# Mappings @ EBI

**James McLaughlin, PhD**
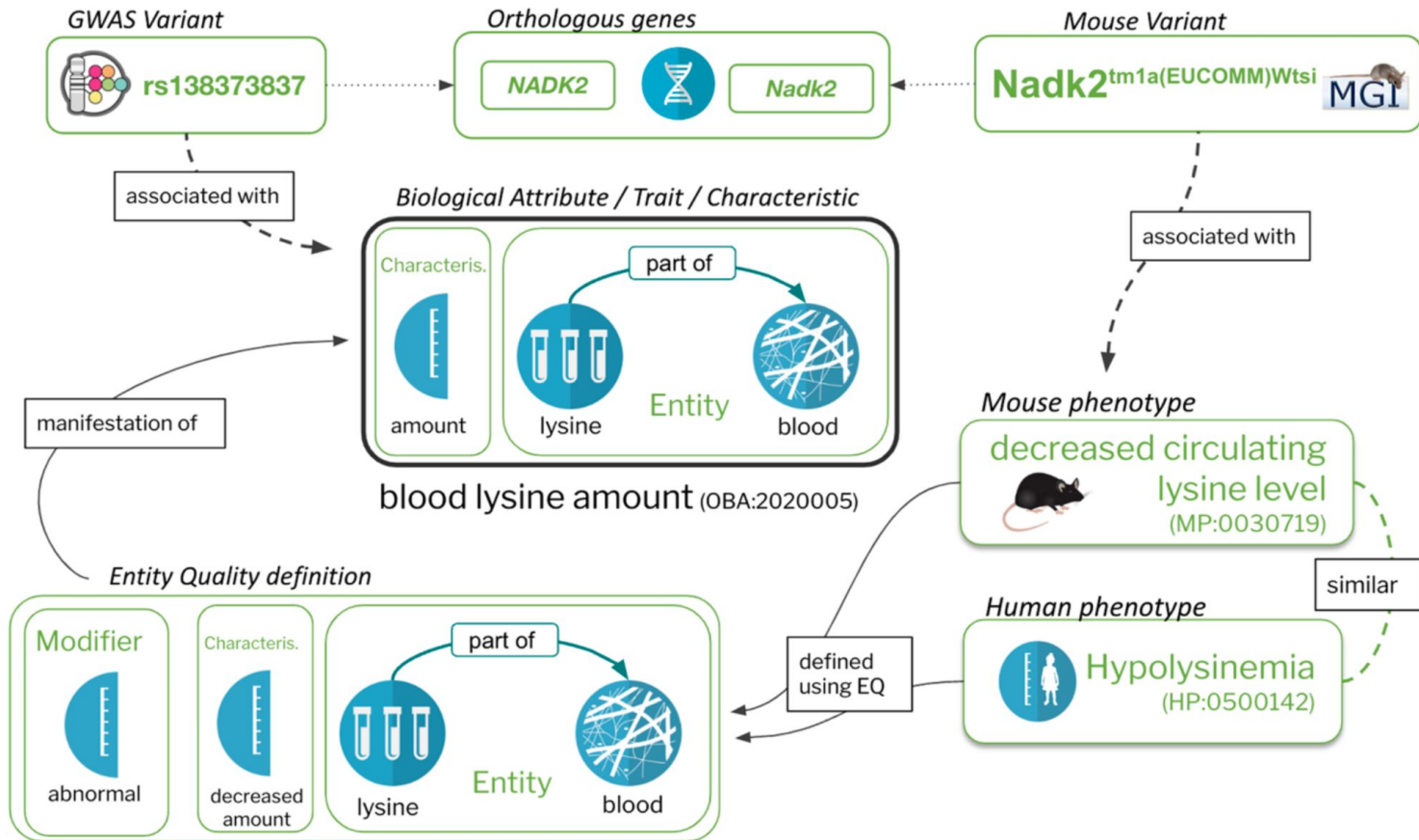
**Ontology Project Lead**
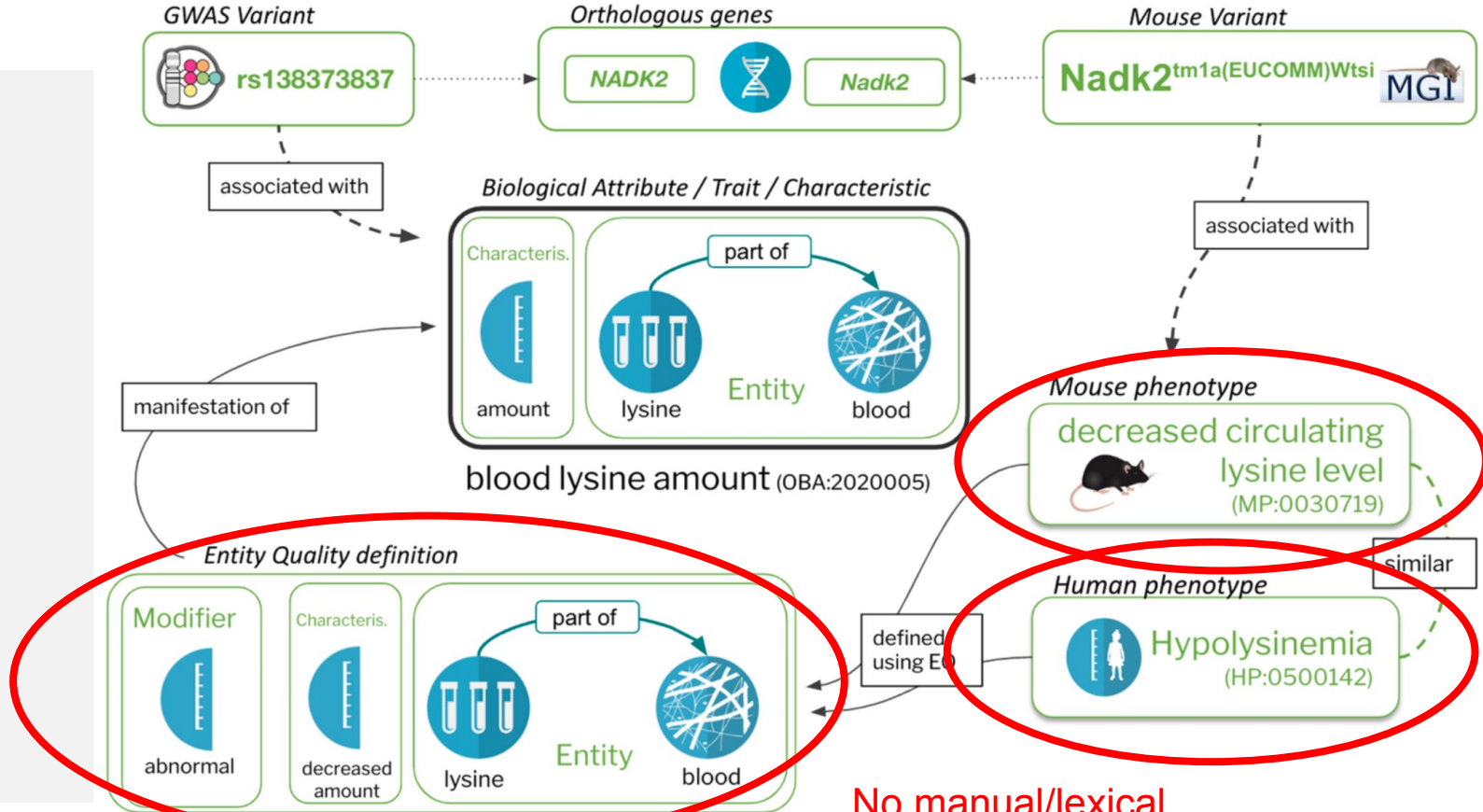
EMBL-EBI

# EMBL-EBI

- We host and manage biomedical DBs including Ensembl, UniProt, GWAS Catalog

- Our data is annotated using ontology terms
    - Our application ontology: the **Experimental Factor Ontology (EFO)** imports terms from domain ontologies such as the Monarch disease ontology (MONDO) and the Gene Ontology (GO)

- We also provide web services for working with ontologies:
    - The **Ontology Lookup Service (OLS)**
    - The **Ontology Xref Service (OxO)**
    - **ZOOMA**
    - **+ OntoString** (beta)

# Mapping use cases @ EBI

- There is LOTS of direct overlap in the ontology ecosystem, e.g. DOID and MONDO are both popular disease ontologies
  - Literally different terms which represent **exactly the same thing**, and to make datasets comparable (e.g. align a patient study with known info about a disease) we need to align them.

- Also less direct overlap e.g. disease-phenotype
  - A phenotype is an observable trait, and a disease is *generally* considered to be a collection of phenotypes
  - Sometimes there is a direct disease->phenotype mapping, sometimes a lot more nuanced.

- Mappings can be scientifically interesting!
  - E.g. mapping between mouse phenotypes in **MP** and human phenotypes in **HP** can make a mouse and human study comparable

EMBL-EBI

# Mapping using semantic definition equivalence

# Explicitly defined mappings

- Ontologies often define mappings **(especially using the "oboInOwl:hasDbXref" predicate)**
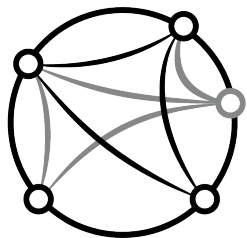
# Explicitly defined mappings

- We are gradually trying to externalise mappings rather than store them in ontologies
  - Because you might want to make mappings between vocabularies you did not develop
  - And mappings can be subjective: the same vocabularies can be mapped in different ways



**SSSOM**: a simple standard for storing mappings in TSV files

**Mapping-commons:** an organisation on github for people who love mappings (come and join us!)

EMBL-EBI

# sssom

SIMPLE STANDARD FOR SHARING
ONTOLOGY MAPPINGS

## A Simple Standard for Sharing Ontological Mappings (SSSOM) ⟜

Nicolas Matentzoglu, James P Balhoff, Susan M Bello, Chris Bizon, Matthew Brush, Tiffany J Callahan, Christopher G Chute, William D Duncan, Chris T Evelo, Davera Gabriel ... Show more

- TSV standard with standard metadata elements to describe mapping

| subject_id | predicate_id | object_id | match_type | subject_label | object_label |
|---|---|---|---|---|---|
| HP:0009124 | skos:exactMatch | MP:0000003 | Lexical | Abnormal adipose tissue morphology | abnormal adipose tissue morphology |
| HP:0008551 | skos:exactMatch | MP:0000018 | Lexical | Microtia | small ears |
| HP:0000411 | skos:exactMatch | MP:0000021 | Lexical | Protruding ear | prominent ears |

Mapping predicates include owl:EquivalentClass; OBO xref
Other columns include confidence score; mapping tool

https://github.com/mapping-commons/SSSOM

EMBL-EBI

# But how do we define the mappings?!

- **Manual curation:** an expert in the subject matter looks at the sets of terms and asserts that they map to each other

- **Lexical matching:** String similarity of labels

- **Semantic similarity:** e.g. Jaccard; what is the semantic overlap of the terms

- And everything in between, e.g. combining these methods with human supervision to make a semi-automated mapping pipeline

- **Next: LLMs ?!**

EMBL-EBI

# Types of mappings

- **Clear <u>correspondence</u>**
  - E.g. skos:exactMatch, owl:equivalentClass

- **Fuzzy <u>correspondence</u>**
  - E.g. gene <codes for> protein
  - Glucose <is measured by> glucose level measurement

- **<u>Associations</u> however are different from mappings**
  - E.g. student <enrolled in> university
  - Phone <manufactured in> china

EMBL-EBI

# Slightly different: <u>**String to term**</u> mappings

- Mappings can be **term to term**….
  - **DOID:162** (cancer) —> **MONDO:0004992** (cancer)
  - (Like all of the examples we have just seen)

- Another class of mappings: <u>**string to term**</u>
  - *"Type 2 diabetes"* **—>** **EFO:0001360** "type II diabetes mellitus"

  - Important for
    - extracting ontology terms from bodies of text
    - Importing datasets that use free text rather than terms

EMBL-EBI

# String to term mapping in SSSOM



sssom

SIMPLE STANDARD FOR SHARING
ONTOLOGY MAPPINGS

```
literal mapping:
  description: Represents an individual mapping between a a literal and an entity
  slots:
  - literal
  - literal_datatype
  - predicate_id
  - predicate_label
  - predicate_modifier
  - object_id
  - object_label
  - object_category
  - mapping_justification
  - author_id
  - author_label
  - reviewer_id
  - reviewer_label
  - creator_id
  - creator_label
  - license
  - literal_source
  - literal_source_version
  - object_type
  - object_source
  - object_source_version
  - mapping_provider
```

EMBL-EBI

# String to term mappings : ZOOMA



(ZOOMA data also now available in SSSOM)

# String to term mapping : OntoString (beta)



- Semi-automated mapping tool

- Users can upload a list of terms, OntoString will guess (using ZOOMA and OLS) a mapping based on previously manually curated mappings

- If no mapping is available, the user can explicitly map it

- This data (string -> term mappings) will eventually inform ZOOMA

(Migrating to SSSOM.)

EMBL-EBI

# Conclusion

- We at EBI use mappings extensively to link biomedical datasets

- We are trying to move away from explicit mappings to **semantic equivalence** by better defining terms in ontologies

- But we still use lots of methods for mapping e.g. manual curation, lexical matching, semantic similarity
  - And maybe LLMs soon?

- We use mappings for relationships that are <u>correspondences</u> rather than <u>associations</u>

- We are aggressively adopting SSSOM in all of our mapping infrastructure!
  - And please don't forget about string to term mappings