

# Table of Contents

<b>1 Introduction</b>	<b>2</b>
<b>1.1. Abstract, Aims</b>	<b>2</b>
<b>2. Related works/ Literature review</b>	<b>3</b>
<b>3. Design</b>	<b>3</b>
3.1. Semantic Data Ontology:	3
3.2. Human Disease Ontology	5
<b>4. Implementation</b>	<b>14</b>
4.1. Identify the Domain, Range, and Competency Questions	14
4.2. The competency questions	15
4.3. Setup SPARQL Endpoint	16
4.4. Queries	16
4.5. Prefixes	16
4.6. Natural Language Processing, Parsing queries	16
<b>5. Evaluation and Use</b>	<b>23</b>
5.1. Reasoning Task	23
5.2. Queries	23
<b>6. Research Mapping/ Conclusion</b>	<b>28</b>
6.1. Critical Reflection	29
6.2. Future Scope	29
6.3. Compare with alternative approaches	29
6.4. Discussion of Limitations	29
<b>7. References</b>	<b>30</b>
<b>8. Appendix</b>	<b>31</b>

# 1 Introduction

This section sets the stages for the ontology application that can impact healthcare and medical research by creating human disease information. You can find this project on my Github page [https://github.com/sarashahin/Huma\\_Disease\\_Ontology](https://github.com/sarashahin/Huma_Disease_Ontology).

## 1.1. Abstract, Aims

The main human diseases for this ontology are chronic respiratory disease, infectious disease, and acute disease, and these diseases are causing death in the world. One of the problems is risk factors, which are associated with chronic diseases and correlated with people's lifestyles. Finding disease in the early stages could prevent some chronic disease deaths.

A possible solution for this is to recognise behaviours related to risk factors, such as drinking alcohol, smoking, and weight gain, which were mapped by the World Health Organisation[14].

The proposed ontology project can track risk factors associated with diseases. Based on the behaviour risk factors, the App can make good recommendations, for example, a healthy diet or quitting smoking.

However, this ontology app performs as a knowledge-based for information systems; the primary goal is to search for the behaviour related to risk factors for chronic, infectious, and acute diseases.

The method I applied for this ontology that has axioms and semantic principles to create queries and inferences, I defined five question queries that enabled inferences to evaluate the human disease ontology.

I have partially imported a dataset from the Bioportal Biotechnology website and classified it into the expected classes and properties. Additionally, I have used SPARQL queries to answer those competence questions for this project.

As a second aim of this project, I gathered and arranged this information about various diseases, causes, symptoms, and treatments into a simple method(App). This will also assist doctors and researchers to easily find and understand information about various diseases. By doing this I hope to make it easier for people who learn more about diseases, causes, risk factors, and treatments, help researchers find new treatments or cures, and make this useful in improving healthcare and medical research.

## 2. Related works/ Literature review

Chronic diseases are one of the most dangerous diseases in the world leading to death, among 7 of 10.

This finding emphasises is necessary to prevent and treat Chronic diseases, such as diabetes, chronic respiratory and heart diseases[1].

By early changing of a healthy lifestyle can lower risk of chronic diseases and minimise the death rates. Reducing tobacco and alcohol consumption, sport activity, and maintaining a balanced diet are a few examples of these adjustments. Thus, encouraging healthy habits need an understanding of human behaviour related to lifestyle choices[2].

Recent research has found that the pandemic of COVID-19 has had an impact on the people's lifestyle and healthcare, mainly with increased intake of alcohol, tobacco, and fast foods and also decreased sport activity. The recent research by Malta et al. showed that around 58% of young ages with noncommunicable chronic diseases and decreased physical activity while during the pandemic, increased the consumption of frozen food by 53% [7].

The paper describes how to create a consistent OWL ontology for vitamin A and its impact on humans by using Protege Tool. It also provides information about vitamin A effect on human health, risk factors affecting these impacts, and diet control. Additionally, it indicates a group of people at risk of Vitamin A deficiency[11].

The Human Disease Ontology significantly has grown with its user-base and disease content developing greatly since 2018. Use of Disease Ontology terms has been increased over five years, and now linked to more than 1.5 million biomedical data sources. The latest research updates added around 1793 disease terms and improved the structure of Disease ontology and disease classification with thousands of new links between diseases[12].

This paper illustrates a system using an Ontology-based approach that improves diagnosis in human disease as well as the research. The proposed Ontology system features Generic Human Disease Ontology, that contains broad information about organised human disease into four categories: 1- disease type, 2- symptoms, 3- causes including both genetics and environmental factors, 4- treatments such as medicine, chemotherapy and surgery. The mentioned Ontology is especially useful for studying complex disorders caused by multiple factors like manic depressive disorders. Moreover, this is the creation of specific Human disease Ontology to help physician and medical research[13].

## 3. Design

### 3.1. Semantic Data Ontology:

An ontology is a formal and shareable representation of a domain knowledge which is made up of well defined concepts and rules. A list of competency questions(SPQRL) can be used to specify the Ontology scope. Therefore, reasoners can help in the search for answers to questions [3].

Ontology, as illustrated via Ontology Web Language(OWL), is a standard language of the World Wide Web Consortium (W3C) (<https://www.w3.org>).

An OWL ontology is made up of four components: instances, concepts, properties, and restrictions.

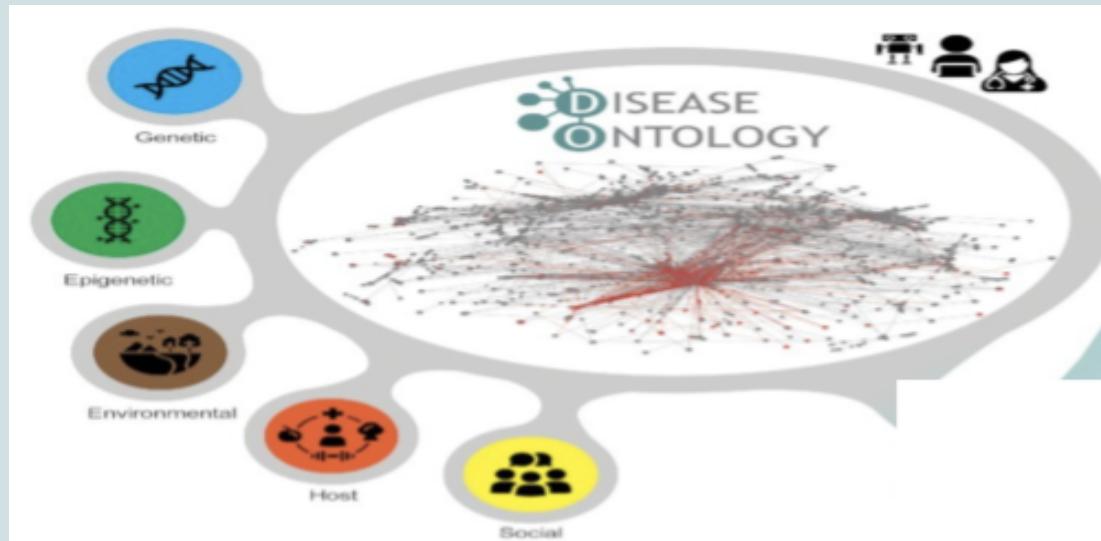
In the ontology definition, concepts are used to define the objects. The individuals of classes defined by instances. Individuals can be associated with other individuals among properties.

Constraints belong to a class where they are used to identify boundaries for the individuals.

SPARQL query language being used for an evaluation of the contents of an ontology. Inferences are related to classes and individuals can be used to identify Semantic Web Rule Language (SWRL)[\[4\]](#).

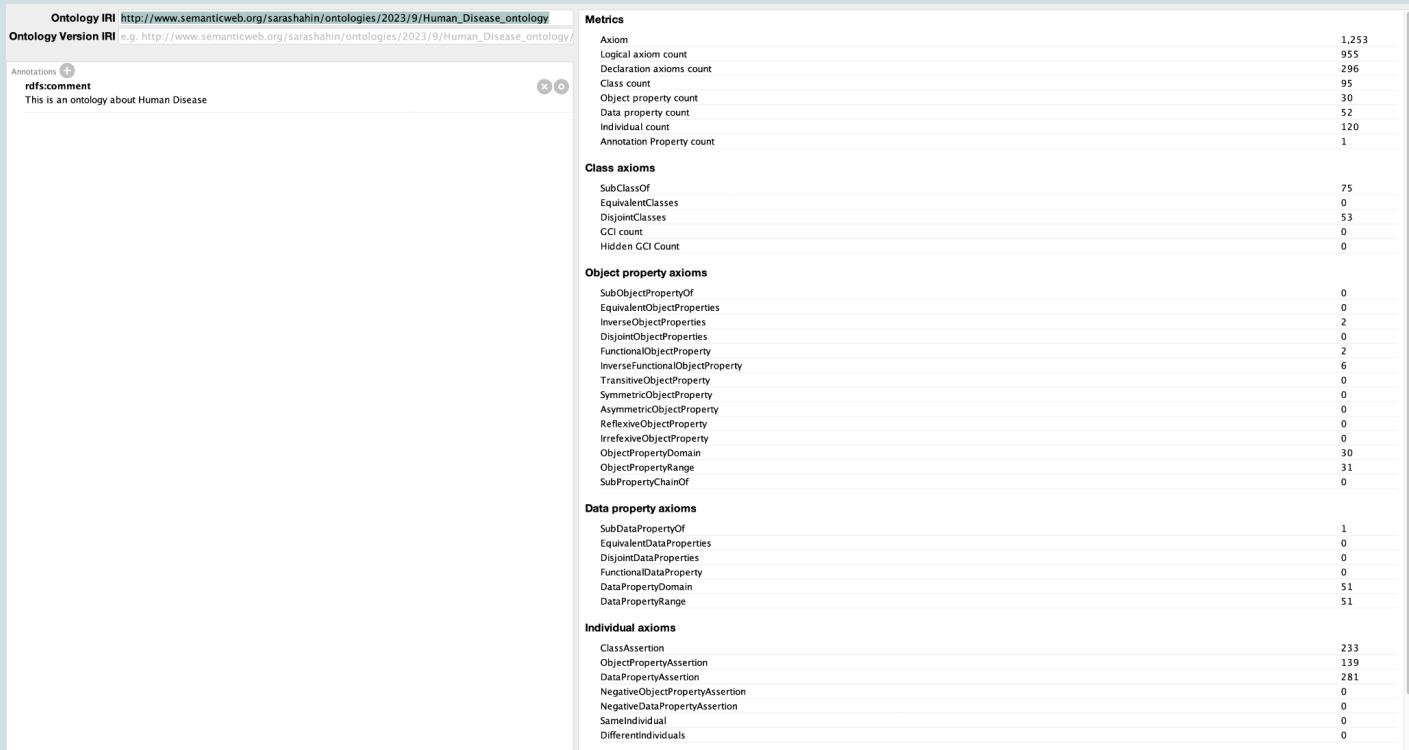
Properties are a binary relationship which is a one-to-one relation between two entities or among itself[\[6\]](#).

As a result, from construction methodologies, ontologies are used showing infer knowledge and domains, in various fields, for example the Web of Things (WoT) and Decision Support Systems (DSS), compared with others[\[5\]](#). Goal of this project to expand an ontology that represents the knowledge for tracking behaviour and providing educational support.



Graphical of Disease[\[8\]](#)

## 3.2. Human Disease Ontology



**Figure 1. Ontology metrics**

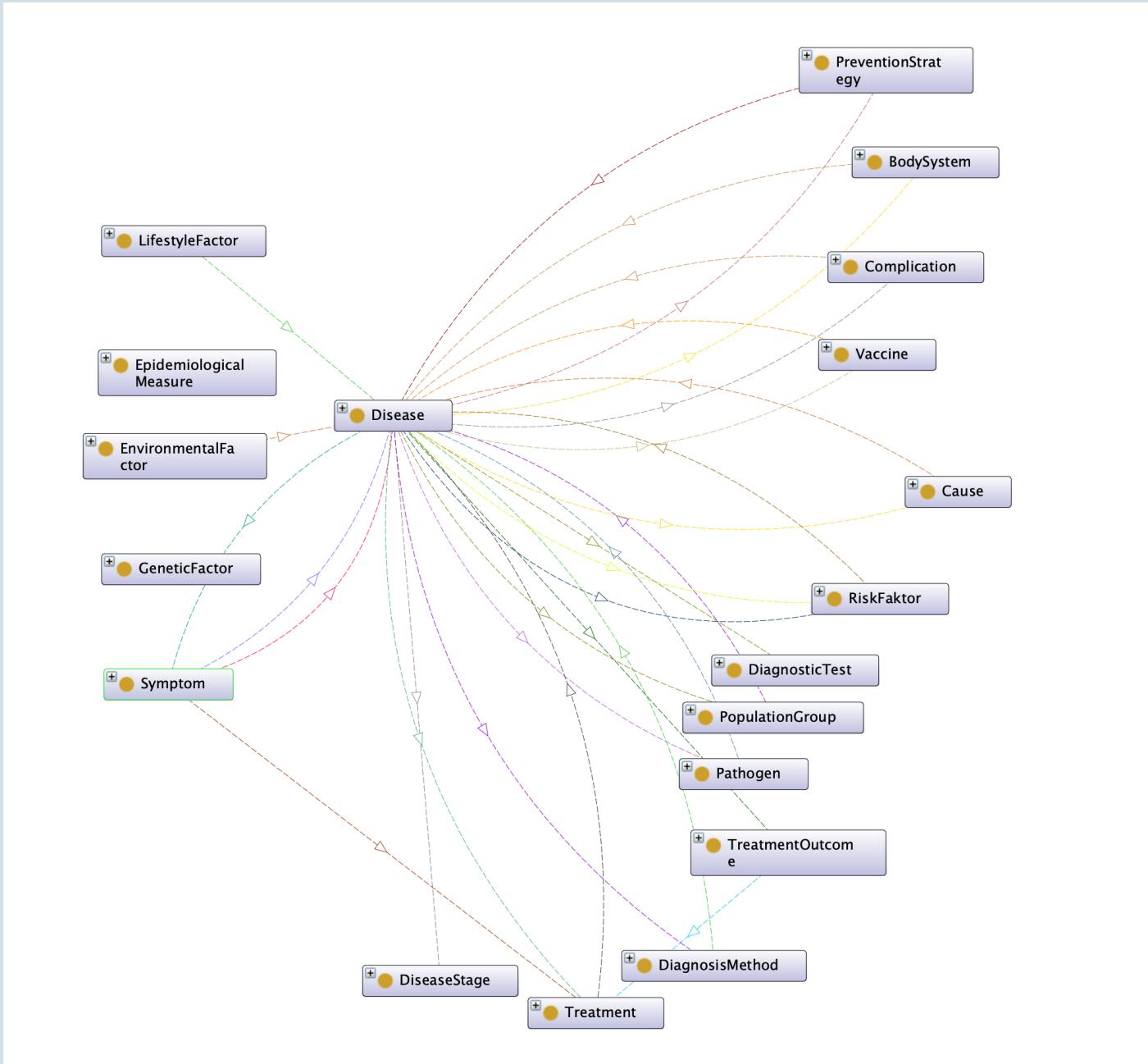
The ontology is implemented by Protege version 5.6.1 in Resource Description Framework (RDF) format. The current version of ontology includes 1253 Axiom, 95 classes, 30 Object properties, 52 Data properties, and 120 Individuals. Axioms are logical expressions that define a concept. The metrics of "Total classes" and "Total subclasses" show the number of classifications in the ontology. Object properties indicate the relationships between instances of two classes. Class and property partially imported from bioportal bio ontology website, initially, I collected data as RDF and imported into Protege then made some modifications(Figure 2)[19].

(<https://bioportal.bioontology.org/ontologies/DOID/?p=summary>)

The figure shows the Human Disease Ontology summary page. The top navigation bar includes tabs for Summary, Classes, Properties, and Notes. Below the tabs is a "Jump to:" input field. To the right is a sidebar containing a hierarchical tree view of ontology categories:

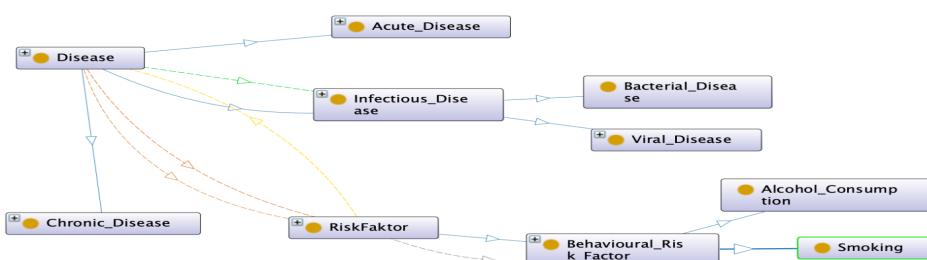
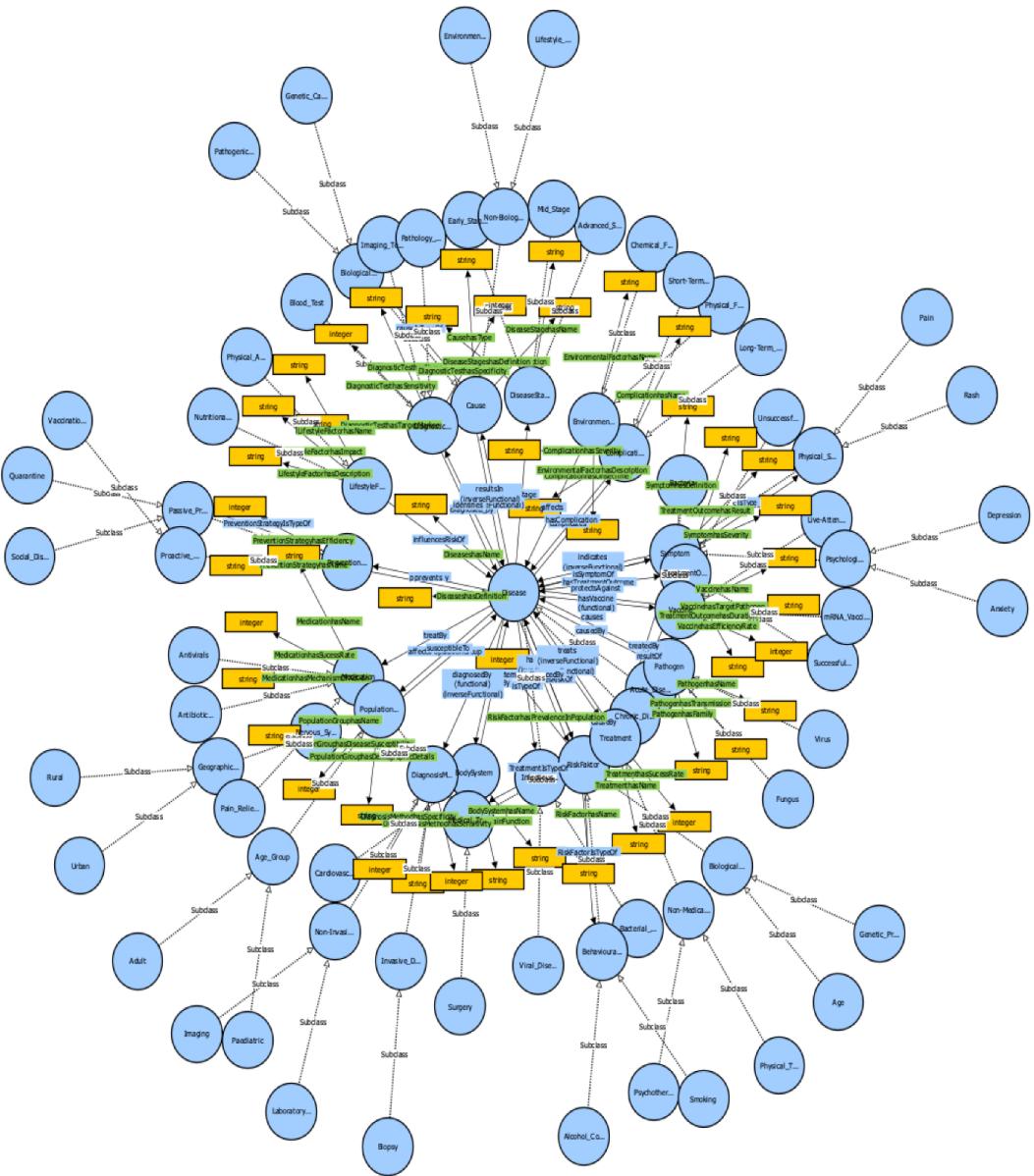
- anatomy
- cell
- chebi
- disease
- disease driver
- evidence
- food material
- inheritance pattern
- ncbitaxon
- omim\_susceptibility
- onset
- phenotype
- sequence
- symptom
- transmission process

**Figure 2. Importing data**



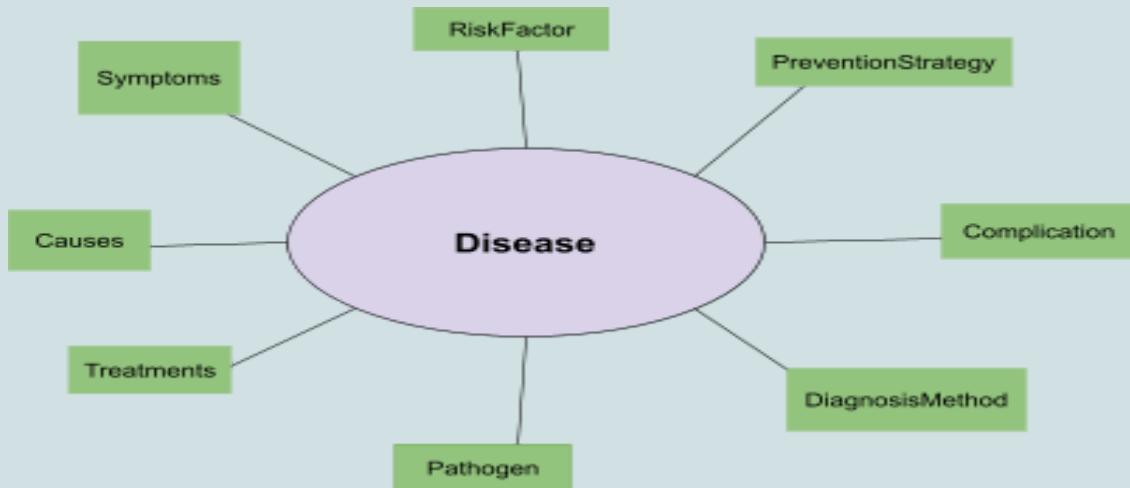
**Figure3. Hierarchical of ontology classes.**

The classes are described in the particular form, and each class can indicate single instance or multi instances. Additionally, this behaviour risk factor is related to chronic diseases. Table 1 describes the main classes.



**Table 1. Main classes of ontology.**

<b>Class</b>	<b>Description/ Subclasses</b>
Disease	Class consisting of subclasses such as Infectious Disease, Chronic Disease and, Acute Disease.
BodySystem	Class represents the body, like nerves or heart.
Cause	Class that includes subclasses like Biological Cause and Non-Biological Cause.
DiagnosisMethod	Classes contain subclasses such as Invasive Diagnosis and Non-Invasive Diagnosis.
DiagnosisTest	Class represents blood tests, such as image tests and pathology tests.
Pathogen	This class consists of bacteria, fungus and viruses.
PreventionStrategy	Class includes subclasses such as Proactive Prevention and Passive Prevention.
RiskFaktor	Class consists of Biological Risk Factor and Behavioural Risk Factor.
Symptom	Class contains subclasses like Physical Symptoms and Psychological Symptoms.
Treatment	Class includes subclasses such as Medical Treatment and Non-Medical Treatment.
LifestyleFactor	Class consists of nutritional choices and physical activity.
Complication	Class includes long term complication or short term complication.



### 3.2.1. Define Relationships and Class Properties:

There are five steps to identify the relationships and properties of the classes.

Figure 3. displays the structure of ontology classes in hierarchical yellow colour. The connections between these classes are as object property and marked in blue, while the attributes of each class are indicated as data property and marked in green. All of this detailed information is available in Protege.

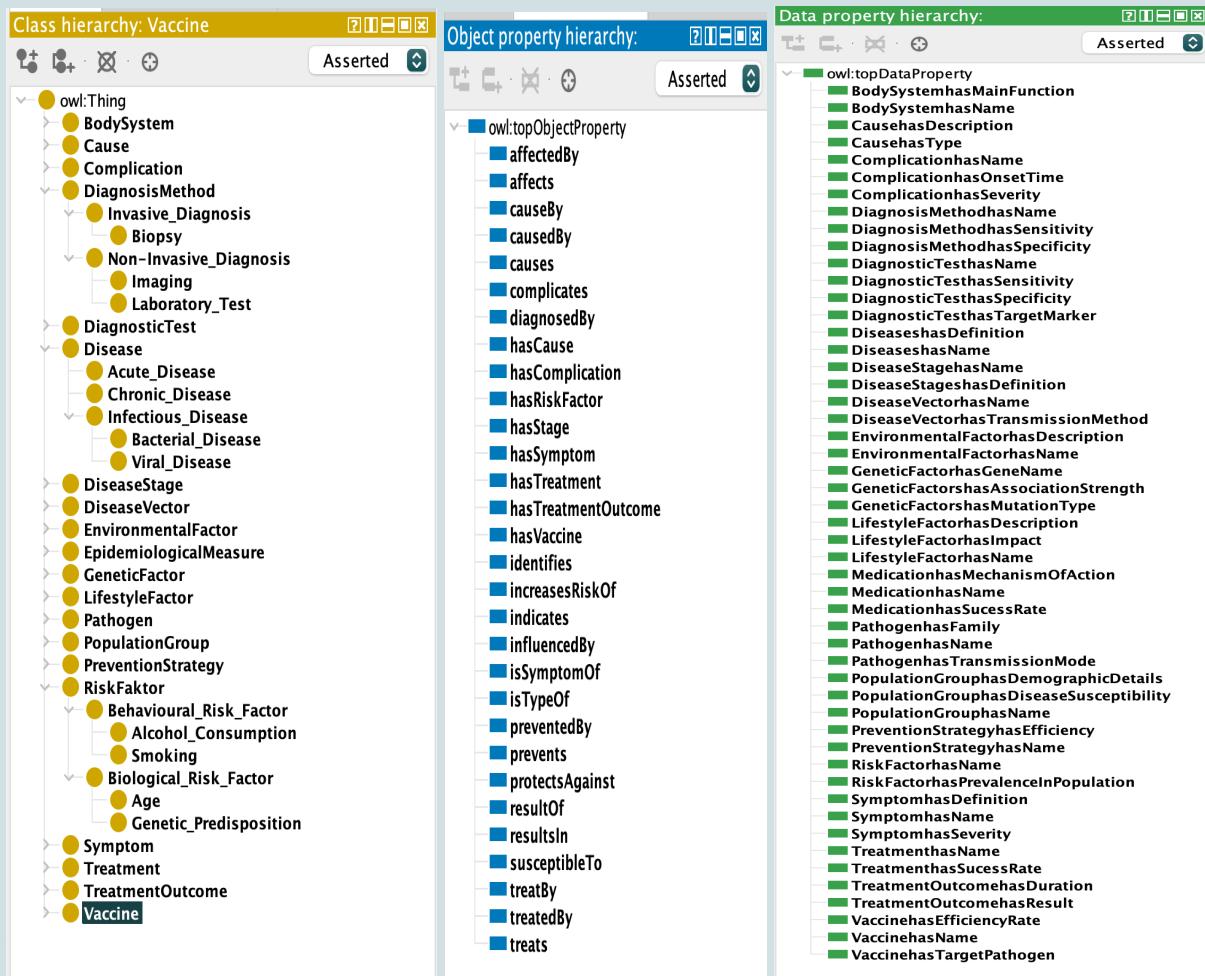
Based on my small project, I have decided to use this format for the data properties, they are clear and follow a pattern, making it easy to understand what each property represents, for example, **DiseasehasDefinition** or **SymptomhasName**.

Figure 3 displays the structured layers of the ontology's classes in yellow, illustrating a clear hierarchy. The connections between these classes are marked in blue (object properties), while the specific characteristics of each class are indicated in green (data properties). All of this detailed information can be accessed through the Protege software.

The screenshot shows the Protege interface with the following tabs open:

- Data properties**: Shows various properties like `diagnoseby`, `hasCause`, etc.
- Annotations: Disease**: Shows annotations for the Disease class.
- Class hierarchy: Disease**: Shows the inheritance chain from `owl:Thing` to `Disease`, listing subclasses like `Acute_Disease`, `Infectious_Disease`, `DiseaseStage`, etc.
- Object properties**: Shows object properties like `hasSymptom`, `hasTreatment`, etc.
- Characteristics** (with tabs for `Characteristic`, `Description`, and `Disjoint With`): Shows characteristics for the `hasSymptom` property, such as Functional, Inverse functional, Transitive, Symmetric, Asymmetric, Reflexive, and Irreflexive.
- Description: DiseasehasName**: Shows data properties for the `DiseasehasName` property, such as `xsd:string` for ranges and `Disease` for domains.
- Disjoint With**: Shows disjointness constraints for the `Disease` class.

The `Disease` class is highlighted in yellow across all tabs, indicating its central role in the ontology.



#### RiskFactorhasName(Figure 4.):

- This is a data property in the human disease ontology and assigned to a `RiskFaktor` entity.
- This represents that I have 30 individual risk factors, and each one of these with a name property set.
- Each risk factor is listed with a triangle which means as individuals.
- The data property `RiskFactorhasName` is showing that it is the type of `xsd:string`, and values are expected to be strings.
- This indicate that the domain of `RiskFactorhasName` is `RiskFaktor` class.

Annotations Usage

**Usage: RiskFactorhasName**

Show:  this  disjoints

Found 30 uses of RiskFactorhasName

- ✓ Aging
  - ◆ Aging RiskFactorhasName "Aging"
- ✓ AirPollution
  - ◆ AirPollution RiskFactorhasName "Air Pollution"
- ✓ Allergies
  - ◆ Allergies RiskFactorhasName "Allergies"
- ✓ CalciumDeficiency
  - ◆ CalciumDeficiency RiskFactorhasName "Calcium Deficiency"
- ✓ DrinkingAlcohol
  - ◆ DrinkingAlcohol RiskFactorhasName "Drinking Alcohol"
- ✓ FamilyHistoryofDiabetes
  - ◆ FamilyHistoryofDiabetes RiskFactorhasName "Family History of Diabetes"
- ✓ FamilyHistoryofGlaucoma
  - ◆ FamilyHistoryofGlaucoma RiskFactorhasName "Family History of Glaucoma"
- ✓ GeneticPredisposition
  - ◆ GeneticPredisposition RiskFactorhasName "Genetic Predisposition"
- ✓ LackofVitamin\_D
  - ◆ LackofVitamin\_D RiskFactorhasName "Lack of Vitamin D"
- ✓ Obesity
  - ◆ Obesity RiskFactorhasName "Obesity"
- ✓ PhysicalInactivity
  - ◆ PhysicalInactivity RiskFactorhasName "Physical Inactivity"
- ✓ RiskFactorhasName
  - ◆ RiskFactorhasName Range: xsd:string
  - ◆ DataProperty: RiskFactorhasName
  - ◆ RiskFactorhasName Domain RiskFaktor
- ✓ WeakenedImmuneSystem
  - ◆ WeakenedImmuneSystem RiskFactorhasName "Weakened Immune System"

**Description: RiskFactorhasName**

Equivalent To +

SubProperty Of +

Domains (intersection) +

**RiskFaktor**

Ranges +

xsd:string

Disjoint With +

Figure 4.

### hasRiskFactor property(Figure 5.):

- This is an object property in the human disease ontology. It is used to link diseases with their associated risk factors.
- There are 28 instances being used.
- For instance, Diabetes is related to FamilyHistoryofDiabetes, DrinkingAlcohol, and Obesity. This represents that having a family history of diabetes, alcohol consumption, and obesity are defined as risk factors for diabetes.
- The domain of hasRiskfactor is Disease, which means that this property is used to describe attributes of disease entities.
- Inverse Functional indicates that for any given risk factor, there can be only one disease related to it.
- The Inverse of shows that hasRiskFactor and increaseRiskOf are semantically related but in opposite directions such as if Diabetes hasRiskFactor Obesity, then Obesity increaseRiskOf Diabetes.

- The Range of hasRiskFactor is RiskFaktor, which indicates that the values of this property are expected to be individuals of the class RiskFaktor.

The screenshot shows the Protégé interface with two main panels. On the left, the 'Usage' tab is selected under 'Annotations'. It displays the 'hasRiskFactor' property with the following details:

- Show:  this  disjoins
- Found 28 uses of hasRiskFactor
  - Asthma
    - Asthma hasRiskFactor AirPollution
    - Asthma hasRiskFactor GeneticPredisposition
  - Diabetes
    - Diabetes hasRiskFactor FamilyHistoryofDiabetes
    - Diabetes hasRiskFactor DrinkingAlcohol
    - Diabetes hasRiskFactor Obesity
  - Glaucoma
    - Glaucoma hasRiskFactor FamilyHistoryofGlaucoma
  - hasRiskFactor**
    - hasRiskFactor Domain Disease
    - InverseFunctional: hasRiskFactor
    - InverseOf increasesRiskOf
    - Range RiskFaktor
    - ObjectProperty: hasRiskFactor
  - increasesRiskOf**
    - hasRiskFactor InverseOf increasesRiskOf
  - Influenza
    - Influenza hasRiskFactor WeakenedImmuneSystem
  - Osteoporosis
    - Osteoporosis hasRiskFactor CalciumDeficiency

Figure 5.

### DrinkingAlcohol(Figure 6.):

- Is listed as an individual instance of the Behavioural\_Risk\_Factor subclass, which means that drinking alcohol is considered a specific behaviour that can increase blood sugar.
- increaseRiskOf Diabetes, this object property(Assertions) links DrinkingAlcohol directed to Diabetes is suggesting that alcohol consumption is a behaviour that can increase the risk of developing diabetes.
- RiskFactorhasName 'Drinking Alcohol', string,: this is data property(Assertions), the name of the risk factor.
- RiskFactorhasPrevalenceInPopulation, 85, integer: this indicates the prevalence of drinking alcohol as a percentage.

The screenshot shows the Protégé interface with the 'Usage' tab selected under 'Annotations'. It displays the 'DrinkingAlcohol' individual with the following details:

- Show:  this  different
- Found 14 uses of DrinkingAlcohol
  - Diabetes
    - Diabetes hasRiskFactor DrinkingAlcohol
  - DrinkingAlcohol
    - DrinkingAlcohol RiskFactorhasName "Drinking Alcohol"
    - DrinkingAlcohol Type RiskFaktor
    - DrinkingAlcohol RiskFactorhasPrevalenceInPopulation 85
    - DrinkingAlcohol Type Behavioural\_Risk\_Factor
    - DrinkingAlcohol increasesRiskOf Diabetes
    - Individual: DrinkingAlcohol

On the right, the 'Description' and 'Property assertions' tabs are selected for the 'DrinkingAlcohol' individual.

**Description: DrinkingAlcohol**

- Types
  - Behavioural\_Risk\_Factor
  - RiskFaktor
- Same Individual As
- Different Individuals

**Property assertions: DrinkingAlcohol**

- Object property assertions
  - increasesRiskOf Diabetes
- Data property assertions
  - RiskFactorhasName "Drinking Alcohol"
  - RiskFactorhasPrevalenceInPopulation 85
- Negative object property assertions
- Negative data property assertions

Figure 6.

### treats (Figure 7.):

- This Object property has been used 18 times to connect treatments to the diseases they are intended to manage.
- Antibiotics are indicated as a treatment for both Influenza and LymeDisease.
- InsulinTherapy is used to treat Diabetes.
- Is characterised as Inverse Functional, showing that each disease is related to a unique treatment.
- The Domain of the treats property is Treatment, which describes attributes of the treatment entities.
- The Range of the ‘treats’ property is Disease, which indicates that the values of this property are expected to be instances of the disease class.

Annotations Usage

Usage: treats

Show:  this  disjoints

Found 18 uses of treats

- Antibiotics
  - Antibiotics treats Influenza
  - Antibiotics treats LymeDisease
- CalciumSupplements
  - CalciumSupplements treats Osteoporosis
- Inhaler
  - Inhaler treats Asthma
- InsulinTherapy
  - InsulinTherapy treats Diabetes

**Characteristics: treats**

Functional  
 Inverse functional  
 Transitive  
 Symmetric  
 Asymmetric  
 Reflexive  
 Irreflexive

**Description: treats**

Equivalent To

SubProperty Of

Inverse Of

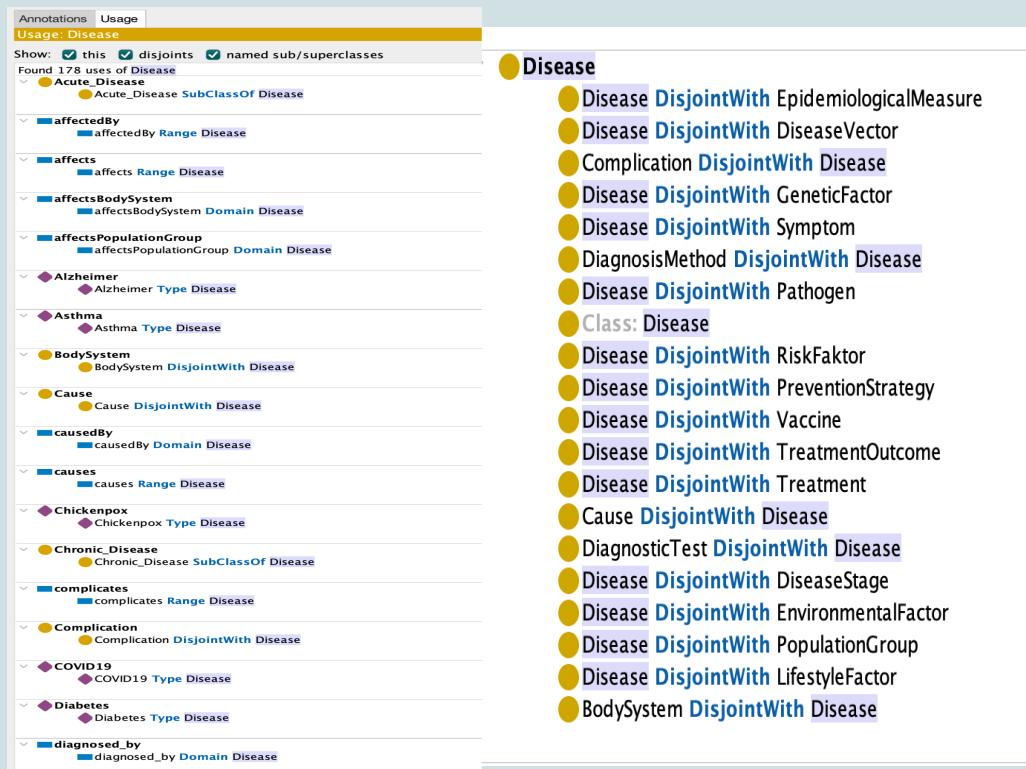
Domains (intersection)  Treatment

Ranges (intersection)  Disease

Disjoint With

SuperProperty Of (Chain)

Figure 7.



## Disease

- Disease DisjointWith EpidemiologicalMeasure
- Disease DisjointWith DiseaseVector
- Complication DisjointWith Disease
- Disease DisjointWith GeneticFactor
- Disease DisjointWith Symptom
- DiagnosisMethod DisjointWith Disease
- Disease DisjointWith Pathogen
- Class: Disease
- Disease DisjointWith RiskFaktor
- Disease DisjointWith PreventionStrategy
- Disease DisjointWith Vaccine
- Disease DisjointWith TreatmentOutcome
- Disease DisjointWith Treatment
- Cause DisjointWith Disease
- DiagnosticTest DisjointWith Disease
- Disease DisjointWith DiseaseStage
- Disease DisjointWith EnvironmentalFactor
- Disease DisjointWith PopulationGroup
- Disease DisjointWith LifestyleFactor
- BodySystem DisjointWith Disease

Figure(8)

## 4. Implementation

### 4.1. Identify the Domain, Range, and Competency Questions

I defined the structure of the knowledge domain of ontology which is for educational support.

I defined the key concepts such as Disease, Symptom, Treatment, RiskFactor, ... ect, and relations(predicate) between them like a Disease hasSymptom or indicates(Figure 9).

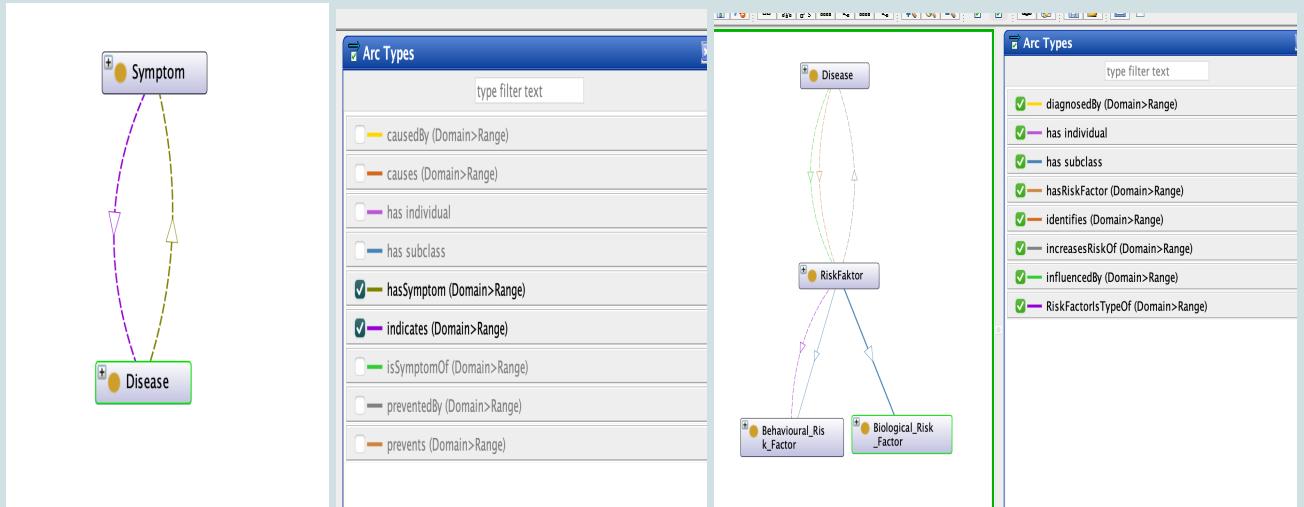
The scope of this ontology is identified, prevent and follow-up of chronic disease.

There are five competency questions that the ontology must answer.

These kinds of questions, apart from proving the ontology, also perform for the evaluation of the model.



(Figure 9)



Figure(10)

#### 4.2. The competency questions are:

- 1- What are the names and counts of symptoms for each disease which has more than 3 symptoms along with the names of those symptoms?
- 2- What are the names of the symptoms, diagnostic tests and treatments associated with LymeDisease?
- 3- list of the diseases and their definitions to handle diseases without a definition provided?
- 4- What are the names and counts of pathogens in each pathogen family along with a list of the names of these pathogens and associated with disease names?
- 5- What are the names and prevalence rates of risk factors for diabetes that the prevalence is more than > 15 percent with the names and efficiencies of prevention strategies used against diabetes?

**4.3. Setup SPARQL Endpoint:** I used a server such as Apache Jena Fuseki that is a SPARQL server, to host the SPARQL. I have loaded the RDF ontology into Fuseki and then allowed me to execute SPARQL queries from the data.

Figure(11)

**4.4. Queries:** With ensuring that the server is running and the RDF ontology loaded, then can run the SPARQL queries to get information from the RDF ontology[18].

<https://www.w3.org/TR/turtle/>

**4.5. Prefixes:** to execute the RDF queries in Apache Jena Fuseki environment , require to have these prefixes.

1. RDF, **PREFIX** `rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>`: is the prefix for RDF syntax namespace. It describes the basic structure of my data, such as Disease, Symptom and Treatment.
2. RDFs, **PREFIX** `rdfs: <http://www.w3.org/2000/01/rdf-schema#>`: is for RDF schema and a semantic extension of RDF.
3. XSD, **PREFIX** `xsd: <http://www.w3.org/2001/XMLSchema#>`: is for XML schema definition, to define data types in RDF like strings, dates, or decimals.
4. hd, **PREFIX** `hd:<http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human\_Disease\_ontology#>`: is a custom prefix which stands for the specific Human Disease Ontology namespace.

## 4.6. Natural Language Processing, Parsing queries:

This python script used spaCy library and a natural language processing library to retrieve information from an ontology about the Human\_Disease by interpreting user queries and formulated a SPARQL query Figure(12). Explain each part individually.

```

[1]: from SPARQLWrapper import SPARQLWrapper, JSON
import spacy
from spacy.matcher import Matcher

# load the spaCy model for english language
nlp = spacy.load("en_core_web_sm")

# initialize a Matcher with the shared vocab
matcher = Matcher(nlp.vocab)

# define patterns for the Matcher
patternsOntology = {
    "risk factor": [{"LOWER": "risk"}, {"LOWER": "factor"}, {"ENT_TYPE": "DISEASE", "OP": "?"}],
    "diagnosis method": [{"LOWER": "diagnosis"}, {"LOWER": "method"}, {"ENT_TYPE": "DISEASE", "OP": "?"}],
    "symptom": [{"LOWER": "symptom"}],
    "treatment": [{"LOWER": "treatment"}],
    "cause": [{"LOWER": "cause"}]
}
for key, pattern in patternsOntology.items():
    matcher.add(key, pattern)

# process text with spaCy and Matcher
def processText(text):
    doc = nlp(text)
    matches = matcher(doc)
    foundTerms = set()
    specificDisease = None
    for matchId, start, end in matches:
        span = doc[start:end]
        # check if span contains a disease entity
        for token in span:
            if token.ent_type_ == "DISEASE":
                specificDisease = token.text
                break
        foundTerms.add(span.text.lower())
    return foundTerms, specificDisease

# create a dynamic SPARQL query
def createSparqlQuery(foundTerms, specificDisease):
    baseQuery = """
PREFIX rdfs <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdf <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hdi <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>
"""

    termMapping = {
        "disease": "hdiDisease",
        "symptom": "hdiSymptom",
        "treatment": "hdiTreatment",
        "cause": "hdiCause",
        "risk factor": "hdiRiskFaktor",
        "diagnosis method": "hdiDiagnosisMethod"
    }

    queryConditions = []
    for term in foundTerms:
        if term in termMapping:
            queryConditions.append(f"?item rdf:type {termMapping[term]}")
    if specificDisease:
        queryConditions.append(f"?item hdi:isRelatedTo '{specificDisease}'")
    if not queryConditions: # default condition if no specific terms are identified
        queryConditions.append(f"?item rdf:type hdiDisease")
    whereClause = " ".join(queryConditions)
    query = f"(baseQuery)SELECT DISTINCT ?item WHERE {{ {whereClause} }} LIMIT 10"
    return query

# run SPARQL query and return results
def executeSparqlQuery(query):
    SPARQL_WRAPPER = ("http://localhost:3030/Human-Disease-Ontology/sparql")
    sparql.setQuery(query)
    sparql.setReturnFormat(JSON)
    try:
        results = sparql.query().convert()
        return results
    except Exception as e:
        print(f"An error occurred: {e}")
        return None

# main function to handle user queries
def main():
    userQuery = input("Enter your query in natural language: ")
    foundTerms, specificDisease = processText(userQuery)
    sparqlQuery = createSparqlQuery(foundTerms, specificDisease)
    print(f"\nSPARQL query to execute:\n{sparqlQuery}")
    results = executeSparqlQuery(sparqlQuery)
    if results:
        results["results"]["bindings"]:
            print("Results found:")
            for result in results["results"]["bindings"]:
                print(result["item"]["value"])
    else:
        print("No results returned or an error occurred.")

if __name__ == "__main__":
    main()

# What are the risk factor for Asthma?
# Give me diagnosis method for Diabetes
# What are some diseases?
# disease has names?
# List all symptoms
# List diagnosis method
# Show treatment.
# Give me cause
# what are the risk factor?
# Give me diagnosis method

```

## Figure(12)

## 1. Importing necessary libraries:

```
from SPARQLWrapper import SPARQLWrapper, JSON
import spacy
from spacy.matcher import Matcher
```

Figure(13)

**SPARQLWrapper:** this is for querying SPARQL endpoints from Python wrapper. It is used for databases to store the data in RDF format such as Human Disease Ontology.

### **Spacy and Matcher:**

spaCy is used for Natural Language Processing(NLP). Matcher is one of the spaCy tools, it is used for matching phrases or word sequences such as Treatment, Symptom and Disease, and is more suited for linguistic context. It acts like a Bag of Words but more advanced than a simple Bag of Words. It allows me to detect multi-words phrases. For example, it checks each word in a sentence if it matches with a pattern I have identified, such as finding a word that is a noun followed by a verb. It is good for grammatical structure. In summary, Matcher tool performs as a filter that it is looking for the patterns which I have defined, can contain base form of words like lemmas or kind of phrase, including named entities. [9]  
<https://spacy.io/api/matcher>.

```
# load the spaCy model for english language
nlp = spacy.load("en_core_web_sm")

# initialize a Matcher with the shared vocab
matcher = Matcher(nlp.vocab)

# define patterns for the Matcher
patternsOntology = {
    "risk factor": [{"LOWER": "risk"}, {"LOWER": "factor"}, {"ENT_TYPE": "DISEASE", "OP": "?"}],
    "diagnosis method": [{"LOWER": "diagnosis"}, {"LOWER": "method"}, {"ENT_TYPE": "DISEASE", "OP": "?"}],
    "disease": [{"LOWER": "disease"}],
    "symptom": [{"LOWER": "symptom"}],
    "treatment": [{"LOWER": "treatment"}],
    "cause": [{"LOWER": "cause"}]
}
for key, pattern in patternsOntology.items():
    matcher.add(key, pattern)
```

Figure(14)

### **2- Implement Matcher:**

First of all, I have loaded a small english model that includes data and algorithms for english text processing.

After that I have initialised Matcher with the vocabulary of the loaded spaCy model.

The reason I have chosen Matcher over regex, using Matcher for text analysing, understanding the relation between words in a phrase and combining named entity recognition[17].  
<https://spacy.io/api/phrasematcher>, <https://spacy.io/api/matcher>

Then, I have defined patterns to match specific words such as 'risk factor' or 'symptom' and these patterns are added to the Matcher tool. In summary, I have setted up the rules in this code for the Matcher to identify specific phrases.

```

# process text with spaCy and Matcher
def processText(text):
    doc = nlp(text)
    matches = matcher(doc)

    foundTerms = set()
    specificDisease = None
    for matchId, start, end in matches:
        span = doc[start:end]
        # check if the span contains a disease entity
        for token in span:
            if token.ent_type_ == "DISEASE":
                specificDisease = token.text
                break
        foundTerms.add(span.text.lower())

    return foundTerms, specificDisease

```

Figure(15)

This function, processText takes text as input to find specific patterns such as disease name or symptom name and defines a specific disease as mentioned. It has used spaCy library for Natural Language Processing and used Matcher tool to find the pattern which matches in the text.

```

# create a dynamic SPARQL query
def createSparql_query(foundTerms, specificDisease):
    baseQuery = """
        PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
        PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
        PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
    """

    termMapping = {
        "disease": "hd:Disease",
        "symptom": "hd:Symptom",
        "treatment": "hd:Treatment",
        "cause": "hd:Cause",
        "risk factor": "hd:RiskFaktor",
        "diagnosis method": "hd:DiagnosisMethod"
    }

    queryConditions = []
    for term in foundTerms:
        if term in termMapping:
            queryConditions.append(f"?item rdf:type {termMapping[term]}")

    if specificDisease:
        queryConditions.append(f'?item hd:diagnosedBy "{specificDisease}"')

    if not queryConditions: # default condition if no specific terms are identified
        queryConditions.append("?item rdf:type hd:Disease")

    whereClause = " . ".join(queryConditions)
    query = f'{baseQuery}SELECT DISTINCT ?item WHERE {{ {whereClause} }} LIMIT 10'

    return query

```

Figure(16)

The creatSparql\_Query function has created dynamically a SPARQL query to retrieve data from Human\_Disease\_Ontology based on user input. termMapping is a dictionary that translates Natural Language terms into the SPARQL endpoint that can be understood like mapping disease.

```

# run SPARQL query and return results
def executeSparql_query(query):
    sparql = SPARQLWrapper("http://localhost:3030/Human-Disease-Ontology/sparql")
    sparql.setQuery(query)
    sparql.setReturnFormat(JSON)
    try:
        results = sparql.query().convert()
        return results
    except Exception as e:
        print(f"An error occurred: {e}")
        return None

```

Figure(17)

Function excuteSparql\_query sends a SPARQL query to the endpoint, fetching the results in JSON format. During the process, handles any potential errors.

```

# main function to handle user queries
def main():
    userQuery = input("Enter your query in natural language: ")
    foundTerms, specificDisease = processText(userQuery)
    sparqlQuery = createSparql_query(foundTerms, specificDisease)
    print(f"SPARQL query to execute:\n{sparqlQuery}")
    results = executeSparql_query(sparqlQuery)
    if results and results["results"]["bindings"]:
        print("Results found:")
        for result in results["results"]["bindings"]:
            print(result["item"]["value"])
    else:
        print("No results returned or an error occurred.")

if __name__ == "__main__":
    main()

```

Figure(18)

The centre of this script of code is the main function. It takes a natural language (NLP) query from the user input, by processing it, extracting the related terms and diseases, then creates and runs a SPARQL query based on this information, and at the end displays the results.

### **The question I have provided for Natural Language query:**

1. Give me a diagnosis method?
2. What are the risk factors?
3. Give me causes ?
4. Show treatment?
5. List all symptoms?
6. What are some diseases?

```

# What are the risk factor for Asthma?
# Give me diagnosis method for Diabetes
# What are some diseases?
# disease has names?
# List all symptom.
# List diagnosis method
# Show treatment.
# Give me cause.
# What are the risk factor?
# Give me diagnosis method

```

Figure(19)

### **Results of this queries:**

```
# Give me diagnosis method

Enter your query in natural language: Give me diagnosis method
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:DiagnosisMethod } LIMIT 10

Results found:
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#CognitiveAssessment
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BloodGlucoseTest
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BloodSugarTest
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BloodTest
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BoneDensityScan
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BoneDensityTest
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#EyeExam
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Spirometry
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#TuberculinSkinTest
```

Figure(20)

Results of the query is a list of diagnosis methods.

```
Enter your query in natural language: What are the risk factor?
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:RiskFaktor } LIMIT 10

Results found:
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Aging
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#AirPollution
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Allergies
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#GeneticPredisposition
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#CalciumDeficiency
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#DrinkingAlcohol
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#FamilyHistoryofDiabetes
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Obesity
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#FamilyHistoryofGlaucoma
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#WeakenedImmuneSystem
```

Figure(21)

Results of the query is a list of risk factors.

```
Enter your query in natural language: Give me cause
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:Cause } LIMIT 10

Results found:
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#InsulinResistance
```

Figure(22)

Results of the query is a list of causes.

```

Enter your query in natural language: Show treatment
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:Treatment } LIMIT 10
Results found:
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Antibiotics
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Acetaminophen
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#CognitiveTherapy
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#AntimalarialDrugs
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Antipyretics
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#AntituberculosisAntibiotics
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Inhaler
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BoneDensityMedication
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#CalciumSupplements
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Rest

```

Figure(23)

Results of the query is a list of treatments.

```

Enter your query in natural language: List all symptom
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:Symptom } LIMIT 10
Results found:
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#MemoryLoss
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Cough
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#ShortnessofBreath
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Hissing
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BlurryVision
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#BodyAches
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Chills
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Fatigue
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Fever
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Headache

```

Figure(24)

Results of the query is a list of symptoms.

```

Enter your query in natural language: What are some diseases?
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:Disease } LIMIT 10
Results found:
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Diabetes
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Alzheimer
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Osteoporosis
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Asthma
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Influenza
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#LymeDisease
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#COVID19
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Heart
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Chickenpox
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Malaria

```

Figure(25)

Results of the query is a list of disease names.

```

# What are the risk factor for Asthma?
# Give me diagnosis method for Diabetes
# What are some diseases?
# disease has names?
# List all symptom.
# List diagnosis method
# Show treatment.
# Give me cause.
# What are the risk factor?
# Give me diagnosis method

Enter your query in natural language: What are some diseases?
userQuery: What are some diseases?
foundTerms : set()
whereClause: ?item rdf:type hd:Disease
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:Disease } LIMIT 10
results: {'head': {'vars': ['?item']}, 'results': {'bindings': [{'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Diabetes'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Alzheimer'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Osteoporosis'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Asthma'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Influenza'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#LymeDisease'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#COVID19'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Heart'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Chickenpox'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Malaria'}}]}
Results found:
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Diabetes
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Alzheimer
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Osteoporosis
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Asthma
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Influenza
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#LymeDisease
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#COVID19
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Heart
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Chickenpox
http://www.semanticweb.org/sarashahn/ontologies/2023/9/Human_Disease_ontology#Malaria'

```

Figure(32)

More detailed.

## 5. Evaluation and Use

### 5.1. Reasoning Task

First, I have validated the ontology model by using an automated reasoning process which was installed in Protege [10].

The reasoner used the ontology's class, object property, data property hierarchies, class assertions, and object property assertions to find all inferences that help to validate the ontology's structure and content.

The screenshot shows the Protege interface with two explanations for the 'hasRiskFactor' property:

- Explanation 1:** Shows the assertion "Diabetes hasRiskFactor DrinkingAlcohol". It includes options to "Show regular justifications" (radio button selected), "Show laconic justifications" (checkbox), and "Limit justifications to" (dropdown menu set to 2). Below the assertion, it says "In NO other justifications" with a question mark icon.
- Explanation 2:** Shows the assertions "hasRiskFactor inverseOf increasesRiskOf" and "DrinkingAlcohol increasesRiskOf Diabetes". It includes options to "Show regular justifications" (radio button selected), "Show laconic justifications" (checkbox), and "Limit justifications to" (dropdown menu set to 2). Below the assertions, it says "In NO other justifications" with a question mark icon.

The screenshot shows the Protege interface with two explanations for the 'hasSymptom' property:

- Explanation 1:** Shows the assertion "Hunger hasSymptom Diabetes". It includes options to "Show regular justifications" (radio button selected), "Show laconic justifications" (checkbox), and "Limit justifications to" (dropdown menu set to 2). Below the assertion, it says "In NO other justifications" with a question mark icon.
- Explanation 2:** Shows the assertion "hasSymptom inverseOf isSymptomOf". It includes options to "Show regular justifications" (radio button selected), "Show laconic justifications" (checkbox), and "Limit justifications to" (dropdown menu set to 2). Below the assertion, it says "In NO other justifications" with a question mark icon.

Figure(26)

**5.2. Queries:** This section is about the SPARQL queries utilised by the application that validate my ontology [16]. Questions have been mentioned in section 4.2.

<https://www.w3.org/TR/2012/PR-sparql11-query-20121108/>,  
[https://www.openresearch.org/wiki/Sparql\\_endpoint/Examples](https://www.openresearch.org/wiki/Sparql_endpoint/Examples)

Q1:

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX hd: <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>
```

```
SELECT ?diseaseName (GROUP_CONCAT(DISTINCT ?symptomName; separator=", ") AS ?symptomNames) (COUNT(DISTINCT ?symptom) AS ?symptomCount)
WHERE {
    ?disease rdf:type hd:Disease ;
        hd:hasSymptom ?symptom ;
        hd:DiseaseshasName ?diseaseName .
    ?symptom hd:SymptomhasName ?symptomName .
}
GROUP BY ?diseaseName
HAVING (COUNT(DISTINCT ?symptom) > 3)
```

/Human-Disease-Ontology/sparql		JSON	Turtle															
<pre>1 v PREFIX rdf: &lt;http://www.w3.org/1999/02/22-rdf-syntax-ns#&gt; 2 v PREFIX rdfs: &lt;http://www.w3.org/2000/01/rdf-schema#&gt; 3 v PREFIX xsd: &lt;http://www.w3.org/2001/XMLSchema#&gt; 4 v PREFIX hd: &lt;http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#&gt; 5 v 6 v #***What are the names and counts of symptoms for each disease that has more than three distinct symptoms with the names of those symptoms? 7 v SELECT ?diseaseName (GROUP_CONCAT(DISTINCT ?symptomName; separator=", ") AS ?symptomNames) (COUNT(DISTINCT ?symptom) AS ?symptomCount) 8 v WHERE { 9 v     ?disease rdf:type hd:Disease ; 10 v         hd:hasSymptom ?symptom ; 11 v         hd:DiseaseshasName ?diseaseName . 12 v     ?symptom hd:SymptomhasName ?symptomName . 13 v } 14 v GROUP BY ?diseaseName 15 v HAVING (COUNT(DISTINCT ?symptom) &gt; 3) 16 v }</pre>																		
<b>Table</b> Response 4 results in 0.056 seconds																		
<table border="1"> <thead> <tr> <th>diseaseName</th> <th>symptomNames</th> <th>symptomCount</th> </tr> </thead> <tbody> <tr> <td>1 Diabetes</td> <td>Blurry vision, Frequent Urination, Increased Thirst, Itchy skin</td> <td>"4"^^&lt;http://www.w3.org/2001/XMLSchema#integer&gt;</td> </tr> <tr> <td>2 Influenza</td> <td>Cough, Body aches, Fever, Headache</td> <td>"4"^^&lt;http://www.w3.org/2001/XMLSchema#integer&gt;</td> </tr> <tr> <td>3 COVID-19</td> <td>Cough, Body aches, Chills, Fatigue, Fever, Headache</td> <td>"6"^^&lt;http://www.w3.org/2001/XMLSchema#integer&gt;</td> </tr> <tr> <td>4 Lyme Disease</td> <td>Chills, Fatigue, Headache, Joint Pain</td> <td>"4"^^&lt;http://www.w3.org/2001/XMLSchema#integer&gt;</td> </tr> </tbody> </table>				diseaseName	symptomNames	symptomCount	1 Diabetes	Blurry vision, Frequent Urination, Increased Thirst, Itchy skin	"4"^^<http://www.w3.org/2001/XMLSchema#integer>	2 Influenza	Cough, Body aches, Fever, Headache	"4"^^<http://www.w3.org/2001/XMLSchema#integer>	3 COVID-19	Cough, Body aches, Chills, Fatigue, Fever, Headache	"6"^^<http://www.w3.org/2001/XMLSchema#integer>	4 Lyme Disease	Chills, Fatigue, Headache, Joint Pain	"4"^^<http://www.w3.org/2001/XMLSchema#integer>
diseaseName	symptomNames	symptomCount																
1 Diabetes	Blurry vision, Frequent Urination, Increased Thirst, Itchy skin	"4"^^<http://www.w3.org/2001/XMLSchema#integer>																
2 Influenza	Cough, Body aches, Fever, Headache	"4"^^<http://www.w3.org/2001/XMLSchema#integer>																
3 COVID-19	Cough, Body aches, Chills, Fatigue, Fever, Headache	"6"^^<http://www.w3.org/2001/XMLSchema#integer>																
4 Lyme Disease	Chills, Fatigue, Headache, Joint Pain	"4"^^<http://www.w3.org/2001/XMLSchema#integer>																
Showing 1 to 4 of 4 entries																		

Figure 27. is designed to have a list of each disease that has more than 3 distinct symptoms.

The results show that the disease's name(groupby), a concatenated string of all the distinct symptom names separated by commas, and total count for each symptom.

Q2:

```

SELECT (STR(?diseaseLabel) AS ?diseaseName)
       (GROUP_CONCAT(DISTINCT ?symptomName; separator=", ") AS ?symptomNames)
       (GROUP_CONCAT(DISTINCT ?treatmentName; separator=", ") AS ?treatmentNames)
       (GROUP_CONCAT(DISTINCT ?diagnosticTestName; separator=", ") AS
?diagnosticTestNames)
WHERE {
  ?disease hd:DiseaseshasName ?diseaseLabel .
  FILTER (STR(?diseaseLabel) = "Lyme Disease")
  ?disease hd:hasSymptom ?symptom .
  ?symptom hd:SymptomhasName ?symptomName .
  ?disease hd:hasTreatment ?treatment .
  ?treatment hd:TreatmenthasName ?treatmentName .
  ?disease hd:diagnosed_by ?diagnosticTest .
  ?diagnosticTest hd:DiagnosticTesthasName ?diagnosticTestName .
}
GROUP BY ?diseaseLabel
  
```

The screenshot shows a SPARQL query interface with the following details:

- Query URL:** /Human-Disease-Ontology/query
- Prefixes:**

```

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX hd: <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>
  
```
- Query Content:**

```

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX hd: <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>

****What are the names of all symptoms, diagnostic test and treatments associated with Lyme Disease?
SELECT (STR(?diseaseLabel) AS ?diseaseName)
       (GROUP_CONCAT(DISTINCT ?symptomName; separator=", ") AS ?symptomNames)
       (GROUP_CONCAT(DISTINCT ?treatmentName; separator=", ") AS ?treatmentNames)
       (GROUP_CONCAT(DISTINCT ?diagnosticTestName; separator=", ") AS ?diagnosticTestNames)
WHERE {
  ?disease hd:DiseaseshasName ?diseaseLabel .
  FILTER (STR(?diseaseLabel) = "Lyme Disease")
  ?disease hd:hasSymptom ?symptom .
  ?symptom hd:SymptomhasName ?symptomName .
  ?disease hd:hasTreatment ?treatment .
  ?treatment hd:TreatmenthasName ?treatmentName .
  ?disease hd:diagnosed_by ?diagnosticTest .
  ?diagnosticTest hd:DiagnosticTesthasName ?diagnosticTestName .
}
GROUP BY ?diseaseLabel
  
```
- Results:**

diseaseName	symptomNames	treatmentNames	diagnosticTestNames
Lyme Disease	Chills, Fatigue, Headache, Joint Pain	Antibiotics	ELISA test

Showing 1 to 1 of 1 entries

Figure 28. This query returns the name of LymeDisease along with a concatenated string of all related to symptoms, diagnostic test and another concatenated string related treatments.

Q3:

```

SELECT ?diseaseName (COALESCE(?definition, "Definition not provided") AS ?diseaseDefinition)
WHERE {
  ?disease hd:DiseaseshasName ?diseaseName .
  
```

```
OPTIONAL { ?disease hd:DiseaseshasDefinition ?definition . }
```

```
}
```

```
LIMIT 5
```

The screenshot shows a SPARQL query interface with the following details:

- Query URL:** /Human-Disease-Ontology/query
- Format:** JSON (selected from a dropdown menu)
- Results View:** Turtle (selected from a dropdown menu)
- Query Text:**

```
1 v PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
2 v PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 v PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
4 v PREFIX hd: <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>
5 v
6 v ****list diseases and their definitions, with a COALESCE to handle diseases without a definition provided:
7 v SELECT ?diseaseName (COALESCE(?definition, "Definition not provided") AS ?diseaseDefinition)
8 v WHERE {
9 v   ?disease hd:DiseaseshasName ?diseaseName .
10 v   OPTIONAL { ?disease hd:DiseaseshasDefinition ?definition . }
11 v }
12 v LIMIT 5
13 v
14 v
15 v
```
- Results Table:**

diseaseName	diseaseDefinition
1 Alzheimer's Disease	A progressive disease that destroys memory and other important mental functions.
2 Asthma	A condition in which a person's airways become inflamed, narrow, and swell, producing extra mucus.
3 COVID-19	A respiratory illness caused by the SARS-CoV-2 virus.
4 Chickenpox	highly contagious disease caused by the varicella-zoster virus (VZV).
5 Diabetes	A metabolic disease characterized by high blood sugar levels over a prolonged period.
- Table Options:** Simple view, Ellipse, Filter query results, Page size
- Page Information:** Showing 1 to 5 of 5 entries

Figure 29. This SPARQL query retrieves the names of diseases and their definitions from RDF data.

Q4:

```
SELECT
```

```
?pathogenFamily  
(GROUP_CONCAT(DISTINCT ?pathogenName; separator=", ") AS ?pathogenNames)  
(GROUP_CONCAT(DISTINCT ?diseaseName; separator=", ") AS ?diseaseNames)  
(COUNT(?pathogen) AS ?pathogenCount)
```

```
WHERE {
```

```
?pathogen hd:PathogenhasFamily ?pathogenFamily ;  
      hd:PathogenhasName ?pathogenName.
```

```
?pathogen hd:causes ?disease.
```

```
?disease hd:DiseaseshasName ?diseaseName.
```

```
}
```

```
GROUP BY ?pathogenFamily
```

/Human-Disease-Ontology/

JSON

Turtle

```

1 v PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
4 PREFIX hd: <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>
5
6 ***What are the names and counts of pathogens within each pathogen family with a list of the names of these pathogens and disease names?
7 SELECT
8 ?pathogenFamily
9 (GROUP_CONCAT(DISTINCT ?pathogenName; separator=", ") AS ?pathogenNames)
10 (GROUP_CONCAT(DISTINCT ?diseaseName; separator=", ") AS ?diseaseNames)
11 (COUNT(?pathogen) AS ?pathogenCount)
12 WHERE {
13 ?pathogen hd:PathogenhasFamily ?pathogenFamily ;
14         hd:PathogenhasName ?pathogenName.
15 ?pathogen hd:causes ?disease.
16 ?disease hd:DiseaseshasName ?diseaseName.
17 }
18 GROUP BY ?pathogenFamily
19
20

```

Table Response 4 results in 0.027 seconds

pathogenFamily	pathogenNames	diseaseNames	pathogenCount
1 Mycobacteriaceae	Mycobacterium tuberculosis	Tuberculosis	"1"^^<http://www.w3.org/2001/XMLSchema#integer>
2 Coronaviridae	Severe acute respiratory syndrome coronavirus 2	COVID-19	"1"^^<http://www.w3.org/2001/XMLSchema#integer>
3 Orthomyxoviridae	Orthomyxoviridae	Influenza	"1"^^<http://www.w3.org/2001/XMLSchema#integer>
4 Spirochaetaceae	Borrelia	Lyme Disease	"1"^^<http://www.w3.org/2001/XMLSchema#integer>

Showing 1 to 4 of 4 entries

Figure 30. This query retrieves a list of pathogen families with the names of the pathogens in each family, disease caused by these pathogens and the total count of unique pathogens in each family.

Q5:

```

SELECT ?riskFactorName ?prevalenceInPercent ?preventionStrategyName ?efficiency
WHERE {
?diabetes rdf:type hd:Disease ;
        hd:DiseaseshasName "Diabetes" ;
        hd:hasRiskFactor ?riskFactor.
?riskFactor hd:RiskFactorhasName ?riskFactorName ;
        hd:RiskFactorhasPrevalenceInPopulation ?prevalenceInPercent.
?preventionStrategy hd:prevents ?diabetes;
        hd:PreventionStrategyhasName ?preventionStrategyName;
        hd:PreventionStrategyhasEfficiency ?efficiency.
FILTER(?prevalenceInPercent > 15)
}

```

The screenshot shows a SPARQL query interface with the following details:

- SPARQL Endpoint:** /Human-Disease-Ontology/sparql
- Content type (SELECT):** JSON
- Content type (GRAPH):** Turtle

```

1 v PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
4 PREFIX hd: <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>
5
6 /***What are the names and prevalence rates of risk factors for diabetes where the prevalence is above 15% with the names and efficiencies of prevention
7   strategies used against diabetes?
8
9  SELECT ?riskFactorName ?prevalenceInPercent ?preventionStrategyName ?efficiency
10 WHERE {
11   ?diabetes rdf:type hd:Disease ;
12   hd:DiseasehasName "Diabetes" ;
13   hd:hasRiskFactor ?riskFactor.
14   ?riskFactor hd:RiskFactorhasName ?riskFactorName ;
15   hd:RiskFactorhasPrevalenceInPopulation ?prevalenceInPercent.
16   ?preventionStrategy hd:prevents ?diabetes;
17   hd:PreventionStrategyhasName ?preventionStrategyName;
18   hd:PreventionStrategyhasEfficiency ?efficiency.
19 }
20

```

**Table Response:** 3 results in 0.05 seconds

riskFactorName	prevalenceInPercent	preventionStrategyName	efficiency
1 Drinking Alcohol	"85"^^<http://www.w3.org/2001/XMLSchema#decimal>	Adopting a Balanced and Healthy Diet	"60"^^<http://www.w3.org/2001/XMLSchema#decimal>
2 Family History of Diabetes	"24"^^<http://www.w3.org/2001/XMLSchema#decimal>	Adopting a Balanced and Healthy Diet	"60"^^<http://www.w3.org/2001/XMLSchema#decimal>
3 Obesity	"18"^^<http://www.w3.org/2001/XMLSchema#integer>	Adopting a Balanced and Healthy Diet	"60"^^<http://www.w3.org/2001/XMLSchema#decimal>

Showing 1 to 3 of 3 entries

Figure 31. This query is important risk factors for diabetes, and also the different strategies used to prevent diabetes, including how effective these strategies are and used FILTER for prevalence greater than 15.

<https://sparqlwrapper.readthedocs.io/en/latest/main.html>

## 6. Research Mapping/ Conclusion

### Refer to Related works section 2.

This report is about Human Disease Ontology. I have started with my aim of understanding OWL by expanding a prototype application in the aim and concept section to create a searchable ontology for human disease.

My project represents an ontology which categorised Human Disease into chronic respiratory, infectious, and acute diseases. The Human Disease Ontology focuses on identifying lifestyle-related issues, like smoking, drinking alcohol, and obesity as early as possible to potentially prevent chronic disease as the cause of death. Therefore, based on behaviour risk factors, and leveraging an application that can make recommendations such as diet changing or quitting smoking.

The project has been using the protege software tool to identify relationships and class properties comprehensible and effectively to make the disease ontology easy to navigate.

The ontology can answer complicated questions about disease, symptoms, treatments, and risk factors by using SPARQL queries and its use as an educational tool and for medical research.

Regardless of small scope, the project has successfully illustrated how semantic data and ontology can be powerful tools to manage large databases in medical research.

## 6.1. Critical Reflection

As my project indicates a well structured approach to build a disease ontology, with significant potential impact on healthcare and medical research. Also by choosing good definitions of properties and classes and the combination of SPARQL queries, that indicates a robust system. Therefore, the project considers the complication of diseases and their complexity of factorial nature, keeping updates and validation against current medical standards would be crucial.

The aim of this project is to make connections in different biomedical data points, and success will be dependent on its ability to remain up to date and medically accurate.

## 6.2. Future Scope

Human Disease Ontology, key areas for developments:

- ❖ Expanding ontology with new symptoms, risk factors and disease while medical research develops.
- ❖ Build personalised healthcare recommendations.
- ❖ Create a platform to contribute and use the ontology for researchers.

## 6.3. Compare with alternative approaches[\[15\]](#)

- ❖ Standard databases like SQL might lack the semantic relationships that ontology provides while my ontology approach has complex relationships between disease, symptoms, and treatments.
- ❖ Statistical models are great for defying trends and make a prediction based on data but might be crucial for understanding disease ontology.
- ❖ Rule based models might be strong for specific tasks but have less flexibility than an ontology model does.

## 6.4. Discussion of Limitations

- ❖ Ontologies can grow complex, fast, and hard to maintain as new data is continuously integrated. Make sure of the accuracy of disease ontology.

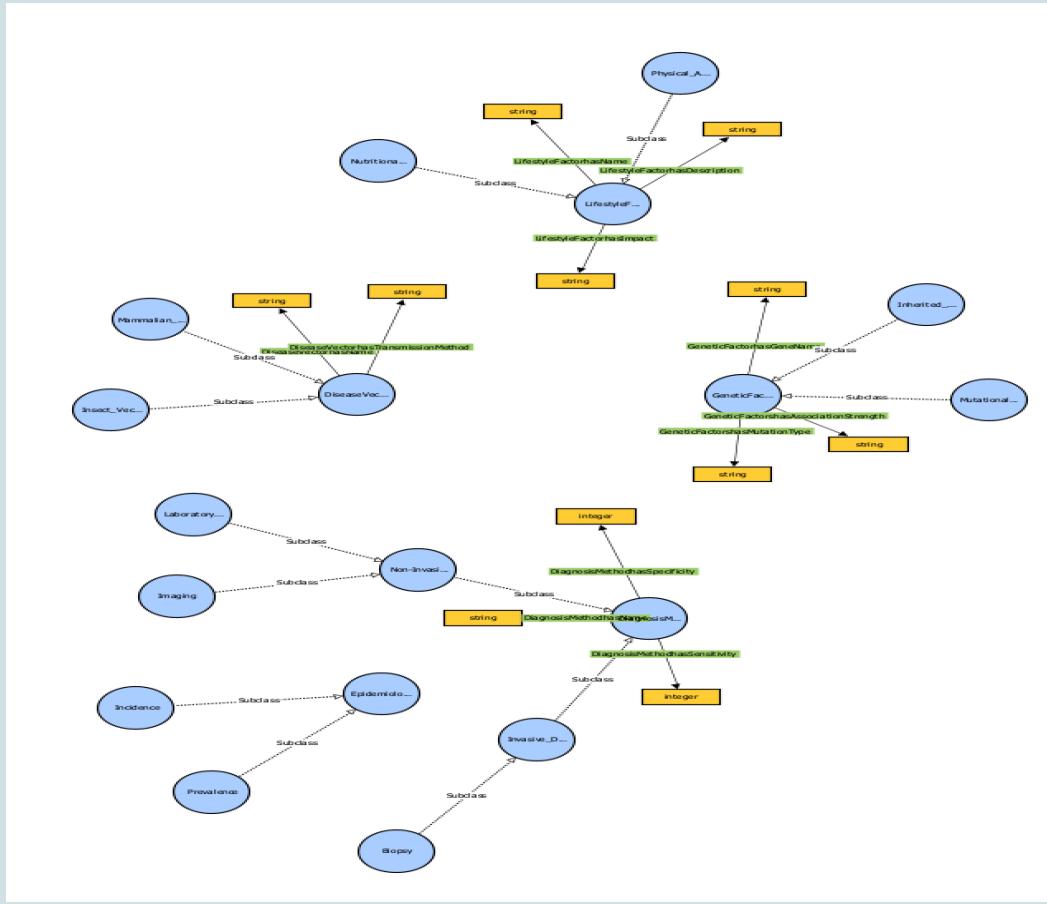
## 7 References

1. WHO. 2021. Strengthening NCD service delivery through UHC benefit package: technical meeting report, Geneva, Switzerland, 14-15 July 2020.
2. James C. Griffiths, Jan De Vries, Michael I. McBurney, Suzan Wopereis, Samet Serttas, and Daniel S. Marsman. 2020. Measuring health promotion: translating science into policy. European Journal of Nutrition 59, 2 (01 Sep 2020).
3. Larentis, A.V.; Neto, E.G.d.A.; Barbosa, J.L.V.; Barbosa, D.N.F.; Leithardt, V.R.Q.; Correia, S.D. Ontology-Based Reasoning for Educational Assistance in Noncommunicable Chronic Diseases. Computers 2021, 10, 128.
4. Larentis, A.V.; Neto, E.G.d.A.; Barbosa, J.L.V.; Barbosa, D.N.F.; Leithardt, V.R.Q.; Correia, S.D. Ontology-Based Reasoning for Educational Assistance in Noncommunicable Chronic Diseases. Computers 2021, 10, 128.
5. Blomqvist, E. The Use of Semantic Web Technologies for Decision Support—A Survey. Semant. Web 2014, 5, 177–201.
6. Logical Design, The Morgan Kaufmann Series in Data Management Systems, 2011, Pages 85-108.
7. Malta, D.C.; Gomes, C.S.; Barros, M.; Lima, M.G. Noncommunicable diseases and changes in lifestyles during the COVID-19 pandemic in Brazil. Rev. Bras. Epidemiol. 2021, 24, e210009.
8. Schriml LM, Munro JB, Schor M, et al. The Human Disease Ontology 2022 update. Nucleic Acids Research. 2022 Jan;50(D1):D1255-D1261. DOI: 10.1093/nar/gkab1063. PMID: 34755882; PMCID: PMC8728220.
9. Matcher · spaCy API Documentation
10. Matthew Horridge, Holger Knublauch, Natasha Noy, Alan Rector et al.

[https://protege.stanford.edu/conference/2007/slides/owlTutorial-reasoning\\_dameron.pdf](https://protege.stanford.edu/conference/2007/slides/owlTutorial-reasoning_dameron.pdf)

11. International Journal of Computer Applications (0975 – 8887) Volume 56– No.6, October 2012.
12. Schriml LM, Munro JB, Schor M, Olley D, McCracken C, Felix V, Baron JA, Jackson R, Bello SM, Bearer C, Lichenstein R, Bisordi K, Dialo NC, Giglio M, Greene C. The Human Disease Ontology 2022 update. Nucleic Acids Res. 2022 Jan 7;50(D1):D1255-D1261. doi: 10.1093/nar/gkab1063. PMID: 34755882; PMCID: PMC8728220.
13. M. Hadzic and E. Chang, "Ontology-Based Support for Human Disease Study," Proceedings of the 38th Annual Hawaii International Conference on System Sciences, Big Island, HI, USA, 2005, pp. 143a-143a, doi: 10.1109/HICSS.2005.472.
14. Lucas Pfeiffer Salomão Dias, Henrique Damasceno Vianna, Weslei Heckler, and Jorge Luis Victória Barbosa. 2023. Ontology-Based Reasoning to Classify Behaviors Associated with Chronic Disease Risk Factors. In XIX Brazilian Symposium on Information Systems (SBSI '23), May 29–June 01, 2023, Maceió.
15. Kamran Munir, M. Sheraz Anjum, The use of ontologies for effective knowledge modelling and information retrieval, Applied Computing and Informatics, Volume 14, Issue 2, 2018, Pages 116-126, ISSN 2210-8327, <https://doi.org/10.1016/j.aci.2017.07.003>.
16. <https://www.w3.org/TR/2012/PR-sparql11-query-20121108/>,  
[https://www.openresearch.org/wiki/Sparql\\_endpoint/Examples](https://www.openresearch.org/wiki/Sparql_endpoint/Examples)
17. <https://spacy.io/api/phrasematcher>, <https://spacy.io/api/matcher>
18. <https://www.w3.org/TR/turtle/>
19. <https://bioportal.bioontology.org/ontologies/DOID/?p=summary>

## 8. Appendix



**Characteristics:** `diagnosedBy`

Functional

Inverse function

Transitive

Symmetric

Asymmetric

Reflexive

Irreflexive

**Description:** `diagnosedBy`

`Equivalent To` [+](#)

`SubProperty Of` [+](#)

`Inverse Of` [+](#)

`Domains (intersection)` [+](#)

**Disease**

`Ranges (intersection)` [+](#)

**DiagnosisMethod**

**DiagnosticTest**

`Disjoint With` [+](#)

`SuperProperty Of (Chain)` [+](#)

**Description:** `Antibiotics` [+](#)

**Property assertions:** `Antibiotics`

`Object property assertions` [+](#)

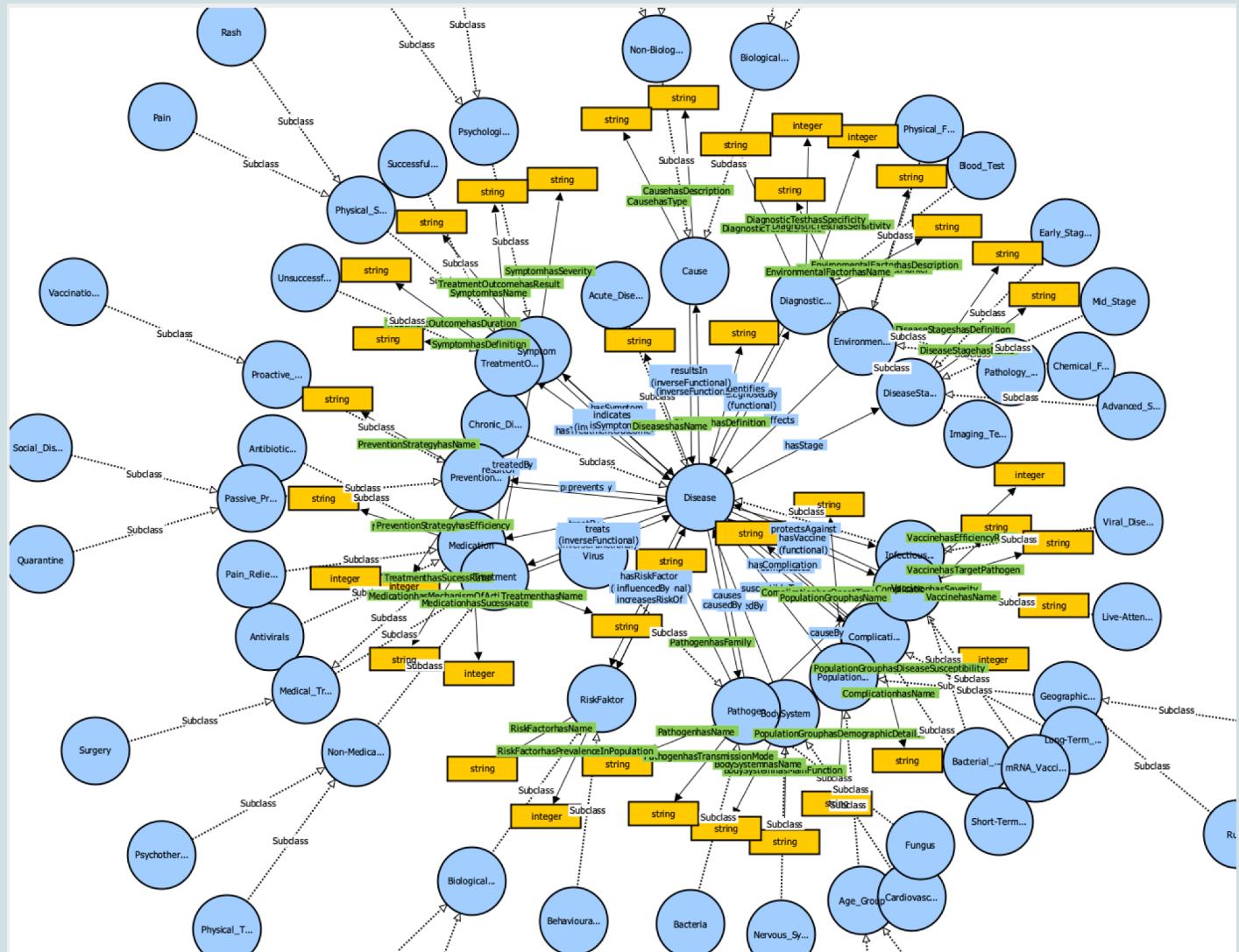
- Medication** [?](#) [@](#) [X](#) [O](#)
- Treatment** [?](#) [@](#) [X](#) [O](#)

`Data property assertions` [+](#)

- TreatmenthasName** "Antibiotics"
- TreatmenthasSuccessRate** 90

`Negative object property assertions` [+](#)

`Negative data property assertions` [+](#)



```

# What are the risk factor for Asthma?
# Give me diagnosis method for Diabetes
# What are some diseases?
# disease has names?
# List all symptom.
# List diagnosis method
# Show treatment.
# Give me cause.
# What are the risk factor?
# Give me diagnosis method

Enter your query in natural language: What are the risk factor?
userQuery: What are the risk factor?
matchId: 13289220008681790430
span: risk factor
token: risk
token: factor
foundTerms : {'risk factor'}
term: risk factor
whereClause: ?item rdf:type hd:RiskFaktor
SPARQL query to execute:

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX hd: <http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#>
SELECT DISTINCT ?item WHERE { ?item rdf:type hd:RiskFaktor } LIMIT 10
results: {'head': {'vars': ['?item']}, 'results': {'bindings': [{'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#Aging'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#AirPollution'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#Allergies'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#GeneticPredisposition'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#CalciumDeficiency'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#DrinkingAlcohol'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#FamilyHistoryofDiabetes'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#Obesity'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#FamilyHistoryofGlaucoma'}}, {'item': {'type': 'uri', 'value': 'http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#WeakenedImmuneSystem'}}]}
Results found:
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#Aging
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#AirPollution
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#Allergies
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#GeneticPredisposition
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#CalciumDeficiency
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#DrinkingAlcohol
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#FamilyHistoryofDiabetes
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#Obesity
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#FamilyHistoryofGlaucoma
http://www.semanticweb.org/sarashahin/ontologies/2023/9/Human_Disease_ontology#WeakenedImmuneSystem

```

