



# ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA

## Gene Ontology Analysis System Report

### Team members

- Elnaz Niloofary
- Luna López-Aparicio Rabasco
- Momchil Petkov
- Alessandra Zoli

**Professor:** Andrea Giovanni Nuzzolese

**Date:** February 2026

# Abstract

This report details the design and implementation of a Gene Ontology (GO) analysis system. The software provides a web-based interface for researchers to upload standard GO data formats (OBO and GAF), perform statistical analysis, and calculate semantic similarity between genes and terms using multiple algorithms (Jaccard, Wu-Palmer, Resnik). The system is built using Python and Flask, featuring a flexible and modular design that allows researchers to easily switch between different analysis algorithms and expand the tool in the future. Key features include a recursive annotation checking, similarity score calculations, searching genes or terms, and pathfinding between ontology terms.

## Introduction

Gene Ontology (GO) is a major bioinformatics initiative to unify the representation of gene and gene product attributes across all species. Analyzing this data requires robust tools capable of handling complex Directed Acyclic Graph (DAG) structures and large annotation datasets.

**Problem Statement:** Researchers need a way to easily explore the relationships between genes and functions, specifically finding common ancestors and checking how similar two genes are based on their function.

**Proposed Solution:** We developed a web application that parses standard OBO (structure) and GAF (annotation) files. The system allows users to search for genes, see term neighborhoods, and compute similarity scores. The backend uses appropriate data structures to handle calculations, while the frontend provides a user-friendly dashboard.

# System Architecture

The software is organized into three main layers, allowing the system to handle data storage, analysis, and user interaction separately.

## 1. Data Handling Layer (The "Brain")

*ontology.py*: This component loads the Gene Ontology structure (the tree/graph of terms). It understands how terms are related (parents/children) and allows the system to find paths between biological terms, like molecular functions, biological processes and cellular components.

*repository.py*: This manages the gene annotation data. It acts like a digital library, allowing the system to quickly look up which genes belong to which terms (and vice versa) and handle the complex associations found in GAF files.

*models.py*: Defines the basic biological objects in the system: a *Gene* (with its symbol and annotations) and a *GO Term* (with its ID, name, and definition).

## 2. Analysis Core (The "Calculator")

*analysis.py*: This contains the mathematical algorithms used to compare terms or genes. It can calculate semantic similarity scores (like Jaccard or Resnik) to determine how functionally similar two genes or terms are.

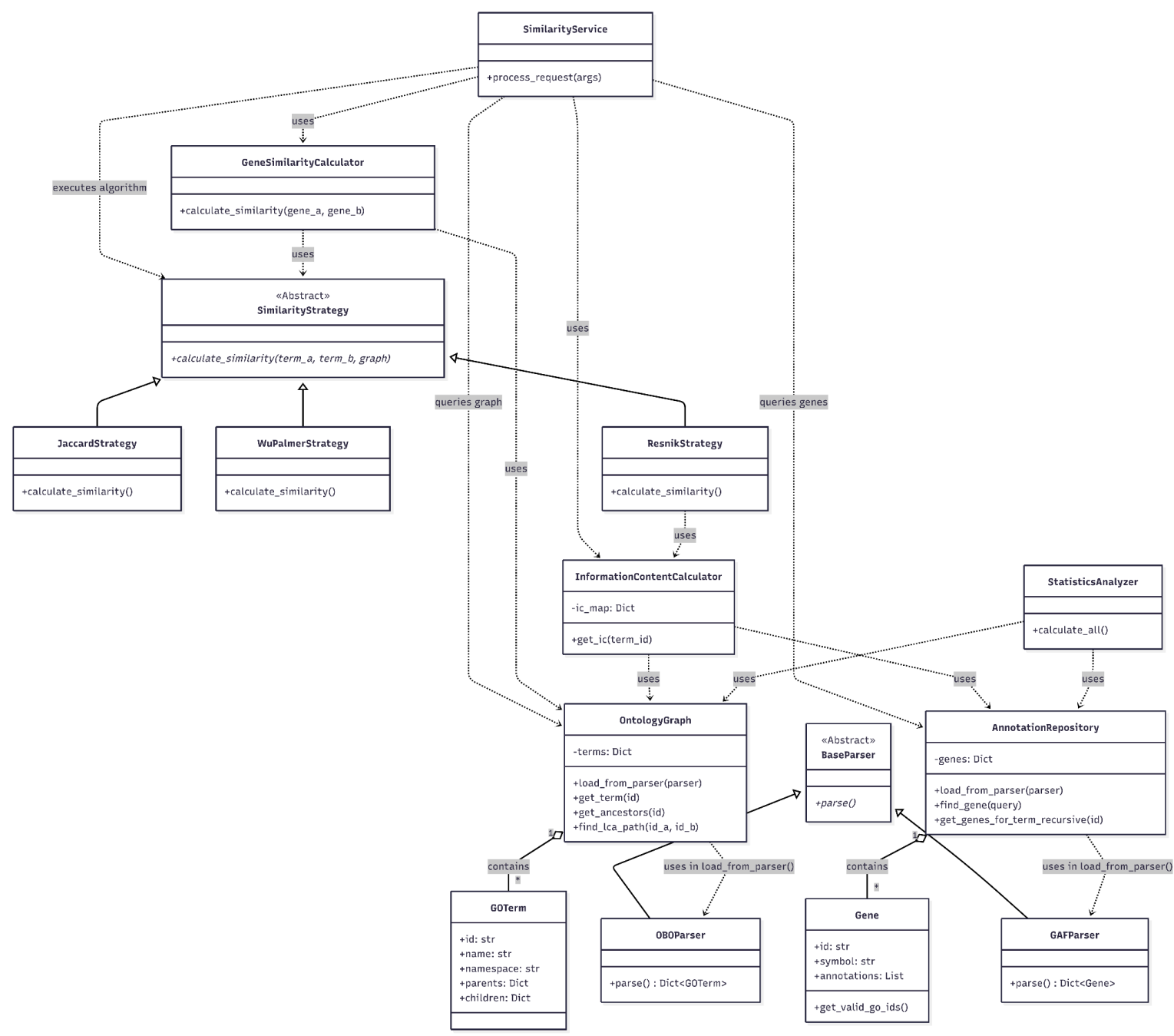
*statistics.py*: This module analyzes the loaded data to generate summary statistics, such as the distribution of annotations across the three GO aspects (Process, Function, Component), top annotated terms and many other metrics.

## 3. User Interface (The "Dashboard")

*main.py*: The central control script that runs the web server. It receives user inputs (like a search query or file upload) and asks the Data and Analysis layers for the results.

*templates*: These provide the visual interface as HTML files, displaying results as interactive tables, heatmaps, and path visualizations in the web browser.

# UML class diagram



# CRC Cards

## Core Data Classes

Class: <a href="#">Gene</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Store biological data (ID, Symbol, Name)		
Store a list of GO annotations		
Filter annotations (e.g., exclude "NOT" relations)		
Count total annotations		

Class: <a href="#">GOTerm</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Store ontology term details (ID, Name, Namespace)		
Maintain a list of parent and child relationships.		
Store metadata like definitions and synonyms.		

## Management & Logic Classes

Class: <a href="#">OntologyGraph</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Manage the entire graph of GOTerm objects	<a href="#">GOTerm</a>	
Load data using a parser	<a href="#">OBOParser</a>	
Calculate term depth and find ancestors or descendants		
Find paths between terms (LCA, Shortest Path)		

Class: <a href="#">AnnotationRepository</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Load data using a parser	<a href="#">Gene</a>	
Manage the collection of Gene objects	<a href="#">GAFParser</a>	
Perform gene searches by symbol or name.	<a href="#">OntologyGraph</a>	
Build a reverse index (Term ID → Genes)		

Class: <a href="#">SimilarityService</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Facade for all similarity operations	<a href="#">OntologyGraph</a>	
Process web requests and format results	<a href="#">AnnotationRepository</a>	
Initialize and manage different SimilarityStrategy instances.	<a href="#">SimilarityStrategy</a>	

## Analysis & Strategy Classes

Class: <a href="#">SimilarityStrategy (Abstract)</a>	Superclass:	Subclasses: <a href="#">JaccardStrategy</a> <a href="#">WuPalmerStrategy</a> <a href="#">ResnikStrategy</a>
Responsibilities	Collaborations	
Define the interface for calculating similarity between two terms	<a href="#">OntologyGraph</a>	
Enforce implementation of <a href="#">calculate_similarity</a>		

Class: <a href="#">JaccardStrategy</a>	Superclass: <a href="#">SimilarityStrategy</a>	Subclasses:
Responsibilities	Collaborations	
Calculate similarity using the Jaccard score	<a href="#">OntologyGraph</a>	
Identify intersection and union of ancestor sets		

Class: <a href="#">WuPalmerStrategy</a>	Superclass: <a href="#">SimilarityStrategy</a>	Subclasses:
Responsibilities	Collaborations	
Calculate similarity using the Wu-Palmer method.	<a href="#">OntologyGraph</a>	
Find the Lowest Common Ancestor (LCA)		
Compute term depths relative to the root		

Class: <a href="#">ResnikStrategy</a>	Superclass: <a href="#">SimilarityStrategy</a>	Subclasses:
Responsibilities	Collaborations	
Calculate semantic similarity using Information Content (IC)	<a href="#">OntologyGraph</a>	
Find the maximum IC among common ancestors	<a href="#">InformationContentCalculator</a>	
Provide a more specific, probability-based measure		

Class: <a href="#">GeneSimilarityCalculator</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Calculate similarity between two Genes	<a href="#">Gene</a>	
Compute the "Best Match Average" score	<a href="#">SimilarityStrategy</a>	
Generate N×N similarity matrices for lists of genes		

Class: <a href="#">InformationContentCalculator</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Calculate the probability of occurrence for every GO term	<a href="#">OntologyGraph</a>	
Compute Information Content (IC) values	<a href="#">AnnotationRepository</a>	
Provide IC lookup for the Resnik algorithm.		

## Utility Classes

Class: <a href="#">BaseParser</a>	Superclass:	Subclasses: <a href="#">OBOParser</a> <a href="#">GAFParse</a>
Responsibilities	Collaborations	
Enforce a common interface (parse method) for all file readers	<a href="#">GOTerm</a>	
Handle basic file validation (check if file path exists)		
Act as a parent class for OBOParser and GAFPaser		

Class: <a href="#">OBOParser</a>	Superclass: <a href="#">BaseParser</a>	Subclasses:
Responsibilities	Collaborations	
Open and read raw text files line-by-line	<a href="#">GOTerm</a>	
Parse OBO syntax format files		
Create GOTerm objects from text data		

Class: <a href="#">GAFParser</a>	Superclass: <a href="#">BaseParser</a>	Subclasses:
Responsibilities	Collaborations	
Open and read raw text files line-by-line	<a href="#">Gene</a>	
Parse GAF tabular format files		
Create Gene objects from text data		

Class: <a href="#">StatisticsAnalyzer</a>	Superclass:	Subclasses:
Responsibilities	Collaborations	
Aggregate data from the graph and repository	<a href="#">AnnotationRepository</a>	
Calculate global stats (e.g., Average Annotations per Gene)	<a href="#">OntologyGraph</a>	
Identify top-ranked genes and terms		

# OOP Principles

- **Abstraction**

We used the *BaseParser* and *SimilarityStrategy* abstract base classes. This hides the complex implementation details of file reading and mathematical formulas from the main application flow. Also the *OntologyGraph* class encapsulates the complex graph traversal logic (BFS algorithm), term depth calculation and finding term ancestors and descendants. External classes simply call *get\_ancestors()* without needing to know how the graph is traversed internally.

- **Encapsulation**

The *Gene* and *GOTerm* classes demonstrate encapsulation by hiding their internal data. Attributes like *\_\_id*, *\_\_name*, and *\_\_annotations* are made private, preventing direct external modification. Access is controlled through public methods and properties. For instance, the *Gene* class forces all annotation additions through the *add\_annotation()* method, which contains logic to handle duplicates and ensures data consistency.

- **Inheritance**

*JaccardStrategy* and *ResnikStrategy* inherit from the parent *SimilarityStrategy*. This ensures they both implement the required *calculate\_similarity* method, enforcing a strict contract. The same goes for *OBOParser* and *GAFParser* classes, which are inherited from the abstract class *BaseParser*.

- **Polymorphism**

The *GeneSimilarityCalculator* can accept any strategy object. This is possible because concrete classes like *JaccardStrategy*, *WuPalmerStrategy*, and *ResnikStrategy* all inherit from the abstract *SimilarityStrategy* and provide their own implementation of the same method. It calls *.calculate\_similarity()* on the object without knowing which specific algorithm is being used, allowing the algorithm to be swapped at runtime.

# Design Decisions

## 1. Modular Similarity Algorithms

In bioinformatics, there is no single "best" way to measure functional similarity. Some analyses benefit from simple graph overlaps (Jaccard), while others require statistical weighting and checking how rare terms are (Resnik). To address this, we designed the similarity engine to be modular.

Instead of using a single formula, the system treats each algorithm as an interchangeable component or "strategy." This allows researchers to switch between Jaccard, Wu-Palmer, and Resnik methods instantly to compare results. Furthermore, this design future-proofs the software. If a new similarity metric is developed in the scientific community, it can be easily used as a new class without needing to rewrite the core application logic.

## 2. Integration of Pandas and NumPy

*Pandas* is used in *statistics.py* and *repository.py* to handle large tabular data from GAF files. It allows for vectorized operations (like *groupby* and *value\_counts*), which are significantly faster than standard Python loops for generating reports and statistical summaries.

*NumPy* is used in *InformationContentCalculator* to perform element-wise logarithmic calculations on array data, which is essential for the Resnik similarity metric.

## 3. Caching & Optimization

To improve performance, the *OntologyGraph* implements a caching mechanism for term depth. Since calculating the depth of a node in a DAG is recursive and expensive, we store the result after the first calculation (*\_\_depth\_cache*), making subsequent lookups  $O(1)$ .

## **Conclusion**

This project successfully implements a comprehensive tool for Gene Ontology analysis. By strictly adhering to Object-Oriented principles like Encapsulation and Polymorphism, we created a system that is both robust and easy to extend. The modular architecture allows for the easy addition of new features, such as new file parsers or similarity metrics, without disrupting the core functionality. The integration of efficient data structures like Pandas DataFrames ensures the system can handle real-world biological datasets effectively.