



Initiatives to Speed up Data Mining

Ségolène Aymé

ICORD 2014, Ede, Netherlands

9 October 2014

Objectives

➔ Make the most of remarkable advances in the molecular basis of human diseases

↳ dissect the physiological pathways

↳ to improve diagnosis

↳ to develop treatments

➔ Make rare diseases visible in health information systems

↳ to gain insight into them

↳ to access real life data already collected



IRDiRC

INTERNATIONAL
RARE DISEASES RESEARCH
CONSORTIUM

scaphocephaly

Medical History

Presenting Problem: [Handwritten text]

Review of Systems: [Handwritten text]

Past Medical History: [Handwritten text]

Physical Examination: [Handwritten text]

Diagnosis: [Handwritten text]

Treatment: [Handwritten text]

Prognosis: [Handwritten text]

Signature: [Handwritten signature]

Date: [Handwritten date]



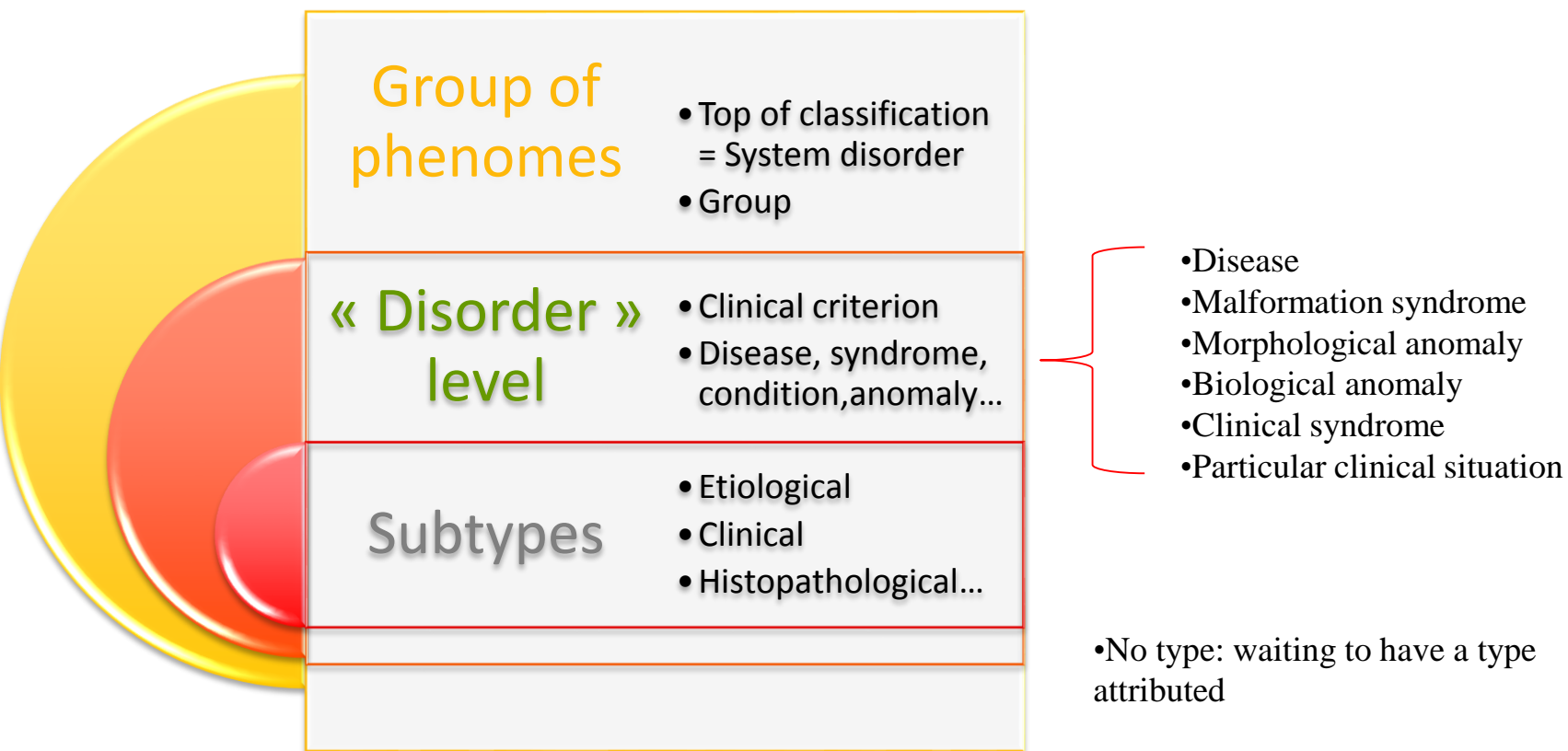
dolichocephaly

[illegible]

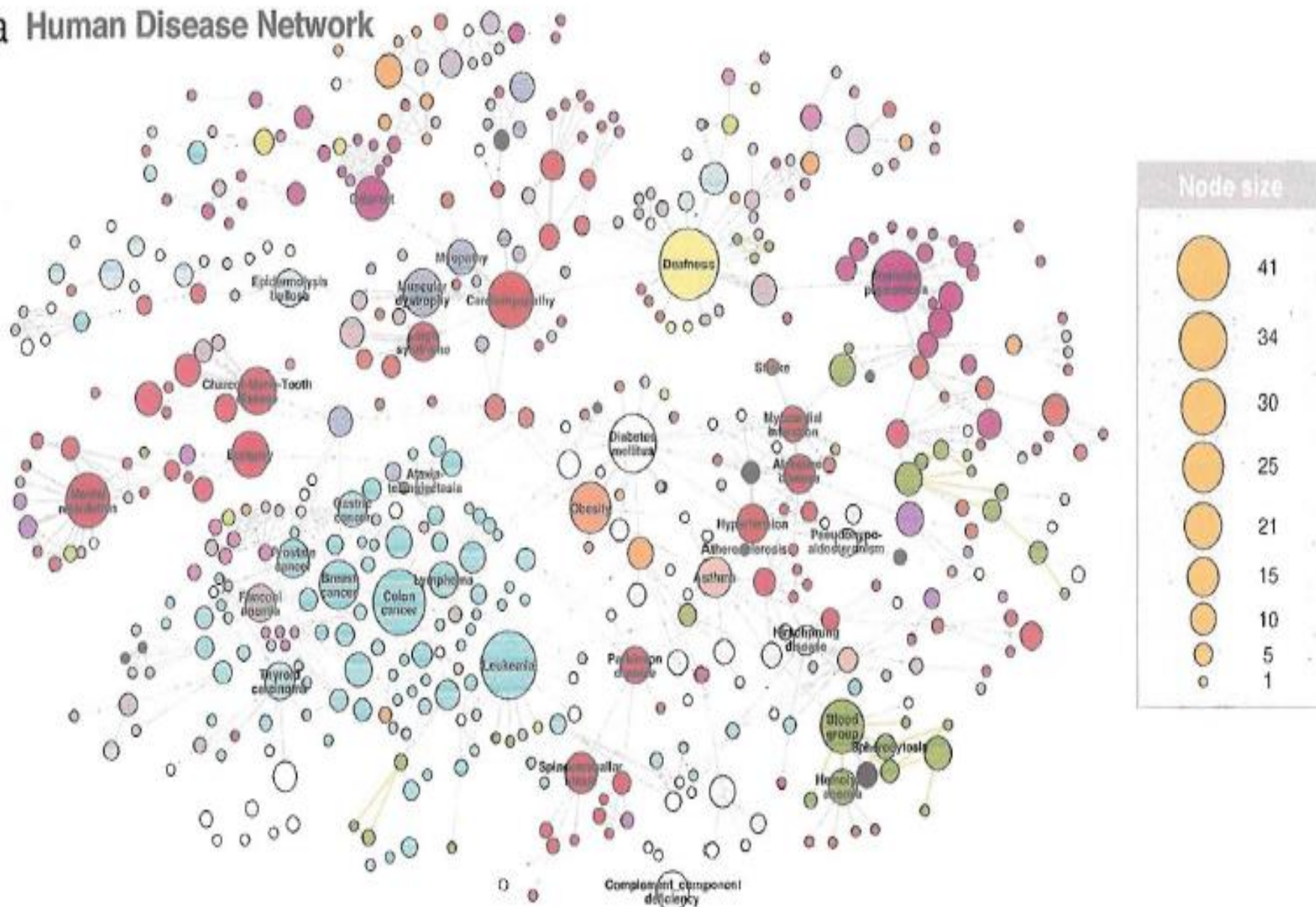
What is the problem ? Computers are not smart enough....

- ▶ Phenotypic descriptions that are very evocative for humans: “myopathic electromyography » or « still walking 25 years after onset”
 - ↪ Unfortunately computers, being significantly dumber than humans, don't quite get it...
- ▶ The following descriptions mean the same thing to you: “generalized amyotrophy”, “generalized muscle atrophy”, “muscular atrophy, generalized” (etc)
 - ↪ But your computer thinks they're completely unrelated

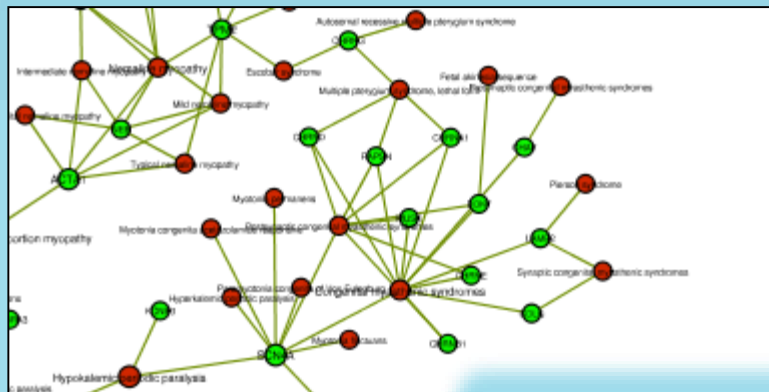
Phenomes: a continuum



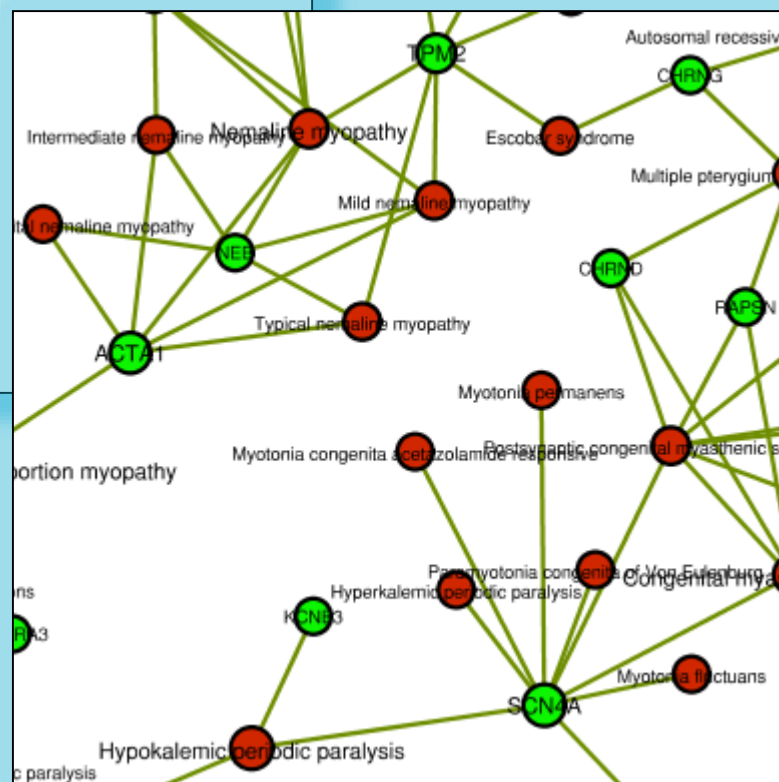
a Human Disease Network



Orphan Diseasome



An Orphan Diseasome permits investigators to explore the orphan disease (OD) or rare disease relationships based on shared genes and shared enriched features (e.g., Gene Ontology Biological Process, Cellular Component, Pathways, Mammalian Phenotype).



The red nodes represent the orphan diseases and the green ones the related genes. A disease is connected to a gene if and only if a mutation which is responsible of the disease has been identified on this gene.

UMLS = Unified Medical Language System

▶ ICD = International Classification of Diseases

↪ Since 1863 by WHO

↪ Used by most countries to code medical activity, mortality data

▶ MeSH = Medical Subject Headings

↪ controlled vocabulary thesaurus used for indexing articles for PubMed by National Library of Medicine (USA)

▶ SnoMed CT = Systematized Nomenclature of Medicine--Clinical Terms

↪ clinical terminology by the International Health Terminology Standards Development Organisation (IHTSDO) in Denmark

↪ Used in the USA and a few other countries

▶ MedDRA = Medical Dictionary for Regulatory Activities

↪ medical terminology to classify adverse event information associated with the use of medical products

↪ by the International Federation of Pharmaceutical Manufacturers and Associations (IFPMA)

Medical History

Examiner: _____ Date: _____

Presenting Problem: _____

History of Present Illness: _____

Review of Systems: _____

Physical Examination: _____

Diagnosis: _____

Treatment: _____

Prognosis: _____

Signature: _____ Date: _____

Free text

[illegible]

Others? Free text?

Menu

Support the Phenomizer, Help

Features

Diseases

Ontology

Dolichocephaly

HP0 is

Feature

HP:0005208

Dolichocephaly

Patient's Features:

HP0

Feature

Modifier

id category: Abnormality of head and neck (1 item)

HP:0005208

Dolichocephaly

observed

id category: Abnormality of the skeletal system (1 item)

HP:0005208

Dolichocephaly

observed

SIMPLE SEARCH

Step 1: Define the rule

Step 2: Select the clinical signs or access the thesaurus

Mandatory

Dolichocephaly scaphocephaly

Delete

☐ Mandatory
 ☐ Optional

Search Thesaurus

☐ Mandatory
 ☐ Optional

Search Thesaurus

☐ Mandatory
 ☐ Optional

Search Thesaurus

☐ Mandatory
 ☐ Optional

Search Thesaurus

Step 3: Search for matching diseases

OK

Orphanet

Each terminology has a purpose– driven approach

- ▶ Indexing health status of individual patients for health management (SnoMED)
 - ↳ Detailed, focus on manifestations and complaints
 - ↳ Adapted to clinical habits
 - ↳ Analytical approach
- ▶ Indexing health status of individual patients for statistical purpose in public health (ICD)
 - ↳ More aggregated, interpreted phenotypic features
 - ↳ Aggregated concepts
 - ↳ Unambiguous to avoid blanks

Purpose-driven approach (2)

- ▶ Indexing health status of individual patients for clinical **research purpose** (HPO / PhenoDB / Elements of morphology)
 - ↪ Highly detailed to fit with the research questions
 - ↪ Specific terminologies developed for disease-specific patient registries
- ▶ Indexing health status of individual patients for retrieving **possible diagnoses** (LDDb,POSSUM,Orphanet)
 - ↪ Agregated concepts
 - ↪ Requires a judgement of clinicians about phenomic expressions that are relevant
 - ↪ Unambiguous to avoid blanks

HOW TO MAKE ALL THESE TERMINOLOGIES INTER-OPERABLE ?

Overview of project progress

- ▶ Sept 2012: start of mappings (Orphanet)
- ▶ EUGT2 – EUCERD workshop (Paris, September 2012)



PhenoDB



HPO



Orphanet

LDDB
Elements of Morphology
POSSUM
SNOMED CT (IHTSDO)

DECIPHER
IRDIRC

Ground Work

- ▶ Orphanet group aligned terms from HPO, PhenoDB, Orphanet terminology, MESH, MedDRA, UMLS, LDDDB, SNOMed-CT, and Elements of Morphology
- ▶ PhenoDB group compared data on use of terms in 3 large scale projects: DDD/DECIPHER, FORGE/Canada using HPO and BHCMG using PhenoDB. Mapped all terms.

Phenotype terminology project

► Aims:

↪ **Map** commonly used clinical terminologies (Orphanet, LDDb, HPO, Elements of morphology, PhenoDB, UMLS, SNOMED-CT, MESH, MedDRA):

- automatic map, expert validation, detection and correction of inconsistencies

↪ **Find common terms** in the terminologies

↪ Produce a **core terminology**

- Common denominator allowing to share/exchange phenotypic data between databases

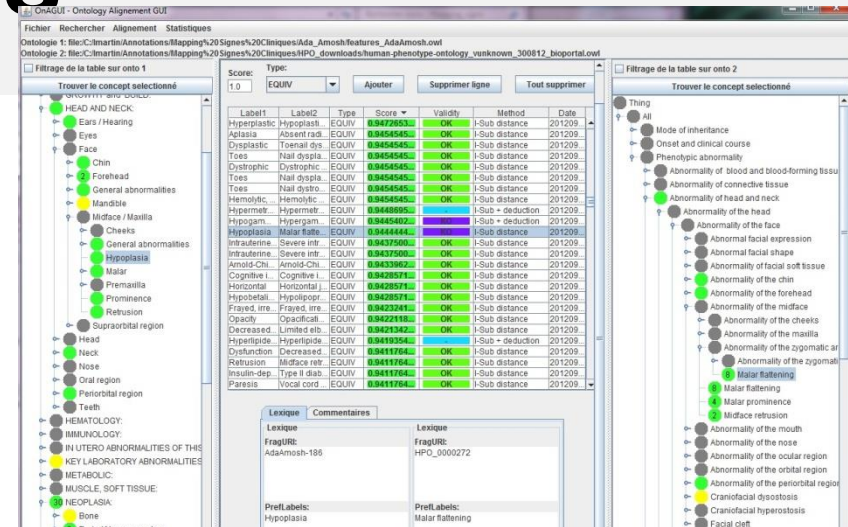


Mapping Terminologies

- ▶ **Orphanet**: 1357 terms (Orphanet database, version 2008)
- ▶ **LDDb**: 1348 dysmorphological terms (Installation CD)
- ▶ **Elements of Morphology**: 423 terms (retrieved manually from publication AJMG, January 2009)
- ▶ **HPO**: 9895 terms (download bioportal, obo format, 30/08/12)
- ▶ **PhenoDB**: 2846 terms (given in obo format, 02/05/2012)
- ▶ **UMLS**: (version 2012AA) (integrating **MeSH**, **MedDra**, **SNOMED**

Tools

- ▶ OnaGUI (INSERM U729): ontology alignment tool
 - Work with file in owl format
 - I-Sub algorithm: detect syntactic similarity
 - Graphical interface to check automatic mappings and manually add ones



- ▶ Metamap (National Library of Medicine): a tool to map biomedical text to the UMLS Metathesaurus

```

Phrase: "total anosmia"
>>>> Phrase
total anosmia
<<<< Phrase
>>>> Candidates
Meta Candidates (1):
      833 C0003126:Anosmia {MDR, MSH, SNOMEDCT} [Sign or Symptom]
<<<< Candidates
  
```

- ▶ Perl scripts: format conversion, launching Metamap, co results...

	A	C	D	E	F	G	H	I	J	K	L	N
	Sign_Id	ORPHA PT	SCORE	CORR	ORPHA SYN	UMLS MATCH	UMLS PT	UMLS CODE	TERMINOLOG	CODE	PT/corr CONCEPT	SEMANTIC TYPE
	10.01.05	Facial cleft	1000		facial cleft	Facial cleft	Cleft face	C0685787	UMLS	C0685787	Cleft face	Congenital Abn
	10.01.05	Facial cleft	1000		facial cleft	Facial cleft	Cleft face	C0685787	SNOMEDCT	92821006	Cleft face (disorder)	Congenital Abn
	11.02.00	Lower lip, ge	770		general abnc Lip	Benign neopla	C0153932	C0153932	UMLS	C0153932	Benign neoplasm of t	Neoplastic Pro
	11.02.00	Lower lip, ge	770		general abnc Lip	Benign neopla	C0153932	C0153932	MDR	10024556	Lip neoplasm benign	Neoplastic Pro
	11.02.00	Lower lip, ge	770		general abnc Lip	Benign neopla	C0153932	C0153932	SNOMEDCT	92185002	Benign neoplasm of li	Neoplastic Pro
	11.02.00	Lower lip, ge	770		general abnc LIP	Lymphoid inte	C0264511	C0264511	UMLS	C0264511	Lymphoid interstitial	Disease or Syn
	11.02.00	Lower lip, ge	770		general abnc LIP	Lymphoid inte	C0264511	C0264511	MDR	10062997	Lymphoid interstitial	Disease or Syn

IRDIRC

INTERNATIONAL
RARE DISEASES RESEARCH
CONSORTIUM

Comparison of mappings and deduction

- ▶ Perl script to compare all the mappings and infer mappings of non-Orphanet terminologies

Eg: Orphanet ID XX mapped to YY in HPO and ZZ in LDDb -> deduction: YY and ZZ should probably map

- ▶ Retrieve HPO mappings versus UMLS, MeSH

	LDDb	El. Morpho	PhenoDB	HPO	UMLS...
Orphanet	E: 1062	E: 416	E: 978	E: 2228	E: 6948
LDDb		D: 275	D: 533	D: 1123	D:2678
El. Morpho			D: 177	D: 716	D: 409
PhenoDB				D: 1045	D:3268
HPO					D: 6307+4800
UMLS...					

Mapping of non-Orphanet terminologies

- ▶ Automatic and inferred mappings were checked by experts

↳ Using OnaGUI for all, except UMLS

✓ Automatic I-Sub: 7.0 + deduction

	El. Morpho	PhenoDB	HPO	UMLS...
LDDDB	D: 257 +23 added	D:528, 92%E A:674, 38%E	D: 1105, 87%E A: 2084, 23%E	D: 2654, 83%E A: 11731
El. Morpho		D:174, 50%E A:189, 74%E	D:393, 93%E A: 436, 16%E	D:405, 84%E A:1248
PhenoDB			D:1018, 91%E A: 4168, 6%E	D: 3222, 82%E A: 18776
HPO				D: 7389 A: 65535
UMLS...				

First list of common terms

- ▶ Present in at least 3 terminologies
- ▶ Definition of rules for nomenclature
- ▶ Addition of terms present in each terminology as synonyms

1	Identifier	Preferred term	Synonyms		
2	T0001	Anomaly of the skull	Skull anomalies	Cranial bones, general abnormalities	Cranium, general
3	T0003	Cranial hyperostosis	Thick skull	Thickened skull	Dense skull
4	T0004	Basilar hyperostosis	Sclerosis of skull base		
5	T0005	Calvarial hyperostosis	Thick calvaria	Thickened calvaria	Dense calvaria
6	T0006	Decreased skull ossification	Poorly ossified skull	Ossification defects of skull	Undermineraliz
7	T0008	Decreased calvarial ossification	Thin calvaria	Absent ossification of calvaria	Thin calvarium
8	T0010	Anomaly of the cranial sutures	Cranial sutures, general abnormalit	Head Sutures anomalies	Abnormality of
9	T0011	Wide cranial sutures	Cranial sutures, wide	Wide cranial sutures (finding)	
10	T0012	Ridged cranial sutures	Cranial sutures, ridged		
11	T0013	Anomaly of the sella turcica	Sella turcica anomaly		
12	T0014	Large sella turcica	Sella turcica, large		
13	T0015	J-shaped sella turcica	Sella turcica, J-shaped	Shoe-shaped sella turcica	
14	T0016	Small sella turcica	Sella turcica, small		
15	T0018	Anomaly of the temporal bone	Abnormality of the temporal bone		
16	T0019	Anomaly of the mastoid process	Mastoids, general abnormalities	Abnormality of the mastoid	
17	T0020	Small foramen magnum	Foramen magnum, small	Foramen magnum stenosis	
18	T0021	Large foramen magnum	Foramen magnum, large		
19	T0022	Delayed pneumatization of the	Delayed pneumatization of mastoids		
20	T0023	Advanced pneumatization of the	Advanced pneumatization of mastoids		

Inter-operability based on mappings

▶ Syntactic:

- ↪ The terms are identical
 - Can be done by machines

▶ Semantic:

- ↪ The concepts are identical
 - Should be done by humans

▶ Structural:

- ↪ The comprehension of the concepts is identical
 - Impossible to maintain

Problems encountered

- ▶ Structure
 - ↳ Criteria for top hierarchy
 - ↳ Granularity
- ▶ Terminology
 - ↳ Synonymy
 - ↳ Ambiguity
 - ↳ Uniqueness
- ▶ Versioning

Granularity

- ▶ Differences between depth degrees between terminologies
- ▶ Even if very granular, none of the terminologies is exhaustive
- ▶ Differences in granularity makes it necessary to qualify the relationships between mapping terms in order to allow interoperability
 - ↪ E = Exact mapping
 - ↪ NTBT = narrower-to-broader
 - ↪ BTNT= broader-to-narrower

International Consortium for Human Phenotype Terminologies (ICHPT)

- ▶ Workshop on 21-22 October 2013 in Boston
- ▶ Participants: HPO, DDD/DECIPHER, FORGE/Canada, PhenoDB, Orphanet, Elements of Morphology, POSSUM, ClinGen, OMIM, Mouse Genome Informatics, SimulConsult
- ▶ Committed to agree on ~2000 high level terms (with definitions and synonyms)
- ▶ Make sure that these are used, mapped, and have behind them an ontology
- ▶ WHO and SNOMed-CT are committed to adopting these core terms

Success!

- ▶ Reviewed 2736 terms appearing 2 or more times in the 6 terminologies in 17 hours
- ▶ 2302 terms chosen, including preferred term
- ▶ Synonyms are clear from the list
- ▶ Definitions are from Elements of Morphology if available, and HPO/Stedman's Medical Dictionary, if not
- ▶ List of terms, mapping to HPO, PhenoDB, Elements of Morphology will be available at <http://ichpt.org> by January 2015.
- ▶ All tools will map to this terminology to allow interoperability among resources

Consensus

- ▶ New tools being developed to allow data sharing of unsolved exomes/genomes will use and/or incorporate this terminology

- ↪ GeneMatcher

- ↪ LOVD

- ↪ PhenomeCentral

- ↪ DECIPHER

Adoption of a core set of >2,300 terms common to all terminologies

Workshop on Terminologies for RD – Paris, 12 September 2012



- ▶ Many terminologies in use to describe phenomes - No interoperability
- ▶ Joint EuroGenTest and EUCERD workshop
- ▶ Organized by Ségolène Aymé
- ▶ Agreement to define a core set of terms common to all terminologies and a methodology
- ▶ Core set identified by cross referencing
 - ↪ HPO
 - ↪ PhenoDB
 - ↪ Orphanet
 - ↪ UMLS: MeSH, MedDRA, SnoMed CT
 - ↪ LDDb
 - ↪ Elements of morphology

Workshop of validation, Boston 21-22 October 2013



- ▶ Workshop supported by HVP and EuroGenTest
- ▶ Organized by Ada Hamosh
- ▶ Expert review of the initial proposal
- ▶ Selection of 2,370 terms
- ▶ Decision to propose them for adoption by all terminologies
- ▶ Establishment of the International Consortium for Human Phenotype Terminologies – ICHPT
- ▶ **Publication on the IRDiRC website** with definitions from
 - ↪ HPO
 - ↪ Elements of morphology

COMPUTERS ARE NOT SMART

FROM A TERMINOLOGY TO AN ONTOLOGY

Why ontologies are needed ?

- ▶ Ontologies are representations of the knowledge in a way which is directly understandable by computers
- ▶ Ontologies allow reasoning
- ▶ Ontologies define the objects AND the relationship between the objects
 - ↳ Is ais part of.... Is a cause of....
- ▶ Anemias
 - ↳ Shistosomias
 - Is a possible cause of anemia but is not an anemia

Standardization of Phenotype Ontologies

Workshop Sympathy, 19 Apr 2013, Dublin

Organized by IRDiRC, supported by the University of Dublin, Forge and EuroGenTest

Conclusion: Adopt HPO & ORDO & cross-reference with OMIM

Orphanet Rare Disease Ontology
Summary Classes Notes Mappings Widgets

Jump To:

- age of onset
- gene
- inheritance
- obsolete_class
- phenome**
- biological anomaly
- clinical subtype
- clinical syndrome
- disease
- etiological subtype
- group of phenome
- histopathological subtype
- malformation syndrome
- morphological anomaly
- particular clinical situation in a disease or syndrome
- point prevalence

Details	Visualization	Notes (0)	Class Mappings (1)
Preferred Name	phenome		
Definitions	A set of phenotypes expressed at totality of all traits of an organism		
ID	http://www.orpha.net/ORDO/Orp		
definition	A set of phenotypes expressed at totality of all traits of an organism		
definition_citation	(Mahner and Kary, 1997)		
label	phenome		
prefixIRI	ORDO:Orphanet_C001		
prefLabel	phenome		
subClassOf	http://www.w3.org/2002/07/owl#		



Human Phenotype Ontology
Summary Classes Notes Mappings Widgets

Jump To:

- All
- Mode of inheritance
- Onset and clinical course
- Phenotypic abnormality
 - Abnormality of blood and blood-forming tissues
 - Abnormality of connective tissue
 - Abnormality of head and neck
 - Abnormality of metabolism/homeostasis
 - Abnormality of prenatal development or birth
 - Abnormality of the abdomen
 - Abnormality of the breast
 - Abnormality of the cardiovascular system
 - Abnormality of the ear
 - Abnormality of the endocrine system
 - Abnormality of the eye
 - Abnormality of the genitourinary system
 - Abnormality of the immune system
 - Abnormality of the integument
 - Abnormality of the musculature

Details	Visualization	Notes (0)	Class Mappings (13)
Preferred Name	All		
Definitions	Root of all terms in the Human Phenotype Ontology		
ID	http://purl.obolibrary.org/obo/HPO		
comment	Root of all terms in the Human Phenotype Ontology		
label	All		
notation	HP:0000001		
prefLabel	All		
subClassOf	http://www.w3.org/2002/07/owl#		

Standardisation of Phenotype Ontologies

Rare Diseases

Orphanet Rare Disease Ontology
Summary Classes Notes Mappings Widgets

Jump To:

- age of onset
- gene
- inheritance
- obsolete_class
- phenome**
 - biological anomaly
 - clinical subtype
 - clinical syndrome
 - disease
 - etiological subtype
 - group of phenome
 - histopathological subtype
 - malformation syndrome
 - morphological anomaly
 - particular clinical situation in a disease or syndrome
 - point prevalence

Details	Visualization	Notes (0)	Class Mappings (1)
Preferred Name	phenome		
Definitions	A set of phenotypes expressed at the totality of all traits of an organism		
ID	http://www.orpha.net/ORDO/Orph		
definition	A set of phenotypes expressed at the totality of all traits of an organism		
definition_citation	(Mahner and Kary, 1997)		
label	phenome		
prefixIRI	ORDO:Orphanet_C001		
prefLabel	phenome		
subClassOf	http://www.w3.org/2002/07/owl#		

bioportal.bioontology.org/ontologies/ORDO

Phenotypic Features

Human Phenotype Ontology
Summary Classes Notes Mappings Widgets

Jump To:

- All
 - Mode of inheritance
 - Onset and clinical course
 - Phenotypic abnormality
 - Abnormality of blood and blood-forming tissues
 - Abnormality of connective tissue
 - Abnormality of head and neck
 - Abnormality of metabolism/homeostasis
 - Abnormality of prenatal development or birth
 - Abnormality of the abdomen
 - Abnormality of the breast
 - Abnormality of the cardiovascular system
 - Abnormality of the ear
 - Abnormality of the endocrine system
 - Abnormality of the eye
 - Abnormality of the genitourinary system
 - Abnormality of the immune system
 - Abnormality of the integument
 - Abnormality of the musculature

Details	Visualization	Notes (0)	Class Mappings (13)
Preferred Name	All		
Definitions	Root of all terms in the Human Phenotype Ontology		
ID	http://purl.obolibrary.org/obo/HP		
comment	Root of all terms in the Human Phenotype Ontology		
label	All		
notation	HP:0000001		
prefLabel	All		
subClassOf	http://www.w3.org/2002/07/owl#		

bioportal.bioontology.org/ontologies/HP

Based on Orphanet multi-hierarchical classification of RD

Genes– diseases relationships

Cross-references:

- For RD nomenclature : OMIM, SNOMED CT, ICD10, MeSH, MedDRA, UMLS
- For genes : OMIM, HGNC, UniProtKB, IUPHAR, ensembl, Reactome

ICHPT

(International Consortium for Human Phenotype Terminologies)

2,307 terms- core terminology

Mapped to:

HPO Elements of Morphology
Orphanet LDDB
SNOMED CT Pheno-DB (OMIM)
MeSH UMLS

Available soon for download at ichpt.org

Orphanet Ontology Browser

- ⊕ obsolete_class
- ⊕ gene
- ⊕ age of onset
- ⊖ phenome
 - ⊕ biological anomaly
 - ⊕ clinical subtype
 - ⊕ clinical syndrome
 - ⊕ disease
 - ⊕ etiological subtype
 - ⊖ group of phenome
 - ⊕ Rare abdominal surgical disease
 - ⊕ Rare allergic disease
 - ⊕ Rare bone disease
 - ⊖ Rare cardiac disease
 - ⊖ Cardiomyopathy
 - ⊕ Arrhythmogenic right ventricular dysplasia
 - ⊕ Dilated cardiomyopathy
 - ⊕ Hypertrophic cardiomyopathy
 - ⊖ Idiopathic giant cell myocarditis
 - ⊕ Restrictive cardiomyopathy
 - ⊕ Unclassified cardiomyopathy
 - ⊕ Coronary artery disease - hyperlipidemia - hypertension - dial
 - ⊕ LMNA-related cardiocutaneous progeria syndrome
 - ⊖ Pericarditis
 - ⊕ Rare cardiac rhythm disease
 - ⊕ Rare cardiac tumor
 - ⊕ Rare circulatory system disease

Legend:

- is a
- develops from
- part of
- other

Help ([hide](#))

Double-click a term to see its children. The ontology browser is populated dynamically. If there are many children for a given term, there may be a small delay while the browser fetches. Click to highlight a term to see any information associated with it. Hover over a term to see its relation with its immediate parent. Root terms will not display any relational information.

Relations

Term Information

ID: [Orphanet:329874](#)

Zoom

Name: Idiopathic giant cell myocarditis

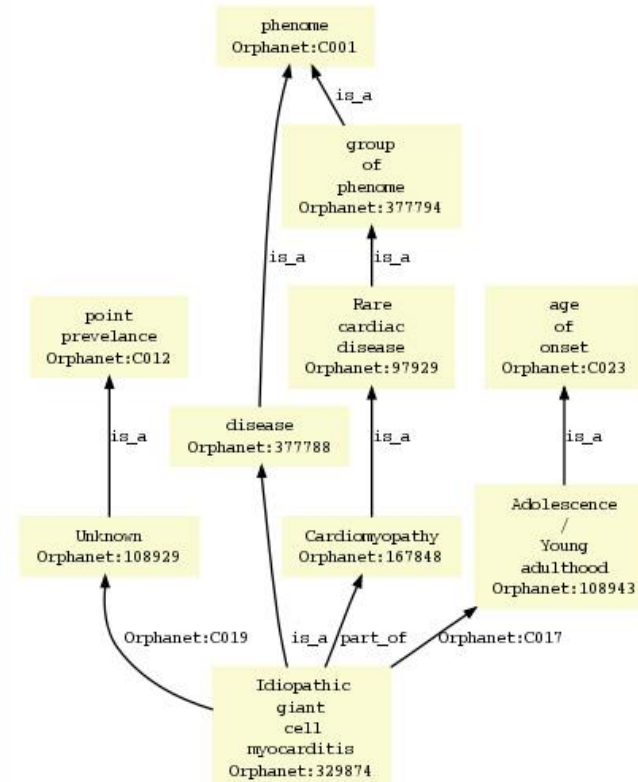
Associated information

alternative:term

IGCM

Term Hierarchy

Paths to Root: ☒ Child relationships: ☐



You can zoom the ontology browser by clicking on a term in the graph.

Consensus at IRDiRC level

- ▶ New tools being developed to allow data sharing of unsolved exomes/genomes should use and/or incorporate this terminology and these ontologies
 - ↳ GeneMatcher
 - ↳ LOVD
 - ↳ PhenomeCentral
 - ↳ DECIPHER
- ▶ ORDO and HPO formally adopted by IRDiRC



IRDiRC

INTERNATIONAL
RARE DISEASES RESEARCH
CONSORTIUM



Please disseminate these tools
to speed up R&D
to the benefit of the patients

Thank you for your attention