# Applications of mOWL for Predicting Gene-Disease Associations

Machine Learning with Biomedical Ontologies
SWAT4HCLS 2023 Conference, Basel, Feb 13-16

**SWAT4LS**          **13 Feb 2023**

Find the causative genes

Patients Clinical Phenotypes **&**
Genomics Sequence Data

Rare/Monogenic Diseases

Mutation

Gene

Inheritance pattern
(dominant or recessive)

Complex Diseases

Gene A

Gene C

Gene B

Gene D

→ Gene variations

Inheritance pattern (complex)

Prediction/Diagnosis
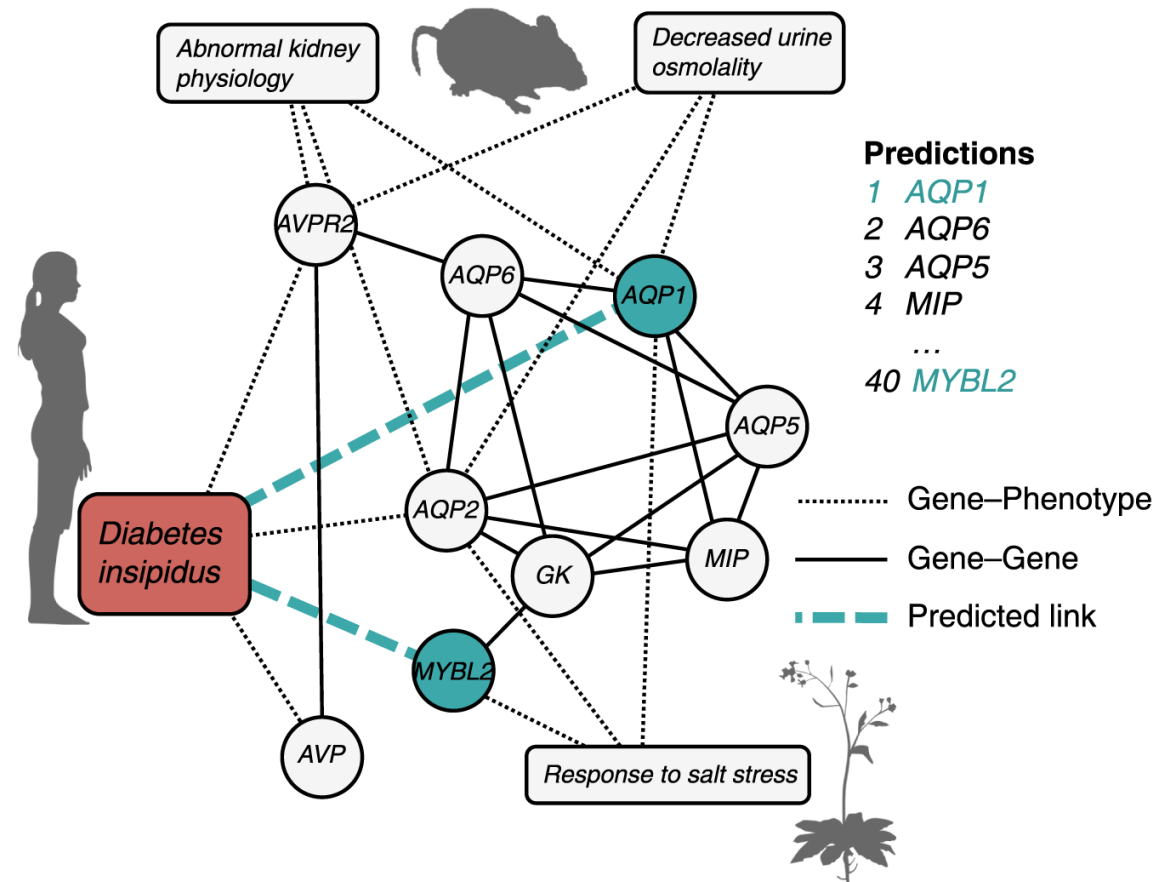
➢ Predicting gene-disease associations based on phenotypic similarity

➢ Diagnosis of disease based on phenotypic similarity

➢ Using the phenotypes of model organism genes and the diseases' phenotypes.

mOWL

Converting knowledge axioms in phenotype ontologies to knowledge graph [1].

[1] Chen, Jun*, Azza Althagafi*, and Robert Hoehndorf. "Predicting candidate genes from phenotypes, functions and anatomical site of expression." Bioinformatics 37.6 (2021): 853-860.

▪ **Using functional and phenotypic characteristics for genes in:**

| Human phenome | Mouse phenome | Functions of the gene products | Gene expression in individual cell types | Anatomical site of expression from the GTEx tissue expression |
|---|---|---|---|---|
| Human Phenotype Ontology (HPO) | Mammalian Phenotype Ontology (MP) | Gene Ontology (GO) | Celltype Ontology (CL) | Uber-anatomy ontology (UBERON) |
| 4,315 genes & 169,281 associations | 13,529 genes & 168,550 associations | 17,786 genes & 208,630 associations | 6,559 genes & 17,149 associations | 20,538 genes & 585,765 associations |

- Generate the representation from the ontology graph (using mOWL).
- Collect features for the gene and disease using different method.

| | |
|---|---|
| **Syntactic embeddings** | Onto2vec |
| | OPA2Vec |
| **Graph-based embeddings** | DL2vec |
| | OWL2Vec* |

Fig1: Onto2Vec Workflow.

Fig2: DL2vec Workflow

- The phenotypes are described using different organism-specific phenotype ontologies.

- Unified Phenotype Ontology (uPheno) include human phenotypes from the Human Phenotype Ontology (HP), → relate mutant model organism phenotypes to human disease-associated phenotypes.

- **Gene annotations:**

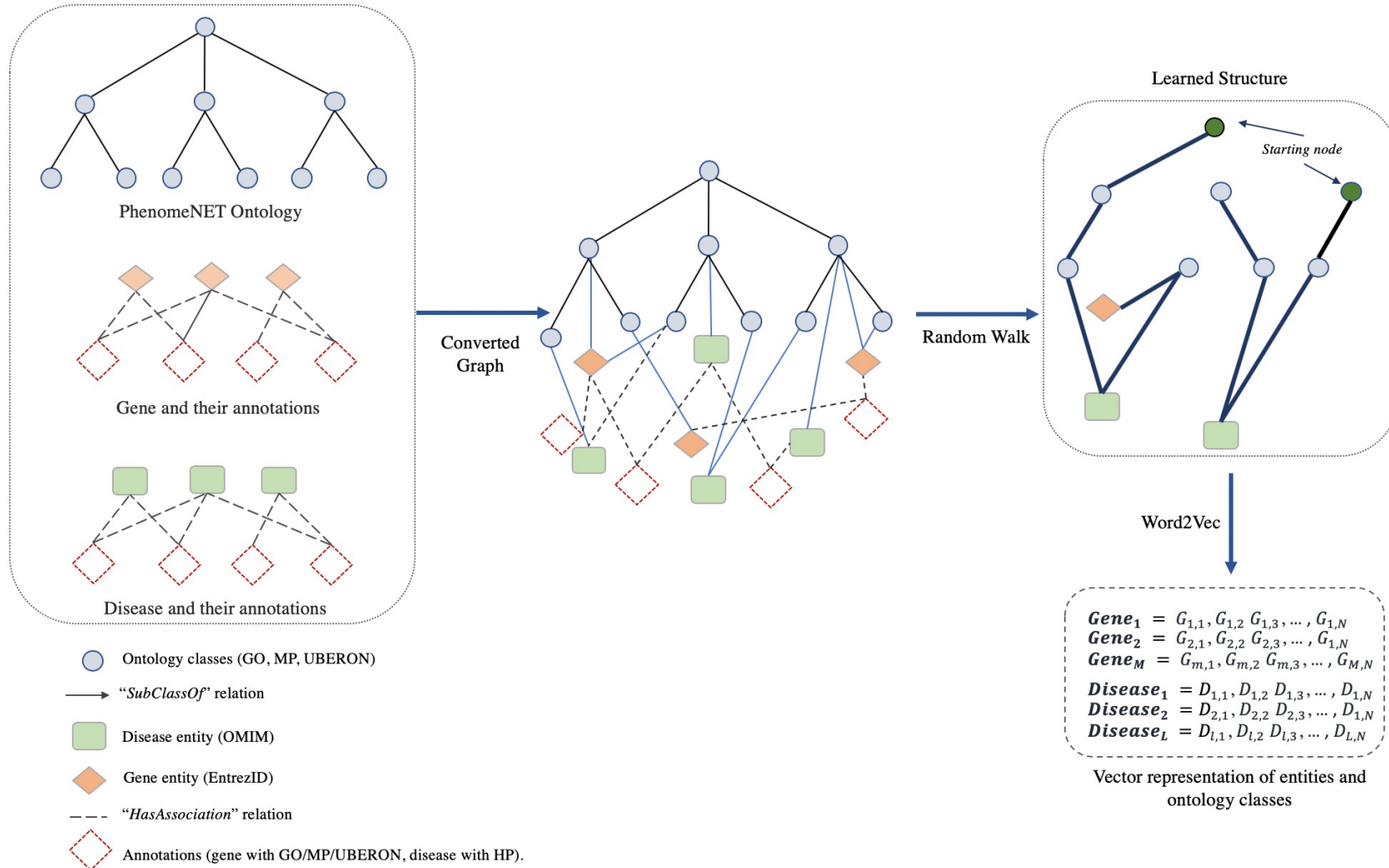  - Mouse/Human Orthology with Phenotype Annotations from HPO database (HMD_HumanPhenotype.rpt).

- **Disease annotations:**

  - obtained from the HPO annotations for rare diseases document (phenotype.hpoa).

- These annotations added to the Unified Phenotype Ontology (uPheno) to build the training ontology.

- mOWL is designed to handle input in OWL format.

  ➢ A mOWL dataset contains 3 ontologies: training, validation and testing.

- Preparee the annoations file (`Genes -> Phenotypes , Diseasese -> Phenotypes`)

- Use mOWL methods to build dataset given an ontology file and the annotations to be inserted to the ontology.

  ➢ Per each row, the first element is the annotated entity and the rest of the elements are the annotating entities (which are the entities in the ontology).

- Use different methods to generate the representation given the annotated ontology file.

mOWL

- **Unsupervised Approach**
  - ➢ Cosine similarity



$$Sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

OMIM:106190

MGI:99418

**Hands On Tutorial**

1. Prepare the dataset (Built-in / your own dataset)

2. Generate the embeddings using different methods

3. Prediction

4. Validation

**Patients Clinical Phenotypes & Genomics Data**

**Can we find the causative variants associated with the phenotypes?**

■ Prioritizing the causative variants.

➢ Determining which variants identified using Whole-exome Sequencing (WES) or Whole-genome sequencing (WGS) are most likely to damage gene function and underlie the disease phenotype.

# Phenotype-based prioritization of candidate genes



Disease

Phenotypes, Functions or site expression

Genes

- Gene-Phenotype annotations
- Gene-Function annotations
- Gene-Celltype annotations
- Gene-Anatomical site of gene expression annotations

- Disease-Phenotypes annotations

➢ Prioritize the candidate genes

$$MP_{V_i} = \text{Max}( (D_{V_i}. G_{1_{V_i}}), (D_{V_i}. G_{2_{V_i}}), \ldots, (D_{V_i}. G_{M_{V_i}}) )$$
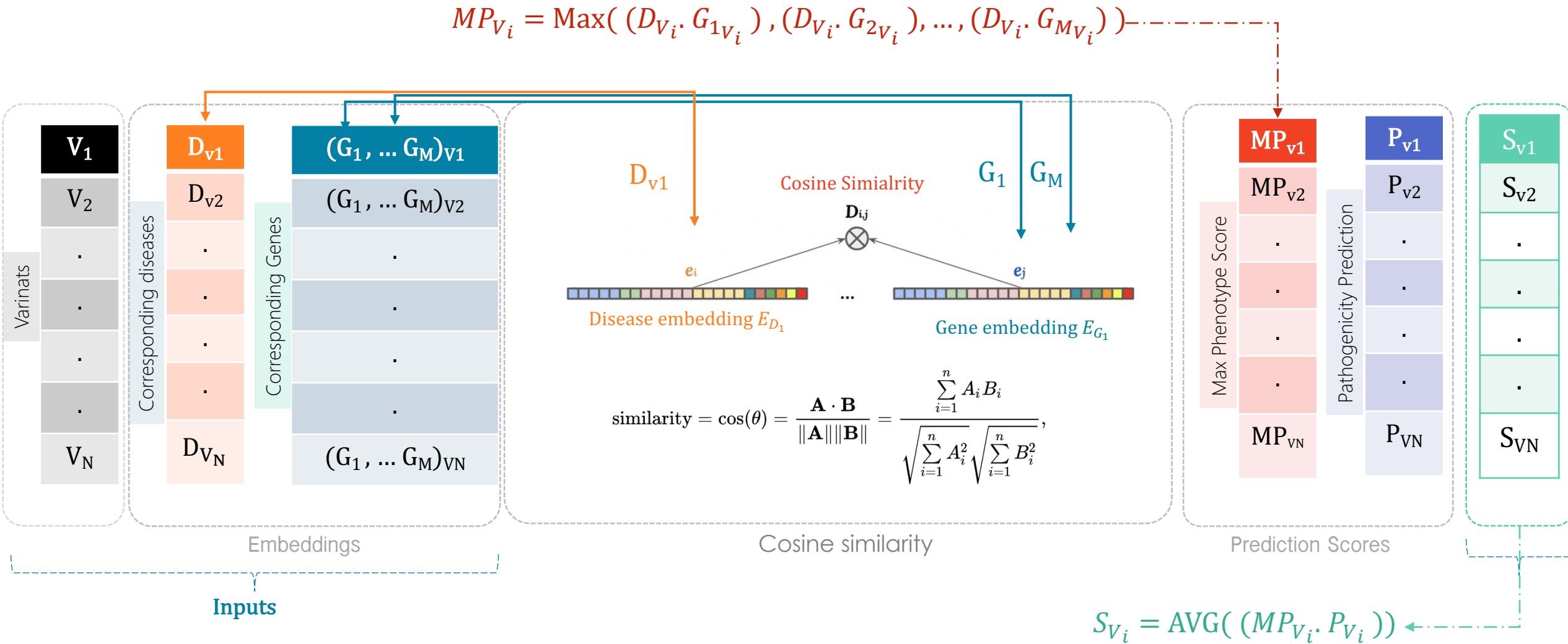
Cosine Similarity

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}},$$

Disease embedding $E_{D_1}$

Gene embedding $E_{G_1}$

Embeddings

Cosine similarity

Prediction Scores

Inputs

$$S_{V_i} = \text{AVG}( (MP_{V_i}. P_{V_i}) )$$

# Resources

- Zhapa-Camacho, Fernando, Maxat Kulmanov, and Robert Hoehndorf. "**mOWL: Python library for machine learning with biomedical ontologies**." Bioinformatics 39.1 (2023): btac811.

- Chen, Jun, Azza Althagafi, and Robert Hoehndorf. "**Predicting candidate genes from phenotypes, functions and anatomical site of expression**." Bioinformatics 37.6 (2021): 853-860.

- Smaili, Fatima Zohra, Xin Gao, and Robert Hoehndorf. "**OPA2Vec: combining formal and informal content of biomedical ontologies to improve similarity-based prediction**." Bioinformatics 35.12 (2019): 2133-2140.

- Smaili, Fatima Zohra, Xin Gao, and Robert Hoehndorf. "**Onto2vec: joint vector-based representation of biological entities and their ontology-based annotations**." Bioinformatics 34.13 (2018): i52-i60.

- Chen, Jiaoyan, et al. "**Owl2vec*: Embedding of owl ontologies**." Machine Learning 110.7 (2021): 1813-1845.

# THANK YOU!