

INSTITUTO DE COMPUTAÇÃO  
UNIVERSIDADE ESTADUAL DE CAMPINAS

**Linked biology technical aspects – linking  
phenotypes and phylogenetic trees**

*E. Miranda*      *A. Santanchè*

Technical Report - IC-14-06 - Relatório Técnico

February - 2014 - Fevereiro

The contents of this report are the sole responsibility of the authors.  
O conteúdo do presente relatório é de única responsabilidade dos autores.

# Linked biology technical aspects – linking phenotypes and phylogenetic trees

Eduardo Miranda \*      André Santanchè †

## Abstract

A large number of studies in biology, including those involving phylogenetic trees reconstruction, result in the production of a huge amount of data – e.g., phenotype descriptions, morphological data matrices, etc. Biologists increasingly face a challenge and opportunity of effectively discovering useful knowledge crossing and comparing several pieces of information, not always linked and integrated. Ontologies are one of the promising choices to address this challenge. However, the existing digital phenotypic descriptions are stored in semi-structured formats, making extensive use of natural language. This technical report is related to a research developed by us [1] to addresses this problem, adding an intermediate step between semi-structured phenotypic descriptions and ontologies. It remodels semi-structured descriptions to a graph abstraction in which the data are linked. Graph transformations subsidize the transition from semi-structured data representation to a more formalized representation with ontologies. The present technical report drills down implementation details of our system. It provides a module to ingest phylogenetic trees and phenotype descriptions – represented in semi-structured formats – into a graph database. Additionally, two approaches to combine distinct data sources are presented and an algorithm to trace changes in phylogenetic traits of trees.

---

\*Institute of Computing – State University of Campinas, 13081-970, Campinas, Brazil. Work partially financed by (CNPq 138197/2011-3), the Microsoft Research FAPESP Virtual Institute (NavScales project), CNPq (MuZOO Project and PRONEX-FAPESP), INCT in Web Science(CNPq 557.128/2009-9) and CAPES, as well as individual grants from CNPq.

†Institute of Computing – State University of Campinas, 13081-970, Campinas, Brazil. Work partially financed by (CNPq 138197/2011-3), the Microsoft Research FAPESP Virtual Institute (NavScales project), CNPq (MuZOO Project and PRONEX-FAPESP), INCT in Web Science(CNPq 557.128/2009-9) and CAPES, as well as individual grants from CNPq.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Basic concepts</b>	<b>3</b>
2.1	Standards for Phenotype Description . . . . .	3
2.2	Life Science Identifiers (LSIDs) . . . . .	5
2.3	The proposed graph data model . . . . .	5
<b>3</b>	<b>System Architecture and Implementation Details</b>	<b>6</b>
3.1	SDD Parser . . . . .	6
3.2	Tree Output . . . . .	8
3.3	Global Names Resolver (GNR) . . . . .	9
3.4	Graph Importer . . . . .	9
3.5	Graph Database . . . . .	10
3.6	Similarity Index . . . . .	12
3.6.1	Practical Implementation of the Similarity Measure . . . . .	13
3.7	Tracing the Evolutionary History . . . . .	14
<b>4</b>	<b>Conclusion</b>	<b>17</b>
<b>A</b>	<b>Demonstration</b>	<b>21</b>
A.1	SDDParser.py . . . . .	21
A.2	TeeOutput.py . . . . .	26
A.3	GlobalNamesResolver.py . . . . .	30
A.4	GNRResultObject.py . . . . .	33
A.5	ITIServices.py . . . . .	35
A.6	CoLServices.py . . . . .	37
A.7	GraphImporter.py . . . . .	40
A.8	SimilarityIndex.py . . . . .	48
A.9	TraceEvolutionaryHistory.py . . . . .	51

## 1 Introduction

In 1859 Charles Darwin published *On the Origin of Species* which is considered the foundation of evolutionary biology. In his book, Darwin set forth the theory of evolution and natural selection. It argues that all life is related and has descended from a common ancestor. The Tree of Life is a metaphor to describe the relationships between living and extinct organisms through their common ancestors. More precisely, it is an abstract form to represent hypotheses about evolutionary relationships, in which all species that have ever existed are taken together with relationships among them, describing their evolutionary lineages. In this abstract representation, the taxa are the leaves of the tree and the internal nodes are common ancestors, or hypothetical taxa.

This huge and complex tree is split into smaller branches, which are investigated separately and then incorporated into the tree. Evolutionary biologists normally work in relatively small chunks of the tree, analyzing a very specific subset of species. A fundamental challenge in this scenario is the creation of a complete evolutionary Tree of Life [2], assembling genomic and morphological data so as to congregate the phylogenetic relationships among all known living or extinct organisms [3, 4, 5]. The integration of these data may contribute to better understand how a morphological trait became organized and evolved over time [6], how organisms interact and how life on Earth came to be.

The main goal of this research is to design and implement a linked biology approach to automatically connect and combine data from independent semi-structured resources of phenotype descriptions and/or phylogenetic trees, exploiting their latent semantics. We propose a graph data model that plays a crucial role, since it is the basis of our linking discovery and combination process. It contributes assisting biologists in the exploration of existing biology assets related to phenotype descriptions and their latent semantics. The present work details algorithms, implementation aspects and the database model related to our research.

The text is organized as follows. Section 2 synthesizes basic concepts necessary for understanding the text. Section 3 discusses implementation details of our system and presents some results. Section 4 presents concluding remarks. In the Appendix the source code is provided with comments explaining its functionalities.

## 2 Basic concepts

In this section, we highlight basic concepts adopted in this text. Subsection 2.1 introduces some key elements of XML formats for phenotype description. Subsection 2.2 we details the Life Science Identifier which is one of the solutions for data interconnection. Subsection 2.3 presents an overview of our proposed graph model.

### 2.1 Standards for Phenotype Description

There is a wide variety of representation formats for phenotype descriptions adopted by information systems and open standards, which represent differently the same information. In [1] we analyze four of them – Xper<sup>2</sup>, SDD, Nexus and NeXML – looking for a common denominator which is the foundation for our graph-based model. SDD, Nexus and NeXML are widely adopted open standards. Xper<sup>2</sup> (<http://lis-upmc.snv.jussieu.fr/lis/>) is a management system adopted by the systematist community, for storing, editing and analyzing phenotype descriptive data. It focuses mainly on taxonomic descriptions, allowing creation, sharing and comparison of identification keys [7, 8]. Xper<sup>2</sup> was developed in the Laboratoire Informatique & Systématique of the University Pierre et Marie Curie and this work is part of a bigger project in collaboration with this lab. Therefore, Xper<sup>2</sup> was adopted for our practical implementation.

In order to transform phenotype observations into digital records and generalize them – e.g., devising general characters and states observed in a genus of monitor lizards – the biologist may use a tool as Xper<sup>2</sup>. Phenotype descriptions can be stored in the Xper<sup>2</sup> native format or can be exported to the SDD open format. The Structured Descriptive Data (SDD) (<http://wiki.tdwg.org/SDD>) is a platform and application-independent XML-based standard developed by the Biodiversity Information Standards (historic acronym: TDWG) for recording and exchanging descriptions of biological and biodiversity data of any type [9]. SDD is adopted by several other phenotype description tools – e.g., Lucid Central (<http://www.lucidcentral.org>) and Linnaeus II (<http://www.eti.uva.nl/>).

We further introduce some key elements of the SDD format, which are recurrent in the formats confronted in [1]. A SDD description comprises, in a single file, a domain schema and its instances. Figure 1 shows a diagram with a fragment of a SDD file containing the description of a varanus lizard. A (C,CS) description in SDD has two main blocks: (i) defines the characters involved and their possible states – Figure 1 top; (ii) describes an Operational Taxonomic Unit (OTU) using the characters defined in (i) – Figure 1 bottom. OTU is a biology term which refers to a given taxon at the rank adopted to the study – e.g., a specimen, a species, a genus etc.

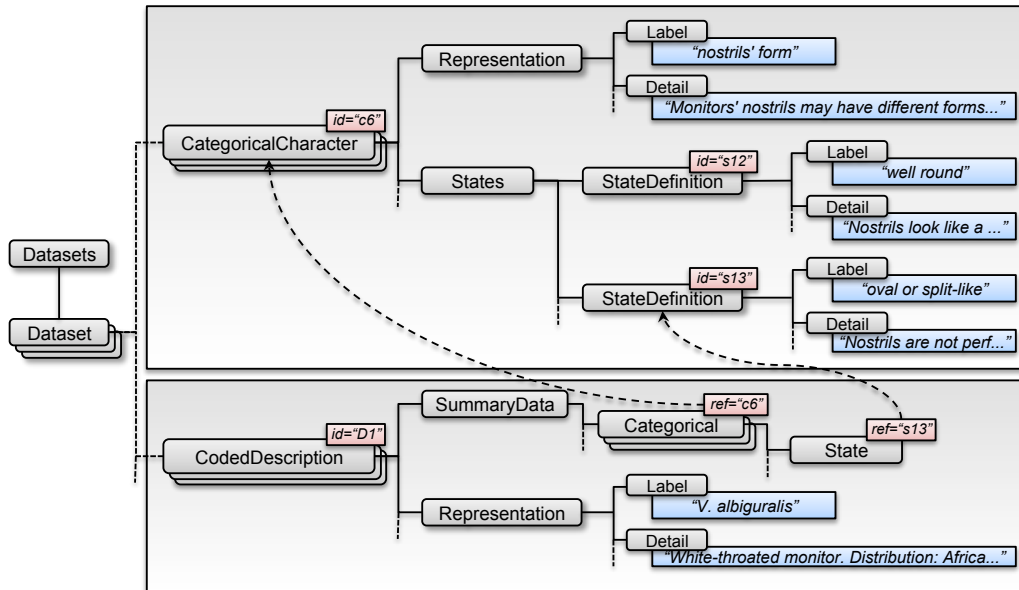


Figure 1: Fragment of SDD Schema with Instances

*<CategoricalCharacter>*s and their *<States>* (shown in Figure 1 top) are primitives to describe an OTU [9]. Each *<CategoricalCharacter>* has its *<Representation>* – comprising a label and a description as plain texts – and a set of *<StateDefinition>* elements with their possible states. *<CategoricalCharacter>* and *<StateDefinition>* elements defined here will be referred throughout the XML document by their ids. The *<CodedDescription>* (Figure 1 bottom) links the described OTU to *States* of each *<CategoricalCharacter>*. It has two essential items: (i) the described OTU, where its name and description are listed in natural language under *<Representation>*; (ii) a set of character and values (*<Categorical>* and *<State>*), which address the characters defined in the previous section through the *ref* attribute. It is possible and usual to assign multiple character-states for a given OTU (i.e. in case of polymorphism). A first integration, problem observed here is that each character or OTU described does not have a global unique identification among documents. Therefore, the description can only be used by the document where it was declared and it is not

possible to guarantee the equivalence of two or more  $\langle \textit{CategoricalCharacters} \rangle$ .

## 2.2 Life Science Identifiers (LSIDs)

One of the problems faced in life science is related to the identification of objects within and across repositories [10]. More precisely, an object may refer to a taxon, gene, anatomical feature, phenotypic description, geographical location etc. Integrating data from different sources is not straightforward and uniquely identifying these objects is undoubtedly a key point for the success of our proposed solution.

During the 18th century, Carolus Linnaeus introduced the binomial nomenclature for naming species that is the basis of modern classification [11]. This system basically concatenates 2 Latin words, where the first part identifies the species genera and the second one the species itself. The binomial nomenclature has been used for the last 250 years [11] and the biological information related to organisms is historically annotated by species names. Hence, the binomial name would appear to be a logical candidate to index information available about species. However, misspelling problems are often encountered [12, 13], moreover, taxonomic names are not unique identifiers [14, 15] because scientists may use (i) similar names to different species (homonyms) or (ii) multiple names for the same specie (synonyms) [10, 16].

Furthermore, each organization has its own means of defining a key, which makes the problem even harder to solve. For example, the species *Aotus ericoides* has the id 11479744 on the Catalogue of Life (CoL), id 42472 on the Australian Plant Name Index (APN), id 643314 on the Encyclopedia of Life (EoL), id 129761-3 on the The International Plant Names Index (IPNI), id 700844 on the Universal Biological Indexer and Organizer (uBio) etc.

In order to address this issue, some organizations – e.g., Universal Biological Indexer and Organizer (uBio), Integrated Taxonomic Information System (ITIS), Catalogue of Life (CoL), The International Plant Names Index (IPNI), National Center for Biotechnology Information (NCBI) etc. – incorporated into their projects the concept of Life Science Identifiers (LSIDs), proposed by the Object Management Group (OMG) (<http://www.omg.org/>). LSID is a persistent, location-independent resource identifier, whose purpose is to uniquely identify biological resources [17]. The persistent property refers to the fact that LSID identifiers are unique, can be assigned to only one object forever and they never expire. The location-independent property specifies that each authority locally creates LSIDs and they are the responsible to guaranteeing the uniqueness of LSIDs.

## 2.3 The proposed graph data model

In this section we will present an overview of our proposed graph model. From the numerous graph data models proposed – see [18, 19, 20] for more details – the *property graph* model was adopted in the present work. In a *property graph*, nodes and relationships can maintain extra metadata as a set of key/value pairs. Moreover, relationships are typed, enabling to create multi-relational networks with heterogeneous sets of edges. Different from single-relational networks, in which edges are of the same type, multi-relational networks are more appropriate to represent complex domain models, due to the variety of relationship types in the same graph [21]. For example: relationships may either represent membership in a social group (family membership) or professional relationships (employer-worker relationship) simultaneously in the same network.

Figure 2 shows our graph data model. The tables below the nodes/edges represent their types and metadata. We mapped the SDD format to the graph model as follows: OTUs are entities (e.g., “*Varanus prasinus*”) and, therefore, were mapped to nodes. A future target of this project is to enrich our model by associating identifiable entities to ontology concepts. One may consider to map Characters and Characters States to key/value pairs, to be related to OTU nodes. However, we decided to map Characters to nodes, in order to unify in the same node equivalent characters observed

in several OTUs and, in a future work, to relate the unified characters with ontologies. Finally, the Character-state makes a semantic bridge (relationship) between OTUs and Characters. Thus, a statement like “*Varanus gouldi ventral pattern is randomly scattered dark spots*” is represented in our model as *Varanus gouldi* (node)  $\rightarrow$  *randomly scattered dark spots* (edge)  $\rightarrow$  *ventral pattern* (node).

Our model comprises, in a single place, phenotype descriptions and phylogenetic trees. For this reason a new node called HTU (Hypothetical Taxonomic Unit) is present in this model. HTUs are internal nodes in phylogenetic trees that represent an inferred ancestral organism. HTUs are hypothetical common ancestors of OTUs nodes and, therefore, can only be connected to themselves (HTU  $\rightarrow$  HTU) or to OTUs (HTU  $\rightarrow$  OTU). For the sake of modeling simplicity, only the *TreeEdge* relationship is allowed between HTU  $\rightarrow$  HTU and HTU  $\rightarrow$  OTU. Finally, there is also a character-state relationship between HTU nodes and character nodes that are strictly created by some algorithms.

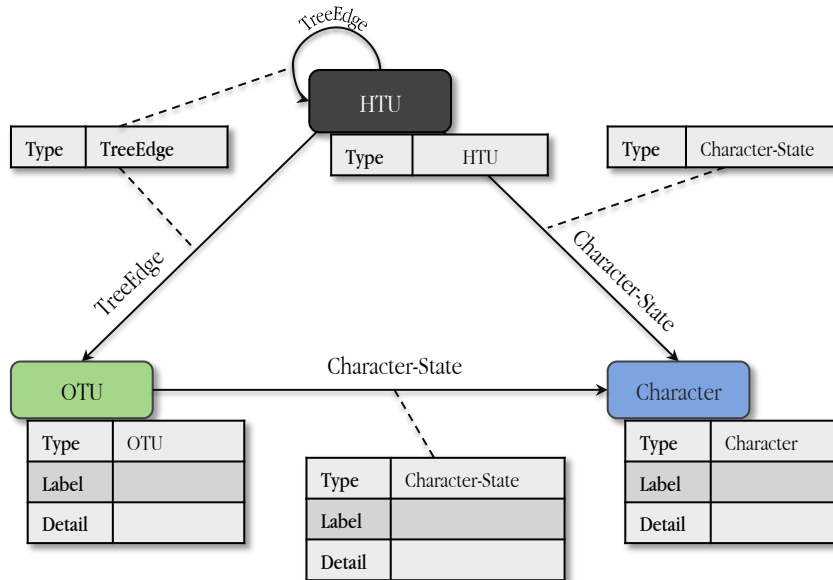


Figure 2: Property Graph Model

### 3 System Architecture and Implementation Details

In this section, we analyze the system architecture and its implementation details, in order to present its main functionalities and operational features. The text is presented progressively. The core functionalities are shown in the first subsections and the algorithms are presented later.

We have developed our platform on top of the Neo4j graph database (<http://www.neo4j.org/>), mainly due to its widespread adoption. Our implementation uses the Python programming language and Py2neo (<http://book.py2neo.org/>), which is an interface connecting Python and Neo4j via REST API. The adopted query language was *Cypher*, which is a declarative graph query language.

#### 3.1 SDD Parser

Our SDD Parser has all functionalities to parse an SDD file (for implementation details see Appendix A.1) using the Python *xml.dom.minidom*, which is a minimal implementation of the Document Ob-

ject Model interface. Listing 1 shows an SDD fragment of a Varanus knowledge base <sup>1</sup>, of which Figure 1 is a simplified abstraction. In addition, all main SDD structures presented in Figure 1 and Listing 1 – Representation, StateDefinition, CategoricalCharacter, Categorical and CodedDescription – were processed to produce our graph.

Listing 1: Varanus.sdd.xml

```

1 <Characters>
2   ...
3   <CategoricalCharacter id="c6">
4     <Representation>
5       <Label>nostrils ' form</Label>
6       <Detail>Monitors ' nostrils mayhave different forms.&lt;br>Look
          at the head in side view or dorsal view in order to
          appreciate this characteristic.</Detail>
7       <MediaObject ref="m40"/>
8     </Representation>
9     <States>
10    <StateDefinition id="s12">
11      <Representation>
12        <Label>well round</Label>
13        <Detail>Nostrils look like a quite perfect circle.</Detail>
14      </Representation>
15    </StateDefinition>
16    <StateDefinition id="s13">
17      <Representation>
18        <Label>oval or split-like</Label>
19        <Detail>Nostrils are not perfectly round: they are oval or they
          present a split-like form.</Detail>
20      </Representation>
21    </StateDefinition>
22  </States>
23 </CategoricalCharacter>
24   ...
25 </Characters>
26   ...
27 <CodedDescriptions>
28 <CodedDescription id="D1">
29 <Representation>
30 <Label>V. albiguralis</Label>
31 <Detail>White-throated monitor&lt;br>&lt;br>Distribution :
          Africa (West and South).&lt;br>&lt;br>CITES: appendix II
          .</Detail>
32 <MediaObject ref="m1"/>
33 </Representation>
34 <SummaryData>
35   ...
36 <Categorical ref="c6">
```

<sup>1</sup>Knowledge base of the genus Varanus  
([http://lis-upmc.snv.jussieu.fr/xper2/infosXper2Bases/details\\_base.php?id\\_base=86](http://lis-upmc.snv.jussieu.fr/xper2/infosXper2Bases/details_base.php?id_base=86))



```

37     <State ref="s13"/>
38     </Categorical>
39     ...
40     </SummaryData>
41     </CodedDescription>
42 </CodedDescriptions>

```

### 3.2 Tree Output

The present work also draws upon phylogenetic trees generated from LisBeth (<http://lis-upmc.snv.jussieu.fr/lis/>). LisBeth is a cladistics software for phylogenetics and biogeography [22] that implements the three-item analysis (3ia) method of phylogenetic inference [23]. It minimizes the conflictual relationships within a set of characters, or maximizes the compatible relationships so as to reconstruct one or several optimal tree(s). We implemented a *TreeOutput* class, which abstracts the functions of interacting with LisBeth output files (for implementation details see Appendix A.2). Listing 2 displays two fragments of a LisBeth output file, focusing in the elements processed in this work, i.e. taxons with their ids and the retained tree – newick tree which is a way to represent a tree in computer-readable form, using parentheses and commas. The *TreeOutput* main function combines the retained tree with the taxon names, retrieved in previous steps, and returns a root node to a tree that represents the retained tree. In this new tree, the internal nodes are renamed to *HTU* and the leaf nodes to its respective taxon names (see Figure 3).

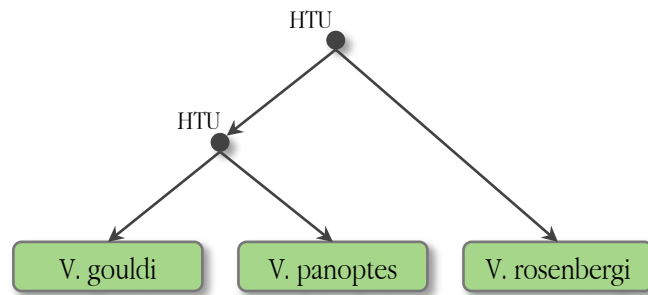


Figure 3: Retained Tree Example

Listing 2: LisBethOutput.3iz

```

1     ...
2     <D02>
3     ...
4     Taxa (3) :
5     .           3   V. gouldii
6     .           7   V. panoptes
7     .          12   V. rosenbergi
8     <F02>
9     ...
10    <D06>
11    ...
12

```

```

13 Retained trees : 1
14 .      1: ((3 7) 12)
15 <F06>
16 ...

```

### 3.3 Global Names Resolver (GNR)

In order to find a valid LSID, we adopted the Global Names Resolver (GNR) web service (<http://resolver.globalnames.org/>) that executes exact or fuzzy matching against canonical forms of scientific names in 170 distinct data sources. The Canonical form (cf) is the simplest, most complete and unambiguous form of a name. The Canonical form of scientific names consists of the genus and species – when applied – with no authorship, rank, nomenclatural annotation or subgenus.

Our system used three of the six types of matching offered by the GNR resolver: (i) exact matching; (ii) exact matching of canonical forms – this process reduces a given name to its canonical form and checks it for an exact match; (iii) fuzzy matching of canonical forms – uses a modified version of the TaxaMatch algorithm [13] and intends to work around misspellings errors. It does a fuzzy match of the canonical form of a given name – even with mistakes – against spellings considered correct. The GNR resolver reports the matching quality (“*confidence score*”) for each match. The other three remaining matching types are: (iv) exact matching of specific parts of names, (v) fuzzy matching of specific parts of names and (vi) exact matching of genus part of names. They were not adopted because we focused in complete names in their canonical form.

Our algorithm extracts all plain text taxon entities present in the SDD file and, for each one, it uses the GNR to transform the taxon name to its canonical form. Only those taxons with confidence score above of 0.988 are considered. After that, the algorithm makes use of the GNR resolver to search for its LSID (for implementation details see Appendix A.3) – only exact matches are considered. The GNR results have the output field “*local id*” which, in the case of uBio, is the LSID. Moreover, we prioritized the uBio LSID, since it indexes and organizes until now more than 11 million names. But there are cases in which the GNR resolver does not retrieve any result from the uBio. In these cases, the algorithm makes use of the Integrated Taxonomic Information System (ITIS) web services (<http://www.usgovxml.com/DataService.aspx?ds=ITIS>), in order to obtain the LSID (for implementation details see Appendix A.5). ITIS is a reliable taxonomic base for species, with more than 740 thousand common names and scientific names indexed. If none of the services return a valid LSID, we also implemented a class to interact with the CoL web service (<http://www.catalogueoflife.org/col/webservice>), attempting to obtain a valid LSID (for implementation details see Appendix A.6).

### 3.4 Graph Importer

Graph Importer is an object class written in Python that is responsible for coupling the phylogenetic trees and phenotype descriptions into the graph database. The insertion process follows the sequence: (1) Starts parsing the SDD XML file and the LisBeth output file – see Listing 1 and 2 respectively. (2) Creates a taxon node for each taxon present in the SDD file – see Figure 1 bottom, tag <Representation>. In this process, it searches for a valid LSID for each taxon node, using the GNR web service, ITIS web service or CoL web service. If the LSID is not found, it creates a taxon node without LSID. (3) Joins the taxon nodes to the tree structure, extracted from the LisBeth output file. (4) A node is created for each character in the SDD file – see Figure 1 top, tag <Representation>. (5) The taxon nodes are linked to the character nodes by their character-states – see Figure 1 top, tag <States>/<StateDefinition>. It will exist character-state relationships where exists a pair <SummaryData>/<Categorical> and <SummaryData>/<Categorical>/<State> – see

Figure 1 bottom, tag `<SummaryData>`. For implementation details see Appendix A.7. Figure 4 shows a visual representation of the retained tree combined with the taxon nodes provided in Listing 2. The figure shows that the edges depart from taxon nodes toward character nodes.

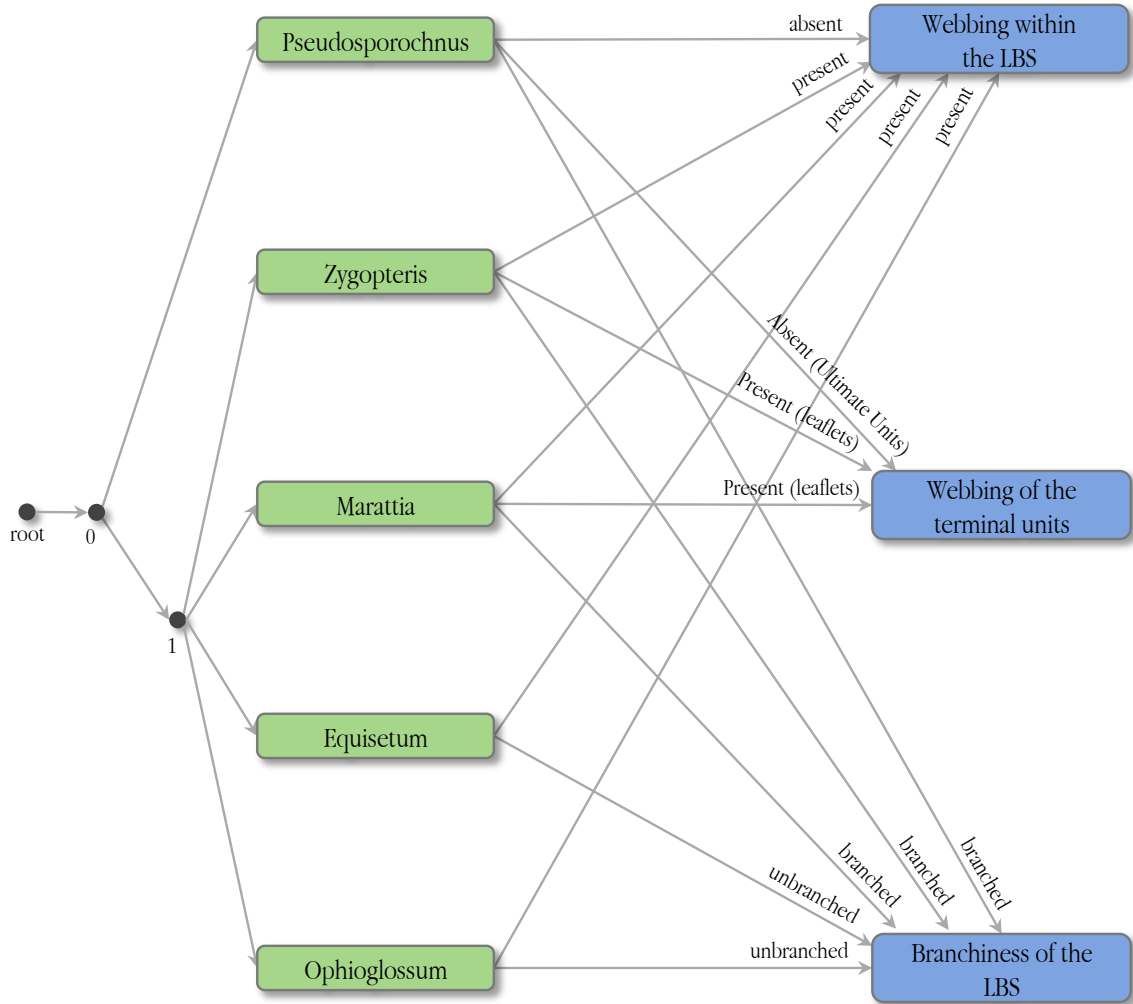


Figure 4: Real Example

### 3.5 Graph Database

We implemented a GraphDB class, which abstracts and centralizes all database operations. We describe each function header, followed by a short description of the main *Cypher* queries used in the system.

```

1 getNodeByLSID( LSID ):
2 // Returns a node for the supplied LSID.
3 START n=node( * )
4 WHERE n.lsid = 'LSID'

```

```

5 RETURN n
6
7 getOutgoingAdjacentNodes( GivenNode ):
8 // Returns all nodes to which the given node points to.
9 START n=node( GivenNode.id )
10 MATCH (n)-->(c)
11 RETURN DISTINCT c
12
13 getIncomingAdjacentNodes( GivenNode ):
14 // Returns all nodes that points to the given node.
15 START n=node( GivenNode.id )
16 MATCH (c)-->(n)
17 RETURN DISTINCT c
18
19 getIncomingAdjacentRelationships( GivenNode ):
20 // Returns all relationships incoming to a given node.
21 START n=node( GivenNode.id )
22 MATCH ()-[r]->(n)
23 RETURN r
24
25 getIncomingAdjacentNodesWithRelationshipInBetween( GivenNode,
26   GivenRelationship ):
27 // Returns all nodes, ordered by their label, that points to a given
28   node with a given relationship in between.
29 START n=node( GivenNode.id )
30 MATCH (c)-[:GivenRelationship.label]->(n)
31 RETURN c
32 ORDER BY c.label
33
34 getOutgoingRelationships( GivenNode ):
35 // Returns all relationships outgoing from a given node.
36 START n=node( GivenNode.id )
37 MATCH (n)-[r]->()
38 RETURN r
39
40 getDistinctRelationshipsInBetween( GivenNodeA, GivenNodeB ):
41 // Returns all distinct relationships that exists between nodes A
42   and B.
43 START a=node( GivenNodeA.id ), b=node( GivenNodeB.id )
44 MATCH (a)-[r]-(b)
45 WITH COLLECT( DISTINCT TYPE( r ) ) as rels
46 RETURN rels
47
48 getDescriptionNodesOfATree( TreeRoot )
49 // Returns all distinct description nodes id, character or character
50   -states depending on the schema, that are conected to a given
51   tree.
52 START root=node( TreeRoot.id )
53 MATCH (root)-[*..]->(d)

```

```

49 WHERE d.type = 'description'
50 RETURN DISTINCT ID(d)
51
52 deleteNodeRelationshipsExceptLabel( GivenNode, RelationshipLabel ):
53 // Deletes all node relationships except for a given relationship
   label.
54 START n=node( GivenNode.id )
55 MATCH n-[r]->()
56 WHERE NOT( r.label = 'RelationshipLabel' ) AND NOT( r.type = '
   TreeEdge' )
57 DELETE r
58
59 deleteRelationshipsTypeFromNode( GivenNode, RelationshipType ):
60 // Deletes all node relationships of a given type.
61 START n=node( GivenNode.id )
62 MATCH n-[r]->()
63 WHERE r.type = 'RelationshipType'
64 DELETE r

```

### 3.6 Similarity Index

We are proposing a heuristic similarity measure that computes the similarity degree between two morphological character descriptions. This measure will represent how closely related they are. The similarity index ( $S_i$ ) is based on 2 weighted aspects. 25% of the index is calculated based on the taxa being described, i.e. it analyzes if two given characters ( $C_1$  and  $C_2$ ) describe the same taxa. The other 75% are based on the meaning of the character-states. It checks if the state labels being used are the same. This heuristic is still a work in progress. The weights assigned to parts of the index are configurable and their values were calibrated based on observations.

Let  $G = (V(G), E(G))$  be a directed graph with vertex-set  $V(G) = \{v_1, \dots, v_n\}$  and edge-set  $E(G) = \{e_1, \dots, e_m\} \subset \{(v_i, v_j) | v_i, v_j \in V(G)\}$ . Let  $C_1, C_2 \in V(G)$  be two distinct vertices of  $G$ . We define the following sets:

$$N_{C_1} = \{v_i \in V(G) \mid (v_i, C_1) \in E(G)\} \quad (1)$$

$$N_{C_2} = \{v_i \in V(G) \mid (v_i, C_2) \in E(G)\} \quad (2)$$

$$S_1 = \frac{|N_{C_1} \cap N_{C_2}|}{\max\{|N_{C_1}|, |N_{C_2}|\}} \quad (3)$$

Let  $f : E(G) \rightarrow \Upsilon$  be a labeling function, where  $\Upsilon$  is a set of labels, and  $f(e) \in \Upsilon$  is the label of edge  $e \in E(G)$ . We define the following sets:

$$L_{C_1} = \{e \mid e = f((v_i, C_1)) \in \Upsilon \text{ and } (v_i, C_1) \in E(G) \text{ and } v_i \in V(G)\} \quad (4)$$

$$L_{C_2} = \{e \mid e = f((v_i, C_2)) \in \Upsilon \text{ and } (v_i, C_2) \in E(G) \text{ and } v_i \in V(G)\} \quad (5)$$

$$S_2 = \frac{|L_{C_1} \cap L_{C_2}|}{\max\{|L_{C_1}|, |L_{C_2}|\}} \quad (6)$$

$$\text{Similarity Index}(S_i) = 0.25 * S_1 + 0.75 * S_2 \quad (7)$$

$S_1$  defines a rate of common OTU vertices with edges for two given characters  $C_1$  and  $C_2$ . The  $S_1$  result lies between 0 (no common OTUs) and 1 (all OTUs are common).  $N_{C_1}$  is the subset

of incoming adjacent vertexes of  $C_1$  and  $N_{C_2}$  is the subset of incoming adjacent vertexes of  $C_2$ . Incoming adjacent vertexes of both  $C_1$  and  $C_2$  are always OTU vertexes, as shown in Figure 2.  $S_2$  defines a rate of common labels of the incoming edges (character-states) for the characters  $C_1$  and  $C_2$ . The  $S_2$  result also lies between 0 (no common character-states) and 1 (all character states are common).  $L_{C_1}$  and  $L_{C_2}$  are the subset of incoming adjacent edge labels (character-states) of  $C_1$  and  $C_2$  respectively.

It is important to note that the character labels of  $C_1$  and  $C_2$  are not being taken into account in the  $S_i$  formula. This intends to avoid weighting in favor of two identical textual characters that do not have the same meaning, and to avoid weighting against two textual characters that are identical but do not have the same meaning. In practice, this will make the solution independent of the label and applicable for both presented scenarios (same label but different meanings and different labels and same meaning). Additionally, the symmetric property of equality is satisfied.

### 3.6.1 Practical Implementation of the Similarity Measure

Our system is able to draw a chart as illustrated in Figure 5, whose algorithm is inspired by the hierarchical edge bundling example (<http://mbostock.github.io/d3/talk/20111116/bundle.html>) of D3.js (<http://d3js.org/>) library. D3.js is a JavaScript library for manipulating documents and it has a wide variety of powerful visualization components. In the case of the hierarchical edge bundling example, it is necessary to provide only a “name” for each node and, inside a related “imports” sentence, the node name to where an edge must be created to. Listing 3 shows the JSON file that encodes the data used to generate Figure 5 (for implementation details see Appendix A.8).

Listing 3: RealExample.json

```

1  [
2  { "name": "root.Cauline cladotaxy" , "imports": ["root.Cauline
3    cladotaxy" , "root.Phyllotaxy" ]},
4  { "name": "root.Protoxylem position within the cauline stele" , "
5    imports": ["root.Protoxylem position within the cauline stele" ]},
6  { "name": "root.Organotaxy of the LBS" , "imports": ["root.Cauline
7    cladotaxy" , "root.Phyllotaxy" ]},
8  { "name": "root.Xylem configuration in the leaflets" , "imports": []}
9  ,
10 { "name": "root.Planation" , "imports": []},
11 { "name": "root.Development of the foliar organ" , "imports": []},
12 { "name": "root.Phyllotaxy" , "imports": []},
13 { "name": "root.Xylem configuration in the rachis" , "imports": []},
14 { "name": "root.Cauline cladotaxy" , "imports": []},
15 { "name": "root.Protoxylem position within the cauline stele" , "
16   imports": []},
17 { "name": "root.Xylem configuration in the rachis" , "imports": []},
18 { "name": "root.Extent of the planation" , "imports": []},
19 { "name": "root.Presence of planated parts within the LBS" , "imports
20   ": []},
21 { "name": "root.Xylem configuration in the leaflets" , "imports": []}
22 ,
23 { "name": "root.Development of the LBS" , "imports": ["root.
24   Development of the foliar organ" ]}
25 ]

```

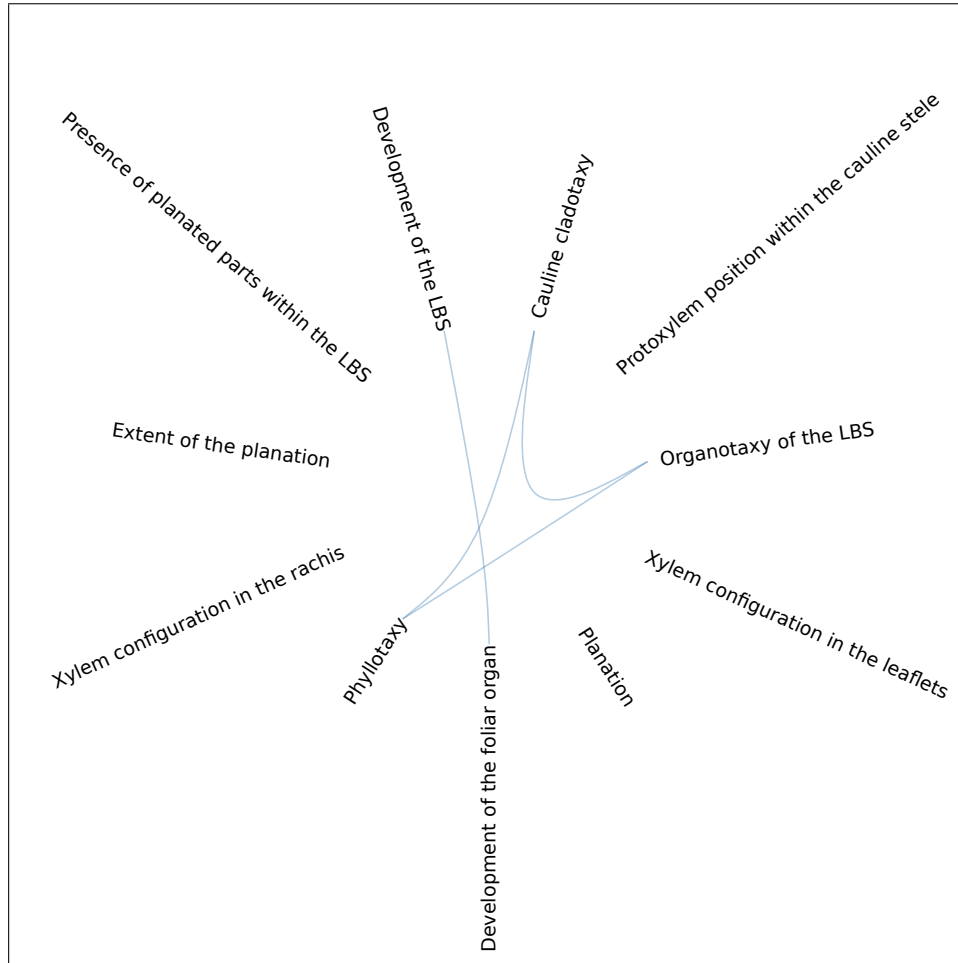


Figure 5: Practical Implementation

### 3.7 Tracing the Evolutionary History

The *TraceEvolutionaryHistory* class abstracts an important algorithm that traces a phylogenetic history of traits changes (for implementation details see Appendix A.9). This algorithm was built on top of our graph data model. It searches in a given tree for traits (characters) that might be the “responsible” for a tree branching, in which branching is considered as any division from a particular ancestor. For example, Figure 4 has two Hypothetical Taxonomic Units (HTU), in which the least nested one after the root has the *Pseudosporochnus* node and another HTU node as children. A typical question that motivated us to create such an algorithm was: What differentiates *Pseudosporochnus* from the other nodes?

The algorithm is divided into two recursive methods that are invoked in sequence. The first one *BottomUpAggregation* starts from a given point in the tree and goes down until it reaches Operational Taxonomic Unit (OTU) nodes. At this point, the method retrieves all outgoing relationships from the OTU node and starts going back towards the root. While the method is traversing internal HTU nodes (*current<sub>HTU</sub>*) from the leaves back towards the root, it performs an union operation with the outgoing relationships of all children nodes – one occurrence for each type of relationship –

and then, for each type of relationship of the resulting union, the method creates an edge departing from the current HTU ( $current_{HTU}$ ) towards the original ending point of the relationship. In the end, the method returns all relationships outgoing from all nodes, including the intermediary HTU nodes ( $current_{HTU}$ ). Figure 6 shows the result of *BottomUpAggregation* method being applied on the graph of Figure 4.

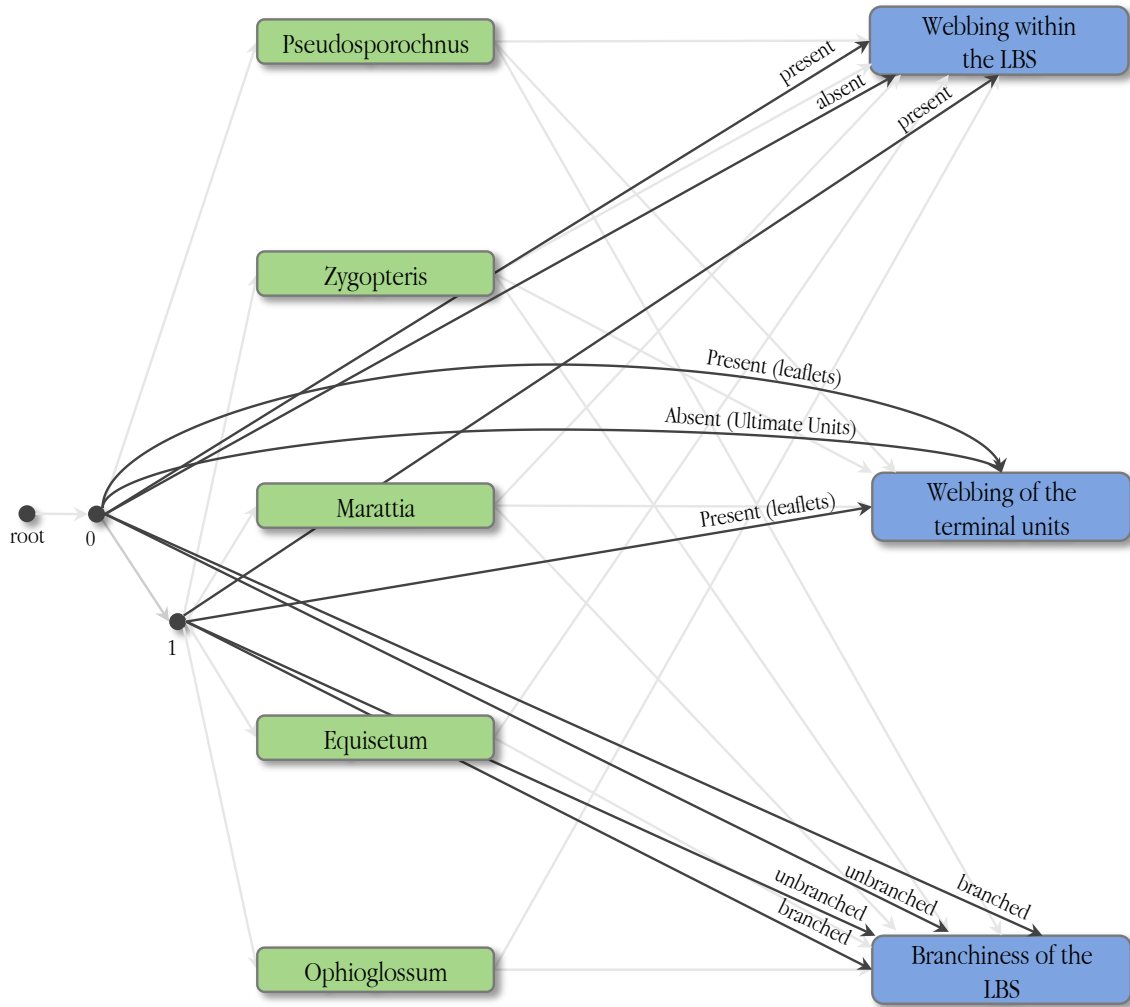


Figure 6: Bottom Up Aggregation

The second part of the algorithm is called *TopDownRefining*. This method is triggered after the *BottomUpAggregation* method, going to the same starting node provided in the *BottomUpAggregation* method. It starts from a given node ( $node_n$ ) traversing down the tree and, in every HTU it reaches, it subtracts the set of character-states that starts in its children nodes ( $node_{children}$ ) and points to a given character ( $node_{character}$ ), from the set of character-states starting from itself ( $node_n$ ) pointing to the same character node ( $node_{character}$ ).

For example, in Figure 6, consider the least nested node ( $node_0$ ), just after the root and linked to the *Webbing within the LBS* character node ( $node_{webbingLBS}$ ). There are two edges connecting the



$node_0$  and the  $node_{webbingLBS}$  with values *present* and *absent*. The *present* edge comes from the most nested part of the tree, composed of the nodes *Zygopteris*, *Marattia*, *Equisetum* and *Ophioglossum*, nested by node 1 ( $node_1$ ) – see Figure 4. The *absent* comes from *Pseudosporochnus* node – see Figure 4.

When the algorithm reaches  $node_0$  it will subtracts the set of character-states (edges) outgoing from *Pseudosporochnus* toward  $node_{webbingLBS}$  from the set of outgoing character-states (edges) outgoing from  $node_0$  toward  $node_{webbingLBS}$ . This set subtraction will be  $\{present, absent\} - \{absent\} = \{present\}$ . If the set subtraction result is not empty, it creates an edge called “*EvolvedTrait*” from itself ( $node_0$ ) toward the character ( $node_{webbingLBS}$ ) as shown in Figure 7.

Also, the algorithm will subtracts the set of character-states (edges) outgoing from  $node_1$  toward  $node_{webbingLBS}$  from the set of character-states (edges) outgoing from  $node_0$  toward  $node_{webbingLBS}$ . This set subtraction will also not be empty ( $\{present, absent\} - \{present\} = \{absent\}$ ) but the “*EvolvedTrait*” edge is created only once between  $node_0$  and  $node_{webbingLBS}$ .

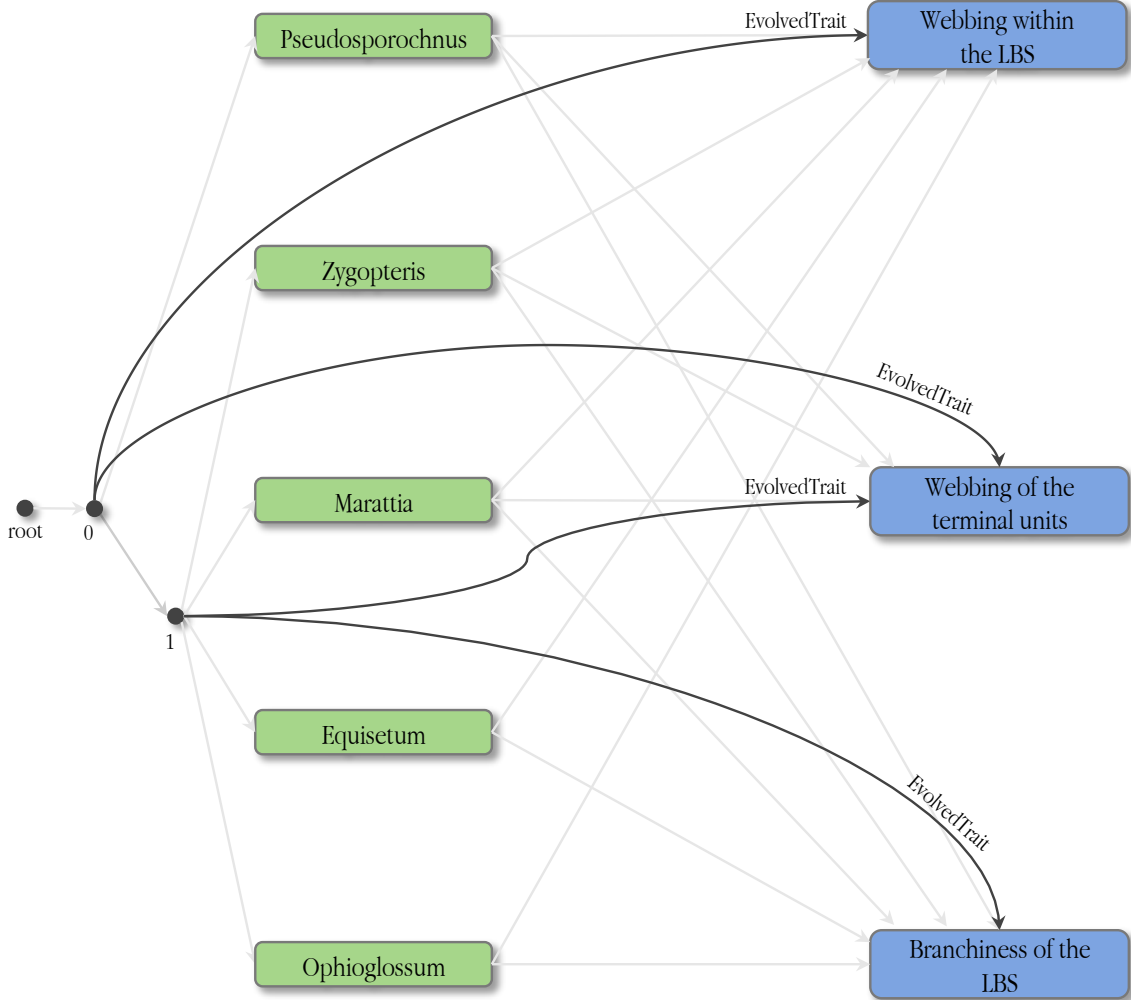


Figure 7: Top Down Refining

In a second iteration, the algorithm will reach  $node_1$  and it will individually subtracts  $node_1$  children nodes (*Zygoteris*, *Marattia*, *Equisetum* and *Ophioglossum*) outgoing character-states toward  $node_{webbingLBS}$  from the set of character-states outgoing from  $node_1$  toward  $node_{webbingLBS}$ . All those set subtractions will be  $\{present\} - \{present\} = \emptyset$ . In such a case (empty set,  $\emptyset$ ), no “*EvolvedTrait*” edge is created, as can be seen in Figure 7.

Finally there is a visual tool that presents to the user the tree structure with all characters flagged with the “*EvolvedTrait*” edge, i.e. the characters that the algorithm “suspect” of being responsible for the branching. Figure 8 is a screenshot of our visual tool.

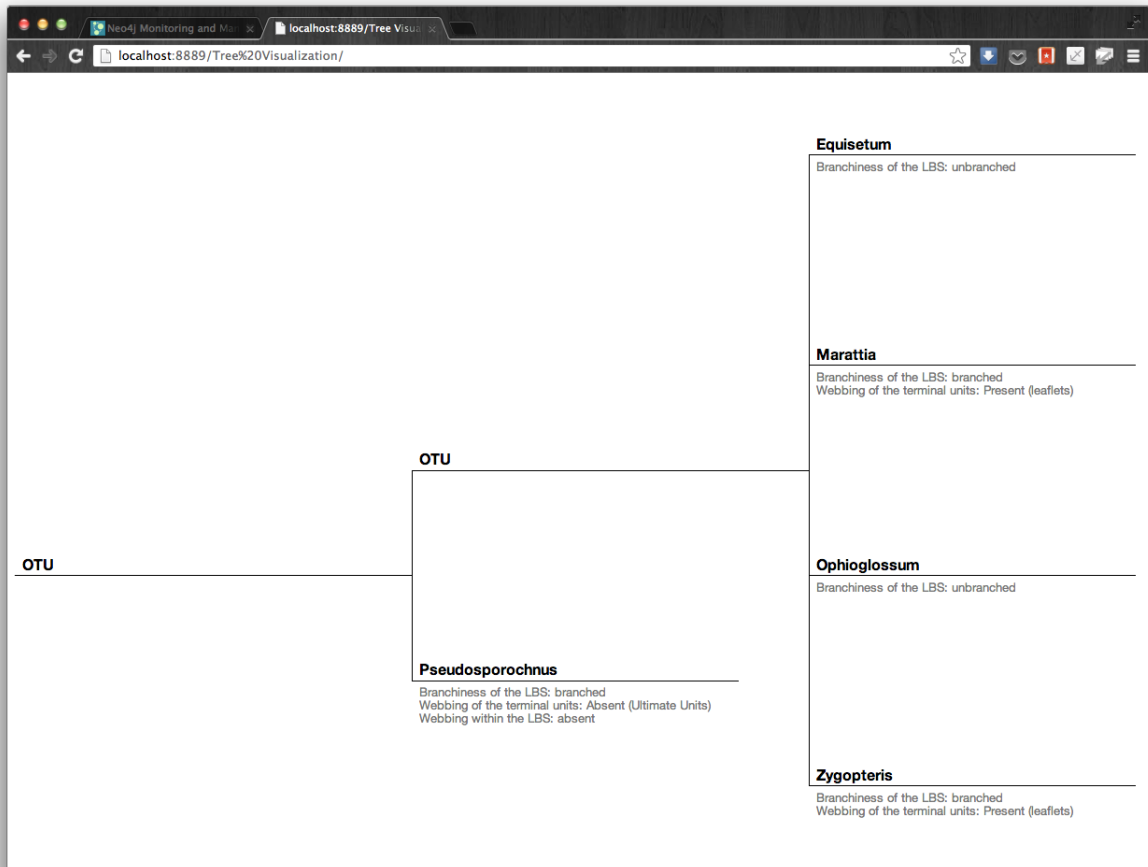


Figure 8: Evolved Traits Visualization

## 4 Conclusion

In this technical report we showed the main functionalities and operational features of the system. We mapped the SDD format to the graph model, remodeling semi-structured descriptions to a graph abstraction, in which the data are linked enabling coupling phylogenetic trees and phenotype descriptions. We drilled down the interconnection process through LSID unification, showing the

required steps to obtain a valid LSID and implementation details of the services used in this process. We presented details regarding a visualization tool implemented on top of the D3.js, to visualize our proposed similarity measure. Such a solution will not only help discovering characters similarity, but will be very important in the next stage of this project, which is the mapping from the graph towards ontologies. Furthermore, an algorithm to trace the phylogenetic history of traits changes has been shown. Finally, *Cypher* database queries and the main classes and methods of the system were provided with detailed comments for each method.

## References

- [1] Miranda, E., Santanchè, A.: Unifying phenotypes to support semantic descriptions. VI Brazilian Conference on Ontological Research (Ontobras) (09 2013)
- [2] Parr, C.S., Guralnick, R., Cellinese, N., Page, R.D.: Evolutionary informatics: unifying knowledge about the diversity of life. *Trends in ecology & evolution* **27**(2) (2012) 94–103
- [3] Ciccarelli, F.D., Doerks, T., Von Mering, C., Creevey, C.J., Snel, B., Bork, P.: Toward automatic reconstruction of a highly resolved tree of life. *Science* **311**(5765) (2006) 1283–1287
- [4] Delsuc, F., Brinkmann, H., Philippe, H.: Phylogenomics and the reconstruction of the tree of life. *Nature Reviews Genetics* **6**(5) (2005) 361–375
- [5] Miller, M.A., Pfeiffer, W., Schwartz, T.: Creating the cypress science gateway for inference of large phylogenetic trees. In: Gateway Computing Environments Workshop (GCE), 2010, IEEE (2010) 1–8
- [6] Mabee, P.M.: Integrating evolution and development: the need for bioinformatics in evo-devo. *BioScience* **56**(4) (2006) 301–309
- [7] Ung, V., Causse, F., Vignes Lebbe, R.: Xper<sup>2</sup>: managing descriptive data from their collection to e-monographs. (2010)
- [8] Ung, V., Dubus, G., Zaragieta-Bagils, R., Vignes-Lebbe, R.: Xper2: introducing e-taxonomy. *Bioinformatics* **26**(5) (2010) 703–704
- [9] Hagedorn, G.: Structuring Descriptive Data of Organisms – Requirement Analysis and Information Models. PhD thesis, Universität Bayreuth, Fakultät für Biologie, Chemie und Geowissenschaften (11 2007)
- [10] Page, R.: Biodiversity informatics: the challenge of linking data and the role of shared identifiers. *Briefings in Bioinformatics* **9**(5) (2008) 345–354
- [11] Godfray, H., et al.: Challenges for taxonomy. *Nature* **417**(6884) (2002) 17–19
- [12] Adler, P.H., Crosskey, R.W.: World blackflies (diptera: Simuliidae): a comprehensive revision of the taxonomic and geographical inventory [2013] (2013) Accessed on July 08 2013.
- [13] Rees, T.: Taxamatch, a "fuzzy" matching algorithm for taxon names, and potential applications in taxonomic databases. In Weitzman, A., Belbin, L., eds.: Provisional Abstracts of the 2008 Annual Conference of the Taxonomic Databases Working Group, Fremantle, Australia, Biodiversity Information Standards (TDWG) and the Missouri Botanical Garden (2008)
- [14] Kennedy, J., Kukla, R., Paterson, T.: Scientific names are ambiguous as identifiers for biological taxa: Their context and definition are required for accurate data integration. In: 2nd Intl. Workshop on Data Integration in the Life Sciences (DILS). LNCS 3615 (July 2005) 80–95
- [15] Patterson, D., Cooper, J., Kirk, P., Pyle, R., Remsen, D.: Names are key to the big new biology. *Trends in ecology & evolution* **25**(12) (2010) 686–691
- [16] Bisby, F.: The quiet revolution: biodiversity informatics and the internet. *Science* **289**(5488) (2000) 2309–2312
- [17] Clark, T., Martin, S., Liefeld, T.: Globally distributed object identification for biological knowledgebases. *Briefings in bioinformatics* **5**(1) (2004) 59–70

- [18] Angles, R.: A comparison of current graph database models. In: Data Engineering Workshops (ICDEW), 2012 IEEE 28th International Conference on. (2012) 171–177
- [19] Angles, R., Gutierrez, C.: Survey of graph database models. ACM Computing Surveys (CSUR) **40**(1) (2008) 1
- [20] Robinson, I., Webber, J., Eifrem, E.: Graph Databases. O'Reilly Media, Inc. (2013)
- [21] Rodriguez, M.A., Shinavier, J.: Exposing multi-relational networks to single-relational network analysis algorithms. Journal of Informetrics **4**(1) (2010) 29 – 41
- [22] Bagils, R.Z., Ung, V., Grand, A., Vignes-Lebbe, R., Cao, N., Ducasse, J.: Lisbeth: New cladistics for phylogenetics and biogeography. Comptes Rendus Palevol **11**(8) (2012) 563 – 566
- [23] Nelson, G., Platnick, N.I.: Three-taxon statements: A more precise use of parsimony? Cladistics **7**(4) (1991) 351–366

## A Demonstration

In this section we present the source code of the system, according to the graph data model presented in previous sections. The code is modularized in files and each file has a class with methods, all with comments explaining their functionality.

### A.1 SDDParser.py

```

1  import os,sys
2
3  from xml.dom import minidom
4  from collections import OrderedDict
5
6  from Representation import *
7  from StateDefinition import *
8  from CategoricalCharacter import *
9  from Categorical import *
10 from CodedDescription import *
11
12 class SDDParser:
13
14     def __init__(self, SDDFile):
15
16         self.CategoricalCharacters = self.__parseCategoricalCharacter(
17             SDDFile )
18         self.CodedDescriptions = self.__parseCodedDescription( SDDFile )
19
20     def __parseRepresentation(self, Repr) :
21         """
22         Representation is a plain text label and description block found
23         inside CategoricalCharacter, StateDefinition and
24         CodedDescription blocks.
25         Args: A XML Representation block and its content.
26         Returns: A SDD Representation object.
27         """
28
29         label = ''
30         detail = ''
31
32         if Repr :
33
34             if 0 < Repr.getElementsByTagName('Label').length :
35                 label = Repr.getElementsByTagName('Label')[0].childNodes[0].
36                     nodeValue.strip()
37
38             if 0 < Repr.getElementsByTagName('Detail').length :
39                 detail = Repr.getElementsByTagName('Detail')[0].childNodes[0].
40                     nodeValue.strip()

```

```

36
37     return Representation( label, detail )
38
39
40 def __parseStateDefinitions(self, StateDefinitions):
41     """
42     StateDefinition has its own id and a Representation block. It is
43     define inside the CategoricalCharacter/States block in which
44     the States groups together all possible states (StateDefinition
45     ) observed at a given Categorical Character.
46     Args: All XML StateDefinition blocks of a particular
47     CategoricalCharacter/States block.
48     Returns: A dictionary of StateDefinition object.
49     """
50
51     # Dictionary with all state definition nodes
52     SStateDefinitionsDictionary = {}
53
54     for State in StateDefinitions :
55
56         Id = State.getAttributeNode('id').nodeValue
57
58         Repr = State.getElementsByTagName('Representation')[0]
59
60         Representation = self.__parseRepresentation( Repr )
61
62         # Add node to Dictionary
63         SStateDefinitionsDictionary[Id] = StateDefinition( Id ,
64             Representation )
65
66     return SStateDefinitionsDictionary
67
68
69 def __parseStates(self, States):
70     """
71     State is define inside CodedDescription/SummaryData/Categorical
72     and it links a taxon CategoricalCharacter to its possible
73     States through the ref parameters.
74     Args: All XML State blocks of a particular CodedDescription/
75     SummaryData/Categorical block.
76     Returns: An array with StateDefinitions references.
77     """
78
79     # Array with references to StateDefinitions
80     StatesDictionary = []
81
82     for state in States :
83
84         ref = state.getAttributeNode('ref').nodeValue

```

```

77     StatesDictionary.append( ref )
78
79     return StatesDictionary
80
81
82 def __parseSummaryData(self, Categoricals):
83     """
84     Categorical is a reference to a CategoricalCharacter object and is
85     composed by a list of references to possible states that a
86     given taxon can take.
87     Args: All XML Categorical blocks of a particular CodedDescription/
88     SummaryData block.
89     Returns: A dictionary of Categorical objects.
90     """
91
92     # Dictionary of Categorical objects
93     SummaryDataDictionary = {}
94
95     for c in Categoricals :
96
97         ref = c.getAttributeNode('ref').nodeValue
98
99         s = c.getElementsByTagName('State')
100
101         States = self.__parseStates( s )
102
103         SummaryDataDictionary[ref] = Categorical( ref , States )
104
105     return SummaryDataDictionary
106
107
108 def __parseCategoricalCharacter(self, SDDFile):
109     """
110     CategoricalCharacter has its own id, a Representation block and a
111     States block.
112     Args: A SDD file name.
113     Returns: A dictionary with all CategoricalCharacters objects in
114     the given file.
115     """
116
117     CC = SDDFile.getElementsByTagName('CategoricalCharacter')
118
119     # Dictionary with all CategoricalCharacters objects
120     CategoricalCharacters = {}
121
122     for Character in CC:
123
124         Id = Character.getAttributeNode('id').nodeValue

```



```

121     States = Character.getElementsByTagName('StateDefinition')
122     Repr = Character.getElementsByTagName('Representation')[0]
123
124     Representation = self.__parseRepresentation( Repr )
125     SStateDefinitionsDictionary = self.__parseStateDefinitions(
126         States )
127
128     CategoricalCharacters[ Id ] = CategoricalCharacter( Id ,
129         SStateDefinitionsDictionary, Representation )
130
131     return CategoricalCharacters
132
133 def __parseCodedDescription(self, SDDFile):
134     """
135     CodedDescription has its own id, a Representation block and a
136     SummaryData block.
137     Args: A SDD file name.
138     Returns: A dictionary with all CodedDescription objects in the
139     given file.
140     """
141
142     CD = SDDFile.getElementsByTagName('CodedDescription')
143
144     # Dictionary with all CodedDescriptions objects
145     CodedDescriptions = {}
146
147     for Description in CD:
148
149         Id = Description.getAttributeNode('id').nodeValue
150
151         SD = Description.getElementsByTagName('Categorical')
152         Repr = Description.getElementsByTagName('Representation')[0]
153
154         Representation = self.__parseRepresentation( Repr )
155         SummaryDataDictionary = self.__parseSummaryData( SD )
156
157         CodedDescriptions[ Id ] = CodedDescription( Id ,
158             SummaryDataDictionary, Representation )
159
160     return CodedDescriptions
161
162 def getAllSates(self):
163     """
164     Returns a dictionary of all 'StateDefinitions' elements.
165     """
166
167     States = {}

```

```
165
166     for key, CategoricalCharacter in self.CategoricalCharacters.
167         iteritems():
168         States.update( CategoricalCharacter.States )
169
170     OrderedStates = OrderedDict( sorted( States.items() ) )
171
172     return OrderedStates
173
174
175 def getAllTaxons(self):
176     """
177     Returns a list of all taxons elements.
178     """
179
180     Taxons = []
181
182     for key, CodedDescription in self.CodedDescriptions.iteritems():
183
184         Taxons.append( CodedDescription.Representation )
185
186     return Taxons
187
188
189 def getAllCharacters(self):
190     """
191     Returns a dictionary of all 'CategoricalCharacter' elements.
192     """
193
194     Characters = {}
195
196     for key, CategoricalCharacter in self.CategoricalCharacters.
197         iteritems():
198
199         Characters[ CategoricalCharacter.id ] = ( CategoricalCharacter.
200             Representation )
201
202     OrderedCharacters = OrderedDict( sorted( Characters.items() ) )
203
204     return OrderedCharacters
```

## A.2 TeeOutput.py

```

1  import re
2  import shlex
3  import mmap
4  import sys
5
6  from TreeNode import *
7  from NodeTypes import *
8
9  class TreeOutput:
10
11     def __init__(self, _TreeOutputFile ):
12
13         self.TreeOutputFile = _TreeOutputFile
14
15
16     def __parseNewickTree( self, NewickTree , parentNode ) :
17         """
18         Newick tree format (New Hampshire tree format) is a way of
19         representing trees in computer-readable form using parentheses
20         and commas.
21         Args:
22         NewickTree: A NewickTree string. For example: ((((((12 18) 22)
23             13) 3) 7) 30) (23 25))
24         parentNode: A node to where NewickTree tree will be attached to.
25         """
26
27         opened = False
28         substring = NewickTree
29
30         i = j = begin = end = 0
31
32         for c in NewickTree :
33
34             if c == '(' :
35                 i += 1
36
37                 if not opened :
38                     begin = j
39
40                 opened = True
41
42             elif c == ')' :
43                 i -= 1
44
45                 if opened and i == 0 :
46                     # (opened and i == 0) means that opening round bracket '(' and
47                     # the corresponding closing round bracket ')' was found.

```

```

44     # It will recursively call _parseNewickTree with brackets
        content. Also, it will remove parentheses block and content
        from NewickTree.
45
46     opened = False
47
48     childrenWithBrackets = NewickTree[ begin : j + 1 ]
49     childrenWithNoBrackets = NewickTree[ begin + 1 : j ]
50
51     child = TreeNode(None)
52     parentNode.appendChild( child )
53
54     self._parseNewickTree( childrenWithNoBrackets , child )
55
56     substring = substring.replace( childrenWithBrackets, "" )
57
58     j += 1
59
60     if "(" not in substring:
61         # When this condition is satisfied, it means that substring will
            only have leaves nodes or it is empty.
62
63         my_splitter = shlex.shlex(substring, posix = True )
64         my_splitter.whitespace += ',,'
65         my_splitter.whitespace_split = True
66
67         for n in my_splitter :
68             parentNode.appendChild( TreeNode(n) )
69
70
71     def getNewickTree( self ) :
72         """
73         This method looks into the file in search of the Newick Tree and
            returns it.
74         The process is pretty straightforward:
75         1. Set the file's current position to the o occurrence of '
            Retained trees'
76         2. Reads this line and discards it
77         3. Reads the next line, which supposedly should contain the
            Newick Tree
78         4. Get the Newick Tree
79         """
80
81         _file = open( self.TreeOutputFile )
82         memorymap = mmap.mmap( _file.fileno(), 0, access = mmap.
            ACCESS_READ )
83
84         RetainedTreesPosition = memorymap.find("Retained trees")
85         memorymap.seek( RetainedTreesPosition )

```

```

86 | memorymap.readline()
87 | FirstRetainedTreeLine = memorymap.readline()
88 | memorymap.close()
89 |
90 | # First occurrence of ')'
91 | begin = FirstRetainedTreeLine.find(')')
92 |
93 | # Last occurrence of '('
94 | end = FirstRetainedTreeLine.rfind('(')
95 |
96 | NewickTree = FirstRetainedTreeLine[ begin : end + 1 ]
97 |
98 | return NewickTree
99 |
100 |
101 | def getTaxons( self ) :
102 |     """
103 |     Get all taxons listed right bellow 'Taxa (# taxons)' inside <D02
104 |     > block and return all those taxons.
105 |     """
106 |
107 |     _file = open( self.TreeOutputFile )
108 |     memorymap = mmap.mmap( _file.fileno(), 0, access = mmap.
109 |         ACCESS_READ )
110 |
111 |     BlockBegin = memorymap.find("<D02>")
112 |     BlockEnd = memorymap.find("<F02>")
113 |
114 |     TaxaPosition = memorymap.find( "Taxa" , BlockBegin , BlockEnd )
115 |
116 |     memorymap.seek( TaxaPosition )
117 |     TaxaLine = memorymap.readline()
118 |     TotalTaxa = int( re.search( re.escape( '(' ) + "(.*?)" + re.escape
119 |         ( ')' ) , TaxaLine ).group(1) )
120 |
121 |     TaxonsDictionary = {}
122 |
123 |     for i in range( TotalTaxa ) :
124 |
125 |         line = memorymap.readline()
126 |
127 |         index = line[ 1 : 21 ].strip()
128 |         taxon = line[ 22 : ].strip()
129 |
130 |         TaxonsDictionary[ index ] = taxon
131 |
132 |     memorymap.close()
133 |
134 |     return TaxonsDictionary

```

```
132
133
134 def __RenameTreeNodes(self, subTree , TaxonsDictionary ):
135     """
136     In a rooted phylogenetic tree, each node is called a taxonomic
137     unit. Internal nodes are generally called hypothetical
138     taxonomic units (HTUs) as they cannot be directly observed.
139     Args:
140     subTree: Is a branch of the tree.
141     TaxonsDictionary: A list of taxons present in the 3iz file.
142     """
143     if subTree.nodes:
144         subTree.value = str( NodeType.HTU )
145
146         for n in subTree.nodes:
147             self.__RenameTreeNodes( n, TaxonsDictionary )
148
149     else:
150         subTree.value = TaxonsDictionary[ subTree.value ]
151
152
153
154 def getTaxonsTreeStructure( self ):
155     """
156     It parse the NewickTree string into a tree structure with
157     Hypothetical Taxonomic Units as internal nodes and the correct
158     Taxon name as the leaves.
159     """
160     NewickTree = self.getNewickTree()
161     TaxonsDictionary = self.getTaxons()
162
163     root = TreeNode(None)
164     self.__parseNewickTree( NewickTree , root )
165
166     self.__RenameTreeNodes( root , TaxonsDictionary )
167
168     return root
```

### A.3 GlobalNamesResolver.py

```

1  from bs4 import BeautifulSoup
2  from GNRResultObject import *
3  import urllib2
4  from enumerator import *
5
6  class GlobalNamesResolver:
7
8  def __init__(self) :
9      self.url = 'http://resolver.globalnames.org/name_resolvers.xml?
          names='
10
11     # Names Data Sources <http://resolver.globalnames.org/data_sources
          >
12     # ID  Source
13     # 169 uBio NameBank
14     #   1 Catalogue of Life
15     #   3 ITIS
16     self.DataSources = enum( CatalogueOfLife = 1, ITIS = 3,
          uBioNameBank = 169 )
17
18     self.DataSourceIds = [self.DataSources.CatalogueOfLife, self.
          DataSources.ITIS, self.DataSources.uBioNameBank ]
19
20
21 def getResultsObjects( self, ScientificName ):
22
23     ScientificName = ScientificName.replace(' ', '%20')
24
25     url = self.url + ScientificName
26
27     if len( self.DataSourceIds ) > 0 :
28         url = url + '&data_source_ids='
29
30         for _id in self.DataSourceIds :
31             url = url + str( _id ) + '|'
32
33     try:
34         GNRServiceUrlResponse = urllib2.urlopen( url ).read()
35
36     except urllib2.HTTPError, e:
37         print "HTTP error: %d" % e.code
38     except urllib2.URLError, e:
39         print "Network error: %s" % e.reason.args[1]
40
41     SoupGNRRResponse = BeautifulSoup( GNRServiceUrlResponse )
42
43     results = SoupGNRRResponse.findAll('result')
```

```

44
45 GNRResultObjects = []
46
47 for result in results :
48
49     DataSourceId      = result.find('data-source-id', {'type': 'integer
50         '})
51     DataSourceTitle  = result.find('data-source-title')
52     gniUUID          = result.find('gni-uuid')
53     NameString       = result.find('name-string')
54     CanonicalForm    = result.find('canonical-form')
55     TaxonId          = result.find('taxon-id')
56     LocalId          = result.find('local-id')
57     MatchType        = result.find('match-type', {'type': 'integer'})
58     Prescore         = result.find('prescore')
59     Score            = result.find('score', {'type': 'float'})
60
61     DataSourceId      = DataSourceId.contents[0]      if DataSourceId
62         else ""
63     DataSourceTitle  = DataSourceTitle.contents[0]   if DataSourceTitle
64         else ""
65     gniUUID          = gniUUID.contents[0]           if gniUUID
66         else ""
67     NameString       = NameString.contents[0]       if NameString
68         else ""
69     CanonicalForm    = CanonicalForm.contents[0]    if CanonicalForm
70         else ""
71     TaxonId          = TaxonId.contents[0]           if TaxonId
72         else ""
73     LocalId          = LocalId.contents[0]           if LocalId
74         else ""
75     MatchType        = MatchType.contents[0]        if MatchType
76         else ""
77     Prescore         = Prescore.contents[0]          if Prescore
78         else ""
79     Score            = Score.contents[0]             if Score
80         else ""
81
82     obj = GNRResultObject( DataSourceId, DataSourceTitle, gniUUID,
83         NameString, CanonicalForm, TaxonId, LocalId, MatchType ,
84         Prescore , Score )
85
86     GNRResultObjects.append( obj )
87
88 return GNRResultObjects
89
90 def getCanonicalForm( self, ScientificName ) :
91     """

```



```

80 Returns the canonical forms of a given scientific name.
81 """
82
83 objects = self.getResultsObjects( ScientificName )
84
85 CanonicalForms = set( [ ] )
86
87 for obj in objects :
88
89     match = int(obj.MatchType)
90
91     # 1 - Exact match
92     # 2 - Exact match by canonical form
93     # 3 - Fuzzy match by canonical form
94     if match == 1 or match == 2 or match == 3 :
95
96         if 0.988 <= float(obj.MatchType) :
97
98             # Add canonical form to the set
99             CanonicalForms = CanonicalForms | set( [ obj.CanonicalForm ] )
100
101         if 1 == len( CanonicalForms ):
102
103             return sorted(CanonicalForms)[0]
104
105         return None
106
107
108 def getLSIDFromCanonicalForm( self, CanonicalForm ) :
109     """
110     Returns the LSID of a given Canonical Form. Only uBio NameBank
111     LSID are retrieved and still only if a exact match occur.
112     """
113
114     ResultsObjects = self.getResultsObjects( CanonicalForm )
115
116     for obj in ResultsObjects :
117
118         if int(obj.MatchType) == 1:
119
120             if int(obj.DataSourceId) == self.DataSources.uBioNameBank:
121
122                 return obj.LocalId
123
124     return None

```

## A.4 GNRResultObject.py

```

1 class GNRResultObject:
2
3     def __init__(self, _DataSourceId, _DataSourceTitle, _gniUUID,
4                 _NameString, _CanonicalForm, _TaxonId, _LocalId, _MatchType,
5                 _Prescore, _Score ):
6
7         # The id of the data source where a name was found.
8         self.DataSourceId = _DataSourceId
9
10        # The data source title where a name was found.
11        self.DataSourceTitle = _DataSourceTitle
12
13        # An identifier for the found name string used in Global Names.
14        self.gniUUID = _gniUUID
15
16        # The name string found in this data source.
17        self.NameString = _NameString
18
19        # A "canonical" version of the name generated by the Global Names
20        # parser
21        self.CanonicalForm = _CanonicalForm
22
23        # Tree path to the root if a name string was found within a data
24        # source classification.
25        # self.ClassificationPath
26
27        # self.ClassificationPathRanks
28
29        # Same tree path using taxon_ids
30        # self.ClassificationPathIds
31
32        # An identifier supplied in the source Darwin Core Archive for the
33        # name string record
34        self.TaxonId = _TaxonId
35
36        # Shows id local to the data source (if provided by the data
37        # source manager)
38        self.LocalId = _LocalId
39
40        # Explains how resolver found the name. If the resolver cannot
41        # find names corresponding to the entire queried name string, it
42        # sequentially removes terminal portions of the name string until
43        # a match is found.
44
45        # 1 – Exact match
46        # 2 – Exact match by canonical form of a name
47        # 3 – Fuzzy match by canonical form
48        # 4 – Partial exact match by species part of canonical form

```

```
39 # 5 - Partial fuzzy match by species part of canonical form
40 # 6 - Exact match by genus part of a canonical form
41 self.MatchType = _MatchType
42
43 # Displays points used to calculate the score delimited by '|' —
    "Match points|Author match points|Context points". Negative
    points decrease the final result.
44 self.Prescore = _Prescore
45
46 # A confidence score calculated for the match.
47 # 0.5 means an uncertain result that will require investigation.
48 # Results higher than 0.9 correspond to 'good' matches.
49 # Results between 0.5 and 0.9 should be taken with caution.
50 # Results less than 0.5 are likely poor matches.
51 # The scoring is described in more details on http://resolver.globalnames.org/about
52 self.Score = _Score
```

## A.5 ITIServices.py

```

1  import suds
2
3  class ITIServices:
4
5      url = "http://www.itis.gov/ITISWebService.xml"
6      client = None
7
8      def __init__(self):
9          self.client = suds.client.Client( self.url )
10
11
12     def getTSNfromScientificName(self, ScientificName ):
13         """
14         Taxonomic Serial Number (TSN) which is the primary key for the
15         scientific name. This method returns a TSN if the provided
16         ScientificName is found and None otherwise.
17         """
18
19         self.client.service.searchByScientificName( ScientificName )
20
21         ScientificNamesResponse = self.client.last_received().getChild("
22             soapenv:Envelope").getChild("soapenv:Body").getChild("ns:
23             searchByScientificNameResponse").getChild("ns:return").
24             getChildren("ax21:scientificNames")
25
26         for sn in ScientificNamesResponse:
27
28             tsn = sn.getChild("ax21:tsn")
29
30             if tsn != None :
31                 return tsn.getText()
32
33         return None
34
35     def getLSIDfromTSN(self, tsn ):
36         """
37         Given a TSN this method returns a LSID if found and None otherwise.
38         """
39
40         self.client.service.getLSIDFromTSN( tsn )
41
42         LSID = self.client.last_received().getChild("soapenv:Envelope").
43             getChild("soapenv:Body").getChild("ns:getLSIDFromTSNResponse").
44             getChild("ns:return").getText()
45
46         if LSID:

```

```
41     return LSID
42
43     return None
```

## A.6 CoLServices.py

```

1 from BeautifulSoup import BeautifulSoup
2 import urllib2
3
4 class CoLServices:
5     """
6     This class contains the main methods to interact with the CoL web
7     service.
8     """
9     def getCoLUrl( self, ScientificName ):
10        """
11        This method uses a XML scraping technique to get the URL of the
12        given Scientific Name from the webservice response.
13        """
14        url = 'http://www.catalogueoflife.org/col/webservice?name='
15
16        ScientificName = ScientificName.replace(' ', '%20')
17
18        try:
19            CoLWebServiceUrlResponse = urllib2.urlopen(url + ScientificName).
20                read()
21            except urllib2.HTTPError, e:
22                print "HTTP error: %d" % e.code
23            except urllib2.URLError, e:
24                print "Network error: %s" % e.reason.args[1]
25
26        SoupCoLWebServiceResponse = BeautifulSoup(CoLWebServiceUrlResponse
27            )
28
29        tagresult = SoupCoLWebServiceResponse.findAll('result')
30
31        CoLUrl = tagresult[0].find('url').contents[0]
32
33        if CoLUrl:
34            return CoLUrl
35
36    def getCoLSpecieID( self, ScientificName ):
37        """
38        This method uses a XML scraping technique to get the ID of the
39        given Scientific Name from the webservice response.
40        """
41
42        url = 'http://www.catalogueoflife.org/testcol/webservice?name='
43
44        ScientificName = ScientificName.replace(' ', '%20')

```

```

43
44     try:
45         CoLWebServiceUrlResponse = urllib2.urlopen(url + ScientificName).
            read()
46     except urllib2.HTTPError, e:
47         print "HTTP error: %d" % e.code
48     except urllib2.URLError, e:
49         print "Network error: %s" % e.reason.args[1]
50
51     SoupCoLWebServiceResponse = BeautifulSoup(CoLWebServiceUrlResponse
        )
52
53     result = SoupCoLWebServiceResponse.find('result')
54
55     if result:
56         findID = result.find('id')
57
58         if findID:
59             SpecieID = findID.contents[0]
60
61             if SpecieID:
62                 return SpecieID
63
64     return None
65
66
67 def getLSIDfromSpecieID( self, SpecieID ):
68     """
69     This method uses a HTML screen-scraping technique to get the LSID
70     of the given SpecieID.
71     """
72
73     url = 'http://www.catalogueoflife.org/testcol/details/species/id/'
74
75     try:
76         SpecieDetailsCoLUrlResponse = urllib2.urlopen(url + SpecieID).
            read()
77     except urllib2.HTTPError, e:
78         print "HTTP error: %d" % e.code
79     except urllib2.URLError, e:
80         print "Network error: %s" % e.reason.args[1]
81
82     SoupSpecieDetailsCoLUrlResponse = BeautifulSoup(
83         SpecieDetailsCoLUrlResponse)
84
85     LSID = SoupSpecieDetailsCoLUrlResponse.find('span', {'class': '
        lsid'}).contents[0]
86
87     return LSID

```

```
86
87
88 def getLSIDfromSpecieUrl( self, SpecieUrl ):
89     """
90     This method uses a HTML screen-scraping technique to get the LSID
91     of the given SpecieUrl.
92     """
93     try:
94         SpecieDetailsCoLUrlResponse = urllib2.urlopen(SpecieUrl).read()
95     except urllib2.HTTPError, e:
96         print "HTTP error: %d" % e.code
97     except urllib2.URLError, e:
98         print "Network error: %s" % e.reason.args[1]
99
100     SoupSpecieDetailsCoLUrlResponse = BeautifulSoup(
101         SpecieDetailsCoLUrlResponse)
102     LSID = SoupSpecieDetailsCoLUrlResponse.find('span', {'class': '
103         lsid'}).contents[0]
104     return LSID
```



## A.7 GraphImporter.py

```

1  from py2neo import rest, neo4j, cypher
2
3  from SDDParser import *
4  from TreeOutput import *
5  from GlobalNamesResolver import *
6  from GraphDB import *
7  from NodeTypes import *
8  from RelationshipTypes import *
9  from ITISServices import *
10 from CoLServices import *
11
12 class GraphImporter:
13
14     SDDFilename = None
15     TreeFilename = None
16
17     def __init__(self, _SDDFilename, _TreeFilename, _IgnoreTreeFilename
18                 ):
19
20         self.SDDFilename = _SDDFilename
21         self.TreeFilename = _TreeFilename
22         self.IgnoreTreeFilename = _IgnoreTreeFilename
23
24     def __CreateTaxonsNodes(self, CodedDescriptions ) :
25         """
26         Add to the Graph DB all taxons elements as nodes. In case the
27         taxon node already exists, it uses the node in GraphDB rather
28         than create a new one.
29         Args: CodedDescriptions: A list of all Coded Descriptions elements
30         Returns: A dict mapping keys to the corresponding added nodes.
31         Each tuple is represented as (Taxon Name , node) where the
32         first element of the tuple is the taxon name and the last one
33         is the node itself.
34         Example:
35         {u'Equisetum' : [Node('http://localhost:7474/db/data/node/142')
36                        ],
37          u'Marattia' : [Node('http://localhost:7474/db/data/node/131')],
38          u'Botryopteris': [Node('http://localhost:7474/db/data/node/222')]
39         }
40         """
41
42     gdb = GraphDB()
43     GDBConn, msg = gdb.getPy2neoGraphDatabaseService()
44
45     if GDBConn is not None:

```

```

39
40 # Dictionary for all taxons nodes
41 TaxonsNodes = {}
42
43 GNR = GlobalNamesResolver()
44 ITIS = ITIServices()
45 CoL = CoLServices()
46
47 for key, CodedDescription in CodedDescriptions.iteritems():
48
49     node = None
50
51     taxonName = CodedDescription.Representation.label
52     taxonNameCF = GNR.getCanonicalForm( taxonName )
53
54     lsid = GNR.getLSIDFromCanonicalForm( taxonNameCF )
55
56     if lsid == None :
57         tsn = ITIS.getTSNfromScientificName( taxonNameCF )
58         lsid = ITIS.getLSIDfromTSN( tsn )
59
60         if lsid == None :
61             SpecieID = CoL.getCoLSpecieID( taxonNameCF )
62             lsid = CoL.getLSIDfromSpecieID( SpecieID )
63
64     n = gdb.getNodeByLSID( lsid )
65
66     if n is None:
67
68         # Create taxon node
69         node = GDBConn.create({ 'label' : taxonNameCF ,
70                               'detail' : CodedDescription.Representation.detail ,
71                               'sourceId' : CodedDescription.id ,
72                               'type' : str( NodeType.OTU ) ,
73                               'LSID' : lsid })
74     else :
75         node = n
76
77     # Add node to Dictionary
78     TaxonsNodes[ CodedDescription.Representation.label ] = node
79
80     return TaxonsNodes
81
82 else:
83     print msg
84     return None
85
86
87 def __CreateStateDefinitionNodes( self, StateDefinitions ) :
```

```

88     """
89     Add to the Graph DB all state definition elements as nodes.
90     Args:
91     StateDefinitions: A dictionary of all 'StateDefinitions' elements
92     Returns: A dict mapping keys to the corresponding added nodes. Each
93             tuple is represented as (Id , node) where the first element of
94             the tuple, Id (For example: s54) is the SDD.XML
95             StateDefinition ID and the last one is the node itself.
96     Example:
97     {u's54': [Node('http://localhost:7474/db/data/node/142')],
98         u's43': [Node('http://localhost:7474/db/data/node/131')],
99         u's46': [Node('http://localhost:7474/db/data/node/222')]}
100     """
101
102     gdb = GraphDB()
103     GDBConn, msg = gdb.getPy2neoGraphDatabaseService()
104
105     if GDBConn is not None:
106
107         # Dictionary for all state definition nodes
108         StateDefinitionsNodes = {}
109
110         for key, State in StateDefinitions.iteritems():
111
112             # Create state definition node
113             node = GDBConn.create({ 'label' : State.Representation.label ,
114                                     'detail' : State.Representation.detail ,
115                                     'sourceId' : State.id ,
116                                     'type' : str( NodeTypes.description ) })
117
118             # Add node to Dictionary
119             StateDefinitionsNodes[ State.id ] = node
120
121         return StateDefinitionsNodes
122
123     else:
124         print msg
125         return None
126
127 def __CreateCharacterNodes( self, Characters ) :
128     """
129     Add to the Graph DB all characters elements as nodes.
130     Args:
131     Characters: A dictionary of all 'Characters' elements.
132     Returns: A dict mapping keys to the corresponding added nodes.
133             Each tuple is represented as (Id , node) where the first
134             element of the tuple, Id (For example: c19) is the SDD.XML

```

```

    CategoricalCharacter ID and the last one is the node itself.
Example:
131 {u'c19': [Node('http://localhost:7474/db/data/node/396')],
132      u'c18': [Node('http://localhost:7474/db/data/node/395')],
133      u'c5': [Node('http://localhost:7474/db/data/node/400')]}
134 """
135
136
137 gdb = GraphDB()
138 GDBConn, msg = gdb.getPy2neoGraphDatabaseService()
139
140 if GDBConn is not None:
141
142     # Dictionary for all characters nodes
143     CharactersNodes = {}
144
145     for ID, Character in Characters.iteritems():
146
147         # Create state definition node
148         node = GDBConn.create({ 'label' : Character.label ,
149                                'detail' : Character.detail ,
150                                'sourceId' : ID ,
151                                'type' : str( NodeTypes.description ) })
152
153         # Add node to Dictionary
154         CharactersNodes[ ID ] = node
155
156     return CharactersNodes
157
158 else:
159     print msg
160     return None
161
162
163 def __JoinTaxonsNodesTreeStructureRecursion(self, TaxonsNodes ,
164       subTree, parentNode ):
165
166     gdb = GraphDB()
167     GDBConn, msg = gdb.getPy2neoGraphDatabaseService()
168
169     if GDBConn is not None:
170
171         if subTree.nodes:
172
173             # Create Hypothetical Taxonomic Unit node
174             htunode, = GDBConn.create({ 'label' : str( NodeTypes.HTU ) ,
175                                       'type' : str( NodeTypes.HTU ) })
176
177             # Join Hypothetical Taxonomic Unit node to its parent node

```

```

177     parentNode.create_relationship_to( htuNode , str(
178         RelationshipTypes.TreeEdge ) , { "type" : str(
179             RelationshipTypes.TreeEdge ) } )
180
181     for n in subTree.nodes:
182         self.__JoinTaxonsNodesTreeStructureRecursion( TaxonsNodes , n,
183             htuNode )
184
185     else:
186         # Get Taxonomic Unit (taxon name) already created, passed
187         # through TaxonsNodes dictionary
188         tuNode = TaxonsNodes[ subTree.value ][0]
189
190         # Join Taxonomic Unit node to its parent node
191         parentNode.create_relationship_to( tuNode , str(
192             RelationshipTypes.TreeEdge ) , { "type" : str(
193                 RelationshipTypes.TreeEdge ) } )
194
195     else:
196         print msg
197         return None
198
199 def __JoinTaxonsNodesTreeStructure(self, TaxonsNodes , Tree ):
200     """
201     Join taxons nodes with the Newick tree structure.
202     """
203
204     gdb = GraphDB()
205     GDBConn, msg = gdb.getPy2neoGraphDatabaseService()
206
207     if GDBConn is not None:
208
209         self.__JoinTaxonsNodesTreeStructureRecursion( TaxonsNodes , Tree,
210             gdb.getRootNode() )
211
212     else:
213         print msg
214         return None
215
216 def ImportUsingTaxonCharacterStateSchema( self ) :
217     """
218     Schema : Taxon(Node) -> CategoricalCharacter (Edge) ->
219             StateDefinition (Node)
220     """
221
222     # Parse the SDD-XML file
223     SDDFile = minidom.parse( self.SDDFilename )

```

```

218
219 SDD = SDDParser( SDDFile )
220
221 CategoricalCharacters = SDD.CategoricalCharacters
222 CodedDescriptions = SDD.CodedDescriptions
223
224 # Create Taxons nodes in the Graph DB
225 TaxonsNodes = self.__CreateTaxonsNodes( SDD.CodedDescriptions )
226
227 # Join Taxons nodes in a tree structure
228 treeOutput = TreeOutput( self.TreeFilename )
229 tree = treeOutput.getTaxonsTreeStructure()
230 self.__JoinTaxonsNodesTreeStructure( TaxonsNodes, tree )
231
232 # Create State Definition nodes in the Graph DB
233 StateDefinitionsNodes = self.__CreateStateDefinitionNodes( SDD.
    getAllSates() )
234
235 for key, CodedDescription in CodedDescriptions.iteritems():
236
237     # Check if the given key exists in the dictionary. Otherwise does
    not proceed by creating the relationship
238     if CodedDescription.Representation.label in TaxonsNodes:
239
240         for key, SummaryData in CodedDescription.SummaryData.iteritems():
241             :
242             States = CategoricalCharacters[SummaryData.ref].States
243
244         for StateRef in SummaryData.States :
245
246             # Check if the given key exists in the dictionary. Otherwise
    does not proceed by creating the relationship
247             if StateRef in StateDefinitionsNodes:
248
249                 taxonNode = TaxonsNodes[ CodedDescription.Representation.
    label ][0]
250                 StateDefinitionsNode = StateDefinitionsNodes[ StateRef ][0]
251
252                 CategoricalCharacter = CategoricalCharacters[ SummaryData.
    ref ].Representation
253                 CategoricalCharacterDetail = CategoricalCharacter.detail if
    CategoricalCharacter.detail else ""
254                 relationshipType = CategoricalCharacter.label.replace( ' '
    , '_' )
255
256
257 # Join Taxon nodes to State Definition node using
    CategoricalCharacter.label as relationship

```

```

258     taxonNode.create_relationship_to( StateDefinitionsNode ,
259         relationshipType , { "label" : relationshipType ,
260             "type" : str( RelationshipTypes.
261                 descriptor ) ,
262             "Detail" : CategoricalCharacterDetail }
263         )
264
265 def ImportUsingTaxonStateCharacterSchema( self ) :
266     """
267     Schema : Taxon (Node) -> StateDefinition (Edge) ->
268         CategoricalCharacter (Node)
269     """
270
271     # Parse the SDD-XML file
272     SDDFile = minidom.parse( self.SDDFilename )
273
274     SDD = SDDParser( SDDFile )
275
276     CategoricalCharacters = SDD.CategoricalCharacters
277     CodedDescriptions = SDD.CodedDescriptions
278
279     # Create Taxons nodes in the Graph DB
280     TaxonsNodes = self.__CreateTaxonsNodes( SDD.CodedDescriptions )
281
282     # Join Taxons nodes in a tree structure
283     treeOutput = TreeOutput( self.TreeFilename )
284     tree = treeOutput.getTaxonsTreeStructure()
285     self.__JoinTaxonsNodesTreeStructure( TaxonsNodes, tree )
286
287     # Create Characters nodes in the Graph DB
288     CharactersNodes = self.__CreateCharacterNodes( SDD.
289         getAllCharacters() )
290
291     for key, CodedDescription in CodedDescriptions.iteritems():
292
293         # Check if the given key exists in the dictionary. Otherwise does
294         # not proceed by creating the relationship.
295         if CodedDescription.Representation.label in TaxonsNodes:
296
297             for key, SummaryData in CodedDescription.SummaryData.iteritems():
298                 :
299
300             # Check if the given key exists in the dictionary. Otherwise
301             # does not proceed by creating the relationship.
302             if SummaryData.ref in CharactersNodes:
303
304                 States = CategoricalCharacters[ SummaryData.ref ].States

```

```
299     for StateRef in SummaryData.States :
300
301         taxonNode = TaxonsNodes[ CodedDescription.Representation.
302             label ][0]
303         CharacterNode = CharactersNodes[ SummaryData.ref ][0]
304
305         StateDefinition = States[ StateRef ].Representation
306         StateDefinitionDetail = StateDefinition.detail if
307             StateDefinition.detail else ""
308         relationshipType = StateDefinition.label.replace(' ', '_')
309
310         # Join Taxon nodes to Categorical Character node using
311         StateDefinition.label as relationship
312         taxonNode.create_relationship_to( CharacterNode ,
313             relationshipType , { "label" : relationshipType ,
314                 "type" : str( RelationshipTypes.
315                     descriptor ) ,
316                 "Detail" : StateDefinitionDetail } )
```



## A.8 SimilarityIndex.py

```

1  from __future__ import division
2  import codecs
3  from py2neo import rest, neo4j, cypher
4  from GraphDB import *
5  from NodeAndRelationshipTypes import *
6
7  class SimilarityIndex :
8
9      def CalculateIndex(self, gdb, n1, n2 ) :
10
11         TAaux = gdb.getIncomingAdjacentNodes( n1 )
12         TBaux = gdb.getIncomingAdjacentNodes( n2 )
13
14         TA = []
15         for n in TAaux : TA.append( n[0] )
16
17         TB = []
18         for n in TBaux : TB.append( n[0] )
19
20         setTA = set( TA )
21         setTB = set( TB )
22
23         S1 = len( setTA & setTB ) / max( len( setTA ) , len( setTB ) )
24
25         TE1aux = gdb.getIncomingAdjacentRelationships( n1 )
26         TE2aux = gdb.getIncomingAdjacentRelationships( n2 )
27
28         TE1 = []
29         for r in TE1aux: TE1.append( r[0]["label"] )
30
31         TE2 = []
32         for r in TE2aux: TE2.append( r[0]["label"] )
33
34         setTE1 = set( TE1 )
35         setTE2 = set( TE2 )
36
37         S2 = len( setTE1 & setTE2 ) / max( len( setTE1 ) , len( setTE2 ) )
38
39         SI = ( 0.25 * S1 + 0.75 * S2 )
40
41         return SI
42
43
44     def CompareStudies(self, TreeRootStudyA, TreeRootStudyB,
45         LowerBoundary, JSONFilename ):
46         """

```

```

46 | It calculates the Similarity Index for all characters between two
    | studies taking them two by two. Only SI greater or equal to
    | LowerBoundary are exported into the given Json file.
47 | Args:
48 |   TreeRootStudyA: Study A tree root.
49 |   TreeRootStudyB: Study B tree root.
50 |   LowerBoundary: Lower Boundary condition.
51 |   JSONFilename: Filename where the JSON data should be saved.
52 | """
53 |
54 | gdb = GraphDB()
55 |
56 | rangeA = gdb.getDescriptionNodesOfATree( TreeRootStudyA )
57 | rangeB = gdb.getDescriptionNodesOfATree( TreeRootStudyB )
58 |
59 | Similarity = SimilarityIndex()
60 |
61 | JSON = "["
62 |
63 | for i in rangeA :
64 |
65 |   ni = gdb.getNode( i )
66 |
67 |   JSON = JSON + "\n" + '{' + "\"name\": \"{0}\" , \"imports\": [".
    |       format( "root." + ni["label"] )
68 |
69 |   imports = False
70 |
71 |   for j in rangeB :
72 |
73 |     nj = gdb.getNode( j )
74 |
75 |     SI = Similarity.CalculateIndex( gdb, ni, nj )
76 |
77 |     if LowerBoundary <= SI :
78 |       JSON = JSON + "\"{0}\" , ".format( "root." + nj["label"] )
79 |       imports = True
80 |
81 |   if imports :
82 |     # Remove the last comma
83 |     JSON = JSON[:-2]
84 |
85 |   JSON = JSON + "]},"
86 |
87 | for j in rangeB :
88 |   nj = gdb.getNode( j )
89 |   JSON = JSON + "\n" + '{' + "\"name\": \"{0}\" , \"imports\": [".
    |       format( "root." + nj["label"] ) + '}, '
90 |

```

```
91 | # Remove the last comma
92 | JSON = JSON[:-1]
93 |
94 | JSON = JSON + "\n]"
95 |
96 | text_file = open( JSONFilename, "w" )
97 | text_file.write( JSON )
98 | text_file.close()
```

## A.9 TraceEvolutionaryHistory.py

```

1  import codecs
2
3  from py2neo import rest, neo4j, cypher
4  from GraphDB import *
5  from NodeAndRelationshipTypes import *
6
7  class TraceEvolutionaryHistory:
8
9  def BottomUpAggregation( self, gdb, node ):
10     """
11     This method starts from anywhere in the tree and goes down until
12     reach Operational Taxonomic Unit (OTU) nodes. When it happens,
13     the method basically retrieves all outgoing relationships from
14     the reached OTU node and start going back toward the root. When
15     the method is traversing internal nodes (Hypothetical
16     Taxonomic Units) from the leaves back toward the root it
17     performs an union operation with all children nodes outgoing
18     relationships – i.e., relationships of the same type are
19     ignored – and then for each relationship in the union the
20     method creates a relationship of the same type changing the
21     starting node to itself and the end node remains the same. In
22     the end, the method returns all relationships outgoing from the
23     given node.
24
25     Returns: Outgoing relationships of the given node. In case the
26     given node is an OTU, it returns only the character–states
27     relationships from the given node to character nodes.
28
29     In case the given node is an HTU, the method returns all outgoing
30     relationships resulted from the union of its children nodes
31     outgoing relationships.
32     """
33
34     if node["type"] != NodeTypes.OTU and node["type"] != NodeTypes.
35         description :
36
37         NeighborsNodes = gdb.getOutgoingAdjacentNodes( node )
38
39         relationships = []
40
41         for neighbor in NeighborsNodes:
42
43             rels = self.BottomUpAggregation( gdb, neighbor[0] )
44
45             relationships.append( rels )
46
47     # At this point we have all children nodes relationships. In such
48     a case, we can implement the first part of the algorithm
49     which is duplicate all relationships (union of children nodes

```

```

relationships) in the given node.
29
30 for rels in relationships:
31
32     if rels is not None:
33
34         for rel in rels :
35
36             if rel[0]["type"] == str( RelationshipTypes.descriptor ) :
37
38                 relType = rel[0].type.encode('ascii', 'ignore')
39
40                 startNode = node
41                 endNode   = rel[0].end_node
42
43                 # creating new relationships only where necessary
44                 gdb.getPy2neoGraphDatabaseService()[0].
45                     get_or_create_relationships( ( startNode, relType,
46                         endNode, { "type" : str( RelationshipTypes.descriptor ) }
47                     ) )
48
49
50 return gdb.getOutgoingRelationships( node )
51
52 def TopDownRefining( self, gdb, node ):
53     """
54     This method essentially should be called just after the
55     BottomUpAggregation method passing the same starting node
56     provided in BottomUpAggregation method. It starts from the
57     given node (gn) back down the tree and in every HIU it
58     traverses it compare the character-states starting from itself
59     (gn) and pointing to a given character (chaN) with every
60     character-states that starts in its children nodes (chiN) and
61     points to the same character node (chaN) for all character
62     nodes it (gn) points to. In case the comparison result is not
63     empty – i.e. the set difference between the character-states
64     starting from the given node (gn) and the set of character-
65     states starting from the children node (chiN) is not empty – it
66     creates a edge called 'EvolvedTrait' from itself (gn) to the
67     given character (chaN).
68     """
69
70     if node["type"] != NodeTypes.OTU and node["type"] != NodeTypes.
71         description :
72
73         NeighborNodes = gdb.getOutgoingAdjacentNodes( node )
74
75         tuNeighborNodes = []

```

```

60     descriptionNeighborNodes = []
61
62     for n in NeighborNodes:
63
64         if n[0]["type"] == NodeTypes.HTU or n[0]["type"] == NodeTypes.
            OTU :
65
66             tuNeighborNodes.append( n )
67
68         elif n[0]["type"] == NodeTypes.description :
69
70             descriptionNeighborNodes.append( n )
71
72     for tu in tuNeighborNodes:
73
74         for desc in descriptionNeighborNodes:
75
76             # Set Semantics
77             # http://www.itmaybeahack.com/book/python-2.6/html/p02/
                p02c06_sets.html
78
79             nodeOutgoingRelationshipTypes = set( gdb.
                getDistinctRelationshipsInBetween( node , desc[0] ) )
80             descOutgoingRelationshipTypes = set( gdb.
                getDistinctRelationshipsInBetween( tu[0] , desc[0] ) )
81
82             # diff will have elements that exist in
                nodeOutgoingRelationshipTypes and does not exists in
                descOutgoingRelationshipTypes
83             diff = nodeOutgoingRelationshipTypes -
                descOutgoingRelationshipTypes
84
85             # Removes EvolvedTrait relationship
86             Difference = diff - set( [ str( RelationshipTypes.EvolvedTrait
                ) ] )
87
88             if 0 < len( Difference ) :
89
90                 # Creates a new type of relationship (EvolvedTrait) which are
                the traits that changed from node to htu
91                 startNode = node
92                 endNode = desc[0]
93                 relType = str( RelationshipTypes.EvolvedTrait )
94
95                 gdb.getPy2neoGraphDatabaseService()[0].
                get_or_create_relationships( ( startNode, relType, endNode,
                { "type" : relType } ) )
96
97

```

```

98     gdb.deleteRelationshipsTypeFromNode( node, str( RelationshipTypes
      .descriptor ) )
99
100
101     for tu in tuNeighborNodes:
102         self.TopDownRefining( gdb, tu[0] )
103
104
105     def __JSONencodingRecursion( self, gdb, node, TraitNodes, nesting )
      :
106         """
107         It is part of the JSONencoding method.
108         Args:
109         node: Given node.
110         TraitNodes: Is the list of character nodes that node's parent has
      a 'EvolvedTrait' edge pointing to.
111         nesting: Is the space (padding) on the left.
112         Returns: JSON string.
113         """
114
115     if node["type"] != NodeTypes.OTU and node["type"] != NodeTypes.
      description :
116
117     EvolvedTraitNodes = gdb.
      getIncomingAdjacentNodesWithRelationshipInBetween( node, str(
      RelationshipTypes.EvolvedTrait ) )
118
119     NeighborNodes = gdb.
      getIncomingAdjacentNodesWithRelationshipInBetween( node, str(
      RelationshipTypes.TreeEdge ) )
120
121     json = ''
122     json = json + "\n" + ' '.ljust( nesting ) + "{"
123
124     json = json + "\n" + ' '.ljust( nesting + 2 ) + "\"{0}\" :
      \"{1}\"", ".format( "otu" , NodeTypes.OTU )
125     json = json + "\n" + ' '.ljust( nesting ) + "\"parents\" : ["
126
127     for nn in NeighborNodes:
128
129         result = self.__JSONencodingRecursion( gdb, nn[0],
      EvolvedTraitNodes , nesting + 2 )
130
131         json = json + result + ","
132
133     # Remove the last comma
134     json = json[:-1]
135
136     json = json + "\n" + ' '.ljust( nesting ) + "]"

```

```

137     json = json + "\n" + ' '.ljust( nesting ) + "}"
138
139     return json
140
141     elif node["type"] == NodeTypes.OTU:
142
143         json = ''
144         json = json + "\n" + ' '.ljust( nesting ) + "{"
145         json = json + "\n" + ' '.ljust( nesting + 2 ) + "\"{0}\" :
            \"{1}\"", ".format( "otu" , node["label"] )
146
147         i = 0
148         for trait in TraitNodes:
149
150             descriptions = gdb.getDistinctRelationshipsInBetween ( node ,
                trait[0] )
151
152             for desc in descriptions:
153
154                 json = json + "\n" + ' '.ljust( nesting + 2 ) + "\"{0}{1}\" :
                    \"{2}\"", ".format( RelationshipTypes.descriptor , str(i) ,
                        trait[0]["label"].encode('ascii', 'ignore') )
155                 json = json + "\n" + ' '.ljust( nesting + 2 ) + "\"{0}{1}\" :
                    \"{2}\"", ".format( NodeTypes.description , str(i), desc.
                        encode('ascii', 'ignore').replace("_", " ") )
156                 i = i + 1
157
158             # Remove the last comma
159             json = json[:-1]
160
161             json = json + "\n" + ' '.ljust( nesting ) + "}"
162
163         return json
164
165
166     def JSONEncoding( self, JSONFilename, startNode ):
167         """
168         It exports to a JSON format the tree structure with all characters
            the algorithm flagged with 'EvolvedTrait' edge.
169         Args:
170             JSONFilename: Filename where the JSON data should be saved.
171             startNode: Node from where the data start being collected.
172
173         """
174
175         gdb = GraphDB()
176
177         json = self.__JSONEncodingRecursion( gdb, startNode, [] , 0 )
178

```



```
179 | text_file = open( JSONFilename, "w" )  
180 | text_file.write( json )  
181 | text_file.close()
```