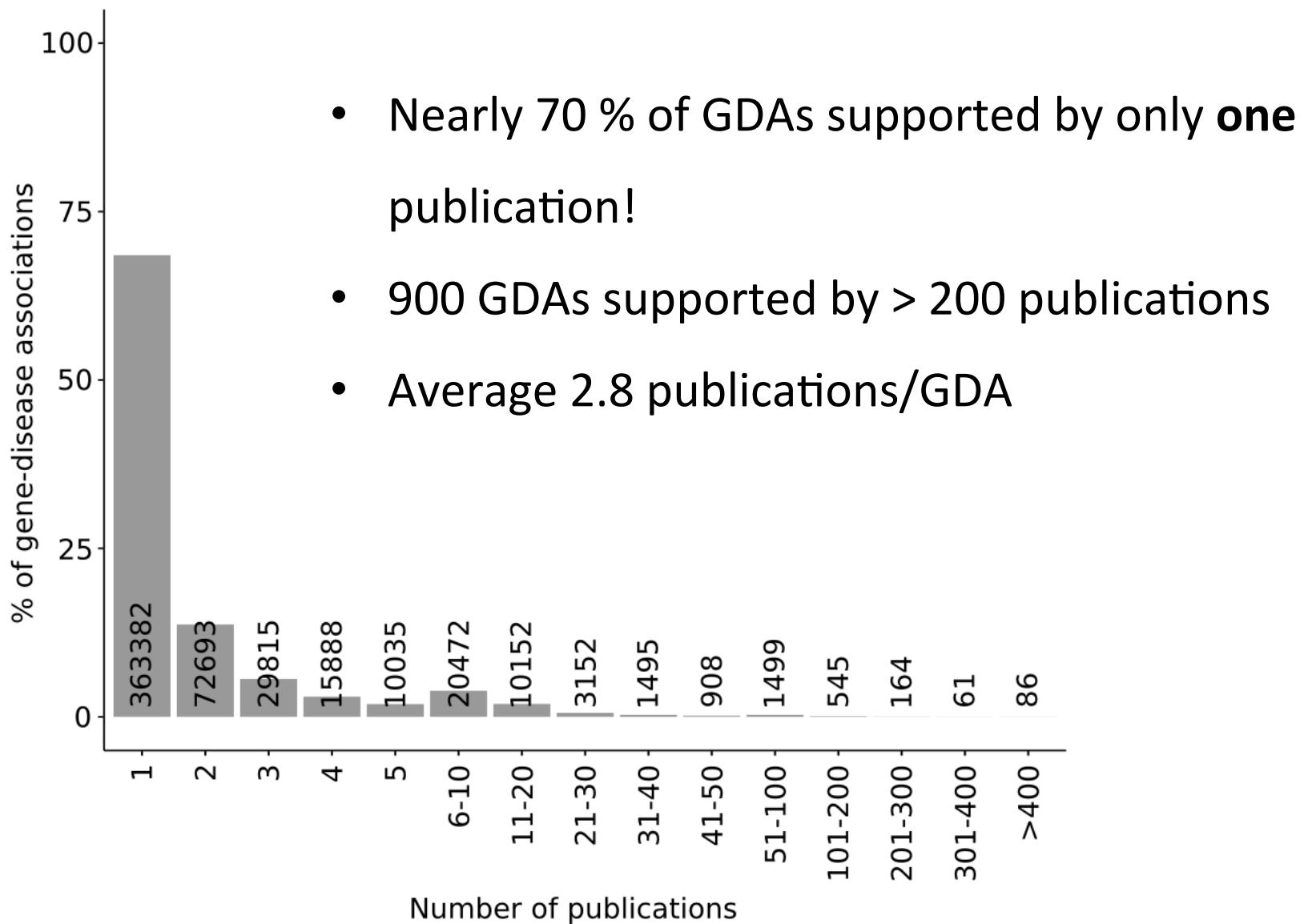
14,777 genes
12,650 diseases
355,976 publications

**530,347
associations**

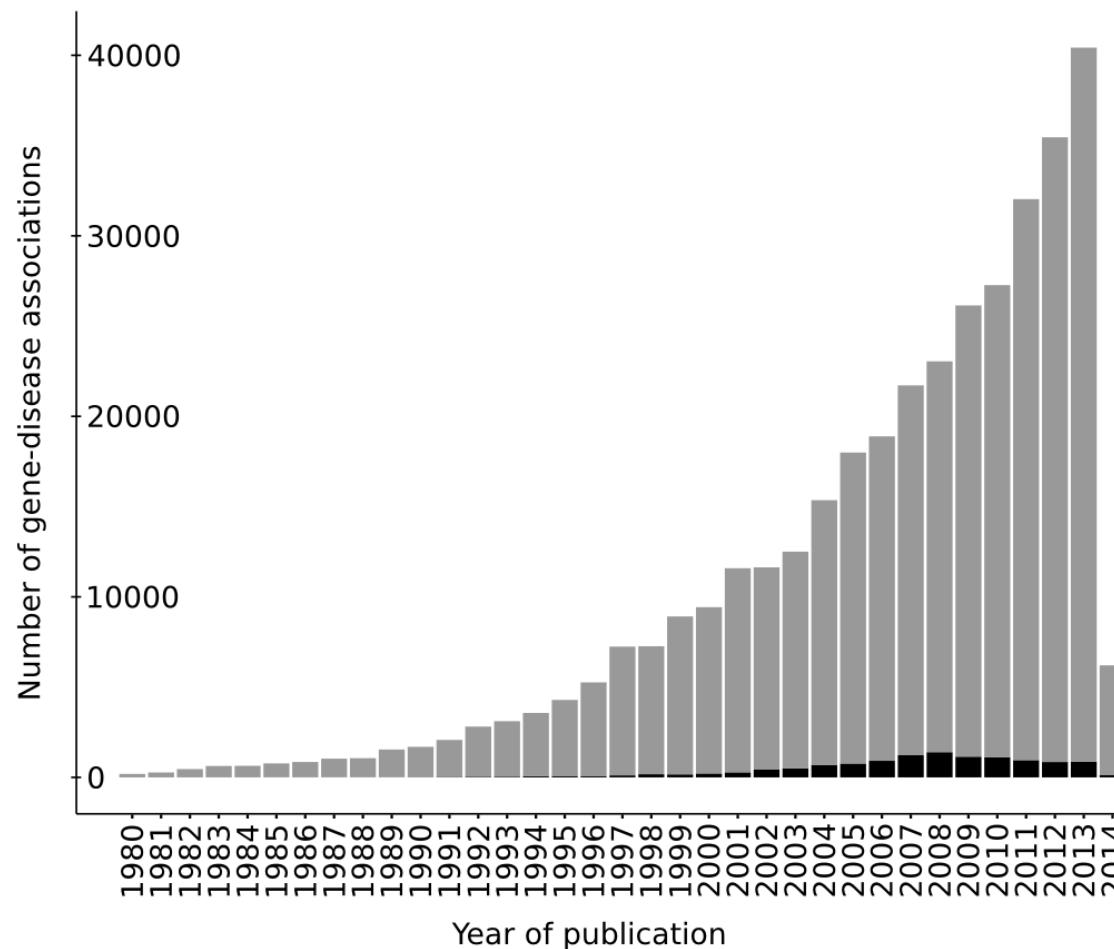
Supervised machine-learning
approach for relation extraction
between genes and diseases



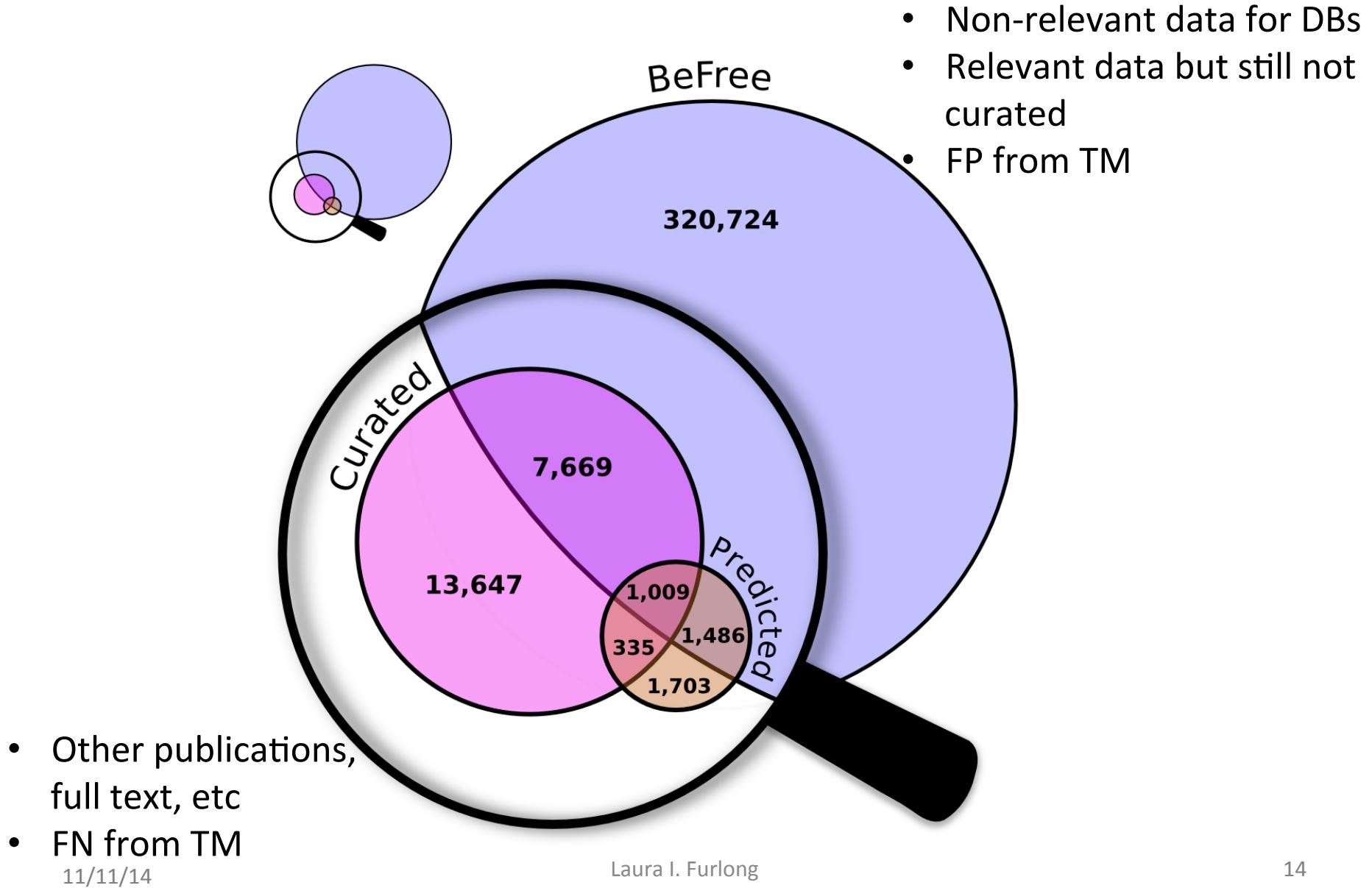
Aprox. 30,000 GDAs



GDAs supported by only one publication



Little overlap between text mined data and data present in curated repositories



- This is not raw data (e.g. from NGS analysis), is data already collated and filtered that have passed peer review before publication
- Text mined from abstracts:
 - Not mining full text, nor supplementary material
 - Not mining tables, figures
- Focus on relations stated on sentences, not handling anaphoras
- We are just looking into a *small fraction* of all the available data!!!

Key points

2. High rate of data generation on GDAs poses challenges to biocuration pipelines

- Need to find alternative strategies to expert curation
- “wisdom of the crowds” approaches (e.g. crowdsourcing)
- ?

Key points

- 3. Data prioritization to support interpretation of data on the genetic determinants of human diseases**

GDAs prioritization

Ranks gene-disease associations based on the supporting **evidence**

$$\text{DisGeNET score} = S_{\text{CURATED}} + S_{\text{PREDICTED}} + S_{\text{LITERATURE}}$$

$$S_{\text{CURATED}} = W_{\text{UNIPROT}} + W_{\text{CTD}}$$

$$S_{\text{PREDICTED}} = W_{\text{Rat}} + W_{\text{Mouse}}$$

$$S_{\text{LITERATURE}} = W_{\text{GAD}} + W_{\text{LHGDN}} + W_{\text{BeFree}}$$

Disease	Gene	Score	UniProt	CTD	Rat	Mouse	Number of publications		
							BeFree	GAD	LHGDN
Wilson's Disease	ATP7B	0.99	0.3	0.3	0.1	0.1	174	31	23
Rett Syndrome	MECP2	0.9	0.3	0.3	0	0.1	438	27	43
Cystic Fibrosis	CFTR	0.9	0.3	0.3	0	0.1	1429	150	78
Obesity	MC4R	0.94	0.3	0.3	0.1	0.1	220	46	0
Alzheimer Disease	APP	0.88	0.3	0.3	0	0.1	1096	18	81

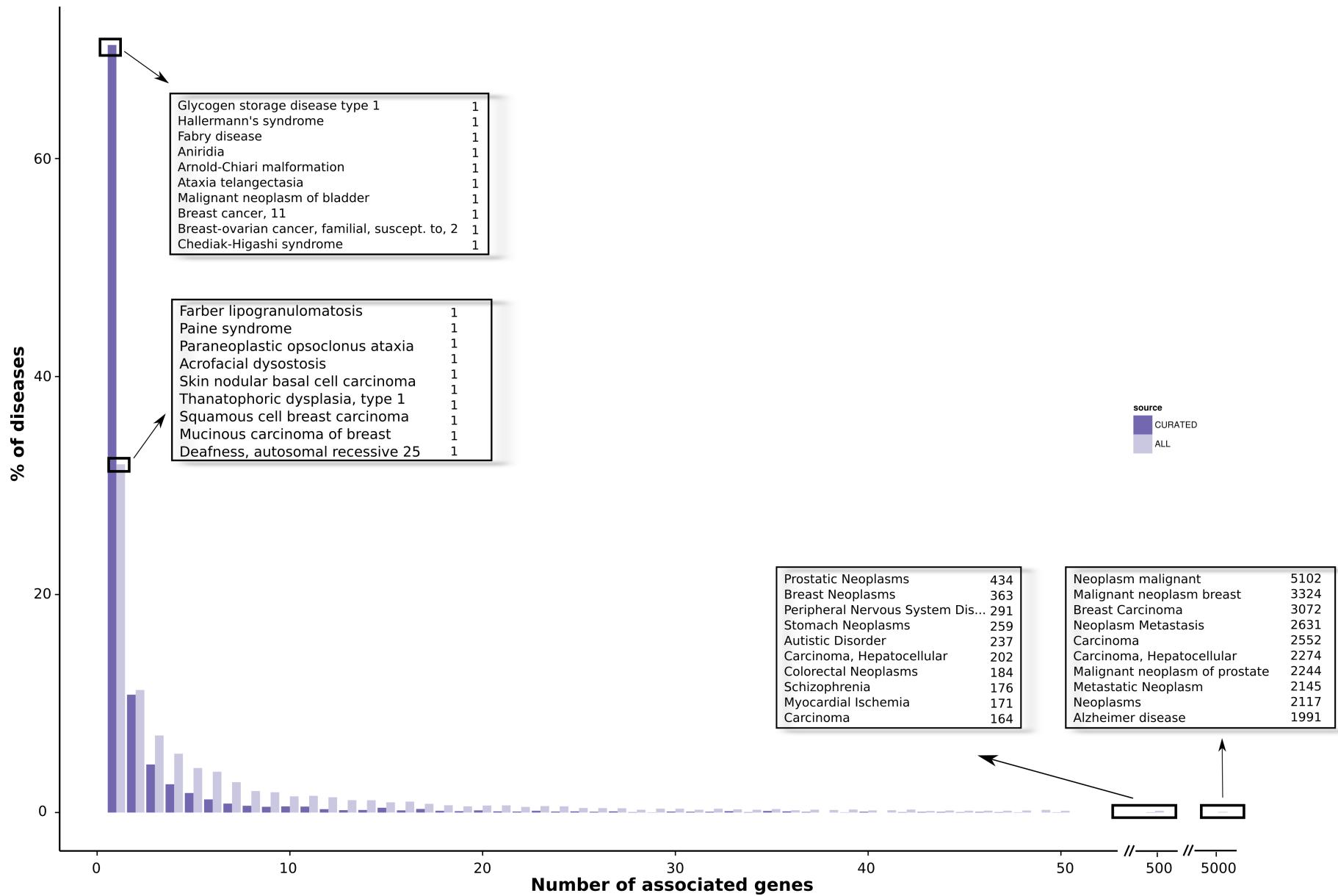
Key points

3. Data prioritization to support interpretation of data on the genetic determinants of human diseases

- Still we have a large dataset with low score (350,000 GDAs)
- Other approaches based on:
 - type of association of the gene to the disease
 - experimental evidence for the GDAs
 - network-based gene-prioritization algorithms
- Take into account contradictory findings
- Different prioritization approaches for different purposes?

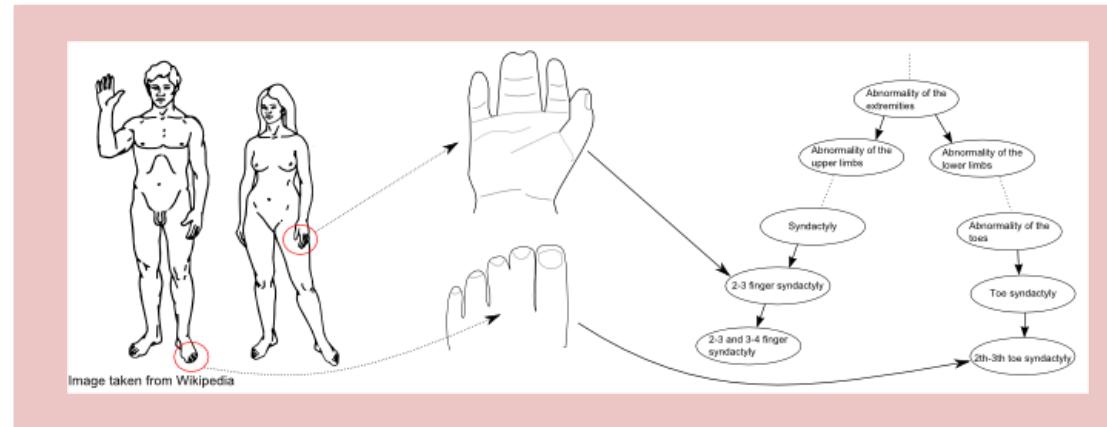
Key points

- 4. Large number of genes associated to some diseases might reflect phenotypic diversity**



Need for deep *phenotyping* of diseases to identify disease subtypes associated to different gene networks

Human Phenotype Ontology (HPO)



<http://www.human-phenotype-ontology.org/>

- Ontology of phenotypic abnormalities of human diseases
- Based on OMIM, Orphanet and DECIPHER

DisGeNET-HPO annotations

33 % of DisGeNET diseases annotated with HPO terms

Mainly come from OMIM diseases (28 % of DisGeNET diseases come from OMIM)

Need for other approaches to annotate the full spectrum of diseases at phenotypic level

Key points

1. Fragmentation of information as a barrier to knowledge on the mechanisms of human diseases
2. High rate of data generation on GDAs poses challenges to biocuration pipelines
3. Data prioritization to support interpretation of data on the genetic determinants of human diseases
4. Large number of genes associated to some diseases might reflect phenotypic diversity

IBI Group

Alba Gutiérrez
Àlex Bravo
Janet Piñero
Núria Queralt-Rosiñach
Miguel A. Mayer
Pablo Carbonell
Laura I. Furlong
Ferran Sanz



Unión Europea

Fondo Europeo
de Desarrollo Regional



RESEARCH
PROGRAMME
ON BIOMEDICAL
INFORMATICS

11/11/

IBI Past members

Montserrat Cases
Michael Rautschka
Anna Bauer-Mehren
Solène Grosdidier



Universitat
Pompeu Fabra
Barcelona



Institut Hospital del Mar
d'Investigacions Mèdiques



26