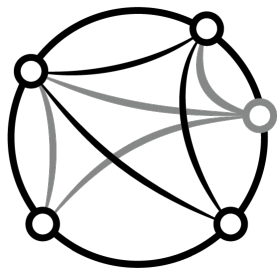


<https://w3id.org/sssom>



sssom

SIMPLE STANDARD FOR SHARING
ONTOLOGY MAPPINGS



FAIR Impact Workshop: SSSOM - a machine actionable model for simple entity mappings

Nicolas Matentzoglou, Why Mappings Matter and how to make them FAIR?
Workshop, 13.04.2023

<https://www.lifewatch.eu/events/why-mappings-matter-and-how-to-make-them-fair-a-fair-impact-workshop/>



* SSSOM can be pronounced "sessom"

What do we mean by “mapping”?

- **Entity mapping:** Determining and documenting the correspondence of an entity in one semantic space to another.
- **Schema mapping:** Determining and documenting the translation rules for converting an entity from one semantic space to another.

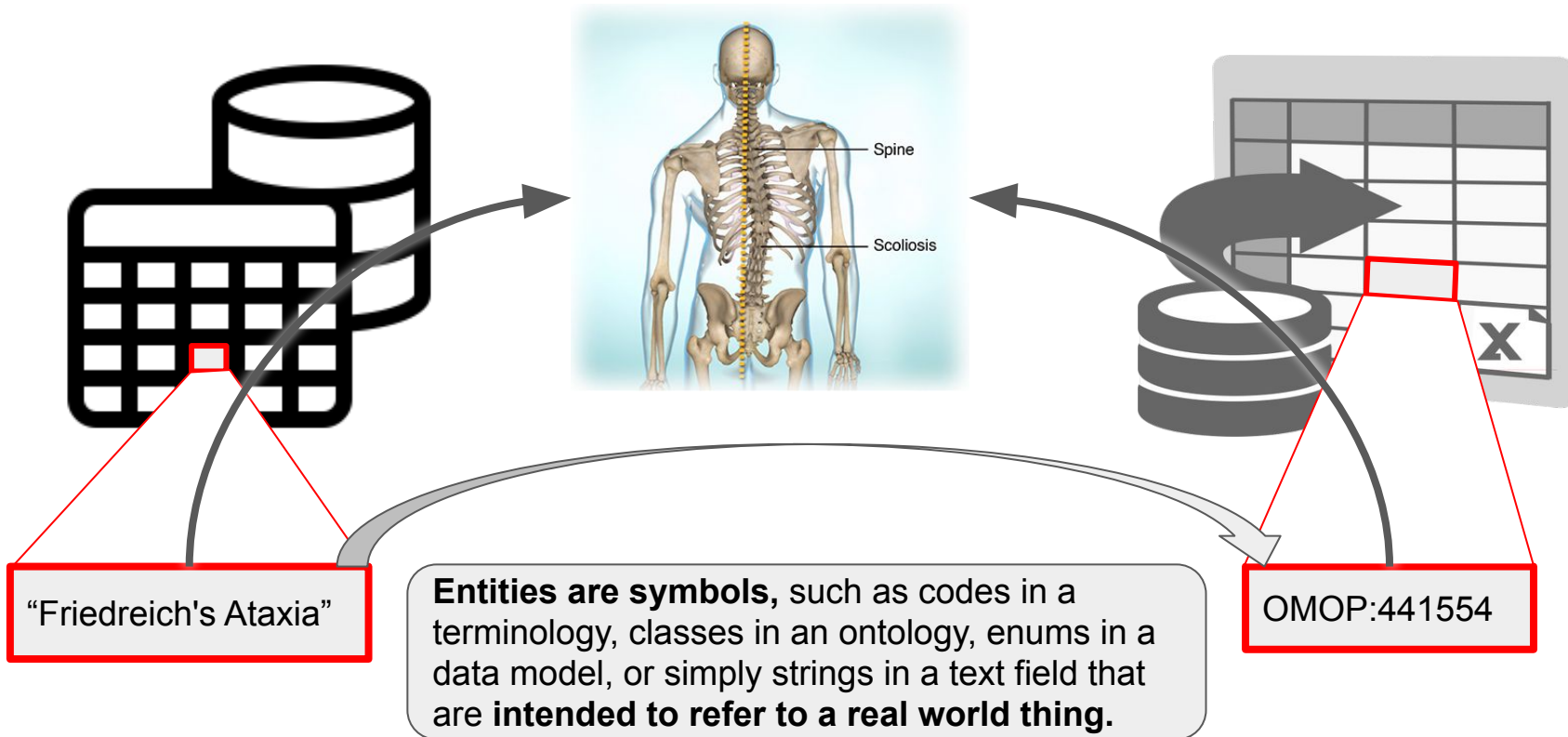
What do we mean by “mapping”?

- **Entity mapping:** Determining and documenting the correspondence of an entity in one semantic space to another.
- ~~**Schema mapping:** Determining and documenting the translation rules for converting an entity from one semantic space to another.~~

Today is all about entity mapping - this is not to say that schema mapping is not super relevant as well for us.

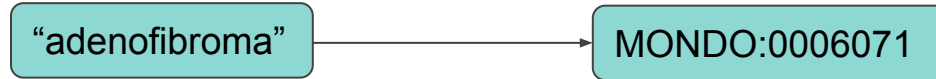


What are entity mappings?

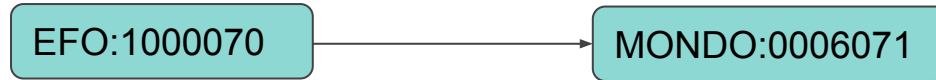


Different types of entity mappings

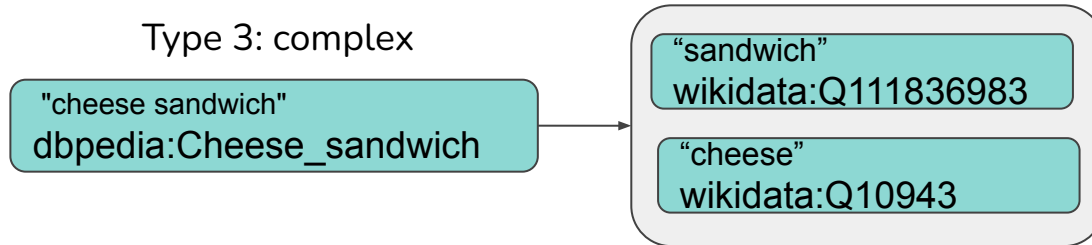
Type 1: string - identifier



Type 2: identifier - identifier



Type 3: complex



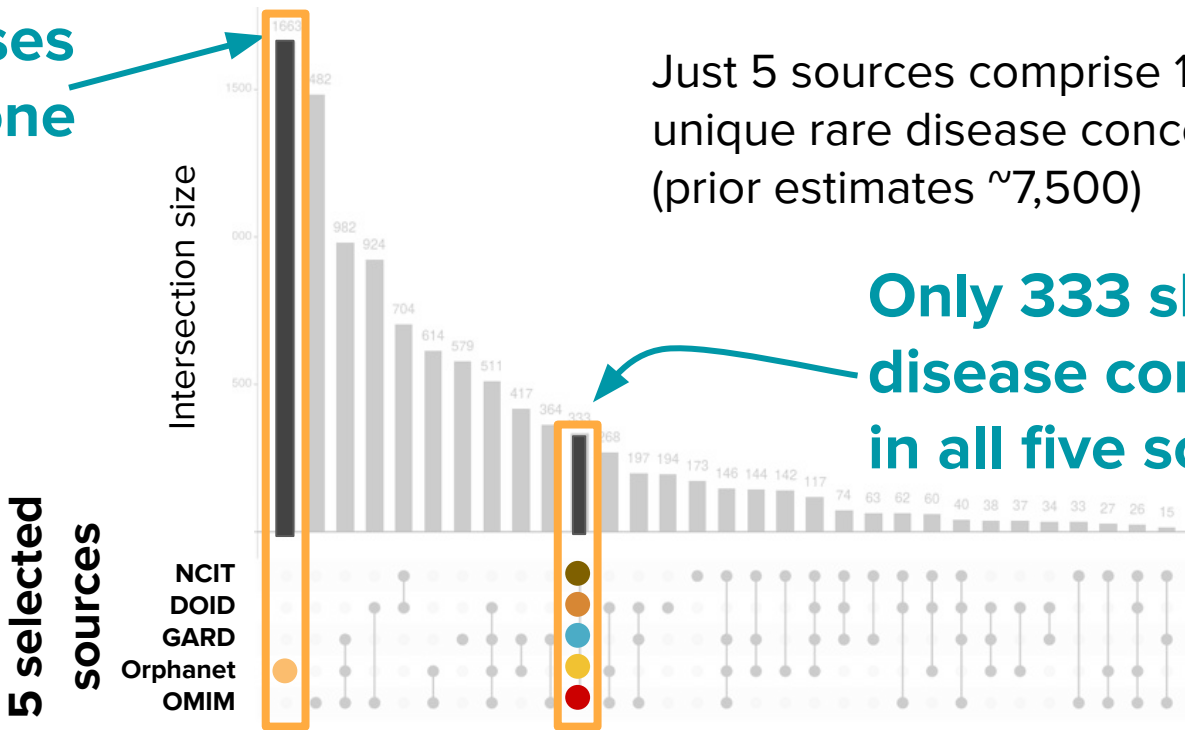
What is the problem?

Mappings frequently lack semantic precision

| ICD10CM Code | ICD10CM Label | Relation | OMOP Label | OMOP ID |
|--------------|---|----------|-------------------|---------|
| A06 | Amebiasis | Maps to | Amebic infection | 438959 |
| D46.A | Refractory anaemia with multi-lineage dysplasia | Maps to | Refractory anemia | 4003185 |

We need to enable analytics across semantic spaces

Many diseases
are in only one
source



Just 5 sources comprise 10,577
unique rare disease concepts
(prior estimates ~7,500)

Only 333 shared
disease concepts
in all five sources



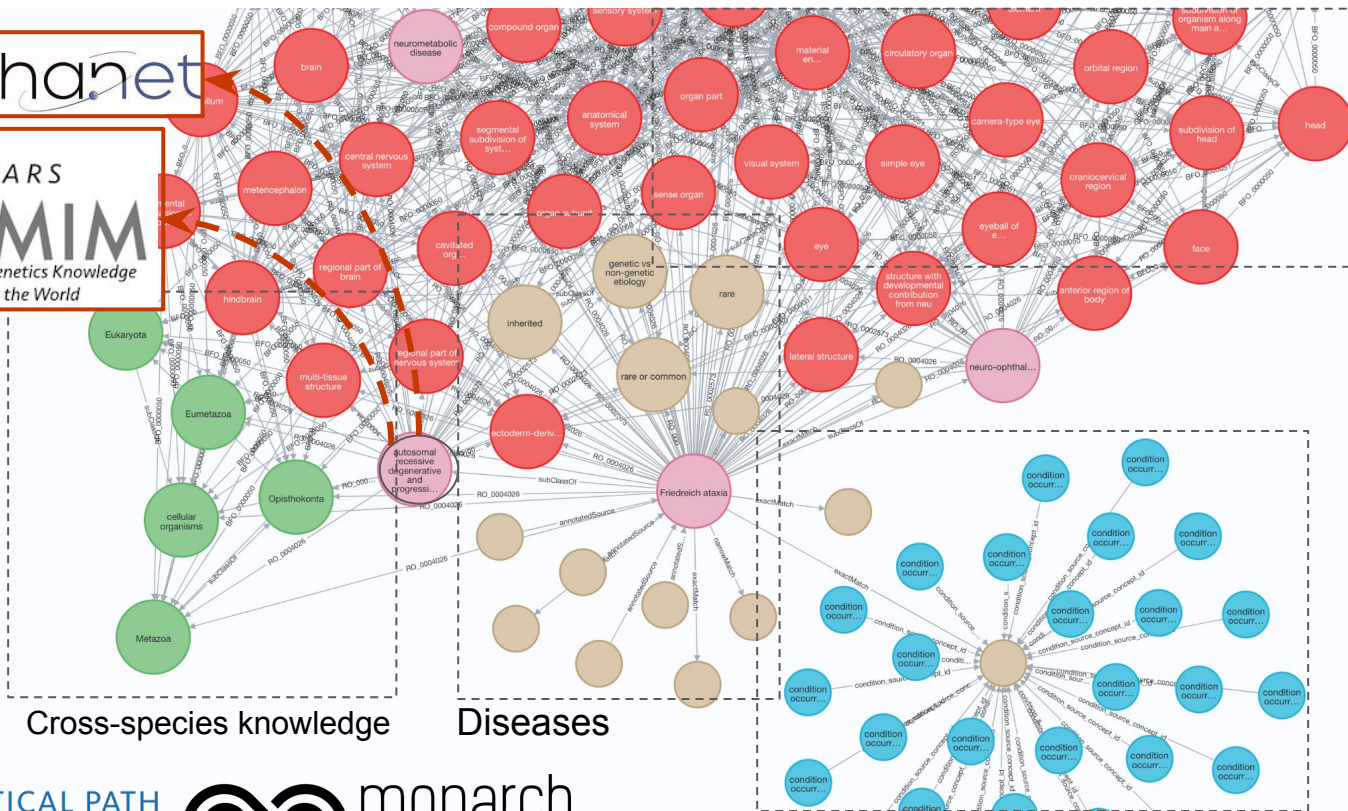
Widely used for
rare disease,
est. 25%
coverage in
SNOMED



We need to integrate genomics and clinical data across diverse semantic spaces

orphanet

5 YEARS
MIM
Human Genetics Knowledge
for the World



Anatomical
reference
models

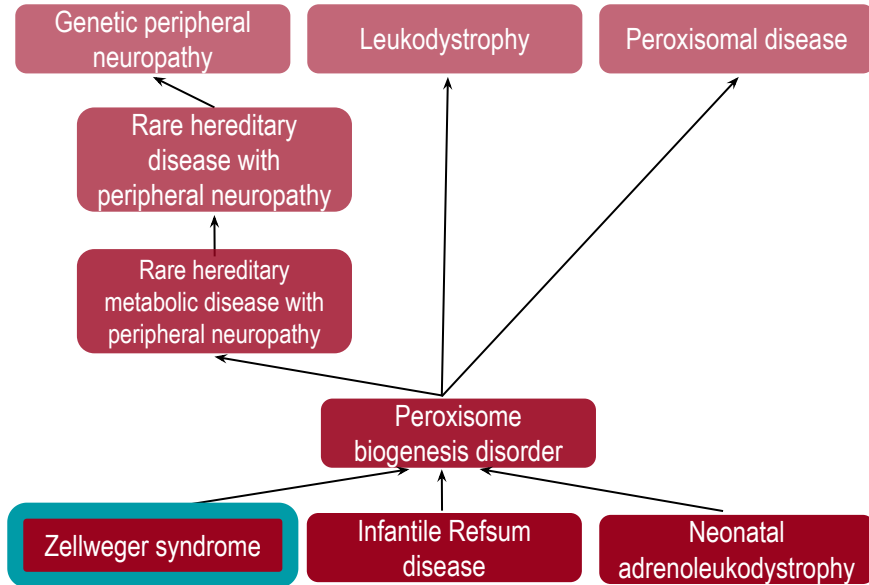
Cross-species knowledge

Diseases

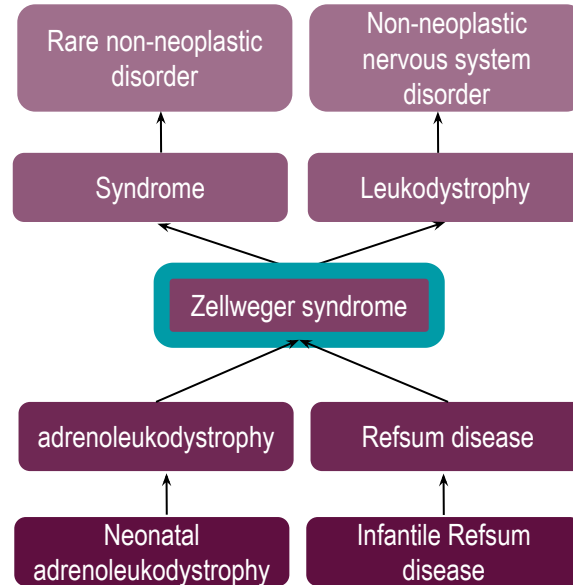
Clinical data
(condition
occurrences)

Curating good mappings is hard - and therefore costly

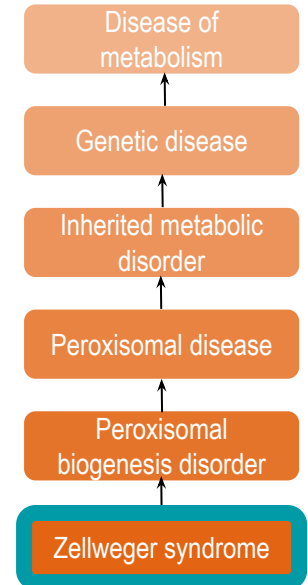
Orphanet



NCIt



DO



We can do better.

- **Share** mappings to avoid creating the same mapping over and over again
- **Enrich** mappings with metadata to enable the combination of mappings from different sources
- **Build** a coordinated decentralised effort of human curators, similar to what we do for (open) ontologies.



Society should not have to fund the incredible duplication of effort we currently have



The SSSOM Metadata Model



Rich YAML schema powered by

```
267 - other
268 - comment
269 mapping:
270   description: Represents an individual
271   slots:
272     - subject_id
273     - subject_label
274     - subject_category
275     - predicate_id
276     - predicate_label
277     - object_id
278     - object_label
279     - object_category
```



Shex shapes for validating rdf



JSON Schema

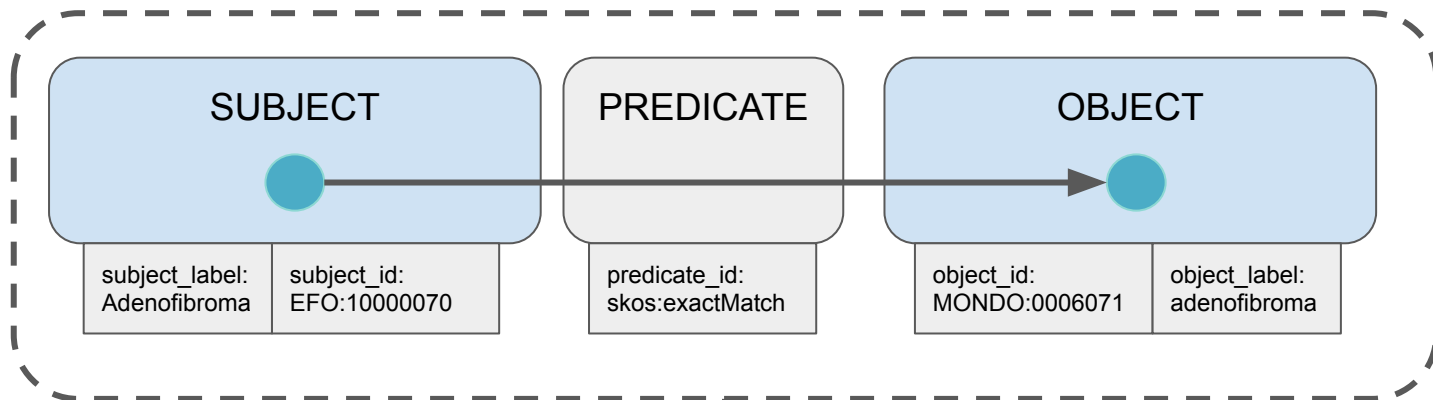


Markdown docs

Browse: <https://w3id.org/sssom/spec>

```
- subject_id
- subject_label
- subject_category
- predicate_id
- predicate_label
- object_id
- object_label
- object_category
- match_type
- creator_id
- creator_label
- license
- subject source
- subject source_version
- object source
- object source version
- mapping provider
- mapping cardinality
- mapping tool
- mapping date
- confidence
- subject match field
- object match_field
- match string
- subject preprocessing
- object preprocessing
- match term type
- semantic_similarity_score
- see also
- other
- comment
```

The anatomy of a semantic entity mapping



JUSTIFICATION

mapping_justification:
semapv:LexicalMatching

subject_match_field: rdfs:label
object_match_field: oio:hasExactSynonym
match_string: adenofibroma
mapping_date: 2022-12-13
reviewer_id: orcid:0000-0002-7356-1779
mapping_tool: wikidata:Q64360017
confidence: 0.8

Example SSSOM TSV file

Can be exported to JSON, RDF, etc.

```
#mapping_set_id: MGI_Full_MP_HPO
#mapping_set_title: All mappings of MP terms to HPO terms generated by MGI
#mapping_set_description: "Consolidated list of all HPO to MP mappings done by MGI...."
```

Provenance and descriptions

```
#creator_id:
# - orcid:0000-0003-4606-0597
# - orcid:0000-0002-6490-7723
# - orcid:0000-0003-2307-1226
# - ror:021sy4w91
# - wikidata:Q1951035
```

<https://bit.ly/ohdsi-sssom-example>

```
#license: https://creativecommons.org/licenses/by/4.0/
```

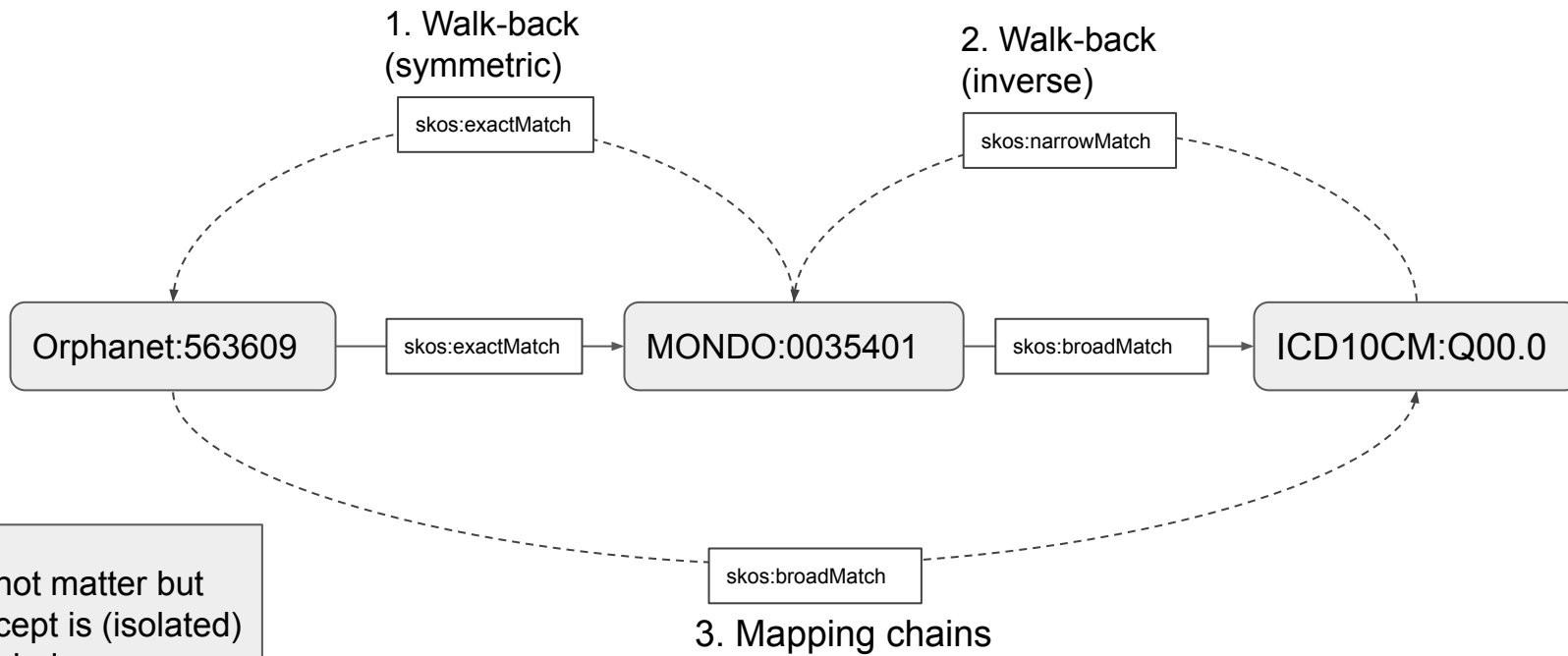
Licensing information in header

```
#object_source: obo:hp
#subject_source: obo:mp
#curie_map:
# HP: http://purl.obolibrary.org/obo/HP_
# MP: http://purl.obolibrary.org/obo/MP_
```

Mapping Table

| object_id | object_label | predicate_id | confidence | subject_id | subject_label | mapping_justification | author_id | mapping_date | comment |
|------------|-------------------|------------------|------------|------------|-----------------|-----------------------|-----------------|--------------|--------------|
| HP:0000016 | Urinary retention | skos:exactMatch | 1 | MP:0003622 | ischuria | semapv:ManualMapp | orcid:0000-0003 | 2022-08-02 | scoliosis |
| HP:0000023 | Inguinal hernia | skos:exactMatch | 1 | MP:0006077 | inguinal hernia | semapv:ManualMapp | orcid:0000-0003 | 2021-05-27 | KidsFirst |
| HP:0000028 | Cryptorchidism | skos:exactMatch | 1 | MP:0002286 | cryptorchism | semapv:ManualMapp | orcid:0000-0003 | 2021-05-27 | KidsFirst |
| HP:0000033 | Ambiguous genital | skos:narrowMatch | 1 | MP:0009198 | abnormal male | semapv:ManualMapp | orcid:0000-0003 | 2022-02-07 | KidsFirst; e |
| HP:0000034 | Hydrocele testis | skos:narrowMatch | 1 | MP:0003623 | hydrocele | semapv:ManualMapp | orcid:0000-0002 | 2021-05-27 | KidsFirst; M |

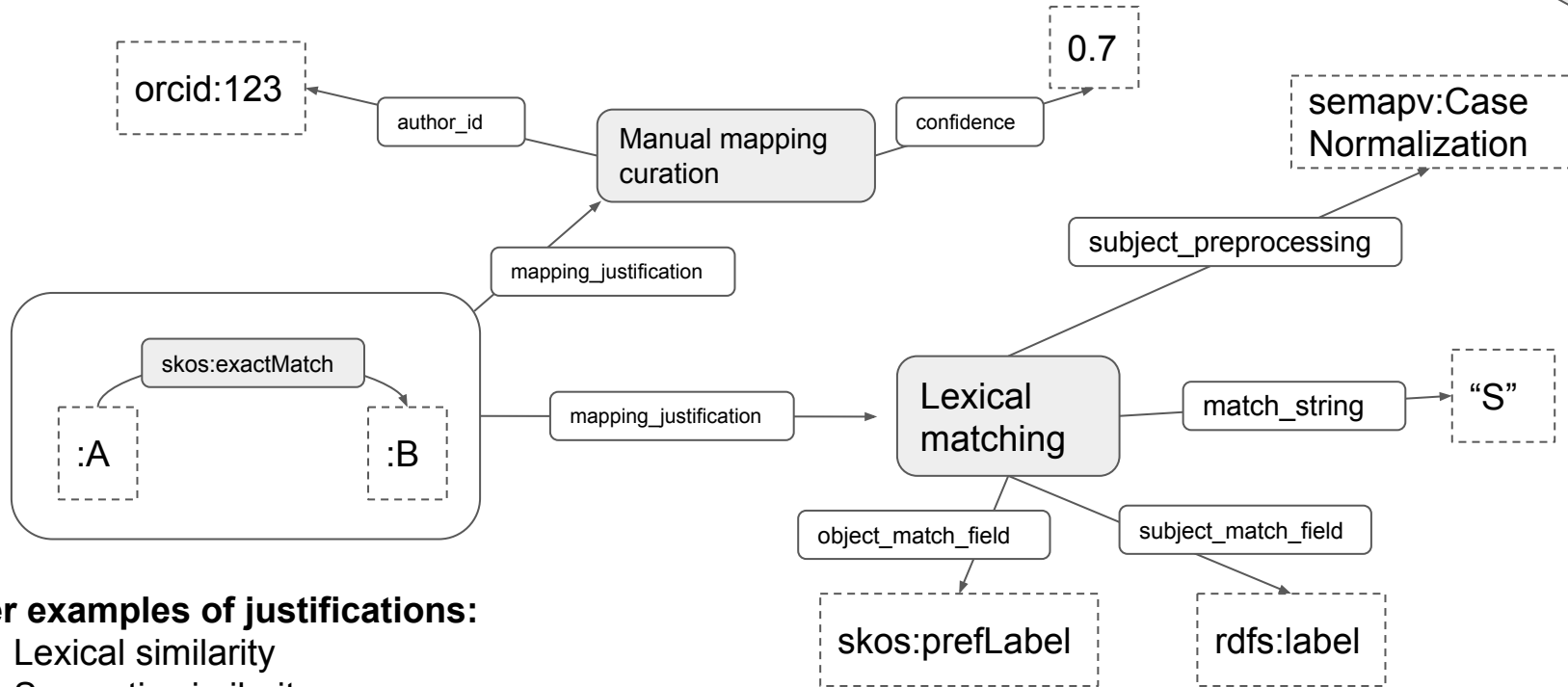
Semantic mapping predicates



It does not matter but
the concept is (isolated)
anencephaly

Mapping justifications

A mapping can have more than one justification!



Other examples of justifications:

- Lexical similarity
- Semantic similarity
- Mapping Chaining

Dereferencable identifiers

- Most metadata elements in SSSOM require the use of “entity references” rather than simple strings
- An “entity reference” should be a globally unique, persistent (and resolvable) identifier (GUPRI)
- The entity reference itself is usually recorded as a Compact URI, or CURIE, which can be resolved to a URI using a special “curie_map”.

Tip of the day: Use dereferencable identifiers **to refer to people**, rather than labels!

semapv:LexicalMatching

<https://w3id.org/semapv/vocab/LexicalMatching>



SSSOM Toolkit and other SSSOM related tools

- SSSOM toolkit (<https://github.com/mapping-commons/sssom-py>)
 - Design philosophy of SSSOM to not require any special tooling
 - Utility methods such as
 - "merge" (to merge two mapping sets)
 - "parse" (to convert a different format, such as EDOAL, into SSSOM)
 - "validate" (to check that a mapping set is legal SSSOM)
 - "filter" command allows to filter a mapping set based on any of its metadata elements
- The Ontology Access Kit (OAK) implements functionality to do basic lexical matching based on term synonyms and extracting SSSOM mappings from ontologies

A first SSSOM aware mapping browser is developed at EBI: Oxo 2



Home | Documentation | About

EFO:0001360
DOID:162
OMIM:180200
MESH:D009202

Search

Showing 10 from a total of 12

Confidence 

Predicate ▾

Download as... ▾

Show 10 ▾

Mapping Provider

- Monarch Initiative
- EBI

Mapping Justification

- Lexical
- Manual Curation

Object Prefix

Previous **1** 2 Next

Subject

EFO:0000400
"diabetes mellitus"
disease

EFO:0000400
"diabetes mellitus"
disease

Predicate

skos:exactMatch
Modifier: None

skos:broadMatch
Modifier: None

Object

MONDO:0005015
"diabetes mellitus"
disease

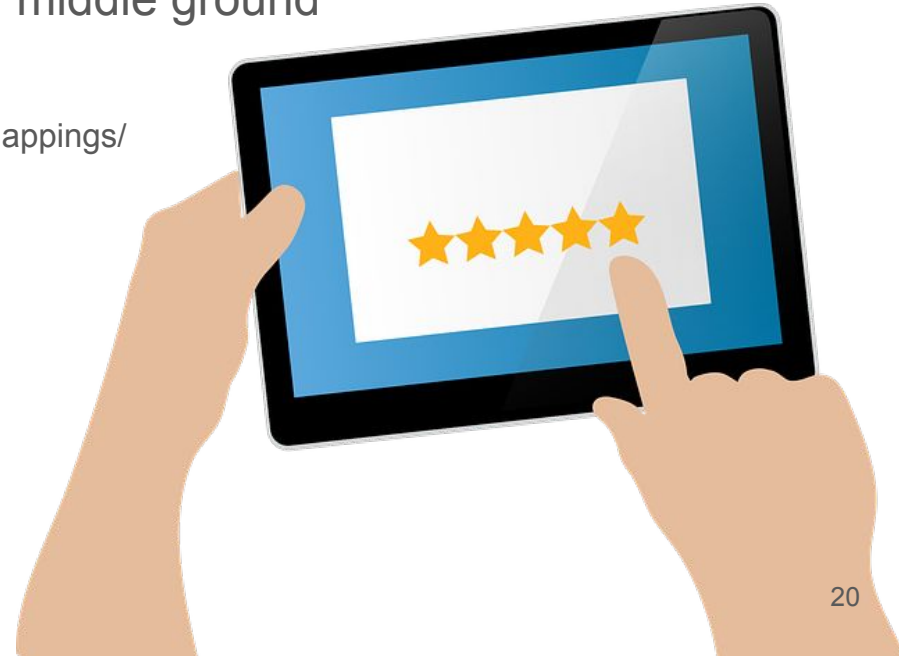
HP:0000819
"diabetes mellitus"
disease



What do we do now? A Five-Star system for mappings

- Standardisation of mappings is costly
- Sometimes, perfect is the enemy of the good enough
- As a community we need to find a good middle ground for what is “good enough”.
- See <https://mapping-commons.github.io/sssom/5star-mappings/>

Please lobby all mapping providers out there to publish mappings using CC-0 or CC-BY licenses under a public URL!



Building FAIR mapping registries: Mapping Commons

mapping-commons / mh_mapping_initiative Public Edit Pins Unv

<> Code Issues 10 Pull requests 1 Discussions Actions Projects

adding scoliosis and KidsFirst mouse model ma

Merged matentzn merged 7 commits into master from MGI_submission on Aug 14

Conversation 1 Commits 7 Checks 0 Files changed 4

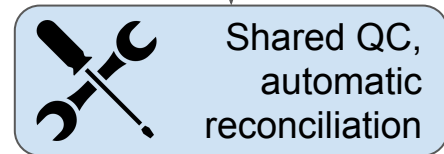
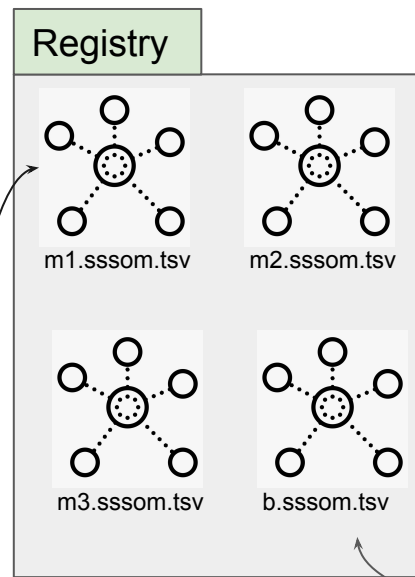
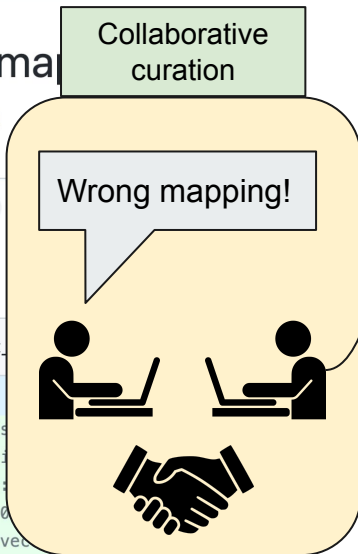
Changes from all commits File filter Conversations

Filter changed files

mappings

- hp_mp_kidsfirst_mgi.sssom....
- hp_mp_scoliosis_mgi.sssom....
- mp_hp_alzheimer_mgi.sso...
- mp_hp_covid19_mgi.sssom.t...

```
353 mappings/hp_mp_kidsfirst
... @@ -0,0 +1,353 @@
1 + #mapping_set_id: MGI_Kids
2 + #mapping_set_title: Mapp
3 + #mapping_set_description:
4 + #creator_id: orcid:0000-0
5 + #license: https://creativ
6 + #object_source: obo:hp
7 + #subject_source: obo:mp
8 + #curie_map:
9 + # HP: http://purl.obolibrary.org/obo/HP
10 + # MP: http://purl.obolibrary.org/obo/MP
```



Acknowledgements

Funding

Phenomics First (NIH / NHGRI
#1RM1HG010860-01): Spec, Mondo integration,
sssom-py CLI

Monarch (NIH / OD #5R24OD011883):
Cross-species mappings, outreach, knowledge
graph integration

Bosch Gift to LBNL: sssom-py IO, testing,
converters, tutorials

DARPA: Young Faculty Award W911NF2010255
(PI: Benjamin M. Gyori)



BOSCH

EMBL-EBI



Community contributions:

<https://w3id.org/sssom>

Core Team, alphabetical order

(<https://github.com/orgs/mapping-commons/teams/sssom-core>)

- Alex H. Wagner (Nationwide Children's Hospital)
- **Anita Caron (EMBL-EBI)**
- Charlie Hoyt (Harvard Medical School)
- Chris Mungall (LBNL)
- **Damien Goutte-Gattat (Flybase)**
- David Osumi-Sutherland (EMBL-EBI)
- **Emily Hartley (C-Path)**
- Ernesto Jimenez-Ruiz (City, Univ. of London)
- Harshad Hegde (LBNL)
- Henriette Harmse (EMBL-EBI)
- Hyeongsik Kim (Bosch)
- Ian Harrow (Pistoia Alliance)
- James McLaughlin (EMBL-EBI)
- Jim Balhoff (RENCI)
- John Graybeal (Stanford)
- Melissa Haendel (CU Anschutz)
- Nicolas Matentzoglou (EMBL-EBI)
- Nicole Vasilevsky (CU Anschutz)
- Núria Queralt Rosinach
- Simon Jupp (SciBite)
- **Sophie Aubin (INRAE)**
- Thomas Liener (Pistoia Alliance)
- Tiffany Callahan (Columbia University)
- Tim Putman (CU Anschutz)
- William Duncan (UFlorida)
- ...many more contributors, see publication

*recent joiners highlighted in bold