

# GOABL Development Log

## Introduction

---

### Part I: Construct Dataset from GEO Sample Data

A species-specified dataset of *Shewanella Oneidensis*.

A dataset of gene expression in *Shewanella Oneidensis* is constructed

#### Collecting Data Sample from GEO

40 gene expression series and 426 samples of *Shewanella Oneidensis* are collected from Gene Expression Omnibus (GEO) database.

#### Obtaining Expression Matrix from GEO Sample

##### Extracting Concept from Natural Language Description

Large Language Models (LLMs) are employed in the preliminary implement. In detail, we used OpenAI's API of GPT 3.5-Turbo to understand natural language-based description of GEO experiment datasets, and summarize them to corresponding concept IDs in GO.

A fundamental risk underlies in the method above: great ambiguity and uninterpretability of LLMs is introduced to the dataset. Therefore, we are exploring better solutions.

#### Unifying Expression Level Data and Integrating Samples as Dataset

**TODO:** Ratio and Signal intensity?

#### Unify Expression Data with Importance Metric

##### non-supervised importance and top 20 relative expression with DEG

A differential expression gene (DEG) analysis is applied to the origin dataset (obtained in last section) to construct a dataset with relative expression level. The DEG was performed with python package PyDESeq2.

Values of first GSM sample in each GSE series is setted as the reference, and relative expression in other GSM samples are measured. In the final dataset, the top 20 most expressed genes (ranked by p-value) of each GSM sample is selected as data labels.

---

## Part II: Embedding Method

In a traditional method, the concept name-described instances are mapped to vectors to train a learner model.

### Ontology Embedding

The concepts of Gene Ontology is mapped to vectors with `owl2vec_star`, which preserves information of knowledge graph structure, word vector and literal meaning of concept names.

### Processing Multiple Concepts

For instances with multiple concepts as input, it is mapped to the mean value of the concept embeddings.

### Obtaining Learner Model with Concept Embeddings

---

## Part III: Knowledge Graph Preprocessing: Remembering

Knowledge graphs (KGs) like GO are often in large scale, containing much information that is irrelevant to specific learning task and a certain degree of noise. To perform tractable reasoning and abductive learning on such large-scale KG, we adopted a method, ABL-KG (Huang et al. 2023) to mine logic rules and filter out irrelevant information from large-scale KGs.

### Structure of OWL Knowledge Graph File

#### Subgraph Extraction

#### Rule Mining

#### Remember Algorithm

---

**Algorithm 1:** Remembering

---

**Result:** Write here the result**Input :** Write here the input**Output:** Write here the output

```
1  $R_{new} \leftarrow \{r \in \mathcal{KB} \mid \text{Sig}(r) \cap \Sigma' \neq \emptyset\};$   
2 while  $t < T$  and  $r_{new} \neq \emptyset$  do  
3    $R_{res} \leftarrow \emptyset;$   
4 end
```

---

## **Part IV: Abductive Learning**

**Building Learner**

**Building Reasoner**

**Bridging Learner and Reasoner**