

Abstracts of Dissertations

Institute of Information and
Communication Technologies

BULGARIAN ACADEMY OF
SCIENCES



2 / 2020



The Open Biodiversity
Knowledge Management
System in Scholarly

Publishing

Viktor Senderov

OpenBiodiv: отворена
система за управление на
знанието за биологичното
разнообразие

Виктор Синдеров

Автореферати на дисертации

Институт по информационни и
комуникационни технологии

БЪЛГАРСКА АКАДЕМИЯ НА НАУКИТЕ

ISSN: 1314-6351

Пореаницата „Авореферати на дисертации на Института по информационни и комуникационни технологии при Българската академия на науките“ представя в електронен формат автореферати на дисертации за получаване на научната степен „Доктор на науките“ или на образователната и научната степен „Доктор“, защитени в Института по информационни и комуникационни технологии при Българската академия на науките. Представените трудове отразяват нови научни и научно-приложни приноси в редица области на информационните и комуникационните технологии като Компютърни мрежи и архитектури, Паралелни алгоритми, Научни пресмятания, Лингвистично моделиране, Математически методи за обработка на сензорна информация, Информационни технологии в сигурността, Технологии за управление и обработка на знания, Грид-технологии и приложения, Оптимизация и вземане на решения, Обработка на сигнали и разпознаване на образи, Интелигентни системи, Информационни процеси и системи, Вградени интелигентни технологии, Йерархични системи, Комуникационни системи и услуги и др.

Редактори

Геннадий Агре

Институт по информационни и комуникационни технологии, Българска академия на науките
E-mail: agre@iinf.bas.bg

Райна Георгиева

Институт по информационни и комуникационни технологии, Българска академия на науките
E-mail: rayna@parallel.bas.bg

Даниела Борисова

Институт по информационни и комуникационни технологии, Българска академия на науките
E-mail: dborissova@iit.bas.bg

Настоящото издание е обект на авторско право. Всички права са запазени при превод, разпечатване, използване на илюстрации, цитирания, разпространение, възпроизвеждане на микрофилми или по други начини, както и съхранение в бази от данни на всички или част от материалите в настоящето издание. Копирането на изданието или на част от съдържанието му е разрешено само със съгласието на авторите и/или редакторите

*The series **Abstracts of Dissertations of the Institute of Information and Communication Technologies at the Bulgarian Academy of Sciences** presents in an electronic format the abstracts of Doctor of Sciences and PhD dissertations defended in the Institute of Information and Communication Technologies at the Bulgarian Academy of Sciences. The studies provide new original results in such areas of Information and Communication Technologies as Computer Networks and Architectures, Parallel Algorithms, Scientific Computations, Linguistic Modelling, Mathematical Methods for Sensor Data Processing, Information Technologies for Security, Technologies for Knowledge management and processing, Grid Technologies and Applications, Optimization and Decision Making, Signal Processing and Pattern Recognition, Information Processing and Systems, Intelligent Systems, Embedded Intelligent Technologies, Hierarchical Systems, Communication Systems and Services, etc.*

Editors

Gennady Agre

Institute of Information and Communication Technologies, Bulgarian Academy of Sciences
E-mail: agre@iinf.bas.bg

Rayna Georgieva

Institute of Information and Communication Technologies, Bulgarian Academy of Sciences
E-mail: rayna@parallel.bas.bg

Daniela Borissova

Institute of Information and Communication Technologies, Bulgarian Academy of Sciences
E-mail: dborissova@iit.bas.bg

This work is subjected to copyright. All rights are reserved, whether the whole or part of the materials is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in other ways, and storage in data banks. Duplication of this work or part thereof is only permitted under the provisions of the authors and/or editor.



Abstract of PhD Thesis

The Open Biodiversity Knowledge Management System in Scholarly Publishing

Viktor Senderov

Supervisor Prof. Lyubomir Penev
Scientific Advisor: Assoc. Prof. Kiril Simov

Approved by Supervising Committee:

Prof. Kalinka Kaloyanova
Prof. Georgi Markov
Prof. Maria Nisheva-Pavlova
Assoc. Prof. Svetla Boycheva
Assoc. Prof. Gennady Agre



Introduction

Importance of the topic

The desire for an integrated information system serving the needs of the biodiversity community dates at least as far back as 1985 when the Taxonomy Database Working Group (TDWG)—later renamed to Biodiversity Informatics Standards but retaining the abbreviation TDWG—was established¹. In 1999, the Global Biodiversity Information Facility (GBIF) was created after the Organization for Economic Cooperation and Development (OECD) had arrived at the conclusion that “an international mechanism is needed to make biodiversity data and information accessible worldwide” (*What is GBIF?*). The Bouchout declaration (*Bouchout Declaration 2014*) crowned the results of the European Union–funded project *pro-iBiosphere* that lasted from 2012 to 2014 and was dedicated to the task of creating an integrated biodiversity information system. The Bouchout declaration proposes to make scholarly biodiversity knowledge freely available as Linked Open Data (LOD). A parallel process in the U.S.A. started even earlier with the establishment of the Global Names Architecture, GNA (Patterson et al., 2010; Pyle, 2016b).

In 2014, the Horizon 2020 BIG4 consortium was formed between academia and industry dedicated to advancing biodiversity science. The project’s mission statement reads “BIG4—Biosystematics, Informatics and Genetics of the big 4 insect groups: training tomorrow’s researchers and entrepreneurs” (University of Copenhagen et al., 2014). An important member of the consortium is the academic publishing house and software company, Pensoft Publishers. Pensoft publishes several dozen well-known open access taxonomic journals² and, as a signatory of the Bouchout declaration, was a prime candidate to push the vision for an Open Biodiversity Knowledge Management System (OBKMS) forward. The presented Ph.D. project is based at Pensoft Publishers and at the Institute of Information and Communication Technology (IICT) of the Bulgarian Academy of Sciences with the goal to follow through *pro-iBiosphere*’s vision.

Previous work

Due to the interdisciplinary nature of the thesis, this section will focus on two areas: (a) knowledge bases and Linked Open Data and (b) biodiversity publishing.

¹A webpage with the history of TDWG dating back to 1985 can be viewed under <http://old.tdwg.org/past-meetings/>; however, a lot of the links are unfortunately broken and the page needs some maintenance.

²For example, ZooKeys, PhytoKeys, MycoKeys, and Biodiversity Data Journal (BDJ).

Knowledge bases and Linked Open Data

We shall start by first introducing *knowledge bases* and *knowledge-based systems*. We use the two terms interchangeably but tend to write the longer variant, knowledge-based system, when we want to emphasize aspects of the knowledge base that are not related to the underlying facts store (database).

It is useful to form one's concept of knowledge-based systems both by looking at explicit definitions and by looking at several examples of knowledge bases in practice. The term was already being widely discussed by the 1980's (Jarke et al., 1989) and early nineties (Harris et al., 1993) and was understood to mean the utilization of ideas from both database management systems (DBMS) and artificial intelligence (AI) to create a type of computer system called *knowledge base management system* (KBMS). Harris et al., 1993 writes that the characteristics of a knowledge base management system are that it contains "prestored rules and facts from which useful inferences and conclusions may be drawn by an inference engine." We should note that the phrase "prestored rules" comes from the time of first-generation AI systems that were rule-based. Recently, there has been progress in incorporating statistical techniques into databases (Mansinghka et al., 2015); however, in this project we are working with the classical rule-based definition. In other words, a knowledge base is, in our understanding, a suitable database tightly integrated with a logic layer.

Another relatively more recent development in knowledge-based systems has been the application of the Linked Data principles (Heath and Bizer, 2011). In fact, most existing knowledge bases emphasize the community aspects of making data more interconnected and reusable. Examples include Freebase (Bollacker et al., 2008), which was recently incorporated in WikiData (Vrandečić and Krötzsch, 2014; Pellissier Tanon et al., 2016), DBPedia (Auer et al., 2007), as well as Wolfram|Alpha (*Wolfram|Alpha, Making the world's knowledge computable*) and the Google Knowledge Graph (Singhal, 2012). What these systems have in common is that an emphasis is placed not only on the logic layer allowing inference but on a unified information space: these systems act as nexus integrating information from multiple places and they follow to various degrees the principles of Linked Open Data (LOD).

Linked Open Data (Heath and Bizer, 2011) is a concept of the Semantic Web (Berners-Lee et al., 2001), which, when applied properly, ensures that data published on the Web is reusable, discoverable, and most importantly ensures that pieces of data published by different entities can work together. We will discuss the Linked Data principles and their application to OpenBiodiv in detail in Chapter 3.

Leveraging these developments modern knowledge bases place a bigger emphasis on interlinking data rather than on developing a complex inference machinery. There has been critique of the idea of bundling logic in the database layer as such bundling leads to increased complexity (Barrasa, 2017). The critique can be summarized with two points. First, bundling the logic near the data (especially when it is excessive for the task at hand) can lead to drastic performance decreases³. Second, the developing of new techniques (e.g. machine learning) can make the existing deep logic layer obsolete. Our view is that data is the commodity which is much more valuable, and the inference strategy (be it a rule-based logic layer, or a statistical machine learning technique) can be replaced as computational science moves forward. These ideas lead to an interesting conundrum in the choice of a database technology discussed in the subsequent sections.

³ We will compare the performance of the stronger Web Ontology Language (OWL) logic layer with a weaker RDF Schema (RDFS) logic layer in Chapter 3. Resource Description Format (RDF) is a data model for storing statements about things discussed later.

Finally, a knowledge-based system ultimately needs to include user-interface components (UI's) and application programming interfaces (API's) or an application layer. These serve as the point-of-contact between human and machine, or machine and machine and are crucial to the success of any such system.

Biodiversity publishing

In the biomedical domain there are well-established efforts to extract information and discover knowledge from literature (e.g. Rebholz-Schuhmann et al., 2005; Momtchev et al., 2009; Williams et al., 2012). The biodiversity domain, and in particular biological systematics and taxonomy (from here on in this thesis referred to as *taxonomy*), is also moving in the direction of semantization of its research outputs (Agosti, 2006; Patterson et al., 2006; Kennedy et al., 2005; Penev et al., 2010a; Tzitzikas et al., 2013). The publishing domain has been modeled through the Semantic Publishing and Referencing Ontologies, SPAR Ontologies (Peroni, 2014). The SPAR Ontologies are a collection of ontologies incorporating, amongst others, FaBiO, the FRBR-aligned Bibliographic Ontology (Peroni and Shotton, 2012), and DoCO, the Document Component Ontology (Constantin et al., 2016). The SPAR Ontologies provide a set of classes and properties for the description of general-purpose journal articles, their components, and related publishing resources. Taxonomic articles and their components, on the other hand, have been modeled through the TaxPub XML Document Type Definition (DTD)—also referred to loosely as XML schema—and the Treatment Ontologies (Catapano, 2010). While TaxPub is the XML-schema of taxonomic publishing for several important taxonomic journals (e.g. ZooKeys, PhytoKeys, Biodiversity Data Journal), the Treatment Ontologies are still in development and have served as a conceptual template for OpenBiodiv-O (discussed in Chapter 2).

Taxonomic nomenclature is a discipline with a very long tradition. It transitioned to its modern form with the publication of the Linnaean System (Linnaeus, 1758). Already by the beginning of the last century, there were hundreds of taxonomic terms in usage (Witteveen, 2015). At present the naming of organismal groups is governed by the International Code of Zoological Nomenclature, ICZN (International Commission on Zoological Nomenclature, 1999) and by the International Code of Nomenclature for algae, fungi, and plants, Melbourne Code (**mcneill_international_2012**). Due to their complexity (e.g. ICZN has 18 chapters and 3 appendices), it proved challenging to create a top-down ontology of biological nomenclature. Example attempts include the relatively complete NOMEN ontology (Dmitriev and Yoder, 2017) and the somewhat less complete Taxonomic Nomenclatural Status Terms, TNSS⁴.

There are several projects that are aimed at modeling the broader biodiversity domain conceptually. Darwin Semantic Web, Darwin-SW (Baskauf and Webb, 2016) adapts the previously existing Darwin Core (DwC) terms (Wieczorek et al., 2012) as Resource Description Framework (RDF). These models deal primarily with organismal occurrence data.

Modeling and formalization of the strictly taxonomic domain has been discussed by Berendsohn (Berendsohn, 1995) and later, e.g., in (Franz and Peet, 2009; Sterner and Franz, 2017). Noteworthy efforts are the XML-based Taxonomic Concept Transfer Schema (Taxonomic Names and Concepts Interest Group, 2006) and a now defunct Taxon Concept ontology. Very recently, the TDWG community has attempted to resurrect the Taxon Concept ontology with the Taxonomic Names and Concepts Interest

⁴Even though it is unknown to the authors whether TNSS was published in peer-reviewed literature, remnants of it can still be found on GitHub, e.g. under https://github.com/pensoft/OpenBiodiv/blob/master/ontology/contrib/taxonomic_nomenclatural_status_terms.owl.

Group. The group discussions can be accessed under <https://github.com/tdwg/tnc>. Interestingly the **very first GitHub issue** discussed OpenBiodiv-O and the possibility of its adoption as a TDWG standard.

By the time the OpenBiodiv project started in June 2015, a number of articles had been previously published on the topics of linking data and sharing identifiers in the biodiversity knowledge space (Page, 2008), unifying phylogenetic knowledge (Parr et al., 2012), taxonomic names and their relation to the Semantic Web (Page, 2006; Patterson et al., 2010), and aggregating and tagging biodiversity research (Mindell et al., 2011). Some partial discussion of OBKMS was to be found in the science blog *iPhylo* (Page, 2014, 2015). The legal aspects of the OBKMS had been discussed by Egloff et al., 2014.

Furthermore, several tools and systems that deal with the integration of biodiversity and biodiversity data had been developed by different groups. Some of the most important ones are UBio, Global Names, BioGuid, BioNames, Pensoft Taxon Profile, and the Plazi Treatment Repository⁵.

Key findings

The key findings from the papers cited in the previous paragraphs can be summarized as follows:

1. Biodiversity science deals with disparate types of data: taxonomic, biogeographic, phylogenetic, visual, descriptive, and others. These data are siloed in unlinked data repositories.
2. Biodiversity databases need a universal system of naming concepts due to the inefficiencies of Linnaean names for modern taxonomy. Taxonomic concept labels have been proposed as a human-readable solution and stable globally unique identifiers of taxonomic concepts had been proposed as a machine-readable solution.
3. There is a base of digitized semi-structured biodiversity information online with appropriate licenses waiting to be integrated as a knowledge base.

Goal and objectives

Given the huge international interest in OBKMS, this dissertation started the OpenBiodiv project, the goal of which is to contribute to OBKMS by creating an open knowledge-based system of biodiversity information extracted from scholarly literature. In order to complete the system, the following objectives need to be achieved:

Objective 1: Architecture. Formally define OpenBiodiv as a knowledge-based system and create its integrated software architecture.

Objective 2: Ontology. Study the domain of biodiversity informatics and biodiversity publishing and develop an ontology allowing data integration from diverse sources.

⁵UBio: <http://ubio.org/>; Global Names: <http://globalnames.org/>; BioGuid: <http://bioguid.org/>; BioNames: <http://bionames.org/>; Pensoft Taxon Profile: <http://ptp.pensoft.eu/>; Plazi Treatment Repository: <http://plazi.org/wiki/>.

Objective 3: Linked open dataset. Create a Linked Open Dataset (LOD) on the basis of published taxonomic articles using the ontology defined in Objective 2.

Objective 4: Library. Develop methods for converting taxonomic publications into the semantic model of the ontology in order to support Objective 3.

Objective 5: Workflows. Develop practical workflows for continuously converting taxonomic data into taxonomic publications and thus updating the LOD dataset.

Objective 6: Web portal. Create a web-portal and example applications on top of the knowledge base.

Methodology

This dissertation has a methods and tools orientation: i.e. its goal is not the testing of particular scientific hypothesis but rather the theoretical design and practical implementation of a knowledge-management system. In this section I shall outline the "meta-choices" that I have made—such as what programming and database paradigms to use—before the design and implementation phase.

Choice of database paradigm for OpenBiodiv

We specify OpenBiodiv as a knowledge-based system with a focus on structuring and interlinking biodiversity data. Two of the possible database technologies that fit this requirement are semantic graph databases (triple stores) such as GraphDB (Ontotext, 2018) and labeled property graphs such as Neo4J (Neo4J Developers, 2012). Semantic graph databases offer a very simple data model: every fact stored in such a database is composed as a triple of *subject*, *predicate*, and *object*. Subjects of triples are always resource identifiers, whereas objects can be other resource identifiers or literal values (e.g. strings, numbers, etc.). Links between resources or between resources and literals are given by the predicates (also specified as identifiers). These links are sometimes referred to as *properties*. Thus, one can visualize a graph whose vertices are the objects or subjects given by resource identifiers or literals and whose edges are predicates.

Semantic graph databases have the unique feature that the logic layer is also expressed as triples stored in the database. This logic layer, known as *ontology*, is not only responsible for drawing conclusions from the data (inference), but also specifies the semantics of how knowledge should be expressed.

Labeled property graphs, on the other hand, offer a freer data model by allowing the edges of the knowledge graph to have properties as well. For example, in a labeled property graph whose vertices are two cities A and B and are connected by a property-predicate *connected by road*, it is possible to additionally attach the value “500 km” to that property. Thus, we indicate that the length of the road connecting the cities is 500 km.

Note that labeled property graphs are not any more expressive than what can be achieved by triples alone. In fact, complex relationships in a simple triple store can be expressed by making relationships into nodes that have properties on their own. This process is known as *reification*. For example, the two cities *A* and *B* can connect to a further vertex, *R* indicating the road. *R* will then have three properties: *start*, *end*, and *length*. The value (object) of *start* will be *A*, of *end* will be *B*, and of *length* will be the literal “500 km.”

TABLE 1: Differences between semantic graph databases (e.g. GraphDB) and labeled property graphs (e.g. Neo4j).

Criterion	Semantic database	Labeled property graph
Semantics	Stored in the database itself as OWL or RDFS statements. Provides a uniform data space. Requires expert ontologists to extract knowledge.	Formal semantics usually are missing. Quick deployment. Uniform data space harder to achieve.
Inference	Provided by the database itself from its ontology or expressed as SPARQL queries. General purpose, slower.	External to the database. Needs to be written for every specific task. Special purpose. Faster.
Community	Has a rich and mature community of ontologists and knowledge engineers. Lots of domain ontologies. Designed for inter-operability. Standards-driven.	Data models are created ad-hoc by data scientists or programmers for a particular task. Inter-operability requires effort and not of primary concern. Applications-driven.

We have summarized the differences between labeled property graphs and semantic graph databases in Table 1. After careful considerations, we settled on the triple store, i.e. semantic graph database as a choice of database technology. This decision was informed by the wide availability of high-quality ontologies and Resource Description Framework (RDF) data models in our domain (Baskauf and Webb, 2016; Peroni, 2014) and the popularity of the Semantic Web (Berners-Lee et al., 2001) in the community. Furthermore, our base at a publisher was more suited to a standards-driven foundational project as opposed to a particular application.

However, we believe that labeled property graphs are a freer and a more natural data model and are perfectly suited for biodiversity informatics. In particular they provide a much more natural formalism for relationships between taxonomic concepts (discussed in Chapter 2). Also, non-RDF semantic databases such as WikiData are gaining in popularity. Therefore, we believe that the applicability of RDF triple stores for OpenBiodiv should constantly be reevaluated.

Choice of information sources

According to *pro-iBiosphere project final report 2014*, biodiversity and biodiversity-related data have two different “life-cycles.” In the past, after an observation of a living organism had been made, it was recorded on paper and then the observation record was published in paper-based form. In order for biodiversity data to be available to the modern scientist, efforts are made nowadays to digitize those paper-based publications by Plazi Agosti et al., 2007 and the Biodiversity Heritage Library (Miller et al., 2012). For this purpose, several dedicated XML schemas have been developed (see Penev et al., 2011 for a review), of which TaxPub (Catapano, 2010) and TaxonX seem to be the most widely used (Penev et al., 2012). The digitization of publications contains

several steps. After scanning and optical character recognition (OCR), text mining is combined with searching for particular kinds of data. This procedure leaves a trace in the form of marked-up (tagged) elements that can then be extracted and made available for future use and reuse (Miller et al., 2015).

In present day, biodiversity data and publications are mostly “born digital” as semantically Enhanced Publications (EP’s, Claerbout and Karrenbach, 1992; Godtshoven et al., 2009; Shotton, 2009). According to Claerbout and Karrenbach, 1992, “an EP is a publication that is enhanced with research data, extra materials, post publication data and database records. It has an object-based structure with explicit links between the objects. An object can be (part of) an article, a data set, an image, a movie, a comment, a module or a link to information in a database.” Semantically enhanced publications are thus natives of the Web and the Semantic Web unlike their paper-based predecessors.

The act of publishing in a digital, enhanced format, differs from the ground up from a paper-based publication. The main difference is that a digitally-published document can be structured in such a format as to be suitable both for machine processing and to the human eye. In the sphere of biodiversity science, Pensoft journals such as ZooKeys, PhytoKeys, and the Biodiversity Data Journal (BDJ) already function by providing EP’s (Penev et al., 2010b).

Given the fact that Pensoft Publishers’ and Plazi’s publications cover a large part of taxonomic literature both in volume and also in temporal span, and the fact that the publications of those two publishers are available as semantic EP’s, we’ve chosen Pensoft’s journals and Plazi’s treatments as our main sources of information.

Furthermore, we incorporate the taxonomic backbone of GBIF GBIF Secretariat, 2017a as a source for data integration. This is further discussed in Chapter 3.

Choice of development methodology and programming environment

In 2016, based on the outcomes of pro-iBiosphere and on the previous work in the area of biodiversity informatics, we published the Ph.D. plan for this research (Senderov and Penev, 2016). This publication can be considered as the first design specification of OpenBiodiv. However, in the course of developing the system, its design was changed iteratively through a feedback loop from collaborators from the BIG4 project⁶ and various international collaborators. We view this positively and in the spirit of both *open science* and *agile software development* (Beck et al., 2001). This iterative approach differs from the waterfall approach where after a through design phase, the specifications “are frozen” and a lengthy implementation phase.

In recent years, the R programming language has been used widely in the field of data science (R Core Team, 2016). R has a rich library of software packages including such for processing XML (Wickham et al., 2018), for accessing rest API’s (Wickham, 2017), and focuses on open science (Boettiger et al., 2015). The capabilities of R as function-oriented and interpreted language allow the iterative software development approach outlined in the previous paragraph to proceed rapidly. Furthermore, R is widely adopted in the biodiversity informatics community. For this reason, the R software environment was chosen as the main programming environment.

⁶The Ph.D. candidate, Viktor Senderov, is part of the Marie Skłodowska-Curie BIG4 International Training Network: Biosystematics, informatics and genomics of the big 4 insect groups: training tomorrow’s researchers and entrepreneurs.

Open Science and The Semantic Web

After having specified the desired design and given the programming language, R, I would like to discuss some methodologies and frameworks that have been adopted to be more efficient, open, and reproducible.

I believe that OpenBiodiv needs to be addressed from the point of view of *Open Science*. According to Kraker et al., 2011 and to *Was ist Open Science?*, the six principles of open science are: open methodology, open source, open data, open access, open peer review, and open educational resources. It is my belief that the aim of open science is to ensure access to the whole research product: data, discoveries, hypotheses, and so on. This opening-up will ensure that the scientific product is reproducible and verifiable by other scientists (Mietchen, 2014). There is a very high interest in development of processes and instruments enabling reproducibility and verifiability, as can be evidenced for example by a special issue in Nature dedicated to reproducible research (*Challenges in irreproducible research* 2010). Therefore, the source code, data, and publications of OpenBiodiv will be published openly.

Moreover, OpenBiodiv should be thought of as integral part of the Semantic Web (Berners-Lee et al., 2001). The Semantic Web is a vision for the future of the web where not only documents but also data are connected.

Structure of the thesis

So far the *raison d'être* of the system and this thesis and an outline of its goal and objectives have been given in this Introduction. In Chapter 1, a formal specification and design of the desired system as well as an outline of its architecture will be presented; this chapter forms Objective 1. The subsequent chapters discuss the implementation of OpenBiodiv. Chapter 2 gives a conceptualization of the domain of scientific taxonomic publishing formalizes it by introducing the central result of this thesis, the ontology of OpenBiodiv (OpenBiodiv-O) and thus forms Objective 2. Chapter 3 describes the Linked Open Dataset that has been generated based on OpenBiodiv-O and forms Objective 3. Chapter 4 describes in detail the RDF4R software package (an R package for working with RDF), which was used to create the Linked Open Data (OpenBiodiv-LOD) and forms Objective 4. In Chapter 5, two case-studies for importing data into OpenBiodiv from important international repositories are discussed and thus it forms Objective 5. Chapter 6 discusses the website that has been prepared to serve on top of OpenBiodiv-LOD and its applications (Objective 6). In the Conclusion, I will explain how the results have been published and summarize the main results.

Chapter 1

Summary of Chapter 1: Architecture of OpenBiodiv

In this chapter, we provide the architectural blueprint, i.e. the specification and design of OpenBiodiv. We break up OpenBiodiv into components that will be treated in detail in subsequent chapters. We describe how these components inter-operate in order to form the OpenBiodiv knowledge-based system.

1.1 What is OpenBiodiv?

The understanding of OpenBiodiv as a knowledge-based system can be summarized as follows: OpenBiodiv is a database of interconnected biodiversity information together with logic and application layers allowing users to not only query the data but also discover additional facts of relevance implied by the data. The primary sources of information in OpenBiodiv are the journals of the academic publisher Pensoft, taxonomic information from Plazi, and the taxonomic backbone of Global Biodiversity Information Facility (GBIF).

The research problem of OpenBiodiv's architecture can be postulated as designing an open-access semantic RDF graph database, incorporating information stored in Pensoft, Plazi, and GBIF, and allowing the users of the system to ask complicated queries.

OpenBiodiv consists of (1) a semantic graph database, (2) a back-end code base, and (3) a front-end in the form of a web-portal facilitating the access to the underlying knowledge base (Fig. 1.1). OpenBiodiv enables the flow of information between international repositories for biodiversity data to Biodiversity Data Journal (BDJ) and other journals that use the ARPHA-BioDiv toolkit (Penev et al., 2017). As a second step, knowledge is extracted from such journals taking advantage of the TaxPub Document Type Definition (DTD)¹ introduced by Catapano, 2010. Example journals include ZooKeys, Biodiversity Data Journal (BDJ), PhytoKeys, MycoKeys, and so on². At the same time, knowledge is extracted from Plazi TreatmentBank, an archive of legacy biodiversity literature containing over 200 thousand treatments³ and updated every day. Last but not least, these sources are interlinked via GBIF's taxonomic backbone (GBIF Secretariat, 2017a). The extracted knowledge is then stored in a semantic graph database (Fig. 1.2).

¹We will take the liberty and refer to TaxPub as an XML schema in the rest of the chapter.

²The journals can be accessed under https://pensoft.net/browse_journals.

³A treatment is a special section in a biological publication describing and discussion a species or a higher taxon. TreatmentBank is accessible under <https://http://plazi.org/resources/treatmentbank/>.

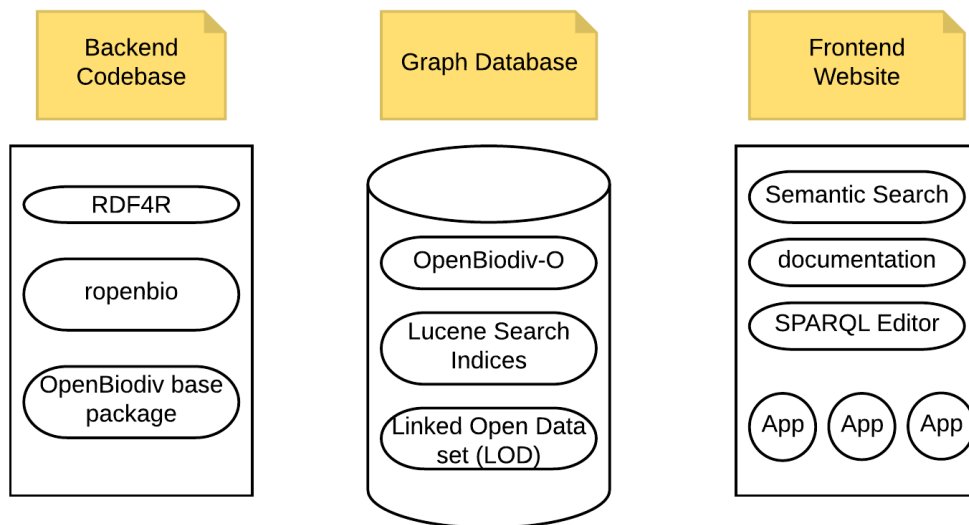


FIGURE 1.1: The components of OpenBiodiv.

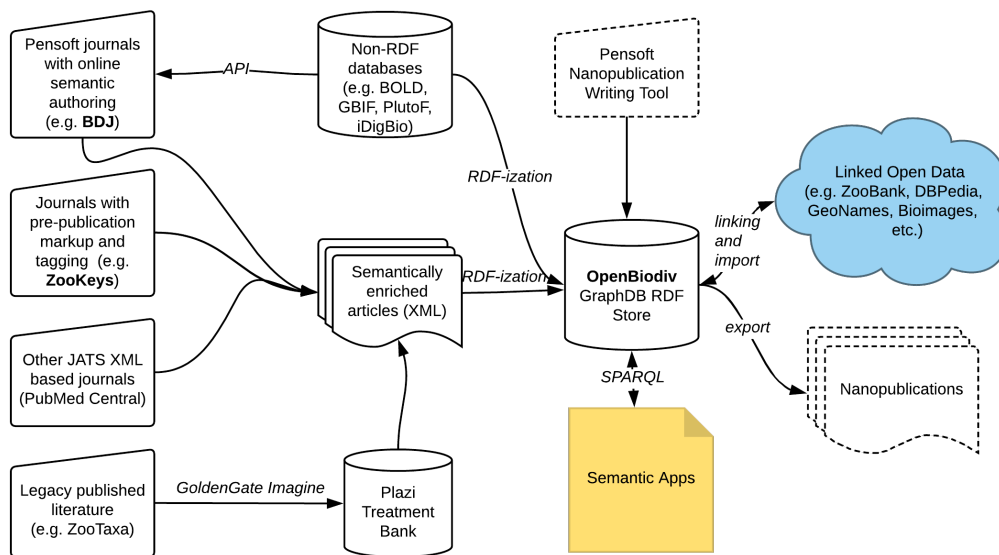


FIGURE 1.2: Flow of information in the biodiversity data space until it reaches the OpenBiodiv semantic database. Dashed lines are components that have not been implemented yet.

1.2 Semantic graph database

A primary output of the OpenBiodiv effort is the creation of a semantic database based on knowledge extracted from the archives of Pensoft and Plazi and GBIF's taxonomic backbone and accessible under <http://graph.openbiodiv.net/>. A discussion of the components of the database follows.

1.2.1 OpenBiodiv ontology (OpenBiodiv-O)

The central result of the OpenBiodiv effort is the creation of a formal domain model of biodiversity publishing, the ontology OpenBiodiv-O (Senderov et al., 2017). The source code of the ontology and accompanying documentation can be accessed under <https://github.com/vsenderov/openbiodiv-o>. A detailed discussion is presented in Chapter 2.

1.2.2 OpenBiodiv Linked Open Dataset (OpenBiodiv-LOD)

Using OpenBiodiv-O and the infrastructure described later in this chapter a dataset incorporating approximately 200 thousand Plazi treatments, five thousand Pensoft articles, as well as GBIF's taxonomic backbone (over a million names) has been created. The dataset is available online through the workbench of the semantic database <http://graph.openbiodiv.net>. It is discussed in detail in Chapter 3.

1.3 Backend

In order to populate a semantic database it is necessary to create the infrastructure that converts raw data (text, images, data tables, etc.) into a structured semantic format allowing the interlinking of resource identifiers and the answering of complex queries. OpenBiodiv creates new infrastructure and extends existing infrastructure for transforming biodiversity scholarly publications into Resource Description Format (RDF) statements with the help of the components described in this section.

1.3.1 RDF4R: R package for working with RDF

One of the greater technical challenges for OpenBiodiv is the transformation of biodiversity information (e.g. taxonomic names, paper metadata, figures, etc.) stored as semi-structured XML into fully-structured semantic knowledge in the form of RDF. In order to solve this challenge, an R package has been developed that enables the creation, manipulation, and submission and retrieval to and from a semantic database of RDF statements. This package is accessible under an open source license on GitHub under <https://github.com/vsenderov/rdf4r>. We describe the package in Chapter 4.

1.3.2 OpenBiodiv Base and ROpenBio

In combination with the RDF4R package, the code-base is completed by one more R package, `ropenbio` and a code-base (OpenBiodiv Base) of scripts and documentation necessary to bootstrap the database. `ropenbio` utilizes the RDF4R package to convert semi-structured XML to RDF. It contains the "mappings" necessary for that conversion. It is available under <https://github.com/pensoft/ropenbio>. OpenBiodiv Base coordinates the invocation of `ropenbio`, contains scripts for the

automatic import of new resources, and other housekeeping details. It is available under <https://github.com/pensoft/openbiodiv>. Their usage to generate the OpenBiodiv-LOD is discussed in Chapter 3.

1.3.3 Workflow for converting ecological metadata to a manuscript

Ecological Metadata Language (EML) is a popular format for describing ecological datasets (Michener et al., 1997). Biodiversity repositories such as GBIF and DataOne make use of this format to describe the datasets that they store. An import pipeline for importing an EML file as a BDJ data paper⁴ has been developed as part of OpenBiodiv (Senderov et al., 2016). We describe this workflow in detail in Chapter 5. To access the pipeline interactively, go to <https://arpha.pensoft.net>, login to the system (registration is free), select “Start a new manuscript,” scroll all the way down to “Import a manuscript,” and follow the necessary steps to upload an EML and use it as a template for your new manuscript.

1.3.4 Workflow for importing specimen data into Biodiversity Data Journal

One of the important types of biodiversity data is occurrence data—data that documents the presence of a properly taxonomically identified organism at a given location and time. Such data is stored at international repositories such as BOLD, GBIF, PlutoF, and iDigBio. In order to facilitate data publishing, as well as to act as an entry point into OpenBiodiv, a pipeline for importing any occurrence record from these databases into a BDJ taxonomic paper has been developed (Senderov et al., 2016). We describe this workflow in detail in Chapter 5. To access the workflow interactively, go to <https://arpha.pensoft.net>, login to the system (registration is free), select "Start a new manuscript," select "Biodiversity Data Journal" as a journal and "Taxonomic Paper" as paper-type and "Create a manuscript." Then, in your new manuscript, expand the "Taxon treatments" section by clicking on the + sign next to it, give a test classification to your treatment (e.g. Animalia), click “Save” and you will be presented with a choice of subsections. Click the “Materials” section on the left to visualize the workflow. Look at the lower-part of the dialog, where “You may place multiple ID’s...”—this is the part where you select external resource identifiers to be imported to your article.

1.4 Frontend

In addition to providing a searchable database endpoint, a website allowing semantic search and containing specific tasks packaged as apps is being developed (<http://openbiodiv.net>). The development of the site extends beyond the scope of the dissertation thesis and is driven by the Pensoft development team. A beta version is already operational Fig. 1.3. A limited discussion is found in Chapter ??.

⁴A data paper (Chavan and Penev, 2011) is a paper in a scholarly (peer-reviewed) journal discussing a scientific dataset.

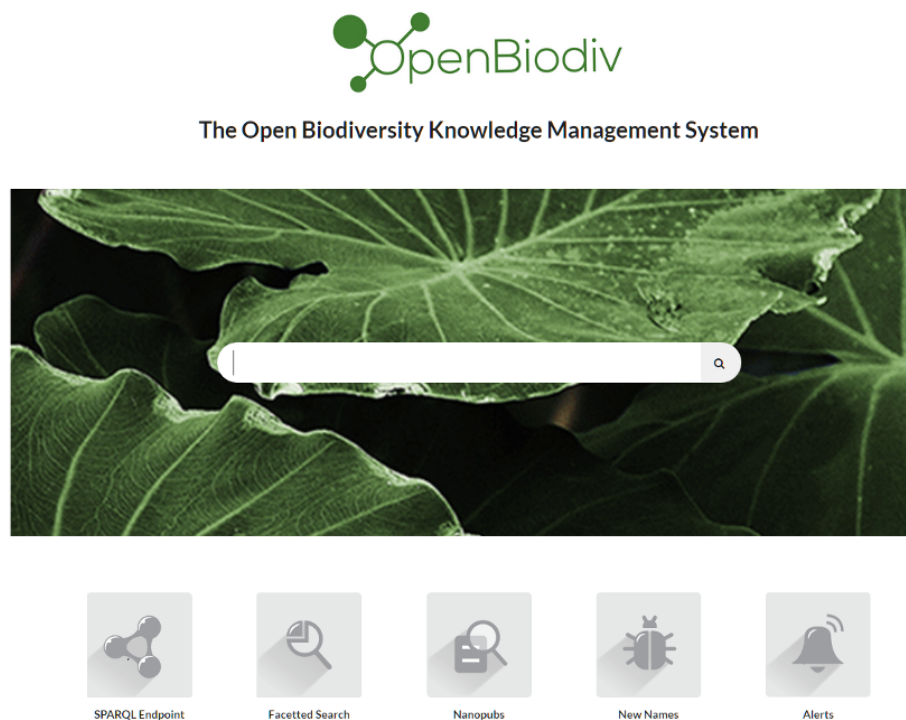


FIGURE 1.3: Beta version of the OpenBiodiv website together with sample app icons.

1.5 IT

The system is deployed on a Debian GNU+Linux virtual machine. GraphDB runs with a 20 GB heap file and with the RDFS-Plus Optimized rule set⁵. Continuous operation is ensured by the automatic execution of scripts from the `run` directory of OpenBiodiv Base.

⁵This is necessitated by the fact that we reached a performance bottleneck the OWL inference. Discussed in Chapter 3

Chapter 2

Summary of Chapter 2: The OpenBiodiv Ontology

OpenBiodiv lifts biodiversity information from scholarly publications and academic databases into a computable semantic form. In this chapter, we introduce OpenBiodiv-O (Senderov et al., 2018), the ontology forming the knowledge and inferencing model of OpenBiodiv. OpenBiodiv-O provides a conceptual model of the structure of a biodiversity publication and the development of related taxonomic concepts. We first introduce the modeled domain in Domain Conceptualization and then formalize it in Results.

By developing an ontology focusing on biological taxonomy, our intent is to provide an ontology that fills in the gaps between ontologies for biodiversity resources such as Darwin-SW and semantic publishing ontologies such as the ontologies comprising the SPAR Ontologies. We take the view that it is advantageous to model the taxonomic process itself rather than any particular state of knowledge.

The source code and documentation are available under the CC BY license¹ from GitHub². We start by introducing the domain of biological taxonomy and the related biodiversity sciences.

2.1 Domain Conceptualization

We give an introduction of the history of modern biological taxonomy starting with Carl Linnaeus (1707-1778) who proposed the modern organism grouping of *kingdoms*, *classes*, *orders*, *genera* and the usage of Latin binomial names in *Systema Naturae* (Linnaeus, 1758). We emphasize that the work of taxonomists to describe and organize biodiversity is far from complete. This informs the creation of our ontology not as a static formalization of the existing biological taxonomy in computer-readable form, but as a formalization of the *scientific process of biological taxonomy*.

We then describe in the detail what the scientific process of biological taxonomy entails. We start by introducing taxonomic concepts and how they are formed. A taxonomic concept is a scientific hypothesis (Deans et al., 2012) that a certain well-defined group of organisms exists in Nature. It is formed by examining specimens and necessarily entails a scientific grouping criterion, often called a species concept (Mallet, 2001; not to be confused with taxonomic concept!). Historically, organisms have been grouped by their appearance (morphological species concept) or reproductive behavior (biological species concept), but recently the focus has shifted towards grouping based on genetic relatedness (phylogenetic and genomic species concepts).

¹Creative Commons Attribution 4.0 International Public License.

²<https://github.com/vsenderov/openbiodiv-o/blob/master/LICENSE.md>

We then describe the ranks of biological taxonomy and how they are regulated by International Codes **mcneill_international_2012**; International Commission on Zoological Nomenclature, 1999). The codes govern lower ranks: species, genus, family, order; higher ranks (e.g. phylum, kingdom, domain, etc.) are however free to be used by researchers as they view fit. This leads to multiple competing viewpoints.

Publishing taxonomic concepts is an integral step in the scientific workflow of every taxonomist. We describe the structure and types of taxonomic publications with a particular emphasis on the Treatment section. A Treatment is the section in a taxonomic publication where a taxonomic concept is circumscribed.

Previous work

We discuss previous efforts made to ontologize scientific publications and biological information. Particularly important are the Semantic Publishing and Referencing Ontologies (SPAR Ontologies, Peroni, 2014) and the TaxPub XML Document Type Definition ((Catapano, 2010) referred to loosely as XML schema). The modeling of biodiversity information is primarily influenced by the Codes (**mcneill_international_2012**; International Commission on Zoological Nomenclature, 1999), that were mentioned in the previous section, and by a variety of standards (e.g. Darwin Core, DwC, Wiczorek et al., 2012), published by the TDWG community.

Finally, we discuss the emerging field of concept taxonomy (Berendsohn, 1995; Franz and Peet, 2009; Sterner and Franz, 2017)—a re-imagination of how the circumscription process in biological taxonomy ought to work.

2.2 Methods

OpenBiodiv-O is expressed in Resource Description Framework (RDF). At the onset of the project, a consideration was made to use RDF in favor of a more complex data model such as Neo4J's (Senderov and Penev, 2016). The choice of RDF was made in order to be able to incorporate the multitude of existing domain ontologies into the overall model.

To develop the conceptualization of the taxonomic process and then the ontology we utilized the following process: (1) domain analysis and identification of important resources and their relationships; (2) analysis of existing data models and ontologies and identification of missing classes and properties for the successful formalization of the domain.

The formal structure of the ontology is specified by employing the RDF Schema (RDFS) and the Web Ontology Language (OWL). It is encoded as a part of a literate programming (Knuth, 1984) document in RMarkdown format titled “OpenBiodiv Ontology and Guide”³. The statements have been extracted from the RMarkdown file via *knitr* and are provided here as an appendix. It is also possible to request the ontology via Curl from its endpoint with the indication of `content-type: application/rdf+xml`. The vocabularies can be found as additional appendices, Taxonomic Statuses and RCC-5, and on the GitHub page⁴.

A dataset (OpenBiodiv-LOD, will be described in detail in the next Chapter) from Pensoft's journals, Plazi's treatments, and GBIF's taxonomic backbone has been generated with OpenBiodiv-O and can be found at the SPARQL Endpoint ⁵. The

³<http://openbiodiv.net/ontology>

⁴<https://github.com/vsenderov/openbiodiv-o>

⁵<http://graph.openbiodiv.net/>

endpoint is also accessible from the website⁶, under “SPARQL Endpoint.” Demos are available as “Saved Queries” from the workbench.

2.3 Results

We understand OpenBiodiv-O to be the *shared formal specification of the conceptualization* (Gruber, 1993; Obitko, 2007; Staab and Studer, 2009) that we have introduced in Background. OpenBiodiv-O describes the structure of this conceptualization, not any particular state of it.

There are several domains in which the modeled resources fall. The first one is the scholarly biodiversity publishing domain. The second domain is that of taxonomic nomenclature. The third domain is that of broader taxonomic (biodiversity) resources (e.g. taxonomic concepts and their relationships, species occurrences, traits). To combine such disparate resources together we rely on SKOS Miles and Bechofer. Unless otherwise noted, the default namespace of the classes and properties for this paper is <http://openbiodiv.net/>. The prefixes discussed here are listed at the beginning of the ontology source code.

2.3.1 Semantic Modeling of the Biodiversity Publishing Domain

We extend the framework of the SPAR Ontologies by introducing a new class for taxonomic articles, its subsections, as well as a new class for the mentioning of a taxonomic name (see next subsection) in an article. These new classes are summarized in Table 2.1.

TABLE 2.1: New biodiversity publishing classes introduced.

Class QName	Comment
<code>:Treatment</code>	section of a taxonomic article
<code>:NomenclatureSection</code>	subsection of Treatment
<code>:NomenclatureHeading</code>	contains a nomenclatural act
<code>:NomenclatureCitationList</code>	list of citations of related concepts
<code>:MaterialsExamined</code>	list of examined specimens
<code>:BiologySection</code>	subsection of Treatment
<code>:DescriptionSection</code>	subsection of Treatment
<code>:TaxonomicKey</code>	section with an identification key
<code>:TaxonomicChecklist</code>	section with a list of taxa for a region
<code>:TaxonomicNameUsage</code>	mention of a taxonomic name

The classes from this subsection are based on the TaxPub XML Document Type Definition (DTD, also referred to loosely as XML schema, Catapano, 2010), on the structure of Biodiversity Data Journal’s taxonomic paper (Smith et al., 2013), and on the Treatment Ontologies (Catapano and Morris, 2016).

Furthermore, we introduce two properties: *contains* (`:contains`) and *mentions* (`:mentions`). *contains* is used to link parts of the article together and *mentions* links parts of the article to other concepts.

A graphical representation of the relationships between instances of the publishing-related classes that OpenBiodiv introduces is to be found in the diagram in Fig. 2.1.

⁶<http://openbiodiv.net/>

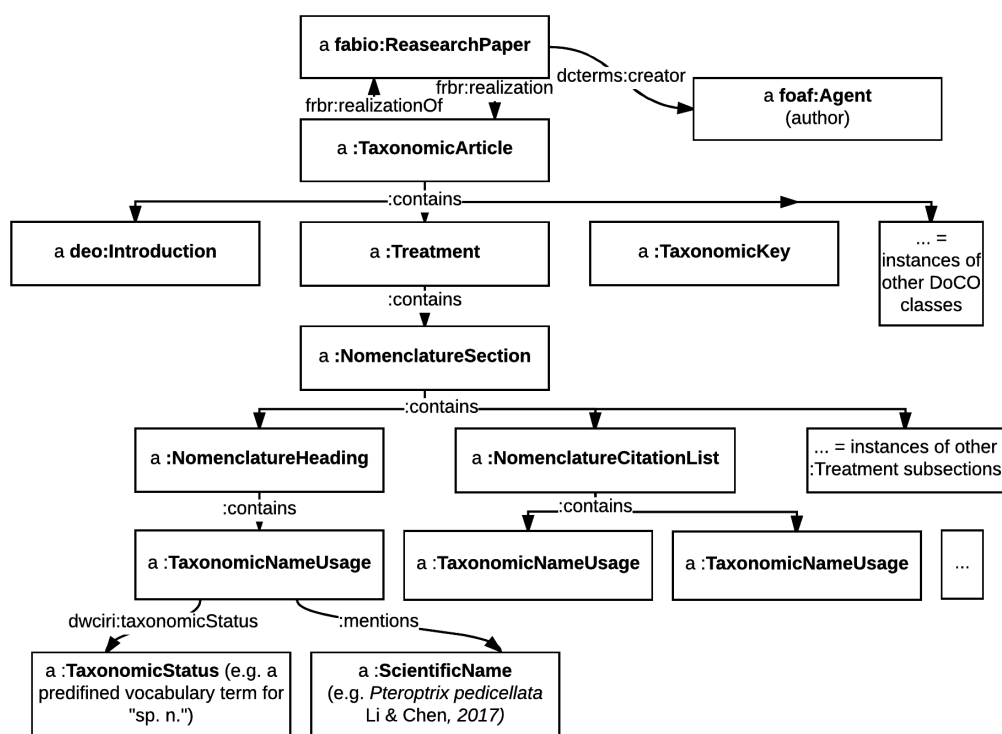


FIGURE 2.1: A graphical representation of the relationships between instances of the publishing-related classes that OpenBiodiv introduces.

Semantics, alignment, and usage

In this section we discuss how the classes and properties that we have introduced align to the Functional Requirements for Bibliographic Records (FRBR) model used by SPAR. In a nutshell taxonomic articles are considered FRBR Expressions of the more abstract FRBR Work that is the intellectual content of the article. Treatments are SPAR discourse elements akin to Introduction, Methods, etc. are also FRBR Expressions. Taxonomic Concepts are their corresponding FRBR Work's.

Figs. 2.2 and 2.3 give example usage in Turtle illustrating these ideas.

2.3.2 Semantic modeling of biological nomenclature

Biological nomenclature is a legacy system with over 200 years of accumulation from before the time of informatics and even from before the time of Darwinian Evolution! It is very hard to model due to complexity and has only partially been covered by the ontologies NOMEN and TNSS (introduced in subsection “Previous work”). With OpenBiodiv-O, I take a bottom-up approach of modeling the use of taxonomic names in articles. Where possible we align OpenBiodiv-O classes to NOMEN.

We have defined the class hierarchy of taxonomic names found in Fig. 2.4. Furthermore, we have introduced the class Taxonomic Name Usage (:TaxonomicNameUsage). Taxonomic name usages have been discussed widely in the community (e.g. in Pyle, 2016a); however, the meaning of term remains vague. The abbreviation TNU is used interchangeably for “taxon name usage” and for “taxonomic name usage.” In

```

:biodiversity-data-journal rdf:type fabio:Journal ;
    skos:prefLabel "Biodiversity Data Journal"@en ;
    skos:altLabel "BDJ"@en ;
    fabio:issn "1314-2836" ;
    fabio:eIssn "1314-2828" ;
    frbr:part :b90f6933-ab5e-4ce1-9379-12de9ef4eaa6 .

<http://dx.doi.org/10.3897/BDJ.1.e953> rdf:type fabio:TaxonomicArticle ;
    skos:prefLabel "10.3897/BDJ.1.e953" ;
    dc:title "Casuarinicola australis Taylor, 2010
(Hemiptera: Triozidae), newly recorded from New Zealand"@en ;
    prism:doi "10.3897/BDJ.1.e953" ;
    dcelements:publisher "Pensoft Publishers"@en ;
    fabio:hasPublicationYear "2013"^^xsd:gYear ;
    prism:publicationDate "2013-9-16"^^xsd:date ;
    dcterms:publisher :pensoft-publishers ;
    frbr:realizationOf :thorpe-2013 .

:thorpe-2013 rdf:type :ResearchPaper ;
    skos:prefLabel "Thorpe 2013"
    skos:altLabel "paper10.3897/BDJ.1.e953" ;
    dcterms:creator :stephen-e-thorpe ;
    prism:keywords "Casuarinicola australis"@en ;
    fabio:hasSubjectTerm :a2ee4929-90dd-4a7a-aa5c-08836f49d549 .

:pensoft-publishers rdf:type :Publisher ;
    skos:prefLabel "Pensoft Publishers"@en .

:stephen-e-thorpe rdf:type foaf:Person ;
    skos:prefLabel "Stephen E. Thorpe" ;
    foaf:firstName "Stephen E." ;
    foaf:surname "Thorpe" ;
    foaf:mbox "stephen_thorpe@yahoo.co.nz" ;
    :affiliation "School of Biological Sciences (Tamaki Campus),
    University of Auckland, Auckland, New Zealand"@en .

:a2ee4929-90dd-4a7a-aa5c-08836f49d549 rdf:type fabio:SubjectTerm ;
    rdfs:label "Casuarinicola australis"@en ;
    skos:inScheme :openbiodiv-subject-terms .

```

FIGURE 2.2: This example shows how to express the metadata of a taxonomic article with the SPAR Ontologies' model and the classes that OpenBiodiv defines. The code is in Turtle.

OpenBiodiv-O, a taxonomic name usage is the mentioning of a taxonomic name in the text, optionally followed by a taxonomic status.

For example, "*Heser stoevi* Deltschev 2016, sp. n." is a taxonomic name usage. The cursive text followed by the author and year of the original species description is the latinized scientific name. The abbreviation "sp. n." stands for the Latin *species novum*, indicating the discovery of a new taxon.

We also introduce the class Taxonomic Concept Label (:TaxonomicConceptLabel). A taxonomic concept label (TCL) is a Linnaean name plus a reference to a publication, where the discussed taxon is circumscribed. The link is via the keyword "sec." (Latin for (*secundum*, Berendsohn, 1995). An example would be "*Andropogon virginicus* var. *tenuispatheus* sec. Blomquist, 1948". Here, Blomquist, 1948 is a valid bibliographic reference to the publication where the concept is circumscribed.

We extracted taxonomic status abbreviations from about 4,000 articles across four taxonomic journals (ZooKeys, Biodiversity Data Journal, PhytoKeys, and MycoKeys) in order to create a taxonomic status vocabulary (see appendices) that covers the eight most common cases (Table 2.2). The Latin abbreviations that have been classified into these classes can be found on the OpenBiodiv-O GitHub page. (See Methods for more details).

```

<http://dx.doi.org/10.3897/BDJ.1.e953>
  :contains :abstract, :casuarinicola-australis-treatment .

:introduction rdf:type deo:Introduction, doco:Section ;
  c4o:hasContent "Casuarinicola australis Taylor, 2010 was described from
  Australia, where it is the most common and widespread member of its
  genus, being widely distributed in New South Wales, Queensland,
  South Australia, Victoria and Western Australia. "

:casuarinicola-australis-treatment rdf:type doco:Section, :Treatment ;
  :contains :casuarinicola-australis-nomenclature ,
            :casuarinicola-australis-materials ,
            :casuarinicola-australis-description ,
            :figure-box-1 ,
            :figure-box-2 .

:casuarinicola-australis-nomenclature rdf:type :NomenclatureSection ;
  :contains :casuarinicola-australis-nomenclature-heading .

:casuarinicola-australis-nomenclature-heading a :NomenclatureHeading ;
  cnt:chars "Casuarinicola australis Taylor, 2010" .

:casuarinicola-australis-materials rdf:type :MaterialsExamined ;
  c4o:hasContent "country: New Zealand;
  verbatimLocality: Mechanics Bay, Auckland City;
  verbatimElevation: 0-5 m;
  verbatimLatitude: 36.8474938105S ;
  verbatimLongitude: 174.7869624545E ;
  eventDate: 6 January 2013;
  sex: 1 male, 1 female;
  recordedBy: Stephen Thorpe;
  institutionCode: Auckland Museum" .

:casuarinicola-australis-description rdf:type :DescriptionSection ;
  c4o:hasContent "On 6 Jan 2013, I examined some Casuarina glauca trees growing
  in the vicinity of Ports of Auckland at Mechanics Bay." .

```

FIGURE 2.3: This examples shows how to express the article structure with the help of `:contains`. The code is in Turtle.

TABLE 2.2: OpenBiodiv Taxonomic Status Vocabulary.

Vocabulary Instance QName	Example Abbrev	Comment
<code>:TaxonomicUncertainty</code>	<i>incertae sedis</i>	Taxonomic Uncertainty
<code>:TaxonDiscovery</code>	<i>sp. n.</i>	Taxonomic Discovery
<code>:ReplacementName</code>	<i>comb. n.</i>	Replacement Name
<code>:UnavailableName</code>	<i>nomen dubium</i>	Unavailable Name
<code>:AvailableName</code>	<i>stat. rev.</i>	Available Name
<code>:TypeSpecimenDesignation</code>	<i>lectotype designation</i>	Type Specimen Designation
<code>:TypeSpeciesDesignation</code>	<i>type species</i>	Type Species Designation
<code>:NewOccurrenceRecord</code>	<i>new country record</i>	New Occurrence Record (for region)

Based on our analysis of taxonomic statuses, we have identified two Code-compliant patterns of relationship between latinized scientific names (Fig. 2.5). The pattern *replacement name*, implemented via the property `:replacementName`, indicates that a certain Linnaean name should be used instead of another Linnaean name. It covers a wide variety of cases in the Codes, such as, for example, the placement of one species taxon in a new genus ("*comb. n.*"), the correction of a name for nomenclatural reasons ("*nomen novum*"), or the application of the Principle of Priority for the discovery of synonyms ("*syn. nov.*", International Commission on Zoological Nomenclature, 2017).

The other pattern is that of *related names* (`:relatedName`). It is a broader pattern, indicating that two names are somehow related. For example, they may be synonyms,

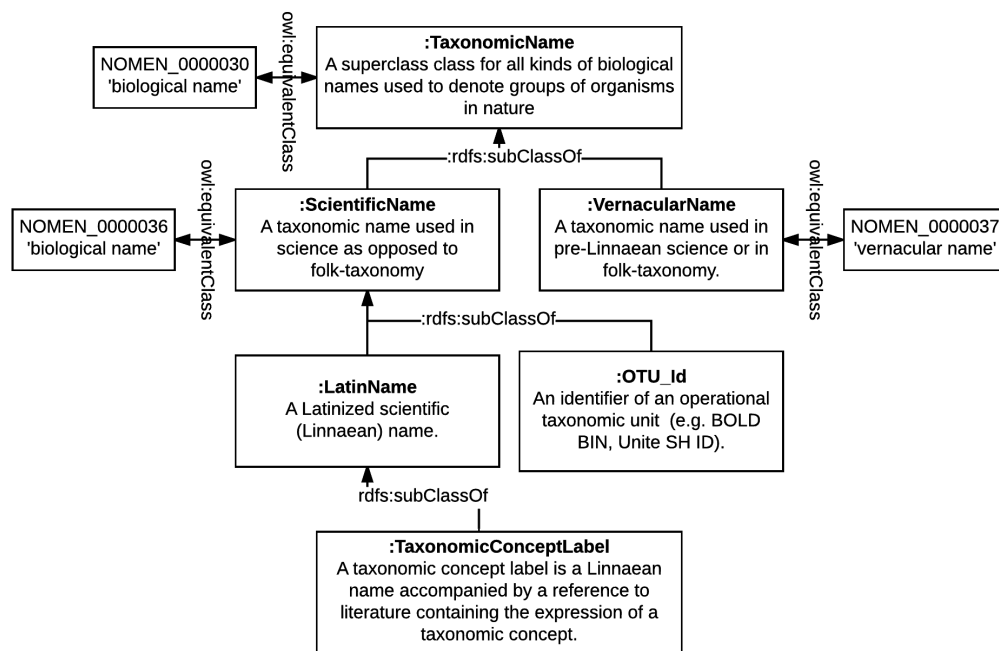


FIGURE 2.4: We created this class hierarchy to accommodate both traditional taxonomic name usages and the usage of taxonomic concept labels and operational taxonomic units.

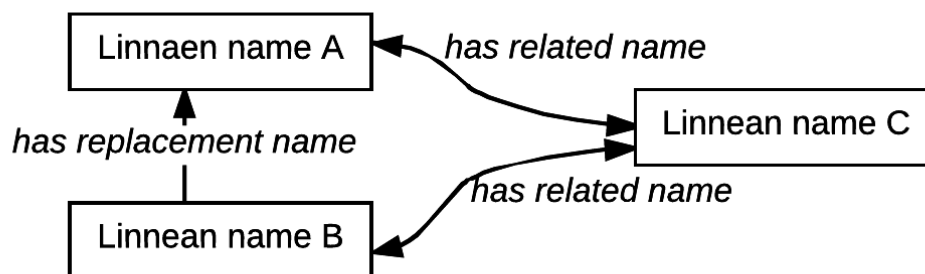


FIGURE 2.5: Chains of *replacement names* can be followed to find the currently used name. *Related name* indicates that two names are related somehow, but not which one is preferable.

with one replacing the other, or they may point to taxonomically related taxonomic concepts. For example, *Harmonia manillana* (Mulsant, 1866) is related to *Caria manillana* Mulsant, 1866 since, as per Poorani and Booth, 2016, a name-bearing type (lectotype) of *Harmonia manillana* (Mulsant, 1866) sec. Poorani Poorani and Booth, 2016 is named *Caria manillana* Mulsant, 1866.

Semantics, alignment and usage

As evident from Fig. 2.4, OpenBiodiv-O taxonomic names are aligned to NOMEN names.

The linking between text and taxonomic names must pass through the intermediary class Taxonomic Name Usage. As parts of the manuscript, taxonomic name usages link document components to taxonomic names. Taxonomic name usages are *contained* in sections such as Treatment, and *mention* a taxonomic name as illustrated in the example in Fig. 2.6.

```

:casuarinicola-australis-nomenclature-heading
  po:contains :casuarinicola-australis-TNU .

:casuarinicola-australis-TNU a :TaxonomicNameUsage ;
  dc:date "2013-9-16"^^xsd:date ;
  cnt:chars "Casuarinicola australis Taylor, 2010" ;
  dwc:genus "Casuarinicola" ;
  dwc:specificEpithet "australis" ;
  dwc:scientificNameAuthorship "Taylor, 2010" ;
  # we can infer the following because we are in the treatment heading
  dwc:nameAccordingToId "doi: 10.3897/BDJ.1.e953" ;
  pkm:mentions :casuarinicola-australis-taylor,
               :casuarinicola-australis-taylor-sec-thorpe-2013 .

:casuarinicola-australis-taylor a :ScientificName ;
  rdfs:label "Casuarinicola australis Taylor, 2010" ;
  dwc:genus "Casuarinicola" ;
  dwc:specificEpithet "australis" ;
  dwc:scientificNameAuthorship "Taylor, 2010" .

:casuarinicola-australis-taylor-sec-thorpe-2013 a :TaxonomicConceptLabel ;
  rdfs:label "Casuarinicola australis Taylor, 2010 sec. Thorpe 2013" ;
  dwc:genus "Casuarinicola" ;
  dwc:specificEpithet "australis" ;
  dwc:scientificNameAuthorship "Taylor, 2010" .
  dwc:nameAccordingToId "doi: 10.3897/BDJ.1.e953" ;
  :nameAccordingTo <http://dx.doi.org/10.3897/BDJ.1.e953> .

```

FIGURE 2.6: This examples shows how taxonomic name usages link document components to taxonomic names. The code is in Turtle.

2.3.3 Semantic Modeling of the Taxonomic Concepts

In OpenBiodiv-O taxonomic names are not the carriers of semantic information about taxa. This task is accomplished by a new class, Taxonomic Concept (`:TaxonomicConcept`). A taxonomic concept is the theory that a taxonomist forms about a taxon in a scholarly biological taxonomic publication and thus always has a taxonomic concept label. We also introduce a more general class, Operational Taxonomic Unit (`:OperationalTaxonomicUnit`) that can be used for all kinds of taxonomic hypotheses, including ones that don't have a proper taxonomic concept label. The class hierarchy has been illustrated in Fig. 2.7.

Taxonomic concepts are related to taxonomic names—including taxonomic concept labels—via the property *has taxonomic name* (`:taxonomicName`) and its sub-properties mimicking in their range the hierarchy of taxonomic names that we introduced earlier. We have defined a property specifically to link taxonomic concepts to taxonomic concept labels, *has taxonomic concept label* (`:taxonomicConceptLabel`). The property hierarchy diagram is shown in Fig. 2.8.

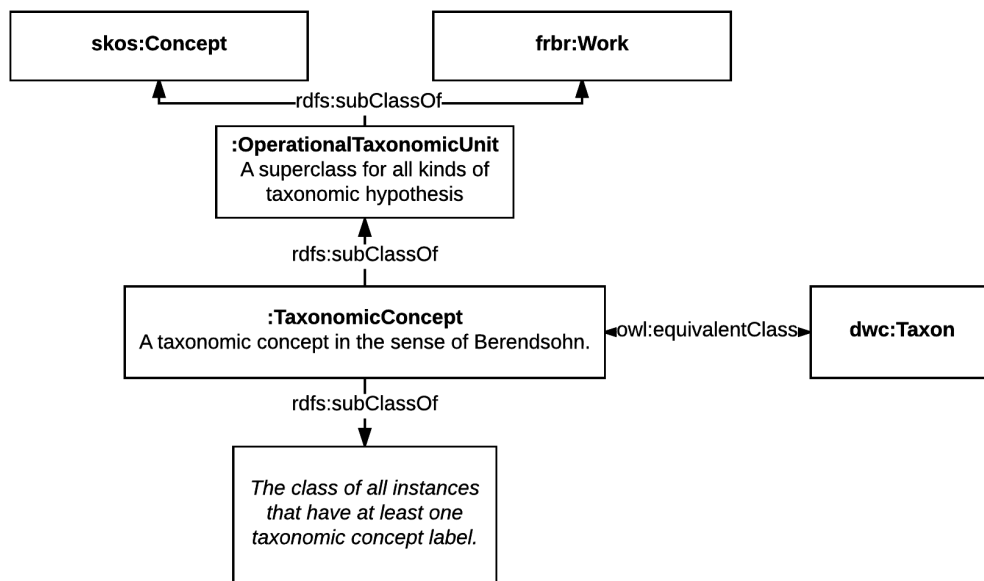


FIGURE 2.7: A taxonomic concept is a `skos:Concept`, a `frbr:Work`, a `dwc:Taxon` and has at least one taxonomic concept label.

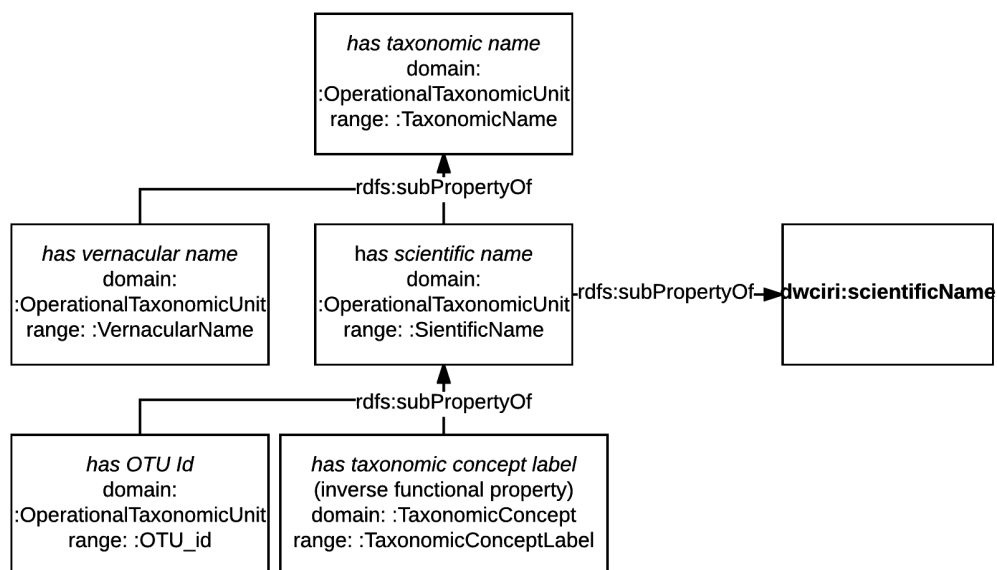


FIGURE 2.8: Property hierarchy is aligned with the taxonomic name class hierarchy and with DarwinCore.

There are two ways to relate taxonomic concepts to each other (Fig. 2.9). As we pointed out earlier, historically taxonomic concepts form the hierarchy known as biological taxonomy. To express such simple semantic relations, it is fully sufficient to use the SKOS semantic vocabulary Miles and Bechofer.

However, these simple relationships are not well suited for machine reasoning.

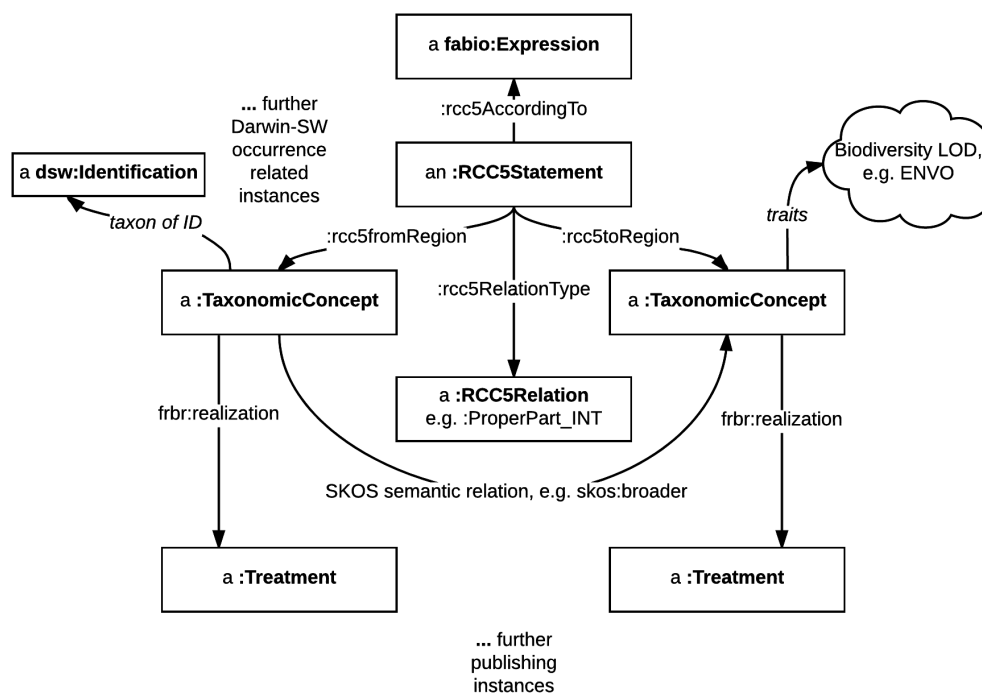


FIGURE 2.9: In order to express an RCC-5 relationship between concepts, create an `:RCC5Statement` and use the corresponding properties to link two taxonomic concepts via it. Further, taxonomic concepts are linked to traits (e.g. ecology in ENVO), occurrences (e.g. Darwin-SW) and realize treatments.

This is why Franz and Peet Franz and Peet, 2009 suggested, building on previous work by e.g. Koperski et al., 2000, to use the RCC-5 language to express relationships between taxonomic concepts. Furthermore, the Euler (Chen et al., 2014) program was developed, which uses Answer Set Programming (ASP) to reason over RCC-5 taxonomic relationships. An answer set reasoner is not part of OpenBiodiv as this task can be accomplished by Euler; however, we have provided an RCC-5 dictionary class (`:RCC5Dictionary`), an RCC-5 relation term class (`:RCC5Relation`), a vocabulary of such terms to express the RCC-5 relationships in RDF (see appendices), as well as a class and properties to express RCC-5 statements (`:RCC5Statement`, `:rcc5Property`, and subproperties).

Semantics and alignment

In this section taxonomic concepts are aligned to DarwinCore (DwC) and a discussion of how taxonomic concepts related to each other either via simple relations (SKOS) and fine-grained (RCC-5) is presented. Also the relationships between biological names and scientific concepts are discussed. We treat instances of our class Taxonomic Concept as functionally equivalent to DwC Taxa. We can now list what types of relationships between names and taxonomic concepts are allowed: (1) The relationship between a taxonomic concept and a name that is not a taxonomic concept label is many-to-many—i.e. one Linnaean name can be a mention of multiple taxonomic concepts, and one taxonomic concept may have multiple Linnaean names. (2) The relationship

between a taxonomic concept and a taxonomic concept label is one-to-many: while a taxonomic concept may have more than one (at least one is needed) labels, every label uniquely identifies a concept. These logical restrictions make taxonomic concept labels into unique identifiers to taxonomic concepts, something that Linnaean names are not.

Usage

In Fig. 2.10, Fig. 2.11, and Fig. 2.12, and Fig 2.13) we provide some useful examples.

```
:concept-casuarinicola-australis-thorpe rdf:type :TaxonomicConcept ;
  :taxonomicConceptLabel :casuarinicola-australis-taylor-sec-thorpe-2013 .

:concept-casuarinicola-taylor rdf:type :TaxonomicConcept ;
  skos:broader concept-thorpe .
```

FIGURE 2.10: We can use SKOS semantic properties to illustrate simple relationships between taxonomic concepts.

```
:statement rdf:type :RCC5Statement ;
  :rcc5FromRegion :concept-casuarinicola-australis-thorpe ;
  :rcc5ToRegion :concept-casuarinicola-taylor ;
  :rcc5AccordingTo <http://dx.doi.org/10.3897/BDJ.1.e953> ;
  :rcc5RelationType :ProperPart_INT .
```

FIGURE 2.11: In order to express an RCC-5 relationship between concepts, create an `:RCC5Statement` and use the corresponding properties to link two taxonomic concepts via it. SKOS relations relate concepts directly.

```
:australian-casuarina-forest rdf:type <http://purl.obolibrary.org/obo/ENVO_01000174> .
:hasHabitat owl:sameAs <http://purl.obolibrary.org/obo/RO_0002303> .
:concept-casuarinicola-australis-thorpe :hasHabitat :australian-casuarina-forest .
```

FIGURE 2.12: We create a shortcut for *has habitat* and instance of the "forest biome" and link them to our taxonomic concept in order to express the fact that specimens of it have been found to live in *Casuarina* trees.

```
:casuarinicola-australis-treatment frbr:realizationOf :concept-casuarinicola-australis-thorpe.
```

FIGURE 2.13: A treatment is the realization of a taxonomic concept.

2.4 Discussion

OpenBiodiv-O is—together with the Treatment Ontologies (Catapano and Morris, 2016)—the first effort to model taxonomic articles as RDF. It introduces classes and properties in the domains of biodiversity publishing and biological taxonomy and

aligns them with the SPAR Ontologies, the Treatment Ontologies, the Open Biomedical Ontologies (OBO), TaxPub, NOMEN, and DarwinCore. We believe this introduction bridges the ontological gap that we had outlined in our aims and allows for the creation of a Linked Open Dataset (LOD) of biodiversity information (biodiversity knowledge graph, Senderov and Penev, 2016; Page, 2016).

Furthermore, this biodiversity knowledge graph, together with this ontology, additional semantic rules, and user software forms the OpenBiodiv system. OpenBiodiv, as any taxonomic information system should, has taxonomic names as a key building block. For any given taxonomic name, the user will be able to rely on two patterns—*replacement name* and *related name*—to get answers to two questions of high importance to the working taxonomist. First: what is the current and historical usage of any given Linnaean name? Second: given a particular name, what other related names ought to be considered in a taxonomic discussion?

In this section we carry out a discussion how the model of OpenBiodiv can be used to store *multiplicity of opinion* about taxonomic relationships and thus democratize the taxonomic process. We further discuss the usefulness of OpenBiodiv to *answer competency questions* from biological taxonomy. These will be touched upon more in the next chapter.

2.5 Conclusions

The chapter provides an informal conceptualization of the taxonomic process and a formalization in OpenBiodiv-O. It introduces classes and properties in the domains of biodiversity publishing and biological systematics and aligns them with the important domain-specific ontologies. By bridging the ontological gap between the publishing and the biodiversity domains, it will enable the creation of Open Biodiversity Knowledge Management System, consisting of (1) the ontology itself; (2) a Linked Open Dataset (LOD) of biodiversity information (biodiversity knowledge graph); and (3) user interface components aimed at searching, browsing and discovering knowledge in big corpora of previously dispersed scholarly publications. Through the usage of taxonomic concepts, we have included mechanisms for democratization of the scholarly process and not forcing a taxonomic opinion on the users.

Chapter 3

Summary of Chapter 3: OpenBiodiv Linked Open Dataset

In Chapter 3 I explore in detail the data sources and their data models.

I, with the help of my support team—see the Acknowledgements in the back—have created a Linked Open Dataset, OpenBiodiv LOD, comprising biodiversity information extracted from Pensoft journals and from Plazi Treatment Bank, and which was integrated with the GBIF Taxonomic Backbone. As ontology, I use the new OpenBiodiv-O developed through the course of the dissertation. I propose to the biodiversity informatics community to use OpenBiodiv LOD as the central point for a biodiversity knowledge graph. OpenBiodiv LOD is an RDF dataset adhering to the principles of Linked Open Data. It is available under <http://graph.openbiodiv.net>, which provides a SPARQL endpoint for it.

OpenBiodiv LOD is a synthetic dataset. It does not contain previously unpublished data. Instead it integrates information previously found in academic journals and databases into one dataset. It also contains extracted, previously inaccessible information from the original datasets in the form of relations. In the next few paragraphs we discuss the sources of information that were combined to form OpenBiodiv LOD and the types of resources that have been extracted, as well as the overall data model. We also discuss the principles of Linked Open Data that tie everything together. The chapter ends with many examples of queries on the dataset and with a technical discussion of how it was generated.

3.1 Data Sources

The data in OpenBiodiv at the time of writing this thesis comes from three major sources: the GBIF Backbone Taxonomy (GBIF Secretariat, 2017b), journal articles published by Pensoft, and Plazi Treatment Bank (Fig. 3.1).

3.1.1 GBIF Backbone Taxonomy

GBIF is the largest international repository of occurrence data, i.e. data about the presence of an organism of a given taxon at a given place and time. GBIF allows its users to do searches on its occurrence data utilizing a taxonomic hierarchy. For example, it is possible to query the database for occurrences of organisms belonging to a specific genus: a search for the beetle genus *Harmonia* sec. GBIF Secretariat, 2017b on 30 June 2018 returned 575,376 results. This search is possible thanks to the GBIF Backbone Taxonomy also known as Nub (GBIF Secretariat, 2017b). Nub is a database organizing taxonomic concepts in a hierarchy covering all names used in occurrence records harvested by GBIF. It is a single synthetic (algorithmically generated) management classification with the goal of covering all names present in

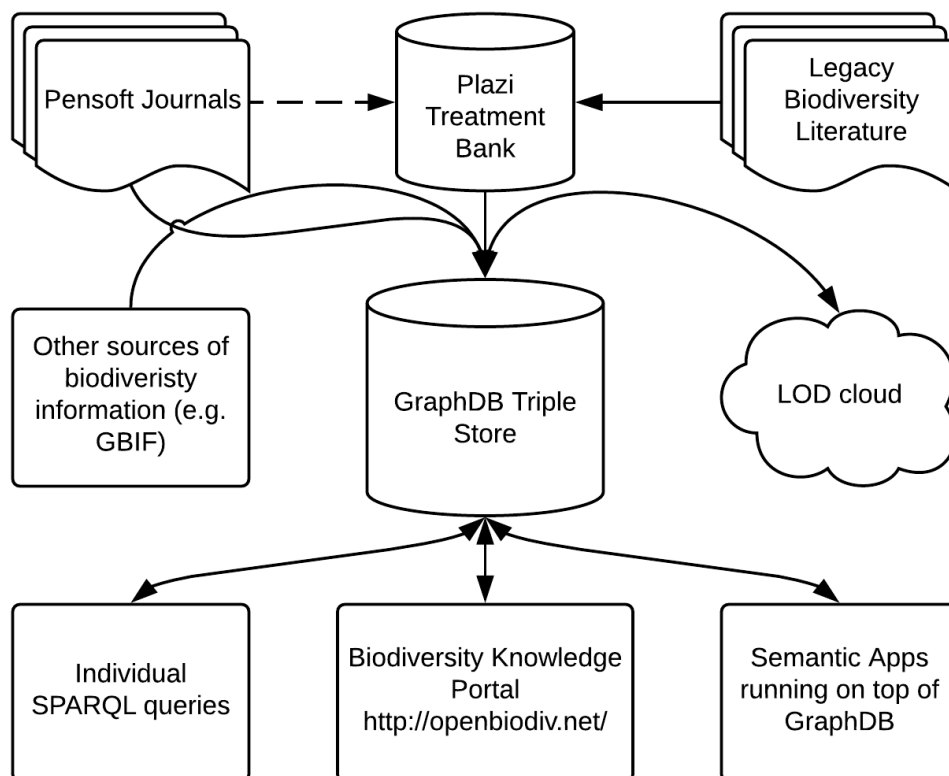


FIGURE 3.1: A simplified version of the OpenBiodiv architecture presented in Chapter 1 focusing on the sources of information.

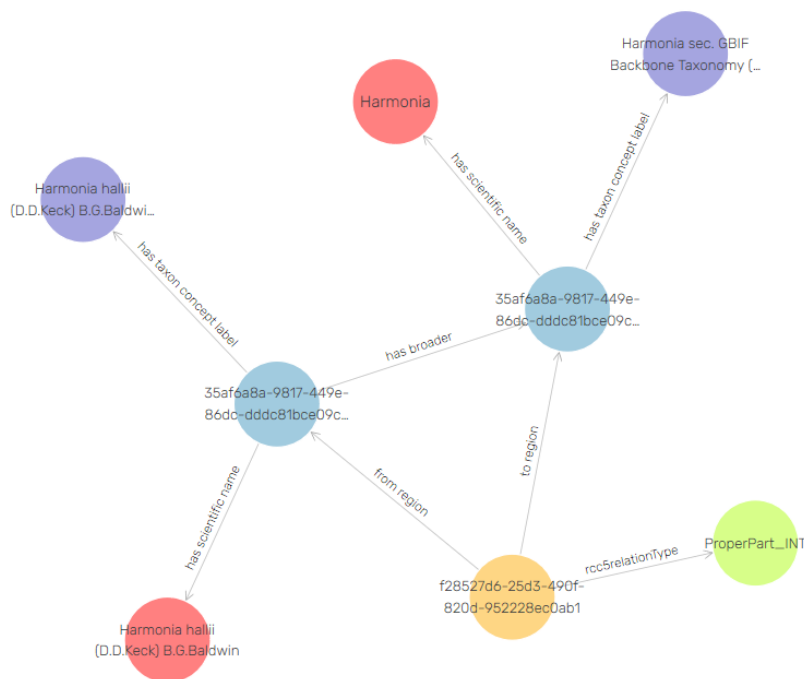


FIGURE 3.2: Illustration of the representation of hierarchical information imported from the GBIF Backbone Taxonomy as two taxonomic concepts, *Harmonia halii* sec. GBIF Secretariat, 2017a and *Harmonia* sec. GBIF Secretariat, 2017a. Each concept has an associated scientific name via *has scientific name*; however, the hierarchical information is not encoded in the names. The hierarchical relationship between *Harmonia halii* sec. GBIF Secretariat, 2017a and *Harmonia* sec. GBIF Secretariat, 2017a is encoded both as SKOS *has broader* and reified via the RCC-5 relationship encoded in f28527d6-25d3-490f-820d-952228ec0ab1.

GBIF's datasets. Thus, the GBIF backbone does not represent an expert consensus on how taxa are hierarchically arranged according to evolutionary criteria in Nature.

Keeping in mind this critique, it is evident how the backbone taxonomy allows GBIF to integrate name based information from diverse sources such as Encyclopedia of Life (EOL), Genbank, or the International Union for Conservation of Nature, and provides a facility for taxonomic searching and browsing.

In order to grant the same capabilities to OpenBiodiv, we have imported Nub as instances of `openbiodiv:TaxonomicConcept` according to OpenBiodiv-O (Fig. 3.2).

3.1.2 Pensoft and Plazi

All valid articles from the journals published by Pensoft listed in Table 3.1 have been converted to RDF and stored in the biodiversity knowledge graph. Additionally, all valid taxonomic treatments from Plazi Treatment Bank have been converted to RDF and stored in the graph as well. Furthermore, the RDF-ization procedure is triggered automatically on a weekly basis and thus the semantic database is always updated with the newest articles published by Pensoft and newest taxonomic treatments extracted

by Plazi. The RDF-ization is made possible by the fact that all Pensoft journals are published as XML according to TaxPub, an extension of the NLM/NCBI journal publishing DTD for taxonomic description (Catapano, 2010) and, similarly, all Plazi treatments follow the TaxonX XML Schema (Penev et al., 2011) (Fig. 3.3).

LISTING 3.1: Taxonomic name usage of the name *P. emarginaticeps* in Taxpub. Name parts are tagged with `tp:taxon-name-part` and the expansion of abbreviations (regularization) is marked up with the attribute `reg`

```
<tp:taxon-name>
  <tp:taxon-name-part taxon-name-part-type="genus" reg="Pristaulacus">
    P.
  </tp:taxon-name-part>
  <tp:taxon-name-part taxon-name-part-type="species" reg="emarginaticeps">
    emarginaticeps
  </tp:taxon-name-part>
  <tp:taxon-name-part taxon-name-part-type="authority">
    Turner 1922
  </tp:taxon-name-part>
</tp:taxon-name>
```

TABLE 3.1: RDF-ized biodiversity journals published by Pensoft.

Journal Name	Submission Style	Number of Articles
ZooKeys	Word document	3829
PhytoKeys	Word document	537
MycoKeys	Word document	127
Biodiversity Data Journal	Web based (ARPHA)	490
Journal of Orthoptera Research	Word document	32

TABLE 3.2: Datatypes marked up in TaxPub and TaxonX articles and the corresponding RDF types of the generated RDF resources. The TaxPub and TaxonX columns contain boolean values indicating whether the information about the datatype is retrieved from files encoded in the corresponding schema.

Datatype	TaxPub	TaxonX	RDF Type
Article metadata	T	T	<code>fabio:JournalArticle</code> and related
Keyword group	T	F	<code>openbiodiv:KeywordGroup</code>
Abstract	T	T	<code>sro:Abstract</code>
Title	T	F	<code>doco:Title</code>
Author	T	T	<code>foaf:Person</code>
Introduction section	T	F	<code>deo:Introduction</code>
Discussion section	T	T	<code>orb:Discussion</code>
Treatment section	T	T	<code>openbiodiv:Treatment</code>
Nomenclature section	T	T	<code>openbiodiv:NomenclatureSection</code>
Materials examined	T	T	<code>openbiodiv:MaterialsExamined</code>
Diagnosis section	T	T	<code>openbiodiv:DiagnosisSection</code>
Distribution section	T	T	<code>openbiodiv:DistributionSection</code>
Taxonomic key	T	T	<code>openbiodiv:TaxonomicKey</code>
Figure	T	T	<code>doco:Figure</code>
Taxonomic name usage	T	T	<code>openbiodiv:TaxonomicNameUsage</code>

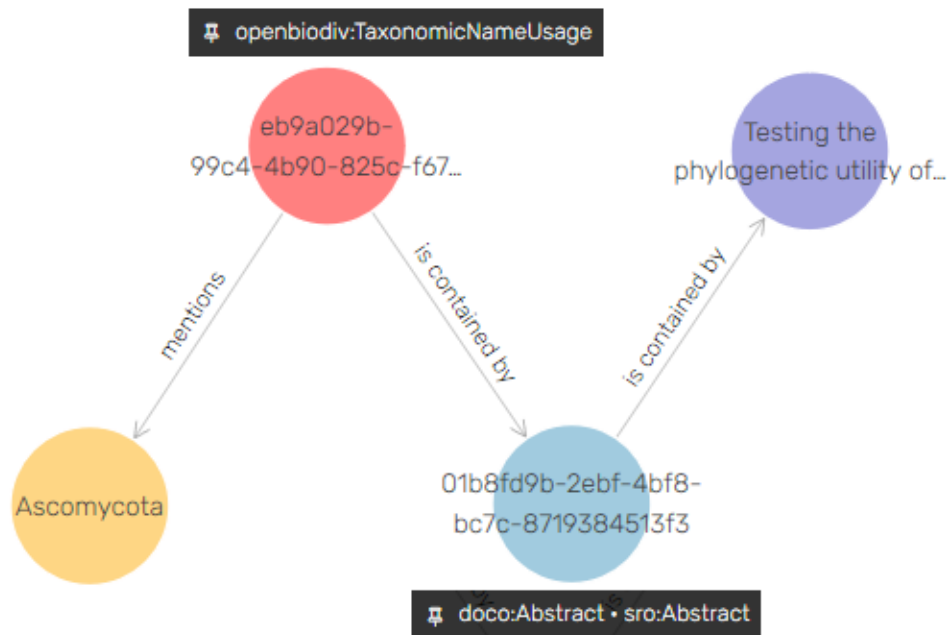


FIGURE 3.3: The taxonomic name usage (`openbiodiv:eb9a029b-99c4-4b90-825c-f670fb88900d`) is linked to the scientific name it mentions, *Ascomycota* and to the part of the article (abstract) that it is contained in.

3.2 Linked Open Data

Linked Open Data (LOD, Heath and Bizer, 2011) is a concept of the Semantic Web (Berners-Lee et al., 2001) applied to ensure that data published on the Web is reusable, discoverable and most importantly to ensure that pieces of data published by different entities can work together. The principles of LOD are the following (Heath and Bizer, 2011)

1. Use URIs as names for things.
2. Use HTTP URIs so people can lookup these things.
3. When someone looks up a URI, provide useful information, using the standards (RDF, SPARQL).
4. Include links to other URIs so they can discover more things.

We have followed these guidelines when creating the OpenBiodiv LOD. We will now discuss each of these points separately.

3.2.1 Usage of URI's as resource identifiers

Every instance in OpenBiodiv LOD is uniquely identifiable by a HTTP URI of the following form: `http://openbiodiv.net/uuid-(suffix)`. All instance identifiers in OpenBiodiv LOD follow this schema. The optional suffix field is assigned only to resources extracted from GBIF.

In this subsection we further discuss how identifiers are assigned to resources extracted from Pensoft and Plazi as well as to the GBIF taxonomic concepts.

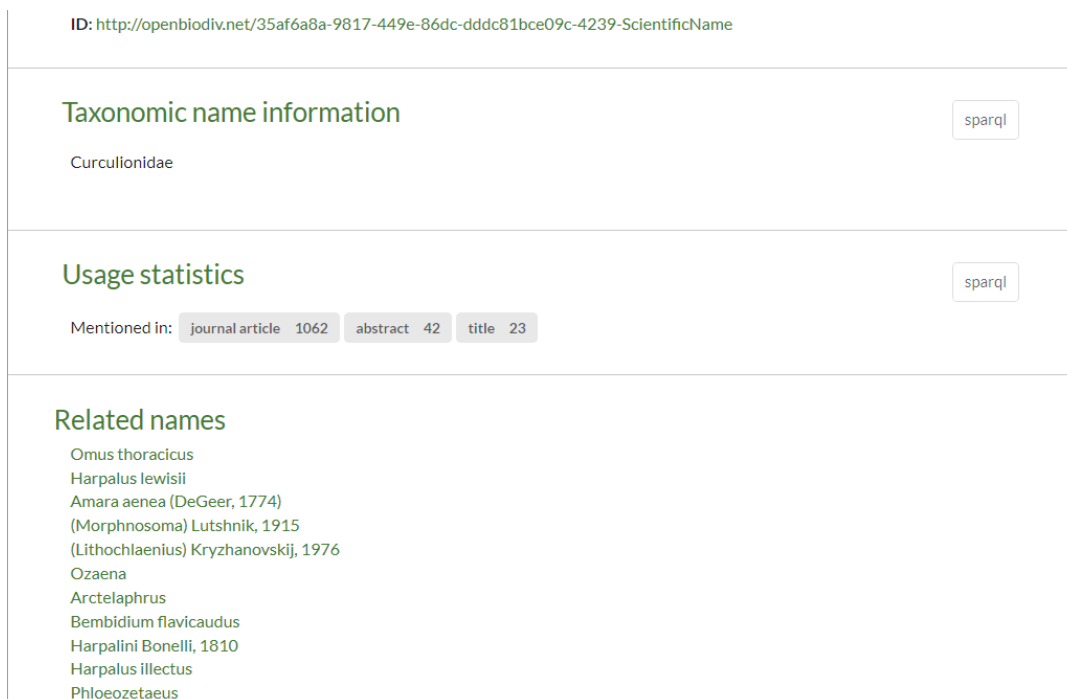


FIGURE 3.4: Visualization.

3.2.2 Usage of HTTP URI's and dereferencing

As per the Linked Data Principles, we use dereferenceable HTTP URIs for our resources. For example, if a web-browser opens <http://openbiodiv.net/35af6a8a-9817-449e-86dc-dddc81bce09c-4239-ScientificName> a web-page is displayed (Fig. 3.4) providing useful information for the name such as where it used and other names are related to it. Also it is possible to request OpenBiodiv resources via Curl with the header `Content-Type: application/rdf+xml` and an RDF representation of the resources is returned.

3.2.3 Linking to other resources

First, all resources in OpenBiodiv form a graph (there are no disconnected parts). The data model is discussed in the next section. Second, taxonomic names are linked to external databases via `dwc:taxonID`. These are strings containing GBIF ID's, ZooBank ID's, LSID's, etc. Unfortunately as HTTP URI's have not gained popularity in the biodiversity informatics community, the only true resource-id-to-resource-id links are within OpenBiodiv itself. However, we hope that the introduction of OpenBiodiv LOD contributes to the amelioration of this situation.

3.3 Data Model

When creating the RDF graph we have conformed to the OpenBiodiv Ontology described in Chapter 2 and well-established community ontologies (Fig. 3.5). In particular, (1) we use the Semantic Publishing and Referencing Ontologies (SPAR, Peroni, 2014) to model entities from publishing such as Journal, Article, Section, Figure, Table, and so on; and (2) we use the DarwinCore (DwC, Wiczorek et al., 2012)

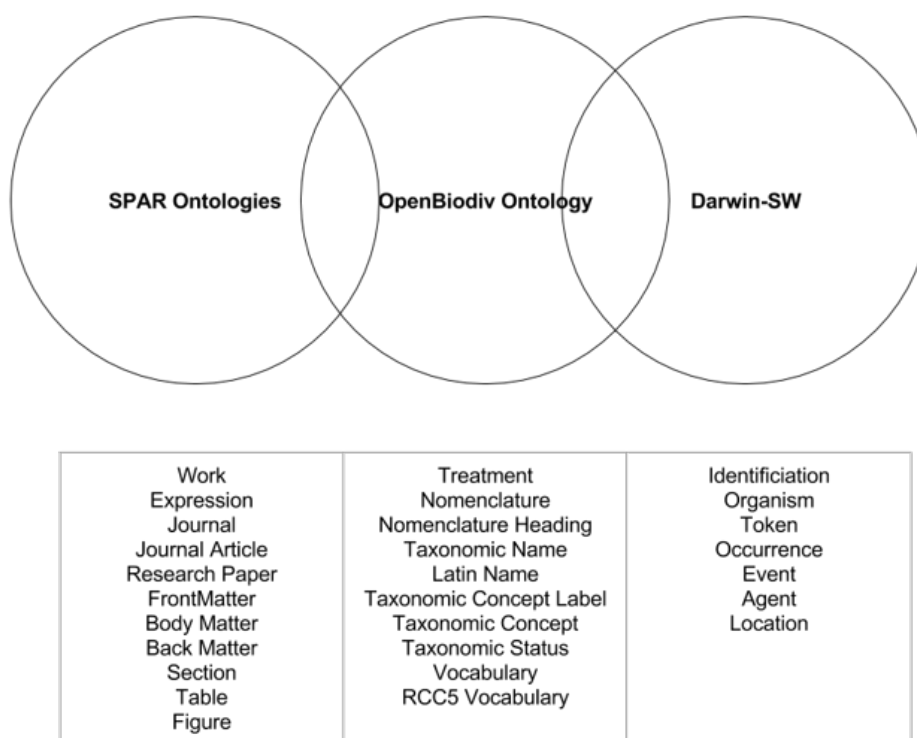


FIGURE 3.5: OpenBiodiv-O is an ontology that links the publishing domain with the biodiversity domain. Major resource types covered by each of the ontology families are given in the box below the Venn diagram. Important resources from the publishing domain are listed in the leftmost column and from biodiversity informatics in the rightmost column. The middle one covers important OpenBiodiv-O resources.

community standard and its extension, the Darwin-SW (Baskauf and Webb, 2016) ontology, to model entities the biodiversity domain.

SPAR provides facilities to deal with the dichotomy between the abstract representation of knowledge through the class `Work` and its concrete representation through the class `Expression`. For example, a `fabio:JournalArticle` can be the realization of a `fabio:ResearchPaper`. On the other hand, the DwC community standard gives a standard way to express properties from taxonomy and biodiversity science and its extension Darwin-SW a way to reify elements of an occurrence instance such as Identification, Organism, Token, and so on. A caveat: the current version of OpenBiodiv-LOD does not store yet occurrence information but all necessary infrastructure is in place to include them in the next release.

3.4 Examples of SPARQL queries

As SPAR, DwC, and OpenBiodiv-O have already been explained elsewhere, we shall illustrate the data model by issuing sample SPARQL queries illuminating aspects of it.

3.4.1 Simple queries

In this section, we give some simple queries. For example, how to search for an author, for a scientific name, etc.

Query the article structure

A unique feature of OpenBiodiv LOD is that articles are broken down into their components (see e.g. Table 3.2 later in this Chapter) and mentions (e.g. taxonomic name usages) connected to the specific part of the article and not just to the article in general. We illustrate how to build queries utilizing this structure.

Query for taxonomic concepts

A key feature of OpenBiodiv-O is that it allows for the separation of taxonomic concepts from scientific names. Scientific names are linked both to the components of an article that mentions them and to taxonomic concepts. To illustrate this, we can create a query uniting information from concepts from the GBIF Backbone Taxonomy with semantics coming from the article structure.

Fuzzy Queries via Lucene

The SPARQL endpoint of OpenBiodiv LOD supports fuzzy matching via a Lucene connector (Ontotext, 2018). In taxonomy, this can be a very useful as due to multiplicity of taxonomic names and the complexities of Latin grammar, one often does not remember the correct spelling of a name. This can lead to no matches in an exact search even though the system may contain information about that name. We illustrate how to do Lucene queries in OpenBiodiv via SPARQL.

3.4.2 Competency question answering via SPARQL

At the end of Chapter 2 I suggested some competency questions that may be answered by OpenBiodiv. In this subsection I show how these can be answered with the help of OpenBiodiv.

Validity of a taxonomic name

Of central importance is the question of whether a given taxonomic name is valid or not. We give the formal criteria on judging the validity of a taxonomic name and translate these into SPARQL.

Investigation of the impact of the lost collections of Museu Nacional

We conclude the discussion of SPARQL queries by showing how OpenBiodiv can be used to assess the impact of the tragically lost collection of the Museu Nacional de Rio de Janeiro (MNRJ).

3.5 Dataset Generation

In the previous section on sources we examined the data formats that each source provides. The inputs are either XML (Pensoft and Plazi) or CSV (GBIF). Thus, the raw data-streams are semi-structured and the dataset generation problem can be thought of as an information retrieval and transformation problem. The input

is encoded in three different data models—DarwinCore CSV (GBIF), TaxPub XML (Pensoft), and TaxonX XML (Plazi). The output of the transformation pipeline is knowledge represented in a fully-structured way according to the ontology.

3.5.1 Obtaining the data

The first step before running any transformation is to obtain the raw inputs. GBIF’s taxonomic backbone is available under

<https://www.gbif.org/dataset/d7dddbf4-2cf0-4f39-9b2a-bb099caae36c>.

There is an RSS feed from which Plazi’s treatments can be downloaded on a daily basis under <http://tb.plazi.org/GgServer/xml.rss.xml>. Each of Pensoft’s journals has a public API endpoint under [http://\[journal_name\].pensoft.net/lib/journal_archive.php](http://[journal_name].pensoft.net/lib/journal_archive.php), where `[journal_name]` ought to be replaced with the name of the Pensoft journal. E.g. `bdj` to make http://bdj.pensoft.net/lib/journal_archive.php.

3.5.2 Tools

In order to carry out the dataset generation we made use of the following tools:

1. RDF4R R package¹, which is described in Chapter 4 and deals with all RDF-related issues such as accessing a triple store, serializing the in-memory resource representations to Turtle files, etc.
2. ROpenBio R package², which implements the data retrieval and transformations described in this chapter.
3. TSV4RDF, which is a PHP library for mapping CSV to RDF developed by Pensoft. It is closed-source and developed outside of the scope of the dissertation and is not discussed in detail.
4. The OpenBiodiv base³, which contains scripts needed for the initialization and updating of the database.

In the rest of the section we describe the transformation from XML as it is implemented in ROpenBio. We do not describe the TSV4RDF transformation of GBIF to RDF as it is a closed source product.

3.5.3 XML to RDF transformation

In order to transform an article represented as an XML document to RDF, we make use of the hierarchical nature of XML and solve the problem recursively with the following Extractor procedure in Algorithm 1. The extractor’s procedure input is an XML node and its output is the RDF corresponding to the XML node. The extractor procedure has three essential steps: atoms extraction, RDF constructions from the extracted atoms, a divide-and-conquer step that recursively calls itself and unites the results. Extraction of a whole article is achieved by calling the Extractor on the root node of the article.

¹RDF4R package on GitHub <https://github.com/vsenderov/rdf4r>

²ROpenBio R package on GitHub <https://github.com/pensoft/ropenbio>

³OpenBiodiv Base <https://github.com/vsenderov/OpenBiodiv>

Algorithm 1 The Extractor procedure

```

1: procedure EXTRACTOR(XML Node  $X$ )
2:    $a \leftarrow$  extract atoms of  $X$                                 ▷ Atoms extraction
3:    $r \leftarrow$  construct RDF from  $a$                                ▷ RDF construction
4:    $C \leftarrow$  find relevant sub-nodes of  $X$                      ▷ Recursively applies itself
5:    $R \leftarrow$  apply Extractor on each  $C_i \in C$ 
6:   return  $r \cup R$ 
7: end procedure

```

Atoms extraction

In this subsection we elaborate on the text-fields of the XML (atoms) are extracted in our framework utilizing the XPATH query language.

RDF Generation

Once the atoms have been extracted they can be put together as RDF. Conceptually, this is straightforward as for each atom we know its type and therefore we know which RDF property to use. The author example is given in Listing 3.2.

LISTING 3.2: RDF snippet of an author. This is a somewhat idealized situation in which the language of the address was available from the article.

```

@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .

:a a foaf:Person ;
  rdfs:label "Aijaz_Ahmad_Wachkoo" .
  :affiliation "Central_Institute_of_Temperate_Horticulture,_Srinagar,_Jammu&Kashmir,_India"@en ;
  foaf:familyName "Wachkoo" ;
  foaf:givenName "Aijaz_Ahmad" .

```

LISTING 3.3: .

```

:2b836ad5-db56-4093-9752-33c9f7892de6 rdf:type fabio:JournalArticle ;
  rdfs:label "Changes_to_publication_requirements_made_at_the_XVIII_International_Botanical_Congress_in_Melbourne_what_does_e-publication_mean_for_you?" ;
  dc:title "Changes_to_publication_requirements_made_at_the_XVIII_International_Botanical_Congress_in_Melbourne_what_does_e-publication_mean_for_you?" ;
  prism:doi "10.3897/mycokeys.1.1961" ;
  dc:publisher "Pensoft_Publishers" ;
  prism:publicationDate "2011-9-14"^^xsd:date ;
  dcterm:publisher openbiodiv:0df76aab-1fcf-4118-8e50-198e830a7bed .
  openbiodiv:151a37ba-a337-4855-8e01-200f5ec0251b rdf:type deo:Introduction ;
  po:isContainedBy openbiodiv:2b836ad5-db56-4093-9752-33c9f7892de6 .
}

```

Divide and conquer

After we have successfully converted the current XML node to RDF, a recursive call to Extractor is made for all nodes that are hierarchically dependent on the current node. For example, the article node contains all the other other nodes such as sections, figures, etc.

Transformation specification

In order for the Extractor to work, therefore, we need to specify an XML schema. The specification includes what XML nodes we are looking for and their location. It then recursively specifies for each node, what sub-nodes we are looking for and their XPATH location relative to their parent node. Finally, for every node we need to give the atom locations and write a constructor. The transformation specification is

done with R6 framework in R. We have specified two schemata that share the same constructors—TaxPub⁴ and TaxonX⁵.

3.5.4 Submission to graph database and post-processing

In the previous section we described how we transform XML documents in TaxPub and TaxonX to RDF statements according to OpenBiodiv-O. In addition, we transform the GBIF backbone taxonomy to RDF according to OpenBiodiv-O with the help of TSV4RDF, a proprietary Pensoft tool. The generated RDF statements are submitted to a repository in a GraphDB instance residing on <http://graph.openbiodiv.net/>. The repository has been initialized with OpenBiodiv-O and the ontologies on which it depends⁶. Finally, after the data has been submitted, update scripts are run to generate further statements from our ontology that have not been encoded in OWL for the updating of scientific name relations.

Update rule for replacement name

We state that a scientific name A replaces a scientific name B , if there exists a taxonomic name usage of A with taxonomic status `:ReplacementName` and B is mentioned by a taxonomic name usage in the nomenclatural citations of the treatment, where the discussed taxonomic name usage of A is in the nomenclature section (Listing 3.4).

LISTING 3.4: Update rule for replacement name.

```
PREFIX dwciri: <http://rs.tdwg.org/dwc/iri/>
PREFIX openbiodiv: <http://openbiodiv.net/>
PREFIX po: <http://www.essepuntato.it/2008/12/pattern#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX pkm: <http://proton.semanticweb.org/protonkm#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>

INSERT
{
  GRAPH <http://openbiodiv.net/Updates>
  {
    ?name2 openbiodiv:replacementName ?name .
  }
}

WHERE {
  ?tnu1 dwciri:taxonomicStatus openbiodiv:ReplacementName ;
    pkm:mentions ?name .
  ?name dwciri:taxonRank ?rank ;
    rdfs:label ?vname .

  ?s po:contains ?tnu .
  ?s po:contains ?citations .
  ?citations rdf:type openbiodiv:NomenclatureCitationsList ;
    po:contains ?tnu2 .
  ?tnu2 rdf:type openbiodiv:TaxonomicNameUsage ;
    pkm:mentions ?name2 .
  ?name2 rdfs:label ?vname2 ;
    dwciri:taxonRank ?rank .
}
```

Update rule for related name

The related names update-rule is similar to the replacement name: two scientific names A and B are considered related if they both mentioned in the nomenclature section of a treatment (Listing 3.5).

LISTING 3.5: Update rule for related name.

```
PREFIX dwciri: <http://rs.tdwg.org/dwc/iri/>
PREFIX : <http://openbiodiv.net/>
PREFIX po: <http://www.essepuntato.it/2008/12/pattern#>
```

⁴<https://github.com/pensoft/ropenbio/blob/redesign/R/taxpub.R>

⁵<https://github.com/pensoft/ropenbio/blob/redesign/R/taxonx.R>

⁶<https://github.com/vsenderov/openbiodiv-o/tree/master/imports>

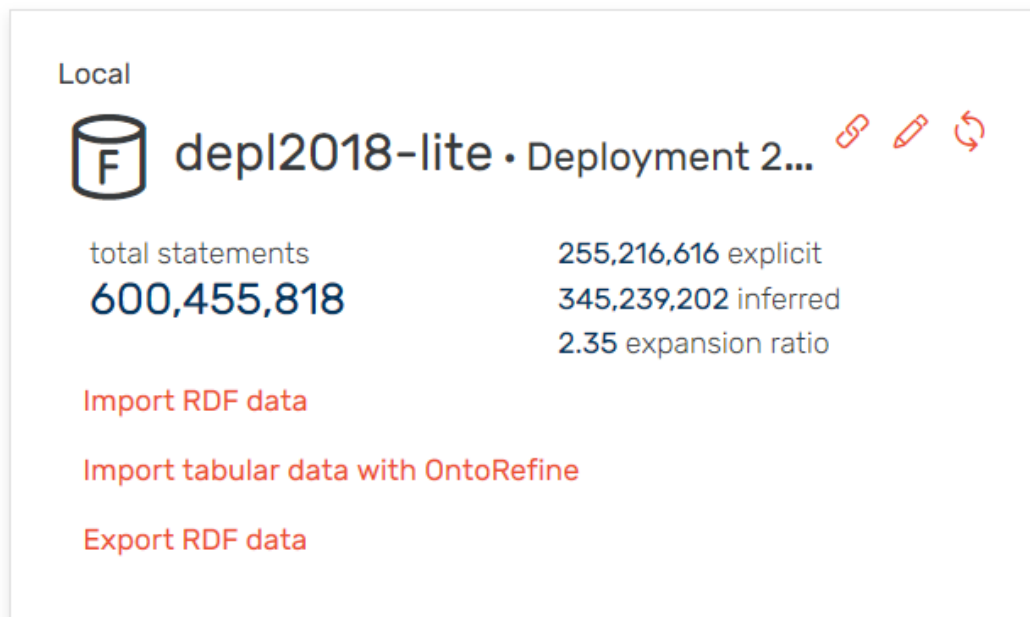


FIGURE 3.6: Statements report from the GraphDB workbench.

```

PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX pkm: <http://proton.semanticweb.org/protonkm#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>

INSERT
{
  GRAPH <http://openbiodiv.net/Updates>
  {
    ?name2 :relatedName ?name .
  }
}

WHERE {
  ?nom_sec rdf:type :NomenclatureSection ;
    :contains ?tnu1 .

  ?tnu1 rdf:type :TaxonomicNameUsage ;
    pkm:mentions ?name.

  ?nom_sec :contains ?tnu2 .

  ?tnu2 rdf:type :TaxonomicNameUsage ;
    pkm:mentions ?name2.

  FILTER(?name != ?name2)
}

```

3.6 Performance degradation analysis

The current iteration of the database holds over 600 million triples (Fig. 3.6). The expansion ratio under the RDFS-Plus (Optimized) ruleset is 2.35, i.e. for each asserted statements we materialize on average 2.35 implicit statements. Under the OWL2-RL ruleset (which contains a full implementation of OWL logic rules), the expansion ratio is about 3.7; however, we encountered significant performance issues using it (Fig. 3.7). Even with the lighter ruleset (RDFS-Plus Optimized), we still see performance degradation with increasing database size. Importing the GBIF backbone taxonomy from file takes about two days under the easier scenario. The subsequent importing of the Pensoft archives takes about two weeks as it is a slower operation requiring not only the time for submission but the time for converting the XML's to RDF.

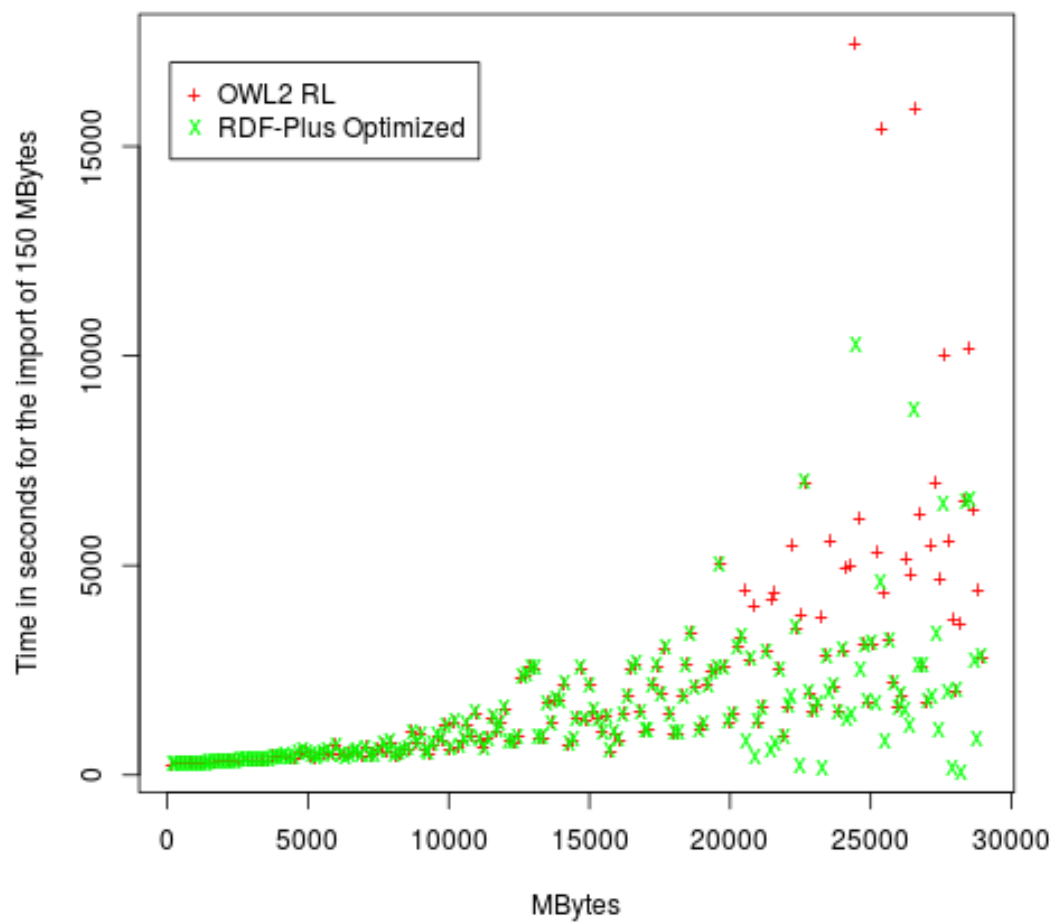
Performance degradation as a function of database size in MBytes

FIGURE 3.7: The graph visualizes the time in seconds needed to import a 150 MB big Turtle data file as a function of the database size. The database size is measured by the adding up the size of the data files that have already been imported.

Chapter 4

Summary of Chapter 4: An R Library for Working with RDF

RDF4R (`rdf4r`) is an R package for working with Resource Description Framework (RDF Working Group, 2014) data. It was developed as part of the OpenBiodiv project but is completely free of any OpenBiodiv-specific code and can be used for generic purposes requiring tools to work with RDF data in the R programming environment (R Core Team, 2016).

4.1 Installation

In this section we describe how to install the RDF4R package. Installation is straightforward and consists of two steps: (1) resolve dependencies and (2) build the package from source using `devtools::install_github`.

4.2 Specification

In this section we present the specifications of RDF4R by detailing the features of the package. Each feature has a dedicated subsection.

4.2.1 Connection to a triple-store

It is possible to establish both basic connections (requiring no password or requiring basic HTTP user-pass authentication) or connection secured with an API access token.

4.2.2 Work with repositories on a triple-store

Once a connection to a triple-store has been established, it is possible to inspect the talk protocol version, view the list of repositories on the database, execute SPARQL Read (SELECT keyword and related) and SPARQL Update (INSERT and related) queries on the database, as well as submit serialized RDF data directly to the database.

4.2.3 Function factories to convert SPARQL queries to R functions

An important feature of RDF4R are its facilities for converting SPARQL queries and the like to R functions.

4.2.4 Work with literals and identifiers

The building blocks of RDF are literals (e.g. strings, numbers, dates, etc.) and resource identifiers. RDF4R provides classes for literals and resource identifiers that are tightly integrated with the other facilities of the package.

4.2.5 Prefix management

Prefixes are managed automatically during serialization by being extracted from the resource identifiers.

4.2.6 Creation and serialization of RDF

The serialization function supports Turtle (and its variant Trig, Bizer and Cyganiak, 2014) and adding new triples.

LISTING 4.1: Using brackets to express RDF blank nodes in Turtle/TriG.

```
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
# :someone knows someone else, who has the name "Bob".
:someone foaf:knows [ foaf:name "Bob" ] .
```

4.2.7 A basic vocabulary of semantic elements

RDF4R has some basic resource identifiers for widely used classes and predicates predefined (e.g. for `rdf:type`, `rdfs:label`, etc.).

4.3 Usage

Here, we explain how to use the package RDF4R by means of examples. In order to fully utilize the package capabilities, one needs to have access to an RDF graph database. We have made available a public endpoint (see next paragraph) to allow the users of the package to experiment. Since write access is enabled, please be considerate and don't issue catastrophic commands.

4.4 Discussion

4.4.1 Related Packages

The closest match to RDF4R is the `rdflib` (Boettiger, 2018). The development of the two packages was simultaneous and independent until `rdflib`'s first official release on Dec 10, 2017. This explains why two closely related R packages for working with RDF exist. After the release of `rdflib` work was started to make both packages compatible with each other. In our opinion, the packages have different design philosophies and are thus complementary.

`rdflib` is a high-level wrapper to `redland` (Jones et al., 2016), which is a low-level wrapper to the C `librdf` (Beckett, 2014), a powerful C library that provides support for RDF. `librdf` provides an in-memory storage model for RDF beyond what is available in RDF4R and also persistent storage working with a number of databases. It enables the user to query RDF objects with SPARQL. Thus, `librdf` can be considered a complete graph database implementation in C.

In our opinion, `redland` is more complex than needed for the purposes of OpenBiodiv. By the onset of the OpenBiodiv project it was available¹; however, we decided not to use it as a decision was made to rely on GraphDB for our storage and querying. Note that RDF4R's main purpose is to provide a convenient R interface for users of GraphDB and similar RDF4J compatible graph databases.

¹But not `rdflib`!

A feature that differentiates `rdflib` from RDF4R is the design philosophy. RDF4R was designed primarily with the Turtle and TriG serializations in mind. This means that RDF4R can work with named graphs, whereas their usage is discouraged or perhaps impossible with `rdflib`², even though `rdflib`'s default format is N-Quads.

Another differentiating feature between RDF4R and `rdflib` is that RDF4R provides facilities for converting SPARQL and related statements to native R functions!

In a future release of RDF4R (2.0) we would like to replace or extend its in-memory model with `rdflib`'s. This is why we would like to make the packages fully compatible and have contributed several patches to `rdflib`³. Thus, it will be possible for the user of RDF4R to retain its syntax and high-level features— constructor factories, functors, etc., and the ability to use named graphs—but benefit from performance increases, stability, and scalability with the `redland/rdflib/librdf` backend.

This will enable the users of the R programming environment to use whichever syntax they prefer and benefit from an efficient storage engine.

4.4.2 Elements of Functional Programming (FP)

In this subsection we discuss how patterns from functional programming were used to create RDF4R.

4.4.3 Elements of Object-Oriented Programming (OOP)

In this subsection we discuss how patterns from object-oriented programming were used to create RDF4R.

²The issue was discussed on the `librdf` GitHub page, <https://github.com/ropensci/rdflib/issues/23>.

³Please, consult the commit history under <https://github.com/ropensci/rdflib>.

Chapter 5

Summary of Chapter 5: Workflows for Biodiversity Data

In this chapter we discuss two automated workflows for exchange of biodiversity data developed as part of OpenBiodiv: (1) automatic import of specimen records into manuscripts, and (2) automatic generation of data paper manuscripts from Ecological Metadata Language (EML) metadata. The workflows were presented at a webinar for the organization iDigBio¹ and published as a paper (Senderov et al., 2016).

The slides from the presentation as well as a PDF of the paper are available from the webinar GitHub page under <https://github.com/vsenderov/idigbio-webinar>.

5.1 Introduction

Information on occurrences of species and information on the specimens that are evidence for these occurrences (specimen records) is stored in different biodiversity databases. These databases expose the information via public REST API's. I focused on the Global Biodiversity Information Facility (GBIF), Barcode of Life Data Systems (BOLD), iDigBio, and PlutoF, and utilized their API's to import occurrence or specimen records directly into a manuscript edited in the ARPHA Writing Tool (AWT).

Furthermore, major ecological and biological databases around the world provide information about their datasets in the form of EML. A workflow was developed for creating data paper manuscripts in AWT from EML files. Such files could be downloaded, for example, from GBIF, DataONE, or the Long-Term Ecological Research Network (LTER Network).

The development of these workflows focuses on two areas: optimizing the workflow of specimen data and optimizing the workflow of dataset metadata. These efforts resulted in the functionality that it is now possible, via a record identifier, to directly import specimen record information from the Global Biodiversity Information Facility (GBIF), Barcode of Life Data Systems (BOLD), iDigBio, or PlutoF into manuscripts in the ARPHA Writing Tool (AWT). No manual copying or retyping is required.

¹Integrated Digitized Biocollections (iDigBio) is a US-based aggregator of biocollections data. They hold regular webinars and workshops aimed at improving biodiversity informatics knowledge, which are attended by collection managers, scientists, and IT personnel. Thus, doing a presentation for iDigBio was an excellent way of making the research and tools-development efforts of OpenBiodiv widely known and getting feedback from the community.

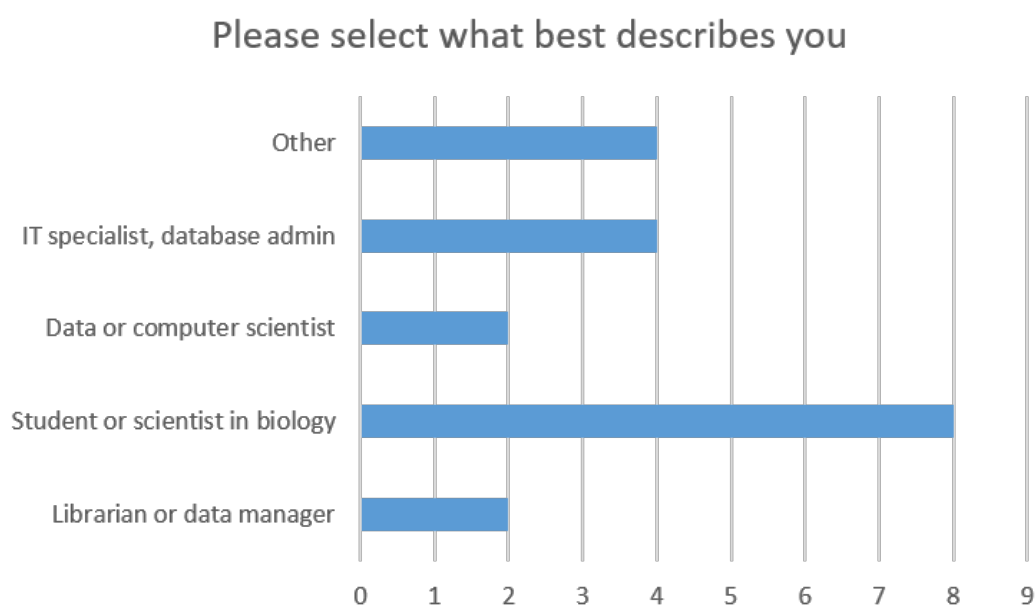


FIGURE 5.1: Poll results about composition of audience during live participation..

5.2 Presentation

A video recording of the presentation is available². More information can be found in the webinar information page³. The slides of the presentation are attached as supplementary files and are deposited in Slideshare⁴.

During the presentation we conducted a poll about the occupation of the attendees, the results of which are summarized in Fig. 5.1. Of the participants who voted, about a half were scientists, mostly biologists, while the remainder were distributed across IT specialists and librarians, with 20% "Other." The other categories might have been administrators, decision-makers, non-biology scientists, collections personnel, educators, etc.

At the end of the presentation, very interesting questions were raised and discussed. For details, see the "Results and discussion" section of this paper.

5.3 Methods

Both workflows discussed rely on three key standards: RESTful API's for the web (Kurtz, 2013), Darwin Core (Wieczorek et al., 2012), and EML (Fegraus et al., 2005).

5.3.1 Development of workflow 1: Automated specimen record import

In this subsection we discuss the development of Workflow 1: Automated specimen record import.

²<http://idigbio.adobeconnect.com/p7sg0aym3e3/>

³<http://www.idigbio.org/content/online-direct-import-specimen-records-idigbio-infrastructure-taxonomic-mar>

⁴<http://www.slideshare.net/ViktorSenderov/online-direct-import-of-specimen-records-from-idigbio-infrastru>

5.3.2 Development of workflow 2: Automated data paper generation

In this subsection we discuss the development of Workflow 1: Automated specimen record import.

5.4 Results and Discussion

5.4.1 Workflow 1: Automated specimen record import into manuscripts developed in the ARPHA Writing Tool

It is now possible to directly import a specimen record as a material citation in an ARPHA Taxonomic Paper from GBIF, BOLD, iDigBio, and PlutoF (Slide 5, as well as Fig. 5.2). The workflow from the user's perspective has been thoroughly described in a blog post; concise stepwise instructions are available via ARPHA's Tips and tricks guidelines. In a nutshell, the process works as follows:

1. At one of the supported data portals (BOLD, GBIF, iDigBio, PlutoF), the author locates the specimen record he/she wants to import into the Materials section of a Taxon treatment (available in the Taxonomic Paper manuscript template).
2. Depending on the portal, the user finds either the occurrence identifier of the specimen, or a database record identifier of the specimen record, and copies that into the respective upload field of the ARPHA system (Fig. 5.3).
3. After the user clicks on "Add," a progress bar is displayed, while the specimens are being uploaded as material citations.
4. The new material citations are rendered in both human- and machine-readable DwC format in the Materials section of the respective Taxon treatment and can be further edited in AWT, or downloaded from there as a CSV file.

Discussion

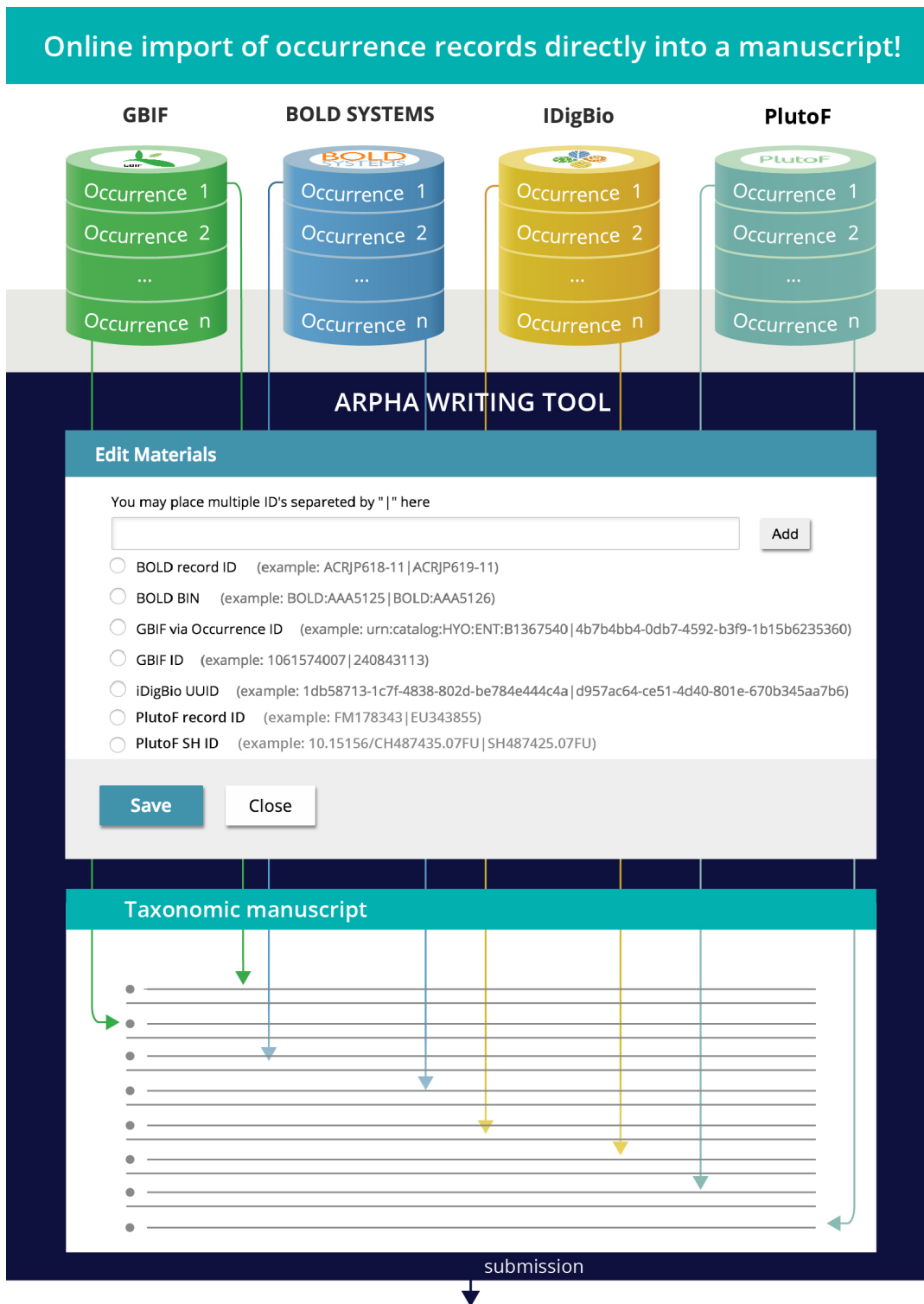
We discuss the availability, or more correctly the lack of persistent unique identifiers (PID's) in the biodiversity informatics space. I furthermore discuss the challenges of importing from our different sources: GBIF, PlutoF, iDigBio, and BOLD. I emphasize how our workflow can be serve as a curation filter for increasing the quality of specimen data via the scientific peer review process.

5.4.2 Workflow 2: Automated data paper manuscript generation from EML metadata in the ARPHA Writing Tool

We have created a workflow that allows authors to automatically create data paper manuscripts from the metadata stored in EML (Fig. 5.4, Fig. 5.5, Fig. 5.6).

Discussion

I discuss the history of data papers and how our implementation greatly improves the availability of data papers to science practitioners. The two workflows presented generated a lively discussion at the end of the presentation, which is summarized in the Chapter.



**Biodiversity
Data Journal**

<http://bdj.pensoft.net>

FIGURE 5.2: This fictionalized workflow presents the flow of information content of biodiversity specimens or biodiversity occurrences from the data portals GBIF, BOLD Systems, iDigBio, and PlutoF, through user-interface elements in AWT to textualized content in a Taxonomic Paper manuscript template intended for publication in the

You may place multiple ID's separated by "|" here

- BOLD record ID (example: ACRJP618-11|ACRJP619-11)
- BOLD BIN (example: BOLD:AAA5125|BOLD:AAA5126)
- GBIF via Occurrence ID (example: urn:catalog:HYO:ENT:B1367540|4b7b4bb4-0db7-4592-b3f9-1b15b6235360)
- GBIF ID (example: 1061574007|240843113)
- iDigBio UUID (example: 1db58713-1c7f-4838-802d-be784e444c4a|d957ac64-ce51-4d40-801e-670b345aa7b6)
- PlutoF Specimen ID (example: AT2000123|TAM0000007)

FIGURE 5.3: User interface of the ARPHA Writing Tool controlling the import of specimen records from external databases.

PENSOFT IPT DATA HOSTING CENTER
free and open access to biodiversity data

Logged in as datascience@pensoft.net [Account](#) [Logout](#) [ENGLISH](#)

[Home](#) [Manage Resources](#) [About](#)

[Summary](#)
[Downloads](#)
[Versions](#)
[Rights](#)
[GBIF Registration](#)
[Keywords](#)
[Contacts](#)
[Geographic Coverage](#)
[Taxonomic Coverage](#)
[Temporal Coverage](#)

[Edit](#)

A checklist to the wasps of Peru (Hymenoptera, Aculeata)
Latest version published by ZooKeys on Feb 17, 2011

The first checklist to the 225 genera and 1169 reported species-group taxa of aculeate wasps of Peru is presented. The list is based on a literature survey and examination of Peruvian entomological collections and include locality references for each taxon. Bibliographic references for the identification of families, genera, and species are provided when available. The occurrence data are published in addition as a downloadable file (doi: 10.3897/zookeys.15.196.app.2.ds, doi: 10.3897/zookeys.15.196.app.3.ds, and 10.3897/zookeys.15.196.app.4.ds) and were uploaded onto GBIF infrastructure simultaneously with the publication process. The following new combinations are proposed: *Ancistroceroides cirrifer* (Zavattari, 1912), *Ancistrocerus epicus* (Zavattari, 1912), and *Stenodynerus corallineipes* (Zavattari, 1912).

[GBIF](#) [DwC-A](#) [EML](#) [RTF](#) [Versions](#) [Rights](#)

FIGURE 5.4: Download of an EML from the GBIF Integrated Publishing Toolkit (IPT).

Biodiversity Data Journal
 Research Ideas and Outcomes

One Ecosystem
 BioDiscovery

Article type

Research ideas	Grant proposals	Brief research outcomes	Early research outcomes
<input type="radio"/> Data Management Plan (Biosciences) <input type="radio"/> Data Management Plan (Generic) <input type="radio"/> Data Management Plan (NSF Generic) <input type="radio"/> PhD Project Plan <input type="radio"/> PhD Project Plan (Free Text) <input type="radio"/> PostDoc Project Plan <input type="radio"/> PostDoc Project Plan (Free Text) <input type="radio"/> Research Idea <input type="radio"/> Small Grant Proposal <input type="radio"/> Small Grant Proposal (Free Text) <input type="radio"/> Software Management Plan	<input type="radio"/> DFG Grant Proposal <input type="radio"/> FP7 Grant Proposal <input type="radio"/> Grant Proposal <input type="radio"/> Grant Proposal (Free Text) <input type="radio"/> H2020 Grant Proposal <input type="radio"/> NIH Grant Proposal <input type="radio"/> NSF Grant Proposal	<input type="radio"/> Conference Abstract <input type="radio"/> Correspondence <input type="radio"/> Ecosystem Inventory <input type="radio"/> Ecosystem Service Mapping <input type="radio"/> Ecosystem Service Models <input type="radio"/> Monitoring Schema <input type="radio"/> Research Poster <input type="radio"/> Research Presentation <input type="radio"/> Single-media Publication	<input type="radio"/> Case Study <input type="radio"/> Case Study (Free Text) <input checked="" type="radio"/> Data Paper (Biosciences) ? <input type="radio"/> Data Paper (Generic) <input checked="" type="radio"/> Forum Paper (Free Text) ? <input type="radio"/> Methods ? <input checked="" type="radio"/> Methods (Free Text) ? <input type="radio"/> Opinion Article <input checked="" type="radio"/> Opinion Article (Free Text) ? <input type="radio"/> Project Report <input type="radio"/> Project Report (Free Text) <input type="radio"/> Questionnaire <input checked="" type="radio"/> Software Description ? <input type="radio"/> Workshop Report
Research outcomes	PhD theses	Editorial matters	
<input checked="" type="radio"/> Alien Species Profile ? <input type="radio"/> Guidelines (Free Text) <input checked="" type="radio"/> Interactive Key ? <input type="radio"/> Policy Brief <input type="radio"/> Policy Brief (Free Text) <input type="radio"/> Replication Study <input checked="" type="radio"/> Research Article ? <input checked="" type="radio"/> Research Article (Free Text) ? <input type="radio"/> Review Article <input type="radio"/> Review Article (Free Text) <input checked="" type="radio"/> Single Taxon Treatment ? <input type="radio"/> Species Conservation Profile <input type="radio"/> Taxonomic Paper <input type="radio"/> Wikipedia Article	<input type="radio"/> PhD Thesis <input type="radio"/> PhD Thesis (Free Text)	<input type="radio"/> Biography <input type="radio"/> Book Review <input type="radio"/> Corrigendum <input type="radio"/> Data Review <input checked="" type="radio"/> Editorial ? <input type="radio"/> Obituary <input type="radio"/> Software Review	

OR

FIGURE 5.5: Selection of the journal and “Data Paper (Biosciences)” template in the ARPHA Writing Tool.

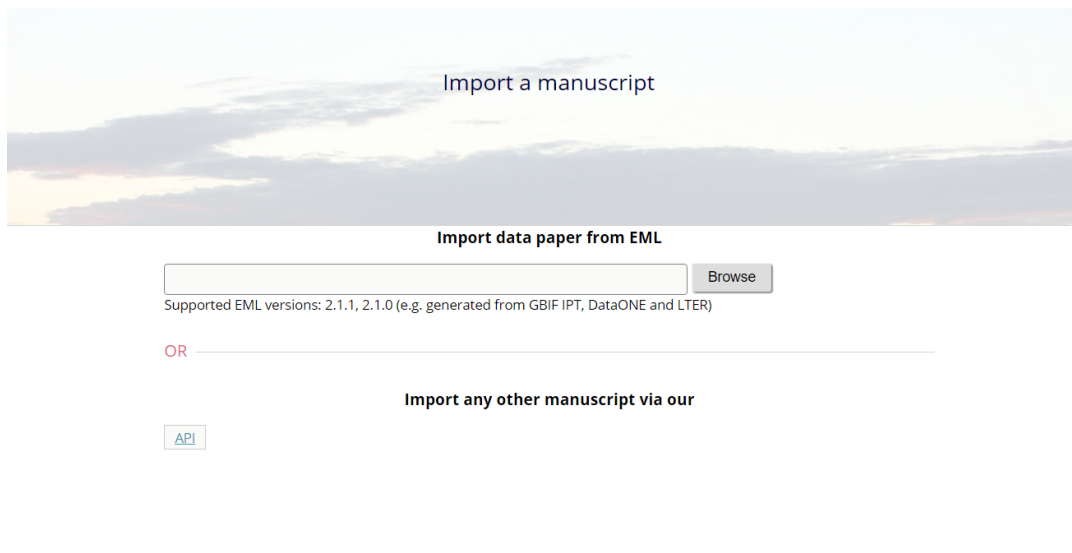


FIGURE 5.6: The user interface field for uploading EML files into ARPHA.

Chapter 6

Summary of Chapter 6: Web portal

Under openbiodiv.net one can reach the main portal giving access to OpenBiodiv resources. This portal was developed by Pensoft to support OpenBiodiv. OpenBiodiv.net presents two visual elements to the user: the search bar and list of application icons in the bottom. Furthermore, under graph.openbiodiv.net (also accessible from the icon SPARQL endpoint) one can reach the OpenBiodiv workbench, a feature of GraphDB that gives web access to the SPARQL endpoint.

These User Interface (UI) features are designed to facilitate the three user types of the system that we envisage:

1. Basic level: uses search bar.
2. Specialist level: uses apps.
3. Power user: uses the work-bench of the system or R.

6.1 Functionality of the system

In this section we discuss how every user-type can use the system.

6.1.1 Basic usage

The basic level of interaction is for users who want a quick look into the system's database; they can be beginners without knowledge of the Semantic Web or of taxonomy, or advanced users with little time or a very basic query. An example of such a user will simply look for an entity (e.g. taxonomic name, person) and would like to retrieve some information about it.

6.1.2 Specialist level

A specialist is someone who has a question of particular taxonomic importance that cannot be answered by a simple name-based look-up. For example, a collection manager at a museum may want to periodically check for articles that make use of their collection in order to justify additional funding to prevent natural disasters. Or a taxonomist interested in a particular region or group may want to stay up to date with published literature fitting those criteria—let's say weevils (Curculionidae) of Arizona, U.S.A.


 Beta
 The Open Biodiversity Knowledge Management System

daniel mietchen Q

Person

ID: <http://openbiodiv.net/2c525459-67b1-427b-9420-844b70c41f03>

Person info

Name: Daniel Mietchen

Affiliations:

- Museum für Naturkunde, Berlin, Germany

Articles



show all articles

Collaborators

Guido Sautter 1	Kevin Richards 1	Pavel Stoev 1	Aleksandra Pawlik 1	Bachir Balech 1	David Eades 1
Donat Agosti 1	Walter G. Berendsohn 1	Adam Brunke 1	Hannes Hettling 1	Robert Hoehndorf 1	Rod Page 1
Nesrine Akkari 1	Tom van Dooren 1	Marko Tsihtinen 1	Quentin John Groom 1	Andreas Plank 1	
Thomas D. Hamann 1	Alan R. Williams 1	Soraya Sierra 1	Thomas Pape 1	Robert A. Morris 1	Gregor Hagedorn 1
Lyubomir Penev 1	Claus Weiland 1	Patricia Kelbert 1	Nicola Nicolson 1	Ayco Holleman 1	Donald Hobern 1
Don Kirkup 1	Teodor Georgiev 1	David King 1	Yuri Lammers 1	Niall Beard 1	Carina Mara de Souza 1
Matthew Blissett 1	Peter Hovenkamp 1	George Gosline 1	Thibaut DeMeulemeester 1	Jeremy A. Miller 1	
Terry Erwin 1	Christian Breninkmeijer 1	Lars Hendrich 1	David Peter Shorthouse 1	Ross Mounce 1	
Henrik Enghoff 1	David Koon-Bong Cheung 1	Michael Balke 1	Serrano Pereira 1		

FIGURE 6.1: Illustration of basic usage of OpenBiodiv to look information about a person.

6.1.3 Power user

The power user is someone with knowledge of the Semantic Web and its technologies (SPARQL, ontologies, etc.). The power user goes to the workbench and executes their queries there, or uses the functionality of the RDF4R package described in Chapter 4 to execute SPARQL directly on the OpenBiodiv endpoint directly from the R environment.

6.2 Implementation

The UI-components of the web portal are developed in the ReactJS JavaScript framework written by Facebook. Server-side processing is done in PHP. This part of OpenBiodiv is not open source and cannot be discussed in detail in the present dissertation effort.

Conclusion

Results

We believe that the presented scientific work fulfills the stated objective and tasks.

Result 1. The central result of the thesis is the creation of a domain conceptualization of biodiversity publishing and a formal ontology OpenBiodiv-O enabling the linking of biodiversity knowledge on the basis of scholarly publications. This result has been described in Chapter 2 and in Senderov et al., 2018 and fulfills Objective 1. The source code of the ontology is available under github.com/pensoft/openbiodiv-o.

Result 2. The second result of the thesis is the creation of the software architecture of the OpenBiodiv system outlined in Chapter 1 and Senderov and Penev, 2016. This result fulfills Objective 2.

Result 3. The third result of the thesis has been the creation of a Linked Open Dataset, OpenBiodiv-LOD, consisting of a transformation to RDF-triples and integration in a single store of information from three major repositories of biodiversity data: the XML sources of biological journals published by Pensoft Publishers, the XML sources of treatments freed by Plazi, and a CSV dump of GBIF's taxonomic backbone. OpenBiodiv-LOD is available under graph.openbiodiv.net and has been described in Chapter 3. This result fulfills Objective 3.

Result 4. In order to create the Linked Open Data, a software package for the R programming environment, RDF4R, was developed. RDF4R enables the manipulation of RDF data within R and facilitates the transformation of scientific publications from a semi-structured XML format to structured semantic RDF. This result has been discussed in Chapter 4 and fulfills Objective 4. The package is available online as free software under github.com/pensoft/rdf4r. Furthermore, additional source code (unoptimized) describing XML schemata of Pensoft and Plazi and working in tandem with RDF4R to convert XML to RDF can be found under github.com/pensoft/ropenbio.

Result 5. The mechanisms to convert semi-structured XML into RDF-triples are complemented by workflows enabling the enrichment of the XML sources of Pensoft journals by data automatically imported from the major international biodiversity data repositories: BOLD, GBIF, iDigBio, as well as PlutoF. Furthermore, it is now possible, thanks to this dissertation effort to automatically create manuscripts from metadata encoded in the Ecological Metadata Language (EML). The discussion of these automated workflows—automatic data paper generation and automatic occurrence record import—is carried out in Chapter 5. It fulfills Objective 5.

Result 6. To complement the creation of OpenBiodiv-LOD, we have developed a website running on top of the knowledge graph openbiodiv.net, containing a semantic search engine and apps. The website is discussed in Chapter 6 and fulfills Objective 6.

Discussion, conclusion, and outlook

OpenBiodiv-O serves as the basis of the Linked Open Data OpenBiodiv-LOD. By developing an ontology focusing on biological taxonomy, we provided an ontology that fills in the gaps between ontologies for biodiversity resources such as Darwin-SW and semantic publishing ontologies such as the ontologies comprising the SPAR Ontologies. Moreover, we take the view that it is advantageous to model the taxonomic process itself rather than any particular state of knowledge. At this stage, the coverage of the ontology of the different types of resources is sufficient to be the basis for creating the LOD. In this sense, it is completed. On the other hand, adding classes and properties for new types of biodiversity data is possible and desirable.

The LOD, similar to the ontology, are already a solid resource for biologists, as they include information from most articles published by Pensoft and Plazi and count over 600 million triplets. Like the ontology, they should be expanded.

Since the RDF4R package was successfully used to create an LOD, it can be considered complete. Like any software package, however, it should be maintained and developed.

The website is still in beta. The functionality that works great is the semantic search engine. For some basic data types there are templates for visualization. However, the site can not be considered complete and most users use the SPARQL search language.

An important conclusion that can be drawn from the work is that it is possible to use a semantic graph for the integration of a large volume of data on biodiversity. We were unexpectedly given the opportunity to illustrate the power of the knowledge graph by analyzing the damage from the tragic fire at the Museu Nacional in Rio de Janeiro. In addition, we have illustrated that it is possible to write relatively simple logical conclusions to check the validity of a taxonomic name.

Due to the large amount of data, we found that although the use of a semantic graph was possible, some of the initially chosen technologies proved to be inapplicable or difficult to apply. We have observed (see Chapter 3) that the practical application of the full logical OWL model is difficult due to performance problems. Instead in the end, we utilized RDFS that is less powerful but faster. Another observation of ours is that although the R programming environment has given us some advantages in rapidly creating the prototype of the system, by increasing the complexity of the program code needed in the real-life system to cover all private cases, a language with dynamic types such as R creates headaches in debugging. At the same time, we were impressed by the powerful functional programming toolkit R provided.

A big difficulty was the disambiguation of resources such as author names or taxonomic names. In the functional design of the RDF4R package we have put modules that allow us to insert a list of functions/rules for disambiguation when searching for an identifier for a given resource. However, we had only limited success with the rule-based disambiguation and for this reason in the production system it was discontinued at the moment.

Considering these and other “lessons,” the future development of the OpenBiodiv project can be outlined in the following not necessarily comprehensive way:

1. As an immediate goal, to expand the LOD and ontology with new data types and new data sources using the existing framework. Such data are e.g. genomic data, occurrence data, (bio-)geographic data, visual data, descriptive data, etc.
2. Look for even closer integration with other existing biodiversity data repositories than GBIF. For example, BioImages, iNaturalist, BOLD, and so on.
3. As a longer-term task to study the transition from a semantic graph to a technology where the inference engine is separated from the data base layer as WikiData or Neo4j. In addition to increased performance, this will give extra flexibility to the project, such as allowing the use of non-RDF-based inference engines such as Euler.
4. Continue developing system software with an even wider application of functional programming and porting it into a functional language like, for example, Haskell or O'CAML.
5. To investigate the problem of disambiguation and related problems for named entity recognition of interesting resources from biodiversity, as well various image recognition tasks, from the point of view of machine learning.
6. Expanding the website with more templates and new applications.

Key scientific and applied contributions

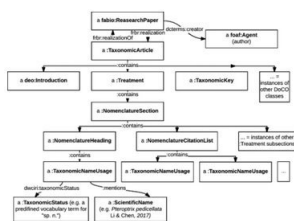
The results discussed in the previous two sections determine the following scientific and applied contributions:

1. Scientific contribution: creating an ontology and a formal model of the field of biodiversity knowledge publication.
2. Applied scientific contribution: analyzing information sources and Creating OpenBiodiv-LOD.
3. Applied scientific Contribution: the implementation of OpenBiodiv software modules.

Our ontology fills the unique niche between bibliographic ontologies such as SPAR and ontologies for biodiversity such as Darwin-SW and as such is undoubtedly of great scientific interest to the biodiversity informatics community. The work has a serious scientific and applied character by providing both a Linked Open Dataset on top of the ontology and software for its users and system developers.

Evaluation of publications

Articles have been published without exception in four international scientific journals: five articles in *Research Ideas and Outcomes*, one article in *ZooKeys* (WoS IF 1.079, Q3 SCOPUS, SJR 0.533), one article in *Biodiversity Data Journal* (WoS SCOPUS, SJR 0.465) and one article in *Journal of Biomedical Semantics* (WoS IF 1.6, Q3 SCOPUS, SJR 0.952). The total number of citations that have been accumulated for the candidate excluding self-citations (cross-citations) is at least 20. The citing articles are given in the list above. The total number of citations that have been



Featured Article: OpenBiodiv-O: ontology of the OpenBiodiv knowledge management system

Extracting information and knowledge discovery from research in the fields of biodiversity and biological taxonomy is a semantic challenge with important applications. In this article, the authors introduce OpenBiodiv-O, an ontology that bridges the gap between scholarly biodiversity publishing and biological taxonomy - enabling the creation of the OpenBiodiv Knowledge Management System.

FIGURE 6.2: The OpenBiodiv-O article is featured on the main web-page of the Journal of Biomedical Semantics..

accumulated including cross-citations and citations of work outside of the scope of the dissertation is at least 48 (Google Scholar).

[1] is an early version of the Introduction as well Chapter 1 and contains work towards Objective 2 (Architecture). The text of publications [2, 3, 5, 6, 7] are not a part of the text of the dissertation one-to-one but contain work towards Objective 5 (Workflows). The ideas presented in these publications have to large degree been incorporated in Chapter 5 whose backbone is formed by [4]; thus Objective 5 (Workflows) is achieved. [7] is published in the peer-reviewed journal ZooKeys with impact factor 1.031 (early 2018). [8] is the most important publication under this dissertation and was published in the high-impact Journal of Biomedical Semantics with impact factor 2.413 (early 2018). [8] makes up the content of Chapter 2 and is the main body of work fulfilling Objective 1 (Ontology). It was a featured article on the home-page of JBS (Fig. 6.2). Chapter 3 and Chapter 4 that form Objectives 3, 4, respectively are currently being prepared as manuscripts in international journals. Furthermore, the software library RDF4R described in Chapter 4 is being submitted to the open source repository rOpenSci¹.

¹“We build software with a community of users and developers, and educate scientists about transparent research practices.” <https://ropensci.org/>

Acknowledgements

This research has been financed through the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 642241. My deep gratitude goes to the European Commission for enabling this wonderful opportunity!

I thank Prof. Lyubomir Penev and Prof. Kiril Simov for the valuable supervision. I also thank the staff and developers at Pensoft Publishers for the support in creating the platform and its popularization; in particular Prof. Pavel Stoev, Teodor Georgiev, Georgi Zhelezov, Iliyana Kuzmova, and Iva Kostadinova. Furthermore, I thank Pensoft's graphic designer, Slavena Peneva, for the help with creating the illustrations for this thesis and in presentations. Last but not least, Margarita Grudova and Elisaveta Taseva for providing valuable administrative support during the elaboration of the thesis.

I thank my colleagues from the Bulgarian Academy of Sciences (Institutes for Information and Communication Technologies and for Biodiversity and Ecosystems Research) for their friendship and advice; in particular Prof. Galya Angelova, Prof. Boyko Georgiev, and Prof. Snejana Grozeva.

I thank my colleagues at the BIG4 training network for the feedback, friendship, and support. In particular Prof. Alexey Solovnikov, but there are too many more names to mention.

I thank my international collaborators for their ideas, reviews, and collaboration on papers. In particular Prof. Nico Franz (Arizona State University), Dr. Daniel Mietchen (National Institutes of Health), Dr. Éamonn Ó Tuama (formerly at GBIF), and Prof. Bob Morris (emeritus UMASS).

I also thank everyone at Plazi for the co-ownership of the vision of the project; in particular, Dr. Donat Agosti, Terry Catapano, and Dr. Guido Sautter.

Last but not least, I would like to acknowledge Ontotext for building the GraphDB database and providing excellent support.



АВТОРЕФЕРАТ НА ДИСЕРТАЦИЯ

за присъждане на образователна и научна степен “доктор” по
научна специалност “Информатика“

**OpenBiodiv: отворена система за управление на
знанието за биологичното разнообразие**

Виктор Синдеров

Ръководител: Проф. Любомир Пенев
Научен консултант: Доц. Кирил Симов

Научно жури:

Проф. Калинка Калоянова
Проф. Георги Марков
Проф. Мария Нишева-Павлова
Доц. Светла Бойчева
Доц. Геннадий Агре



**Институт по информационни и
комуникационни технологии**

**Секция „Лингвистично моделиране и обработка
на знания“**

Това изследване е осъществено благодарение на финансовата подкрепа на Е.С. в рамките на изследователски грант No. 642241 по проект BIG4 на програмата Мария-Склодовоска Кюри, част от “Хоризонт 2020”.

Увод

Актуалност и значимост на темата

Необходимостта от създаване на интегрирана информационна система, обхващаща знание за биологично разнообразие, датира от 1985 г. Тогава е учредена Работна група по таксономични бази данни (TDWG), която в последствие е преименувана на Стандарти за информационни технологии за биологичното разнообразие, но запазва съкращението TDWG¹. През 1999 г. е създаден Глобален институт за информация за биологичното разнообразие (GBIF), след като Организацията за икономическо сътрудничество и развитие (OECD) стига до извода, че е необходим международен механизъм за достъп до данни и информация за биоразнообразието в световен мащаб (*What is GBIF?*). Декларацията Bouchout (*Bouchout Declaration 2014*) ознаменува резултатите от финансираните от Европейския Съюз проект *pro-iBiosphere*, който продължава от 2012 г. до 2014 г. и е посветен на задачата за създаване на интегрирана информационна система за биологичното разнообразие. Декларацията “Бухаут” призовава към свободното предоставяне на научни знания за биологичното разнообразие като Linked Open Data (LOD). Успореден процес в САЩ започва още по-рано със създаването на Архитектура за глобални имена, GNA (Patterson et al., 2010; Pyle, 2016).

През 2014 г., в рамките на европейската програма “Хоризонт 2020”, между академични институции и частни фирми е създаден консорциумът BIG4. BIG4 цели да допринесе за напредъка на науката за биологичното разнообразие. Мисията на проекта е: “BIG4 - Биосистематика, информатика и генетика на големите 4 групи насекоми: обучение на учените и предприемачите от утрешния ден” (University of Copenhagen et al., 2014). Важен член на консорциума е академичното издателство и софтуерна компания “Пенсофт”. Пенсофт публикува няколко десетки добре познати таксономични списания с отворен достъп². Като подписваща страна на декларацията “Бухаут”, Пенсофт е естествен кандидат за изграждането на Отворена система за управление на знанията за биоразнообразието (ОВКМС).

Поради тези причини настоящата дисертация за придобиване на докторска степен е базирана в академичното издателство “Пенсофт” и в Институт по информационни и комуникационни технологии на Българска академия на науките (ИИКТ-БАН).

Преглед на литературата по проблема

Поради интердисциплинарния характер на дисертацията, в рамките на литературната справка ще се прегледат две области: (а) “бази от знания и свързани отворени данни”, а също така и (б) “публикуване на знания и данни в сферата на биоразнообразието”.

¹Историята на TDWG стартира през 1985 г. и може да се види на <http://old.tdwg.org/past-meetings/>. Много от препратките на тази уеб страница за съжаление са невалидни и страницата се нуждае от известна поддръжка.

²Напр. ZooKeys, PhytoKeys, MycoKeys, Biodiversity Data Journal (BDJ).

Бази от знания и свързани отворени данни

Въвеждат се на термините “база от знания (knowledge base)” и “система, използваща знание (knowledge-based system)”. В дисертацията двата термина се използват взаимозаменяемо. Понякога се изписва по-дългият вариант – система, използваща знания/ система от знания – когато се подчертават аспекти на базата от знания, които не са свързани с технологичните аспекти на базата данни – напр. потребителски интерфейс, източници на информация и т.н. Понятието се дефинира посредством изрични дефиниции, а също и чрез разглеждане на реализации на няколко системи от знания на практика. Терминът е широко обсъждан още през осемдесетте години на 20 век (Jarke et al., 1989) и началото на деведесетте години (Harris et al., 1993). Значението, вложено му тогава, е “използване на идеи, както от системите за управление на бази данни (DBMS), така и от изкуствения интелект (ИИ) за създаване на тип компютърна система, наречена *knowledge management system* (KBMS)”. Harris et al., 1993 пише, че характеристиките на KBMS са, че тя съдържа “предписани правила и факти, от които полезни изводи могат да бъдат извлечени от машина за правене на извод (inference engine)”. Трябва да се отбележи, че това тълкуване идва от времето на първото поколение ИИ (rule-based systems). В по-ново време бе постигнат напредък във включването на статистически техники в базите данни (Mansinghka et al., 2015). Независимо от това, в този проект се работи с класическото разбиране за KBMS, основано на логически правила. С други думи под *база от знания* се разбира “подходяща база данни, тясно интегрирана с логически слой, който позволява правенето на логически изводи и придава семантика на данните”.

Друго относително по-ново развитие е приложението на принципите на свързаните отворени данни (Heath and Bizer, 2011). Повечето съществуващи бази от знания наблягат на социални аспекти, които правят данните взаимосвързани и годни за повторно използване. Примери са Freebase (Bollacker et al., 2008), която неотдавна беше включена в WikiData (Vrandečić and Krötzsch, 2014; Pellissier Tanon et al., 2016), DBPedia (Auer et al., 2007), както и Wolfram|Alpha (*Wolfram|Alpha, Making the world's knowledge computable*), а също така и Google Knowledge Graph (Singhal, 2012). Общото между тези системи е, че акцентът се поставя не само върху логическия слой, който позволява изводи, а и върху единното информационно пространство: тези системи действат като интегрират информация от множество източници и следват принципите на свързаните отворени данни, Linked Open Data (LOD). Семантичната мрежа (Semantic Web, Berners-Lee et al., 2001) предлага редица препоръки за LOD (Heath and Bizer, 2011), които, когато се прилагат правилно, гарантират, че публикуваните в Мрежата данни са годни за използване от всички потребители на Мрежата. Принципите на свързаните данни и тяхното приложение към OpenBiodiv се разглеждат в глава 3.

Водени от тези тенденции, модерните бази от знания наблягат повече върху свързването на данни, отколкото върху разработването на сложни механизми за логически изводи. Критика на идеята за интегриране на логически слой в базата данни се състои в наблюдението, че интегрирането води до повишена сложност (Barrasa, 2017). Критиката може да бъде обобщена с две точки. Първо, поставянето на логика близо до данните (особено, когато е прекомерно мощна за задачата) може да доведе до драстично намаляване на производителността³. Второ, нови техники (например машинно обучение) правят съществуващия логически слой излишен.

³Ефективността на логическия слой OWL се сравнява с по-слабата RDF схема (RDFS) в глава 3.

От наша гледна точка, данните са обектът, който е много по-ценен, докато стратегията за вадене на изводи (дали е логически слой, основан на правила, или техника за статистическо машинно обучение) може да бъде подменяна с напредването на компютърните науки. Тези идеи водят до интересна главоблъсканица в избора на технология за база данни, разгледана в следващите раздели.

И накрая, всяка система, базирана на знания, трябва да включва задължително и компоненти на потребителския интерфейс и/или интерфейс за програмиране на приложения (API), както и приложения (apps). Те служат като точка на контакт между човека и машината и са от решаващо значение за успеха на всяка такава система.

Публикуване на знания данни за биоразнообразието

В биомедицинската област отдавна се работи по извличането на информация и откриването на знания от първична литература (напр. Rebholz-Schuhmann et al., 2005; Momtchev et al., 2009; Williams et al., 2012). Областта на биоразнообразието, и по-специално биологичната систематика и таксономия (от тук нататък в дисертацията съкратено наричана *таксономия*), също се движи в посока към семантизация (напр. Agosti, 2006; Patterson et al., 2006; Kennedy et al., 2005; Penev et al., 2010a; Tzitzikas et al., 2013). Академичната издателска дейност е моделирана чрез Онтологичните издателски и референтни онтологии, SPAR Ontologies (Peroni, 2014). Онтологията на SPAR са колекция от онтологии, включващи, наред с другото, библиографската онтология FaBiO (Peroni and Shotton, 2012) и DoCO, онтология за компонентите на даден документ (Constantin et al., 2016). Онтологията на SPAR осигуряват класове и свойства за описанието на академични статии с общо предназначение. Таксономичните статии и техните компоненти, от друга страна, са моделирани чрез TaxPub XML Document Type Definition (DTD) и неформално се нарича “XML схема” (Catarano, 2010). TaxPub е XML схемата на няколко важни таксономични списания (напр. ZooKeys, PhytoKeys, Biodiversity Data Journal) и служи като концептуален шаблон за онтология OpenBiodiv-O (глава 2), разработена в хода на дисертацията.

Таксономичната номенклатура е дисциплина с много дълга традиция. Тя се трансформира в модерната си форма с публикуването на линеевата система (Linnaeus, 1758). Вече към началото на миналия век са били използвани стотици таксономични термини (Witteveen, 2015). Понастоящем именуването на групи от организми се регулира от Международния кодекс за зоологична номенклатура (ICZN, International Commission on Zoological Nomenclature, 2017) и Международния кодекс за номенклатура на водорасли, гъби и растения, известен като Кодекс на Мелбърн (*International code of nomenclature for algae, fungi and plants (Melbourne code) 2012*). Поради тяхната сложност (напр. ICZN има 18 глави и 3 приложения), се оказва непреодолимо до този момент предизвикателство да се създаде непротиворечива онтология базирана на кодексите за биологична номенклатура. Опити за преодоляването на това предизвикателство са сравнително пълната онтология на NOMEN (Dmitriev and Yoder, 2017) и не до там завършените Термини за статута на таксономичните номенклатури, TNSS.

Съществуват няколко проекта, целящи моделирането на по-широката област на биологичното разнообразие. Darwin-SW (Baskauf and Webb, 2016) адаптира предишните термини “DarwinCore” (Wieczorek et al., 2012) като Resource Description Framework (RDF). Двата модела се занимават основно с данни за наличие на организми. Моделирането и формализирането на строго таксономичния домейн е обсъдено от Verendsohn, 1995 и по-късно в напр. Franz and Peet, 2009; Sterner

and Franz, 2017. За отбелязване е и XML-базираната схема за прехвърляне на таксономични концепции (Taxonomic Names and Concepts Interest Group, 2006) и вече невалидната Taxon Concept Ontology. Наскоро в общността на TDWG се поднови интересът към таксономичните концепции, чрез създаването на Група за имена и таксономични концепции. Груповите дискусии могат да бъдат достъпни под <https://github.com/tdwg/tnc>. Интересното е, че в първата дискусия в GitHub се обсъди OpenBiodiv-O и възможността за приемането му като стандарт на TDWG.

През юни 2015 г. се стартира проектът OpenBiodiv. До старта на проекта са публикувани редица важни статии по темите за свързване на данни и споделяне на идентификатори в областта биологичното разнообразие (Page, 2008), относно обединяването на филогенетични знания (Page et al., 2012), и по темата на таксономичните имена и връзката им със семантичната мрежа (Page, 2006; Patterson et al., 2010), както и относно обобщаването на изследванията на биологичното разнообразие (Mindell et al., 2011). Дискусията относно OVKMS се намира в научния блог iPhylo (Page, 2014, 2015). Правните аспекти на OVKMS са обсъдени от Egloff et al., 2014. Освен това няколко системи за интегриране на данни за биоразнообразието са разработени от различни групи. Някои от най-важните са UBio, Глобални имена, BioGuid, BioNames, Профил на таксони на Пенсофт и Плаци⁴.

Основни изводи от литературната справка

Основните изводи от разгледаните източници са обобщени както следва:

1. Биоразнообразието се занимава с различни видове данни: таксономични, биогеографски, филогенетични, визуални, описателни и други. Тези данни са трудно достъпни в не взаимно-свързани хранилища за данни.
2. Базите данни за биологичното разнообразие се нуждаят от универсална система за именуване на таксономични концепции поради недостатъците на линеевите имена за модерната таксономия. Етикетите за таксономични концепции са предложени като решение, разбираемо от човека. А глобално-стабилни уникални идентификатори (GUID) на таксономични концепции – като решение, удобно за обработка в машинен вид.
3. Налице е основа от цифровизирана полу-структурирана информация за биологичното разнообразие в Мрежата с подходящи лицензи, чакащи да бъдат интегрирани като база от знания.

Цел и задачи

Предвид огромния международен научен интерес към отворената система за управления на знания и данни за биоразнообразието, тази дисертация стартира проект OpenBiodiv, който да допринесе за създаването на системата, като се концентрира върху знание, извлечено от научната литература.

Целта на проекта е създаването на формален семантичен модел на областта на публикуването на знание за биологичното разнообразие и приложението на този модел за създаване на система за свързани отворени данни за биоразнообразието.

За да се завърши системата, са формулирани следните задачи:

⁴UBio: <http://ubio.org/>; Глобални имена: <http://globalnames.org/>; BioGuid: <http://bioguid.org/>; BioNames: <http://bionames.org/>; Профил на таксон на Пенсофт: <http://ptp.pensoft.eu/>; Плаци: <http://plazi.org/wiki/>.

Задача 1: Онтология. Изучаване на информатиката на биоразнообразието и публикуването на данни и разработване на онтология, която позволява интегрирането на данни за биологичното разнообразие от различни източници.

Задача 2: Софтуерна архитектура. Формализация на OpenBiodiv като основана на знания система и създаването на интегрираната ѝ софтуерна архитектура.

Задача 3: Свързани отворени данни. Създаване на свързани отворени данни (LOD) въз основа на публикувани таксономични статии, използващи онтологията от Задача 1.

Задача 4: Софтуерна библиотека. Разработване на инструменти за преобразуване на таксономичните публикации в семантичния модел на онтологията с цел подпомагане на Задача 3.

Задача 5: Методи за работа. Разработване на практически методи за работа за непрекъснато преобразуване на таксономичните данни в таксономични публикации и по този начин актуализиране на набора от данни LOD.

Задача 6: Уеб портал. Създаване на уеб портал с примерни приложения в допълнение към базата от знания.

Методология

В този раздел се очертават изборите, направени преди започване на фазата на проектиране и изпълнение на системата. Те включват, но не се ограничават до напр. парадигмите за програмиране и за база данни.

Избор на парадигма на база данни за OpenBiodiv

OpenBiodiv е формулирана като система, основана на знание, с фокус върху структурирането и взаимовръзката на данни за биоразнообразието. Две от възможните технологии за бази данни, с които системата може да се реализира, са семантична бази данни (triple stores), като напр. GraphDB (Ontotext, 2018) и графична⁵ бази данни (labeled property graphs), като напр. Neo4J (Neo4J Developers, 2012). И двата типа бази данни всъщност се основават на модела на математическия граф и могат да бъдат наричани “графични”. Семантичните (графични) бази данни предлагат много прост модел за данни: всеки факт, съхраняван в такава база данни, се представя като тройка от подлог *subject*, сказуемо *predicate* и пряко допълнение *object*. Subjects са винаги идентификатори на ресурси, докато objects могат да бъдат други идентификатори или буквални стойности (литерали, напр. низове, числа и т.н.). Връзките се дават от predicates (посочени също като идентификатори). Тези връзки се наричат предикати или свойства (*properties*). По този начин може да се визуализира граф, чиито върхове са идентификатори или буквални стойности и, чиито ребра са свойства. Семантичните бази данни имат уникалната черта, че логическият слой също се изразява като тройки, съхранени в базата данни. Този логически слой, известен като онтология (*ontology*), не

⁵Тук трябва да се внимава: графична в случая означава базирана на граф (т.е. съвкупност от ребра и върхове), а не: на графика.

само отговаря за извличането на нови факти от данните (извод), но също така формализира и семантиката (смисловия модел) на знанието.

Не-семантичните графични бази данни, наричани просто графични, предлагат по-свободен модел за данни, като позволяват и ребрата на графа от знания да имат етикети. Например, в граф, чиито върхове са два града *A* и *B*, които са свързани посредством свойството *свързани с път*, е възможно допълнително да се приложи стойността “500 км” към това свойство. По този начин посочваме, че дължината на пътя, свързващ градовете, е 500 км. Този разширен модел не е по-изразителен от простия семантичен, разгледан в предишния параграф. Т.е. няма твърдение, което може да бъде изказано с разширения модел, а да не може да бъде изказано посредством семантичния модел. Сложни взаимоотношения в обикновен triple store могат да бъдат изразени чрез преобразуването на сложни свойства в нови ресурси, които имат собствени свойства. Този процес е известен като реификация (*reification*). Например, двата града *A* и *B* могат да се свържат към друг връх, *R*, посочващ пътя. *R* ще има три свойства: *start*, *end* и *length*. Стойността (object) на *start* ще бъде *A*, на *end* ще бъде *B*, а *length* ще бъде литералът “500 км” или числото 500.

Разликите между графичните и семантичните бази данни са обобщени в Таблица 1. След внимателни разсъждения, семантична база данни е избрана на технология за база данни. Това решение е информирано от широката наличност на висококачествени онтологии и RDF модели в областта на биологичното разнообразие (Baskauf and Webb, 2016; Peroni, 2014) и от популярността на семантичната мрежа (Berners-Lee et al., 2001) в научната общност.

Въпреки това смятаме, че несемантичните графи са по-свободен и по-естествен модел на данни и са напълно подходящи за информатиката за биоразнообразието. По-специално, те осигуряват много по-естествен формализъм за изразяване на взаимоотношенията между таксономичните концепции (обсъдени в глава 2). Също така напоследък семантични бази данни, различни от RDF, като WikiData, стават популярни. Ето защо считаме, че приложимостта на RDF за OpenBiodiv трябва постоянно да бъде преоценявана.

Избор на източници на информация

Биоразнообразието и свързаните с биоразнообразието данни имат два различни “цикъла на живот” (*pro-iBiosphere project final report 2014*). В класическия вариант, след като е направено наблюдение на жив организъм, то е записвано в тетрадка и след това бележка за наблюдение е публикувана в научна статия или монография. Не-правителствената организация “Плаци”, както и Biodiversity Heritage Library (Miller et al., 2012), полагат усилия за дигитализиране на публикации, публикувани по този начин на хартиен носител (Agosti et al., 2007). За тази цел са разработени няколко специални XML схеми (вж. Penev et al., 2011 за преглед), от които TaxPub (Catarano, 2010) и TaxonX са най-широко използвани (Penev et al., 2012). Дигитализирането на публикациите съдържа няколко стъпки. След сканиране и оптично разпознаване на символи (OCR), се извършва text-mining. Тази процедура маркира елементи (semantic markup), които могат да бъдат извлечени и предоставени за бъдещо използване и повторно използване (Miller et al., 2015).

В днешно време знание и данни за биоразнообразието и се публикуват в цифров формат като семантично подобрени публикации, enhanced publications (EP, Claerbout and Karrenbach, 1992; Godtsenhoven et al., 2009; Shotton, 2009). Според Claerbout and Karrenbach, 1992, “EP е публикация, която е разширена с данни

ТАБЛИЦА 1: Разлики между бази данни върху семантични графи (напр. GraphDB) и не-семантични графи (напр. Neo4j).

Критерий	Семантичен граф	Не-семантичен граф
Семантика	<p>Съхранява се в самата база данни като OWL или RDFS-изрази.</p> <p>Осигурява единно пространство за данни.</p> <p>Изисква експерти-онтолози да извличат знания.</p>	<p>Формална семантика обикновено липсва.</p> <p>Бързо внедряване.</p> <p>Унифицирано пространство за данни по-трудно постижимо.</p>
Извод	<p>Осигурен от самата база данни от нейната онтология или формулиран посредством SPARQL заявки. С общо предназначение, по-бавен.</p>	<p>Външен за базата данни. Трябва да се напише за всяка конкретна задача. Със специално предназначение.</p> <p>По-бърз.</p>
Общност	<p>Има богата и зряла общност от онтолози и инженери по знания.</p> <p>Много онтологии за различните дисциплини. Стремяща се към стандарти.</p>	<p>Моделите се създават ad-hoc от програмисти за дадена задача.</p> <p>Постигането на интероперативност на данните изисква усилия и не е от първостепенно значение. Стремяща се към работещи приложения.</p>

от изследвания, допълнителни материали и допълнителни данни. Тя има структура, базирана на обекти, с изрични връзки между обектите. Един обект може да бъде (част от) статия, набор от данни, изображение, филм, коментар, модул или връзка към информация в база данни”. По този начин семантично подобрените публикации са “родени в Интернет и в семантичната мрежа” за разлика от техните хартиени предшественици.

Актът за публикуване в цифров, подобрен формат, се различава основно от публикуването в печатен източник. Основната разлика е, че цифрово публикуваният документ може да бъде структуриран в такъв формат, че да е подходящ както за машинна обработка, така и за човешкото око. В сферата на науката за биологичното разнообразие, списания на Пенсофт като ZooKeys, PhytoKeys и Biodiversity Data Journal (BDJ) от години предлагат публикуване на EP (Penev et al., 2010b).

Предвид факта, че публикациите на Пенсофт и Плаци покриват голяма част от таксономичната литература както по обем, така и по хронология, и че публикациите на тези две издателства са достъпни като семантични EP, периодичните издания на Пенсофт и Плаци бяха избрани като основни източници на информация за OpenBiodiv.

Освен това в системата е включен таксономичният гръбнак на GBIF GBIF Secretariat, 2017 като източник за интеграция на данни. Допълнителни разяснителни детайли са дадени в глава 3.

Избор на методология и среда за програмиране

През 2016 г., въз основа на резултатите от pro-iBiosphere и на съществуващи разработки в областта на информатиката за биологичното разнообразие, е публикуван план за докторантурата (Senderov and Penev, 2016). Тази публикация може да се счита за първата спецификация на проекта на OpenBiodiv. Въпреки това, в хода на разработването на системата, дизайнът на системата бе променен итеративно чрез обратна връзка от сътрудници от проекта BIG4 project⁶. Тези промени са в духа на отворената наука (*open science*) и на *agile разработката на софтуер* (Beck et al., 2001). Този итеративен подход се различава от подхода waterfall, където след фазата на проектиране, спецификациите “са замразени” през дълга изпълнителна фаза.

През последните години програмният език R се използва широко в областта на науката за данни data science (R Core Team, 2016). R има богата библиотека от софтуерни пакети, включваща пакети за обработка на XML (Wickham et al., 2018), за достъп до API (Wickham, 2017) и се фокусира върху отворена наука (Voettiger et al., 2015). Възможностите на R като функционално-ориентиран и интерпретиран език облекчават подхода за итеративно разработване на софтуер, очертан в предходния параграф. Освен това R се използва широко в общността на информатиката за биологичното разнообразие. Поради тази причина средата на софтуера R е избрана за основна програмна среда.

Отворена наука и Семантична Мрежа

Следните методологии правят резултатите от изследването отворени и възпроизводими.

⁶Кандидатът Виктор Сендеров е част от Международната мрежа за обучение на “Мария Склодовска-Кюри” BIG4: Биосистематика, информатика и геномика на големите 4 групи насекоми: обучение на утрешните изследователи и предприемачи.

OpenBiodiv трябва да бъде разгледан от гледна точка на *Open Science*. Съгласно Kraker et al., 2011 и *Was ist Open Science?*, шестте принципа на Отворената наука са: отворена методология, отворен код, отворени данни, отворен достъп, отворени рецензии и отворени образователни ресурси. Целта на Отворената наука е да осигури достъп до всички изследователски продукти: данни, открития, хипотези, код и т.н. Това отваряне гарантира, че научният продукт може да бъде възпроизводим и проверим от други учени (Mietchen, 2014). Съществува голям интерес към разработването на процеси и инструменти, които дават възможност за възпроизводимост и проверка. Тези проблеми са разгледани напр. в специален брой в *Nature*, посветена на възпроизводими изследвания (*Challenges in irreproducible research* 2010). Поради това изходният код, данните и публикациите на OpenBiodiv ще бъдат публикувани открито.

Освен това OpenBiodiv се разглежда като неразделна част от Семантичната Мрежа (Berners-Lee et al., 2001). Семантичната мрежа е визия за бъдещето на Мрежата, където са свързани не само документи, но и данни.

Структура на дисертацията

Във Въведение е дадена мотивировка за съществуването на системата OpenBiodiv, както и обобщение на нейните цели и задачи.

В глава 1 е представена формалната спецификацията и дизайнът на системата, както и нейната архитектура; тази глава формира Задача 2, но в логическия ход на дисертацията е хубаво да бъде разгледана първо. Следващите глави обсъждат изпълнението на OpenBiodiv. Глава 2 предлага формална концептуализация на областта на публикуването на данни за биологичното разнообразие. Въвежда се централният научен резултат на дисертацията – онтологията OpenBiodiv (OpenBiodiv-O) и се покрива Задача 1. Глава 3 описва отворените данни, които са генерирани въз основа на OpenBiodiv-O и формира Задача 3. Глава 4 подробно описва софтуерния пакет RDF4R (R пакет за работа с RDF), който бе използван за създаване на Linked Open Data (OpenBiodiv-LOD) и формира Задача 4. В глава 5 се обсъждат два казуса за внасяне на данни в OpenBiodiv от важни международни хранилища, Задача 5. Глава 6 обсъжда уеб-сайта, който се подготвя да служи на на OpenBiodiv-LOD и приложенията му и покрива Задача 6. В Заключение резултатите на дисертацията са обобщени и научните ѝ и научно-приложни приноси са изтъкнати. Разисква се публикуването и популяризирането на резултатите.

Глава 1

Архитектура на OpenBiodiv

Въвеждат се компонентите на OpenBiodiv, които ще бъдат разгледани подробно в следващите глави. Описва се, как взаимодействат тези компоненти, за да се формира базираната на знанието система OpenBiodiv.

1.1 Какво е OpenBiodiv?

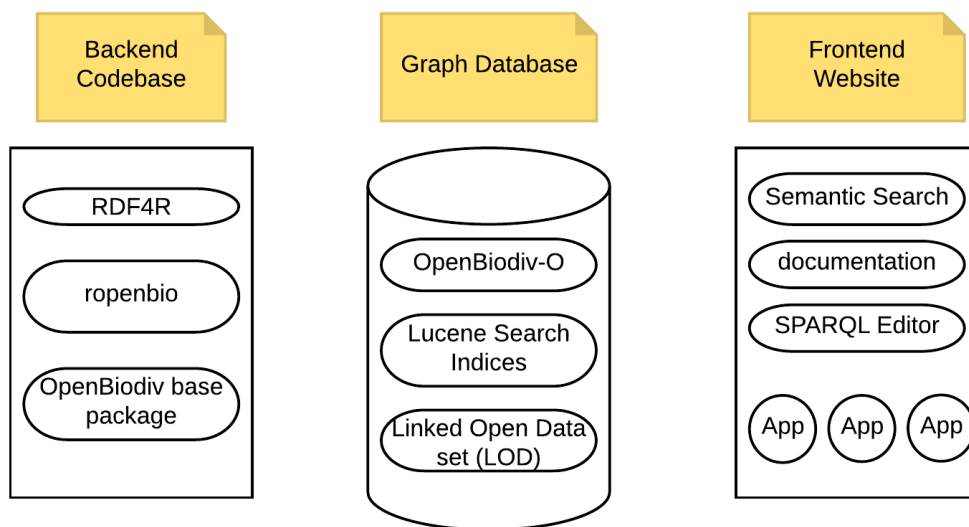
Разбирането за OpenBiodiv като система, базирана на знание, може да се обобщи по следния начин: OpenBiodiv е база данни от взаимосвързана информация за биологичното разнообразие, заедно с логика и приложения, позволяващи на потребителите не само да се допитват до данните, но и да откриват допълнителни факти, свързани с данните. Основните източници на информация в OpenBiodiv са списанията на академичния издател “Пенсофт”, таксономичната информация от Плаци и таксономичният гръбнак на Global Information Biodiversity Facility (GBIF).

Изследователският проблем на архитектурата на OpenBiodiv може да се формулира като проектиране на семантична графична база данни на основата на RDF, която е с отворен достъп и включва информация, предоставяна от Пенсофт, Плаци и GBIF и позволява на потребителите на системата да задават сложни заявки.

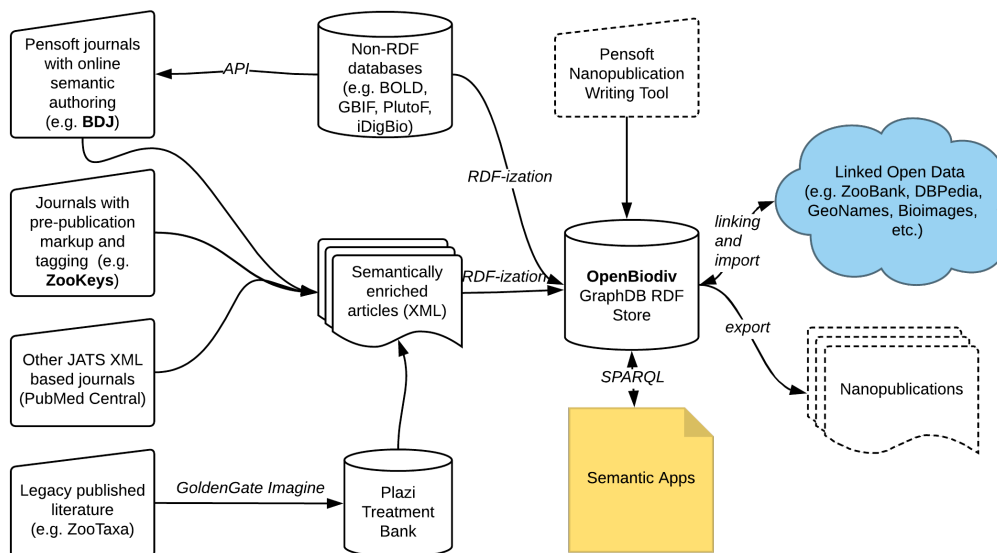
OpenBiodiv се състои от (1) семантична графична база данни, (2) програмен код, осигуряващ функционирането на базата и (3) динамична уеб страница (frontend), улесняваща достъпа до основната база от знания (Фиг. 1.1). OpenBiodiv позволява динамичното вмъкване на данни от хранилища за данни за биологичното разнообразие като BOLD или GBIF. OpenBiodiv извлича от таксономични списания знание (напр. ZooKeys, Biodiversity Data Journal (BDJ), PhytoKeys, MycoKeys и т.н.¹). В същото време знание под формата на факти (triples) се извлича от Plazi TreatmentBank, архив на литература за биологичното разнообразие, съдържащ над 200 хиляди таксономични дискусии² и актуализиран всеки ден. Не на последно място, тези факти са взаимосвързани чрез таксономичния гръбнак на GBIF (GBIF Secretariat, 2017). След това извлеченото знание се съхранява в нашата семантична база данни (Фиг. 1.2).

¹Списанията могат да бъдат достъпени под https://pensoft.net/browse_journals

²Таксономичната дискусия е специален раздел в биологична публикация, която описва и дискутира вид или по-висок таксон. TreatmentBank е достъпен под <https://plazi.org/resources/treatmentbank/>



ФИГУРА 1.1: Компоненти на OpenBiodiv.



ФИГУРА 1.2: Поток на информация в пространството за данни за биологичното разнообразие. Пунктирните линии са компоненти, които все още не са създадени.

1.2 Семантична база данни

Основен резултат от усилията на OpenBiodiv е създаването на семантична база данни, базирана на знания, извлечени от архивите на Пенсофт и Плаци и таксономичния гръбнак на GBIF и достъпни под <http://graph.openbiodiv.net/>. Следва обсъждане на компонентите на базата данни.

Централният резултат от усилията по OpenBiodiv е създаването на формален модел на областта за публикуването на знание за биоразнообразието. Този формален модел е онтологията OpenBiodiv-O (Senderov et al., 2017). Изходният код на онтологията и придружаващата документация могат да бъдат достъпни под <https://github.com/pensoft/openbiodiv-o>. Детайлна дискусия е представена в глава 2.

Използвайки OpenBiodiv-O и инфраструктурата, описана по-нататък в тази глава, са създадени свързани отворени данни, включващи приблизително 200 хиляди записа от Плаци, пет хиляди статии от Пенсофт, както и таксономичния гръбнак на GBIF (над милион биологични имена). Данните са достъпни онлайн чрез работния инструмент на семантичната база данни <http://graph.openbiodiv.net>. Разясняват се подробно в глава 3.

1.3 Backend

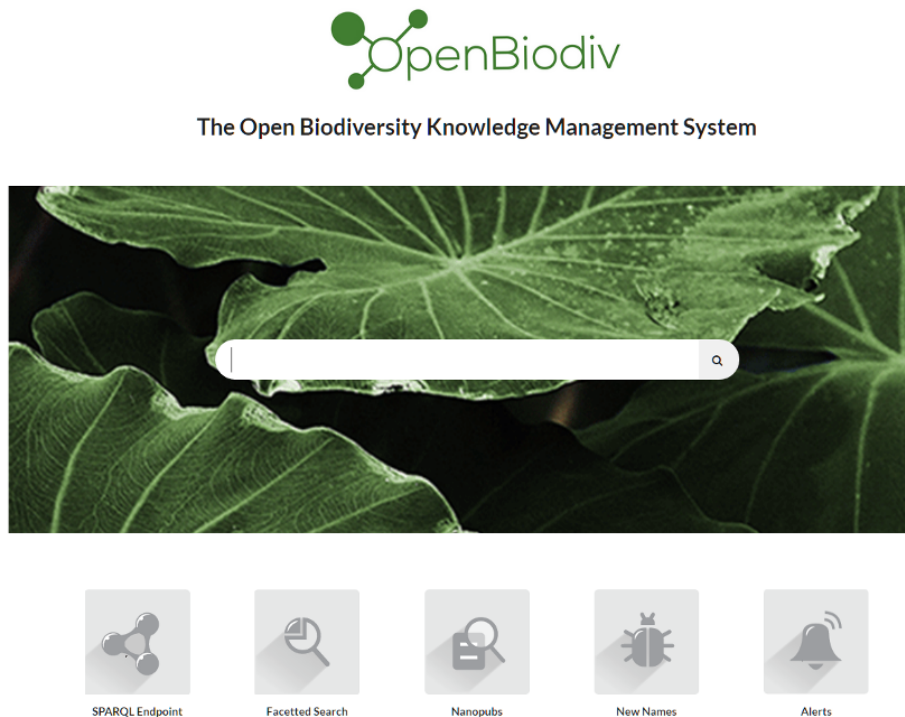
За да се попълва семантичната база данни, е необходимо да се създаде инфраструктура, която преобразува необработени данни (текст, изображения, таблици с данни и т.н.) в структуриран семантичен формат. OpenBiodiv предоставя инфраструктура за трансформиране на научни публикации за биоразнообразието в твърдения под формата на RDF с помощта на инструментите, описани в този раздел.

Едно от по-големите технически предизвикателства за OpenBiodiv е трансформирането на информация за биологичното разнообразие (напр. таксономични имена, метаданни, фигури и т.н.), съхранявани като полу-структуриран XML, в напълно структурирани семантични знания под формата на RDF. За да се реши това предизвикателство, е разработен R пакет, който позволява създаването, манипулирането и записа в семантична база данни на създадения RDF. Този пакет е достъпен под лиценз с отворен код на GitHub под <https://github.com/vsenderov/rdf4r> и е описан в глава 4.

В комбинация с пакета RDF4R, програмният код съдържа още един R пакет, `ropenbio` и базисен програмен код от скриптове и документация, необходими за стартиране на базата данни. `ropenbio` използва пакета RDF4R за преобразуване на полу-структуриран XML в RDF. Той съдържа преобразуванията, необходими за тази реализация. Той е достъпен под <https://github.com/pensoft/ropenbio>. Базисният софтуерен код координира извикването на `ropenbio`, съдържа скриптове за автоматично импортиране на нови ресурси и други подробности. Той е достъпен под <https://github.com/pensoft/openbiodiv>. Генерирането на OpenBiodiv-LOD с помощта на тези пакети е обсъдено в глава 3.

1.4 Работни процеси

Език за екологични метаданни (EML) е популярен формат за описване на екологични данни (Michener et al., 1997). Хранилища на данни за биоразнообразието, като GBIF и DataOne, използват този формат за метаданните, които съхраняват.



ФИГУРА 1.3: Бета версия на потребителския интерфейс.

Автоматичното преобразуване на EML файл в data paper ръкопис от Biodiversity Data Journal³ е възможно с помощта на системата OpenBiodiv (Senderov et al., 2016). Този работен процес е описан подробно в глава 5.

Един от важните видове данни за биологичното разнообразие са данни за наблюдения на организми, occurrence data. Това са данни, които документират наличието на правилно таксономично идентифициран организъм на дадено място и време. Такива данни се съхраняват в международни хранилища като BOLD, GBIF, PlutoF и iDigBio. За да се улесни публикуването на такъв тип данни, е разработен работен процес за импортиране на такива записи от тези бази данни в таксономична статия (taxonomic paper) в списанието Biodiversity Data Journal (Senderov et al., 2016). Този работен процес е описан подробно в глава 5.

1.5 Интерфейс

В допълнение към предоставения endpoint за база данни с възможност за търсене, се разработва уеб сайт, позволяващ семантично търсене и капсулиращ специфични задачи, пакетирани като приложения (<http://openbiodiv.net>). Бета версията вече е в действие. Фиг. 1.3. Дискусия е представена в глава 6.

³Data paper (Chavan and Penev, 2011) е научна статия, обсъждаща научни данни.

1.6 Дискусия

Проектирането на системата е стартирано през втората половина на 2015 г., след разглеждане на различни алтернативи за база данни (Neo4J, GraphDB, WikiBase) и различни технологии за RDF-изация на знание. Освен това е направен избор на източници за информация и типове данни. След анализ на основните модели данни и онтологии, е публикувана спецификацията на системата в Senderov and Penev, 2016 като отворен проект за дисертация (PhD project plan). По време на имплементацията се оказва, че първоначалният план се нуждае от доразвиване, за да отговаря на изменящите се изисквания на системата и на новите предизвикателства. По тази причина през втората и третата година се преминава към модела agile, където спецификацията е разбита на по-малки user-stories, които биват реализирани ad-hoc в рамките на едно- или двумесечни спринтове. Визия за бъдещето на системата е работата по нея да се поеме от agile team, който да се възползва от пълния арсенал на методологията.

Глава 2

Онтология OpenBiodiv-O

OpenBiodiv трансформира информация за биоразнообразието от научни публикации и академични бази данни в семантичен вид. В тази глава е представена OpenBiodiv-O (Senderov et al., 2018) – онтологията, формираща модела на OpenBiodiv за знания и извод. OpenBiodiv-O предоставя концептуален модел на структурата на академична публикация в сферата на биоразнообразието и на съдържаните в нея таксономични концепции. За първи път формално е концептуализирана областта на публикуването на знание за биоразнообразието.

Чрез разработването на онтология, насочена към биологичната таксономия са запълнени празнините между онтологии като Darwin-SW и онтологията за семантично публикуване като онтологията “SPAR”. Счита се, че е предимство да се моделира самият таксономичен процес, а не някакво конкретно състояние на знанието.

Изходният код и документацията са достъпни под лиценза CC BY¹ от GitHub². Следва увод в областта на биологичната таксономия и биоразнообразието.

2.1 Концептуализация на областта

Представена е историята на модерната биологична таксономия, като се започне с Карл Линей (1707-1778), който предлага модерното групиране на организма на *царства, класове, разреди, родове* и използването на латинските биномиални имена в *Systema Naturae* (Linnaeus, 1758). Подчертава се, че работата на учените-таксономи за описване и организиране на биоразнообразието далеч не е пълна. Това информира създаването на новата онтология не като статично формализиране на съществуващата биологична таксономия в компютърно-четима форма, а като формализиране на *научния процес на биологичната таксономия*.

След това се описва подробно, как протича научният процес в биологичната таксономия. Започва се с въвеждането на таксономични концепции и начина, по който се формират. Таксономичната концепция е научна хипотеза (Deans et al., 2012), че определена група от организми съществува в природата. Тя се формира чрез изследване на екземпляри и задължително включва научен критерий, по който те да се групират, често наричан видова концепция (Mallet, 2001). Исторически погледнато, организмите могат да бъдат групирани по външен вид и вътрешно устройство (концепция за морфологични видове) или репродуктивно поведение (биологична видова концепция), но напоследък фокусът се е насочил към групиране въз основа на генетична свързаност (филогенетични и геномни видови концепции).

¹Creative Commons Attribution 4.0 International Public License

²<https://github.com/vsenderov/openbiodiv-o/blob/master/LICENSE.md>

След това се описват единиците на биологичната таксономия и начина, по който те са регламентирани от международните кодекси International Commission on Zoological Nomenclature, 1999; *International code of nomenclature for algae, fungi and plants (Melbourne code) 2012*). Кодексите уреждат по-ниските рангове: видове, род, семейство, разред; по-високите рангове (напр. отдел, царство, домейн и т.н.) могат да бъдат използвани от изследователите, както сметнат за подходящо. Това води до множество конкуриращи се гледни точки.

Публикуването на таксономични концепции е неразделна стъпка в научния работен поток на всеки таксоном. Описва се структурата и типовете таксономични публикации, като се акцентира специално върху секцията “Таксономична дискусия”, раздела в таксономичната публикация, където таксономичната концепция е дефинирана.

Направена е литературна справка, като се разглеждат опитите областта на биологичната таксономия да бъде формално концептуализирана. Интересни са SPAR Ontologies, Peroni, 2014) и TaxPub XML Document Type Definition (Catapano, 2010).

Концептуализацията основно е повлияна освен от практиката и от кодексите (International Commission on Zoological Nomenclature, 1999; *International code of nomenclature for algae, fungi and plants (Melbourne code) 2012*), а така и от стандартите, създадени от TDWG (напр. Darwin Core, DwC, Wiczorek et al., 2012).

Направен е и обзор на областта на концептуалната таксономия (Berendsohn, 1995; Franz and Peet, 2009; Sterner and Franz, 2017), която представлява нова гледна точка за това, как трябва да протича процесът на описване на видове в биологичната таксономия, с оглед на напредъка на информационните технологии.

2.2 Методи

OpenBiodiv-O е изразена посредством RDF чрез използване на RDF Schema (RDFS) и Web Ontology Language (OWL).

За да се разработи онтологията, е използван следният процес: (1) анализ на областта и идентифициране на важните класове обекти и техните взаимоотношения (наричани свойства); (2) анализ на съществуващите информационни модели и онтологии и идентифициране на липсващите класове и свойства за успешно формализиране на областта.

2.3 Резултати

OpenBiodiv-O е споделена формална спецификация на концептуализация на областта на биоразнообразието по смисъла на Gruber, 1993; Obitko, 2007; Staab and Studer, 2009. Тяхното разбиране за онтология е въведено в Background.

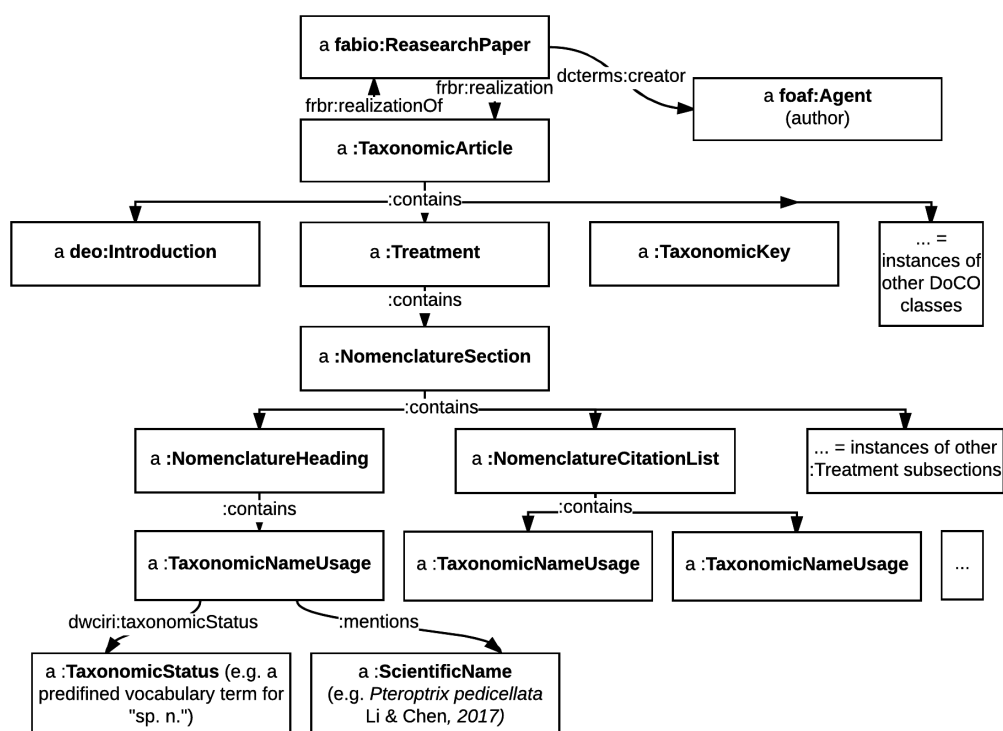
Има няколко домейна, в които моделираните ресурси падат. Първият е научната област за публикуване на данни за биоразнообразието. Втората област е тази на таксономичната номенклатура. Третата област е на по-широката таксономия (например таксономични понятия и техните взаимоотношения, видови събития, черти).

Публикуване на данни. Рамката на онтологията “SPAR” е разширена посредством като въвеждането на нова категория за таксономичните статии, техните подраздели, както и нов клас за споменаване на таксономично име (вж. следващ подраздел). Тези нови класове са обобщени в Таблица 2.1.

ТАБЛИЦА 2.1: Нови класове в областта на публикуването на данни за биоразнообразието.

Class QName	Comment
:Treatment	section of a taxonomic article
:NomenclatureSection	subsection of Treatment
:NomenclatureHeading	contains a nomenclatural act
:NomenclatureCitationList	list of citations of related concepts
:MaterialsExamined	list of examined specimens
:BiologySection	subsection of Treatment
:DescriptionSection	subsection of Treatment
:TaxonomicKey	section with an identification key
:TaxonomicChecklist	section with a list of taxa for a region
:TaxonomicNameUsage	mention of a taxonomic name

Графичното представяне на взаимоотношенията между ресурси на класовете, свързани с публикуването, които OpenBiodiv въвежда, може да се намери в диаграмата на фиг. 2.1.



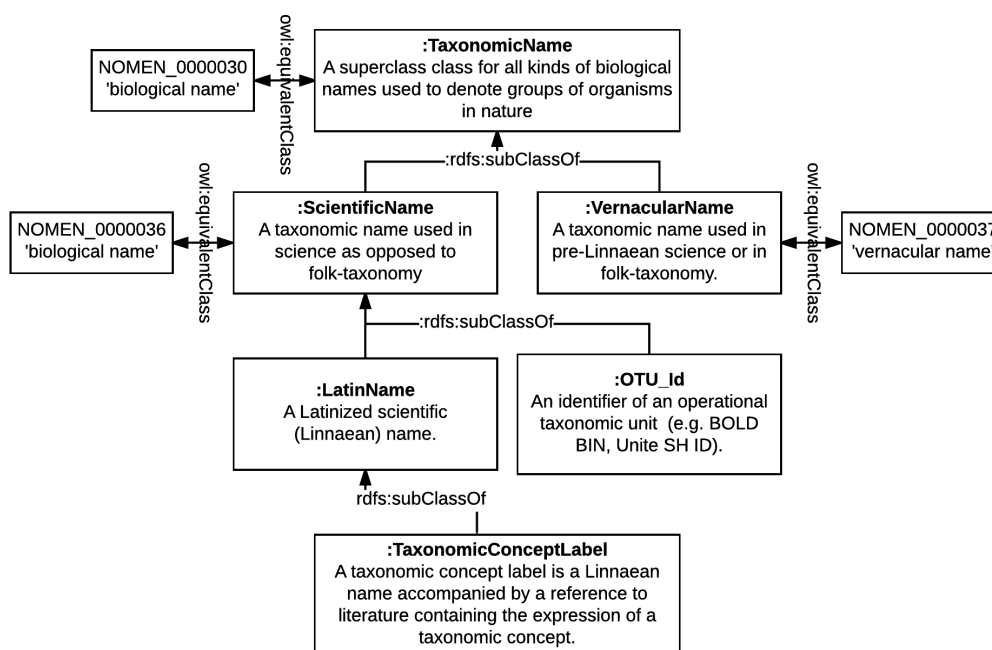
ФИГУРА 2.1: Графична представяне на взаимоотношенията между ресурсите, които OpenBiodiv въвежда за публикуване на данни за биоразнообразието.

Класовете и свойствата, които се въвеждат, са в съответствие с модела на функционалните изисквания за библиографски записи (FRBR), използван от SPAR. Таксономичната статия се третира като конкретен израз/запис, FRBR Expression, на абстрактното понятие работа, FRBR Work, представляващо интелектуалното

съдържание на статията. Таксономичната дискусия се третира подобно на Въведение, Методи, Резултати и т.н., т.е. също е FRBR Expression и DEO discourse element. Таксономична концепция е съответният абстрактен ресурс от клас FRBR Work на дадена таксономична дискусия.

Биологичната номенклатура. Биологичната номенклатура е система с над 200 годишна традиция датираща до преди времето на информатиката и дори до преди времето на Теорията за еволюцията на Дарвин. Много е трудно да се моделира поради сложността си и само частично е обхваната от онтолозиите NOMEN и TNSS (въведени в подраздел “Previous Work”). С OpenBiodiv-O се използва подход "отдолу-нагоре" за моделиране на използването на таксономични имена в статиите. Където е възможно, понятията от OpenBiodiv-O са подравнени към NOMEN.

Йерархия от класовете на таксономичните имена е поместена във фиг. 2.2. Освен това се въвежда taxonomic name usage (TaxonomicNameUsage).



ФИГУРА 2.2: Йерархията, включваща както традиционните таксономични наименования, така и таксономични концептуални етикети и оперативни таксономични единици.

Въвежда се Taxonomic Concept Label (TaxonomicConceptLabel). Етикетът на таксономична концепция (TCL) е линеено име плюс позоваване на публикация, където дискутираният таксон е дефиниран. Връзката се осъществява чрез ключовата дума “sec.” (Латински за (*secundum* Verendsohn, 1995). Напр. *Andropogon virginicus* var. *tenuispatheus* Blomquist, 1948. Тук Blomquist, 1948 е валидна библиографска справка за публикацията, в която концепцията е дефинирана.

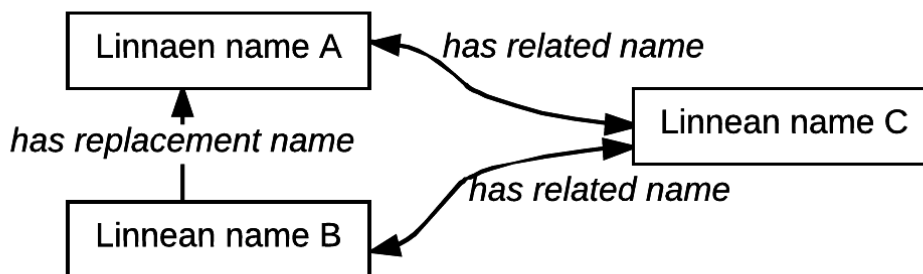
Извадени са съкращения от таксономични термини от около 4000 статии в четири таксономични списания (ZooKeys, Biodiversity Data Journal, PhytoKeys и MusoKeys) и е създаден таксономичен речник на термините, който обхваща осемте най-често срещани случая Таблица 2.2). Латинските съкращения, които

са класифицирани в тези класове, могат да бъдат намерени на страницата на OpenBiodiv-O GitHub. (вж. Методи за повече подробности).

ТАБЛИЦА 2.2: Речник с таксономични термини.

Vocabulary Instance QName	Example Abbrev	Comment
:TaxonomicUncertainty	<i>incertae sedis</i>	Taxonomic Uncertainty
:TaxonDiscovery	<i>sp. n.</i>	Taxonomic Discovery
:ReplacementName	<i>comb. n.</i>	Replacement Name
:UnavailableName	<i>nomen dubium</i>	Unavailable Name
:AvailableName	<i>stat. rev.</i>	Available Name
:TypeSpecimenDesignation	<i>lectotype designation</i>	Type Specimen Designation
:TypeSpeciesDesignation	<i>type species</i>	Type Species Designation
:NewOccurrenceRecord	<i>new country record</i>	New Occurrence Record (for region)

Въз основа на анализ на термините за таксономични статуси, са идентифицирани два модела за съответствие между кодирането на латинизираните научни имена (Фиг. 2.3). Моделът *заместващо име*, изпълняван чрез свойството `replacementName`, показва, че вместо едно латинизирано име, трябва да се използва друго. Този модел обхваща голямо разнообразие от случаи в кодексите, като например поставянето на един вид таксон в нов род (нова комбинация), поправката на наименования поради номенклатурни причини (ново име), или прилагането на Принципа на приоритет за откриването на синоними ("syn nov. International Commission on Zoological Nomenclature, 2017).

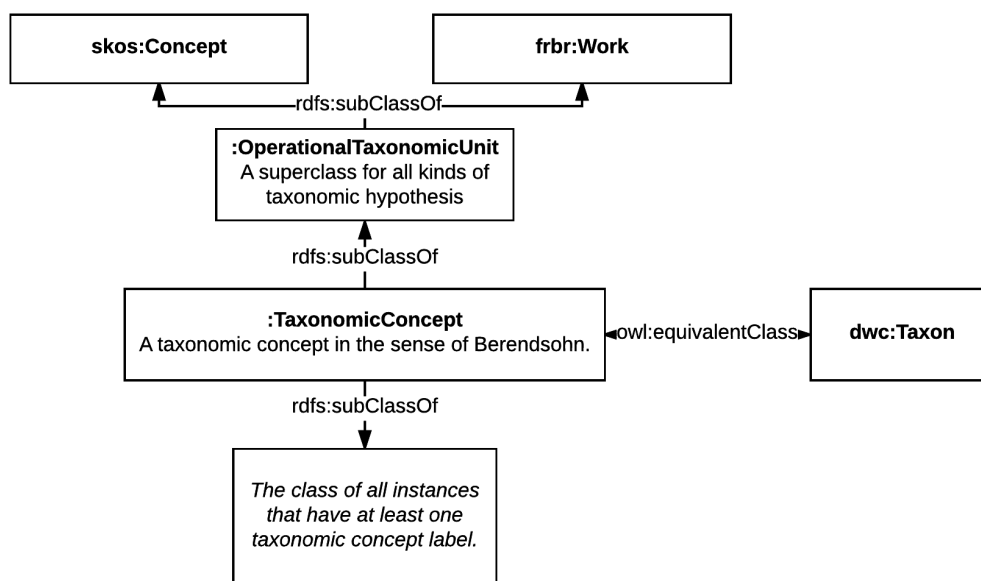


ФИГУРА 2.3: Веригите на *заместващи имена* могат да бъдат проследени, за да се намери използваното понастоящем име. *Свързано име* показва, че две имена са свързани по някакъв начин, но не кой е предпочитан.

Другият модел е този на *свързани имена* (`relatedName`). Това е по-широк модел, който показва, че две имена са някак си свързани. Например, те могат да бъдат синоними, едното да замества другото, или да сочат към таксономично свързани таксономични понятия. Например, *Harmonia manillana* (Mulsant, 1866) е свързано с *Caria manillana* Mulsant, 1866 *Harmonia manillana* (според Poorani and Booth, 2016 лектотипус на *Harmonia manillana* (Mulsant, 1866) sec. Poorani and Booth, 2016 носи името *Caria manillana* Mulsant, 1866).

Както се вижда от фиг. 2.2, таксономичните имена на OpenBiodiv-O са приравнени към NOMEN имена.

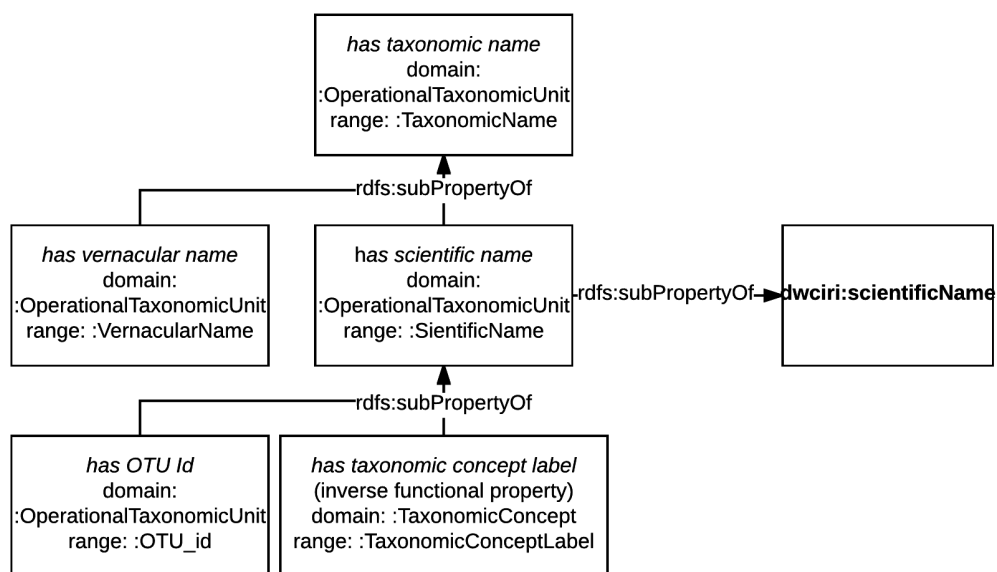
В OpenBiodiv, таксономичните имена не са носители на семантична информация за таксоните. Тази задача се изпълнява от нов клас, Taxonomic Concept (`TaxonomicConcept`). Таксономична концепция е теорията, която таксономът формира около таксон посредством научна таксономична публикация. Тя винаги има етикет, състоящ се от таксономичното име и библиографско позоваване на статията, в която името е описано. Въвежда се и по-общ клас, оперативно таксономично звено (`OperationalTaxonomicUnit`), което може да се използва за всички видове таксономични хипотези, включително такива, които нямат правилен таксономичен етикет. Класовата йерархия е илюстрирана на Фиг. 2.4. Свойствата на таксономичните имена са илюстрирани в Фиг. 2.5. Двата начина за изразяване на взаимоотношения между таксономични концепции са дадени във Фиг. 2.6).



ФИГУРА 2.4: Таксономичната концепция е от клас `skos:Concept`, `frbr:Work`, `dwc:Taxon` и има поне един етикет.

Простите отношения между таксони не са подходящи за машинен извод. Ето защо Franz and Peet, 2009, позовавайки се на Koperski et al., 2000 предлагат да се използва езикът RCC-5, за изразяване на взаимоотношенията между таксономичните понятия. От техни колеги е разработена програмата Euler (Chen et al., 2014), която използва Answer Set Programming (ASP) за разсъждения и извод по таксономичните взаимоотношения, стъпвайки на RCC-5. Изводната машина за RCC-5 не е част от OpenBiodiv, тъй като тази задача може да бъде изпълнена от euler; въпреки това, ние OpenBiodiv-O предоставя RCC-5 терминологичен речник.

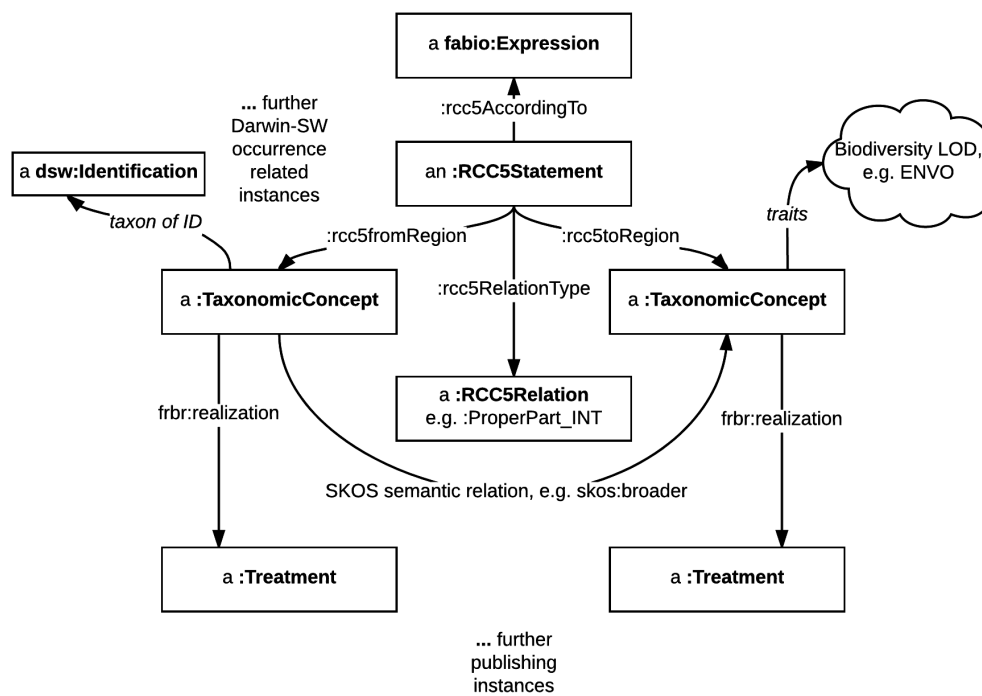
Таксономичните понятия са приведени в съответствие с DarwinCore (DwC) и е проведена дискусия за това как са представени таксономичните понятия, свързани както с прости отношения (SKOS), така и в подходящ за извод вид (RCC-5). Също така се обсъждат взаимоотношенията между биологичните имена и таксономичните концепции. Обсъжда се как OpenBiodiv-O е първият по рода онтологичен модел на таксономична публикация. Тя проправя пътя към създаване на граф от знания за биологичното разнообразие.



ФИГУРА 2.5: Йерархията на свойствата съвпада с йерархията на таксономичните имена и е приравнена към DarwinCore.

2.4 Заключение

Главата предоставя концептуализация на таксономичния процес и формализация в OpenBiodiv-O. Въвеждат се класове и свойства в областта на публикуването на биоразнообразие и биологичната систематика и ги привежда в съответствие с важните онтологии, специфични за този домейн.



ФИГУРА 2.6: За да изразите RCC-5 връзка между концепциите, създайте `RCC5Statement` и използвайте съответните свойства, за да свържете две таксономични концепции през него. Освен това, таксономичните понятия са свързани със качества (например екология в ENVO), събития (например Darwin-SW) и са абстрактният клас на Таксономичните дискусии.

Глава 3

Свързани отворени данни

Създадени са свързани отворени данни OpenBiodiv-LOD, съдържащи информация за биологичното разнообразие, извлечена от списанията на Пенсофт и базата данни на Плаци. Данните са интегрирани с помощта на таксономичния гръбнак на GBIF. Като онтология се използва новата OpenBiodiv-O, разработен в хода на дисертацията. Предлага се на общността на информатиката за биоразнообразието да използва OpenBiodiv-LOD като централна точка на семантичния граф за знания за биоразнообразието. OpenBiodiv-LOD са достъпни под <http://graph.openbiodiv.net>.

OpenBiodiv-LOD е синтетичен набор от данни. Той не съдържа предварително непубликувани данни. Вместо това той интегрира в една база от знания информация, която преди това е била оповестена в академични списания и бази данни. Интеграцията позволява материализация на скрити взаимоотношения. В следващите няколко параграфа ще обсъдим източниците на информация, които са комбинирани в OpenBiodiv-LOD и видовете ресурси, които са извлечени, както и общия модел на данните. Също така се обсъждат принципите на Linked Open Data, които свързват всичко заедно. Главата завършва с много примери за заявки в набора от данни и с техническа дискусия за начина, по който тя е генерирана.

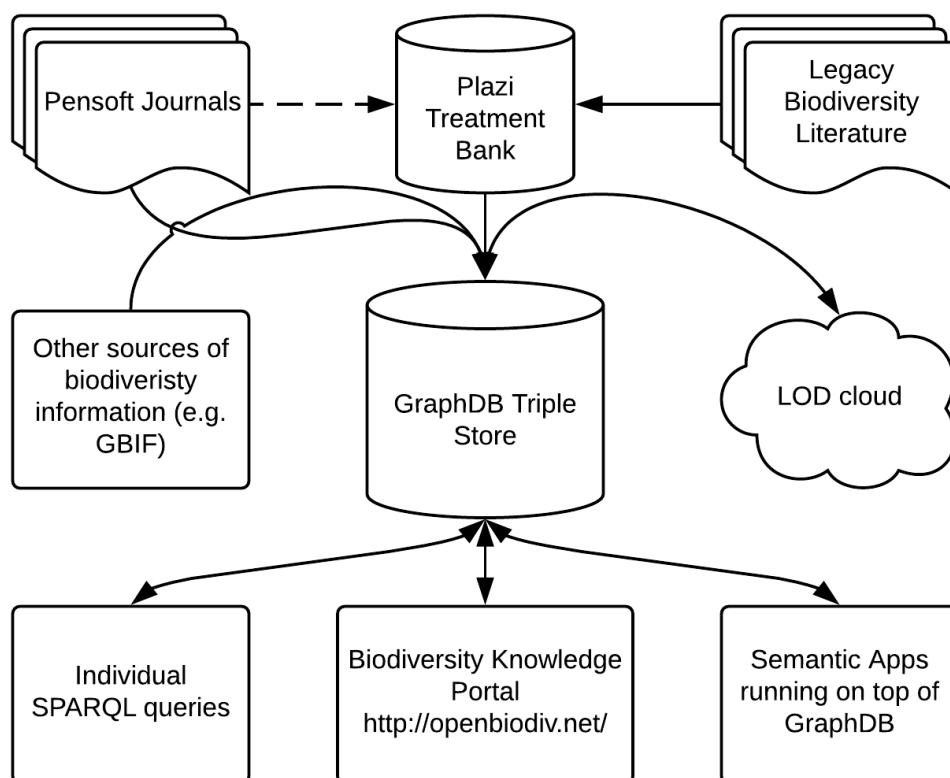
3.1 Източници на данни

Данните в OpenBiodiv към времето на писането на дисертацията идват от три основни източника: таксономичния гръбнак на GBIF (GBIF Secretariat, 2017), и научни статии публикувани от Пенсофт и Плаци (Fig. 3.1).

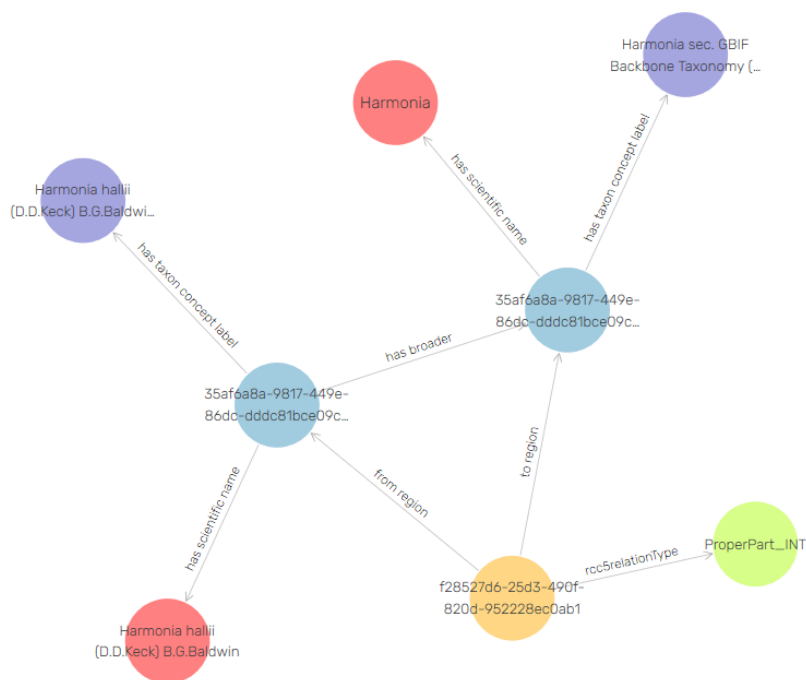
GBIF е най-голямото международно хранилище на данни за наблюдения на организми (occurrence data). GBIF позволява на потребителите да търсят в тяхната система, като използват таксономичната йерархия, Nub (GBIF Secretariat, 2017). Nub е база данни, която организира таксономични концепции в йерархия, обхващаща всички биологични имена, събрани от GBIF. Тя е синтетична, т.е. алгоритмично генерирана класификация. По този начин гръбнакът GBIF не представлява експертен консенсус за това как таксоните са йерархично подредени според еволюционните критерии в природата, но въпреки това е много полезен на практика.

За да псе предоставят същите възможности на OpenBiodiv, концепциите от Nub са импортирани като `openbiodiv:TaxonomicConcept` по OpenBiodiv-O (фиг. 3.2).

Всички валидни статии от таксономичните списания, публикувани от Пенсофт и упоменати в Таблица 3.1 бяха конвертирани в RDF и съхранени в графа от знания на биологичното разнообразие. В допълнение, всички валидни таксономични дискусии (treatments) на Плаци бяха конвертирани в RDF и също съхранени в графа. Процедурата по RDF-изиране се повтаря всяка седмица и



ФИГУРА 3.1: Опростен модел на архитектурата на OpenBiodiv от глава 1 фокусираща върху източниците на информация.



ФИГУРА 3.2: Илюстрация на представянето на йерархична информация, импортирана от GBIF като таксономични концепции в OpenBiodiv.

по този начин семантичната база данни винаги съдържа най-новите статии и таксономични дискусии. RDF-изацията е възможна благодарение на факта, че списанията на Пенсофт публикуват статии в TaxPub XML (Catapano, 2010), докато Плаци публикува дискусиите си в TaxonX (Penev et al., 2011) (Fig. 3.3). И двете схеми са стандартни и общо-достъпни.

ТАБЛИЦА 3.1: Списания ма Пенсофт, които са превърнати в RDF.

Journal Name	Submission Style	Number of Articles
ZooKeys	Word document	3829
PhytoKeys	Word document	537
MycoKeys	Word document	127
Biodiversity Data Journal	Web based (ARPHA)	490
Journal of Orthoptera Research	Word document	32

ТАБЛИЦА 3.2: Типове данни, маркирани в TaxPub and TaxonX и тяхната кореспонденция към RDF типовете, които се използват в OpenBiodiv.

Datatype	TaxPub	TaxonX	RDF Type
Article metadata	T	T	fabio:JournalArticle and related
Keyword group	T	F	openbiodiv:KeywordGroup
Abstract	T	T	sro:Abstract
Title	T	F	doco:Title
Author	T	T	foaf:Person
Introduction section	T	F	deo:Introduction
Discussion section	T	T	orb:Discussion
Treatment section	T	T	openbiodiv:Treatment
Nomenclature section	T	T	openbiodiv:NomenclatureSection
Materials examined	T	T	openbiodiv:MaterialsExamined
Diagnosis section	T	T	openbiodiv:DiagnosisSection
Distribution section	T	T	openbiodiv:DistributionSection
Taxonomic key	T	T	openbiodiv:TaxonomicKey
Figure	T	T	doco:Figure
Taxonomic name usage	T	T	openbiodiv:TaxonomicNameUsage

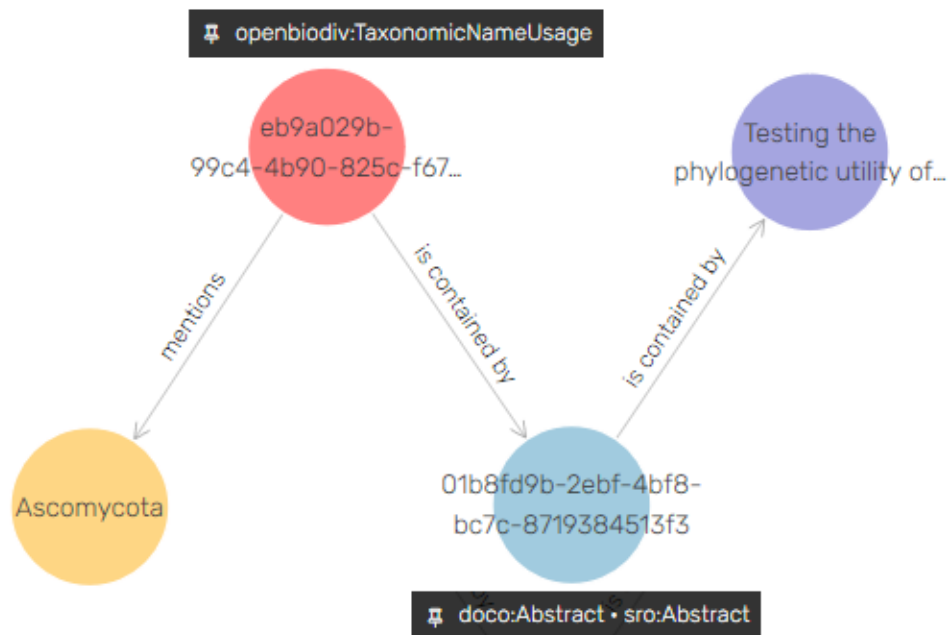
3.2 Свързани отворени данни

Linked Open Data (LOD, Heath and Bizer, 2011) е идея на Семантичната Мрежа (Berners-Lee et al., 2001), целяща да гарантира полезността на данни, публикувани в Мрежата, като улеснява тяхното повторно намиране и употреба от трети лица. В тази секция са разяснени принципите на Свързаните Отворени Данни и тяхното приложение в OpenBiodiv (Heath and Bizer, 2011).

Връзка между модела OpenBiodiv-O, подробно разяснен в глава 2, и допълнителните онтологии е дадена във Fig. 3.5.

3.3 Примерен SPARQL

В тази секция е илюстрирано, как представеният модел за данни, заедно с инкорпорираната в базата данни информация е достатъчен за да се изпълнят някои



ФИГУРА 3.3: Употреба на таксономични имена в OpenBiodiv (taxonomic name usage).

ID: <http://openbiodiv.net/35af6a8a-9817-449e-86dc-dddc81bce09c-4239-ScientificName>

Taxonomic name information sparql

Curculionidae

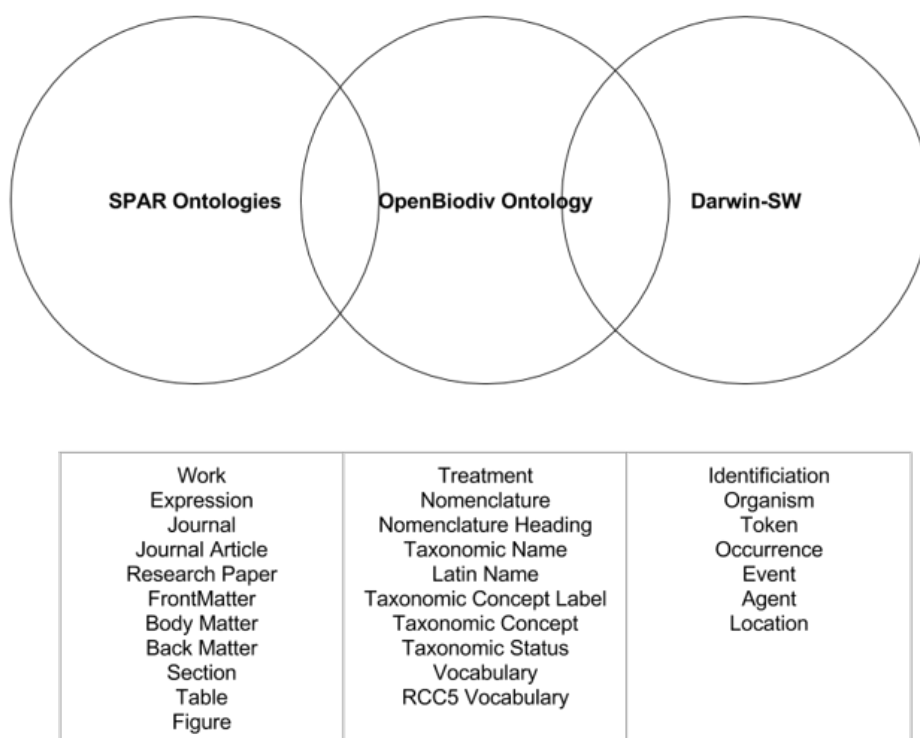
Usage statistics sparql

Mentioned in: journal article 1062 abstract 42 title 23

Related names

- Omus thoracicus
- Harpalus lewisii
- Amara aenea (DeGeer, 1774)
- (Morphnosoma) Lutshnik, 1915
- (Lithochlaenius) Kryzhanovskij, 1976
- Ozaena
- Arctelaphrus
- Bembidium flavicaudus
- Harpalini Bonelli, 1810
- Harpalus illectus
- Phloeozetaeus

ФИГУРА 3.4: Визуализация.



ФИГУРА 3.5: Връзка между OpenBiodiv-O и други публично-достъпни онтологии.

много интересни заявки върху графа за знания.

Ето кратък списък с типовете заявки, които са илюстрирани:

1. Заявка върху структурата на статията: Тъй като в OpenBiodiv LOD статиите са разбити по техните компоненти (вж напр. Таблица 3.2) и споменаванията на таксони (taxonomic name usages) винаги са свързани със специфична част на статията, могат да се правят запитвания, използващи тази структура.
2. Запитване за таксономични концепции
3. Размито търсене с Lucene: използван е Lucene connector (Ontotext, 2018) на GraphDB, за размито търсене, където се допускат правописни грешки.
4. Отговаряне на експертни въпроси: поради машината си за извод, OpenBiodiv може да функционира като експертна система. Напр. за оценка на научната загуба след пожара в Museu Nacional в Бразилия, системата може да бъде запитана относно списък с видовете, чиито екземпляри бяха загубени в пожара в Рио де Жанейро.
5. Проверка на валидността на таксономично име: проверк, дали дадено таксономично име е валидно, т.е. дали е уместно да се употребява в научен контекст, или е било премахнато или заменено от друго име.

3.4 Генериране на данните

В тази секция се разкриват техническите подробности за това, как данните бяха създадени. Използват се следните инструменти

1. Пакет RDF4R
2. Пакет ROpenBio
3. Пакет TSV4RDF разработен от Пенсофт
4. Базисен код на OpenBiodiv

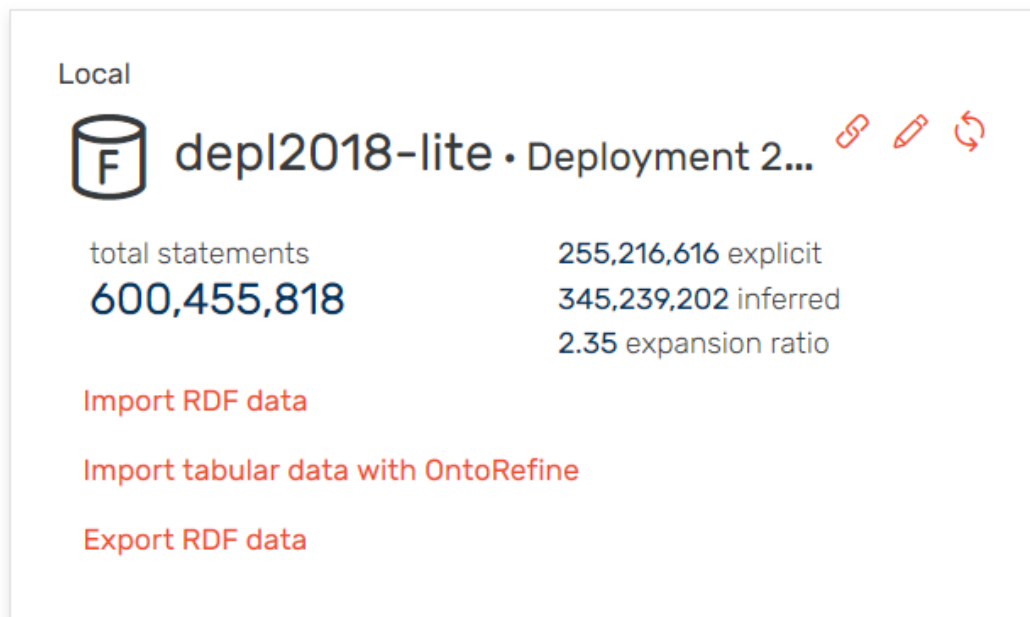
3.4.1 Трансформация от XML в RDF

За трансформация от XML в RDF е използвана йерархичната структура на XML, за да се реши проблемът рекурсивно. Решението е реализирано посредством Алгоритъм 1.

Algorithm 1 Екстрактор

```
1: procedure EXTRACTOR(XML Node  $X$ )
2:    $a \leftarrow$  extract atoms of  $X$                                  $\triangleright$  Atoms extraction
3:    $r \leftarrow$  construct RDF from  $a$                                 $\triangleright$  RDF construction
4:    $C \leftarrow$  find relevant sub-nodes of  $X$                         $\triangleright$  Recursively applies itself
5:    $R \leftarrow$  apply Extractor on each  $C_i \in C$ 
6:   return  $r \cup R$ 
7: end procedure
```

Текстовите полета в XML са наричани атоми. Тяхното извличане се осъществява с помощта на езика XPATH. След извличане на атомите, те могат да бъдат



ФИГУРА 3.6: Брой твърдения.

събрани обратно във формата на RDF. Проблемът за RDF-изирането на цял XML файл е реализиран посредством разбиването му на малки проблеми, за които решението е тривиално.

За да работи Extractor, трябва да се дефинира схемата, в която XML е кодиран. Спецификацията включва, какви XML възли се търсят и къде се намират. След това рекурсивно се указва за всеки възел, кои под-възли се търсят и тяхното местоположение по XPATH спрямо техния родителски възел. И накрая, за всеки възел трябва да се укажат адресите на атомите и да се напише конструктор, който ги сглобява.

Спецификацията се извършва в R6 клас на R. Дефинираме и схемите на Пенсофт и Плаци: TaxPub¹ и TaxonX².

Извлечените RDF твърдения се изпращат в GraphDB, който се намира на адрес <http://graph.openbiodiv.net/>. OpenBiodiv-O автоматично материализира допълнителни твърдения посредством машината за извод. Допълнително се изпълняват следното допълнително правило:

Ъпдейт на заместващите имена: изгражда връзките от тип `replacementName` между вмъкнатите имена.

3.5 Дискусия

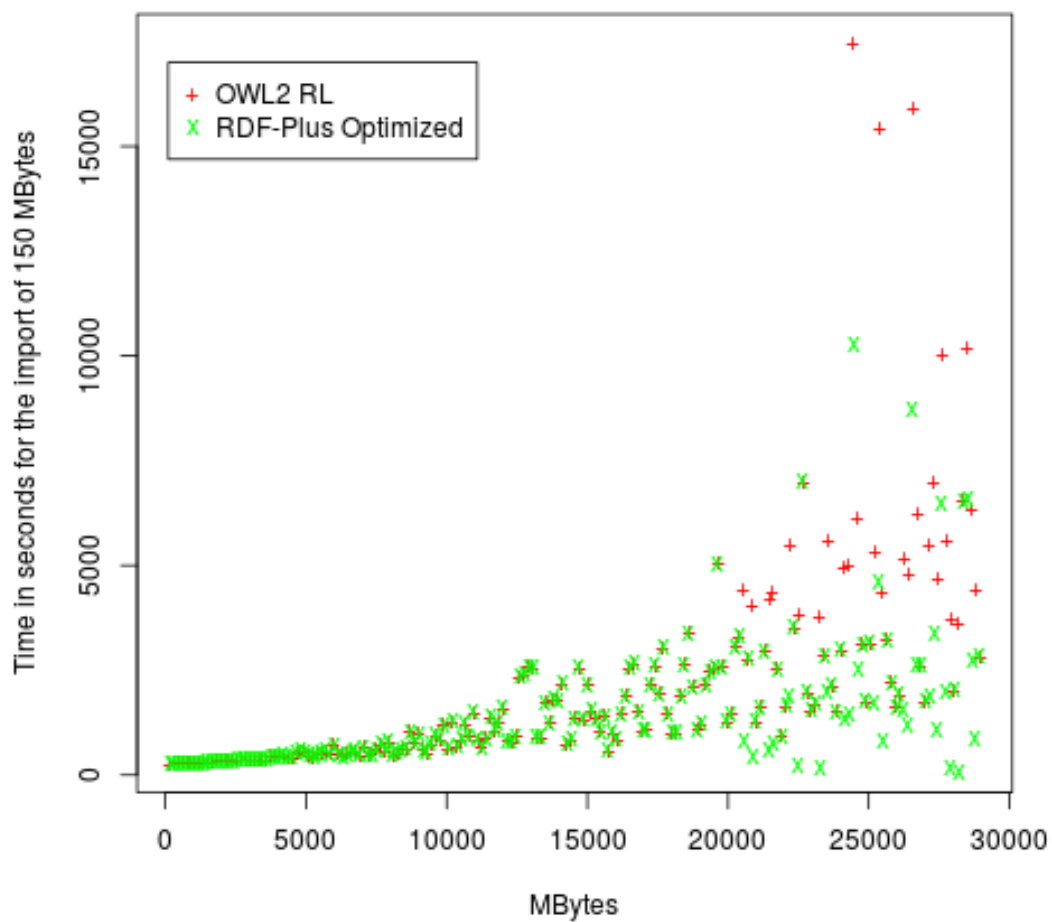
Установено е, че обемът данни, който се обработва, в комбинация с пълна OWL логика води до неприемлива производителност. В тази секция се разглежда този проблем (Фиг. 3.6).

Направени са някои основни изводи и са очертаваме насоки за бъдещо развитие.

¹<https://github.com/pensoft/ropenbio/blob/ redesign/R/taxpub.R>

²<https://github.com/pensoft/ropenbio/blob/ redesign/R/taxonx.R>

Performance degradation as a function of database size in MBytes



ФИГУРА 3.7: Визуализация на деградацията на производителността.

Глава 4

Библиотека за R за работа с RDF

RDF4R (`rdf4r`) е R пакет за работа с RDF. Тя е разработена като част от проекта OpenBiodiv, но е напълно свободна от специфичен за OpenBiodiv код и може да се използва за общи цели, изискващи инструменти за работа с RDF данни в средата за програмиране R (R Core Team, 2016).

В тази глава се разяснява функционалността на пакета.

Основните функции на пакета са:

1. Връзка със семантична база данни
2. Функции за преобразуване на SPARQL заявки в R функции
3. Работа с литерали и идентификатори
4. Работа с префикси
5. Създаване и сериализация на RDF
6. Основна терминология на семантични понятия

Дискутирани са приликите и разликите на RDF4R с подобни пакети като `rdflib`. Обсъждат се, как са използвани модели от функционалното програмиране за създаване на RDF4R. Освен това се обсъждат, как се използват модели от обектно-ориентирано програмиране за създаване на RDF4R.

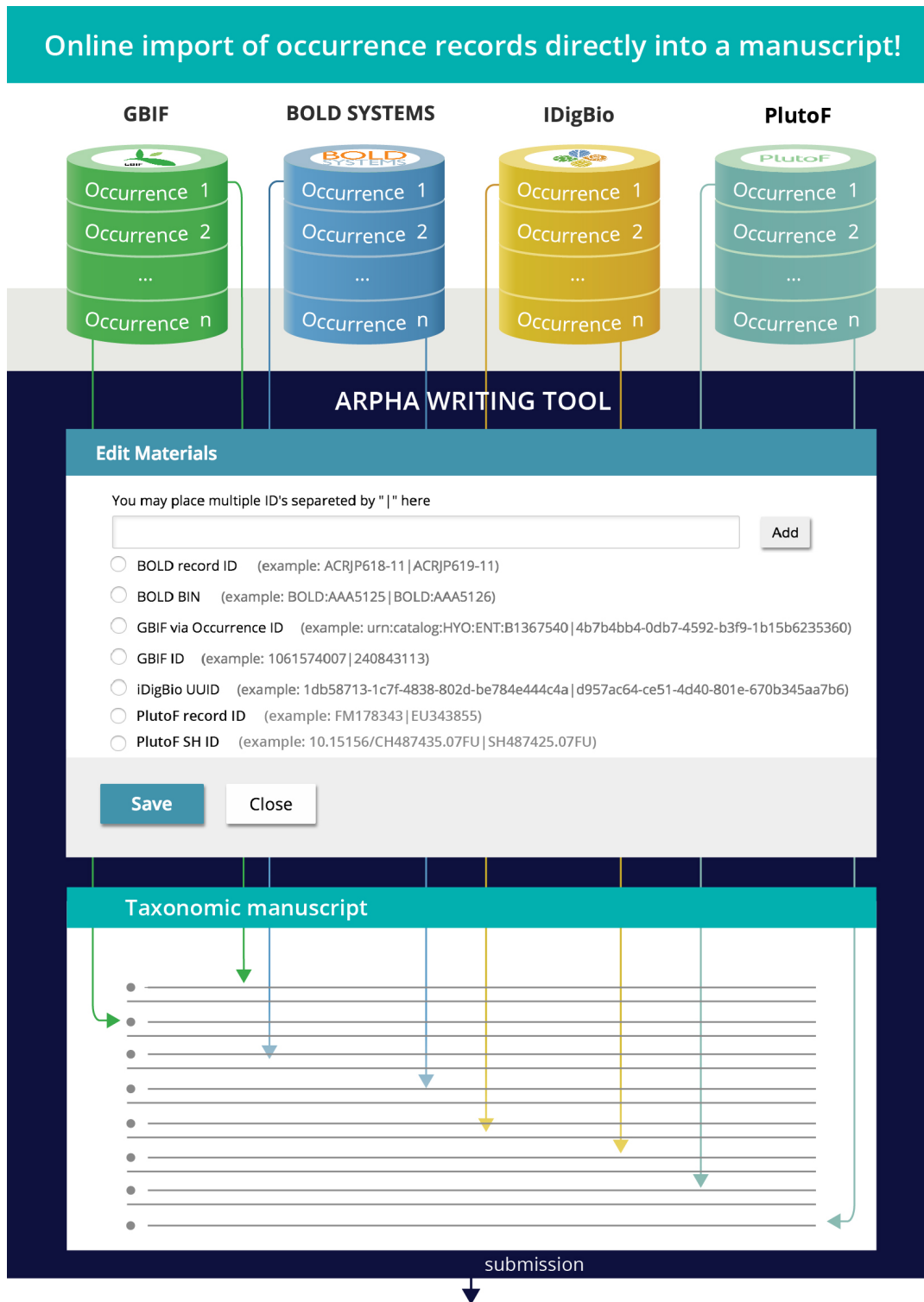
Глава 5

Работни процеси

В тази глава се обсъждат два автоматизирани работни процеса за обмен на данни за биологичното разнообразие, разработени като част от OpenBiodiv: (1) автоматично внасяне на записи за наблюдения на видове (occurrence data) в ръкописи и (2) автоматично генериране на ръкописни материали от екологичните метаданни (EML). Работните потоци са представени на webinar за организацията iDigBio и публикувани като научна статия (Senderov et al., 2016). Слайдовете от презентацията под github.com/vsenderov/idigbio-webinar.

Работен процес 1: Автоматизиран импорт на occurrence data в ръкописи, разработени в ARPHА Writing Tool Разработен е автоматичен импорт на occurrence data в таксономична статия в системата ARPHА от международните бази данни GBIF, BOLD, iDigBio и PlutoF (Фиг. 5.1).

Работен процес 2: Автоматизирано генериране на ръкопис от EML метаданни в ARPHА Writing Tool Създаден е работен процес, който позволява на авторите да създават автоматично ръкописни материали от метаданни, съхранявани в EML (фиг. 5.2).



**Biodiversity
Data Journal**

<http://bdj.pensoft.net>

ФИГУРА 5.1: Работният поток от информационните портали GBIF, BOLD Systems, iDigBio и PlutoF чрез елементи на потребителски интерфейс в AWT.



ФИГУРА 5.2: Полето за потребителски интерфейс за качване на EML файлове в ARPHA.

Глава 6

Уеб портал

Под openbiodiv.net може да се стигне до основния портал, който дава достъп до ресурсите на OpenBiodiv. Този портал е разработен от Пенсофт в подкрепа на OpenBiodiv и представя два визуални елемента на потребителя: лентата за търсене и списък с икони на приложения в долната част. Освен това, под graph.openbiodiv.net (също достъпен от иконката SPARQL) може да се достигне работния плот OpenBiodiv за на SPARQL.

Тези функции на потребителския интерфейс (UI) са предназначени да улеснят трите типа потребители на системата:

1. Основно ниво: използва лентата за търсене.
2. Ниво на специалист: използва приложения.
3. Power user: използва работната плот за SPARQL или R.

The screenshot displays the OpenBiodiv web portal interface. At the top, the logo for OpenBiodiv Beta is shown, with the text 'The Open Biodiversity Knowledge Management System' below it. A search bar at the top center contains the text 'daniel mietchen' and a search icon. Below the search bar, the profile for 'Person' is displayed, with the ID: <http://openbiodiv.net/2c525459-67b1-427b-9420-844b70c41f03>.

The 'Person info' section includes the name 'Daniel Mietchen' and a list of affiliations: 'Museum für Naturkunde, Berlin, Germany'. A 'sparql' button is visible in the top right corner of this section.

The 'Articles' section features a bar chart titled 'Articles per Year' showing the number of articles published in 2011, 2013, and 2014. The x-axis represents the number of articles (0 to 20), and the y-axis represents the year. The chart shows 2 articles in 2011, 1 article in 2013, and 1 article in 2014. A 'show all articles' link is located below the chart. A 'sparql' button is also present in the top right corner of this section.

The 'Collaborators' section displays a list of names, each followed by a '1' indicating the number of collaborations. The names are: Guido Sautter, Kevin Richards, Pavel Stoev, Aleksandra Pawlik, Bachir Balech, David Eades, Donat Agosti, Walter G. Berendsohn, Adam Brunke, Hannes Hettling, Robert Hoehndorf, Rod Page, Nesrine Akkari, Tom van Dooren, Marko Tsihtinen, Quentin John Groom, Andreas Plank, Thomas D. Hamann, Alan R. Williams, Soraya Sierra, Thomas Pape, Robert A. Morris, Gregor Hagedorn, Lyubomir Penev, Claus Weiland, Patricia Kelbert, Nicola Nicolson, Ayco Holleman, Donald Hobern, Don Kirkup, Teodor Georgiev, David King, Yuri Lammers, Niall Beard, Carina Mara de Souza, Matthew Blissett, Peter Hovenkamp, George Gosline, Thibaut DeMeulemeester, Jeremy A. Miller, Terry Erwin, Christian Breninkmeijer, Lars Hendrich, David Peter Shorthouse, Ross Mounce, Henrik Enghoff, David Koon-Bong Cheung, Michael Balke, and Serrano Pereira. A 'sparql' button is located in the top right corner of this section.

ФИГУРА 6.1: Илюстрация на основната употреба на OpenBiodiv за търсене на информация за човек.

Заклучение

Резултати

Считаме, че представеният научен труд изпълнява поставената цел и задачи.

Резултат 1. Основният резултат на дисертацията е създаването на концептуален модел на областта на публикуването на знание за биологичното разнообразие и логическото му формализиране му под формата на онтологията OpenBiodiv-O. Този резултат е подробно дискутиран в глава 2 и в Senderov et al., 2018 и изпълнява Задача 1. Изходният код на онтологията е достъпен под github.com/pensoft/openbiodiv-o.

Резултат 2. Създаването на софтуерната архитектура на OpenBiodiv, очертана в глава 1 и Senderov and Penev, 2016. Този резултат отговаря на Задача 2.

Резултат 3. Създаването на свързани отворени данни, OpenBiodiv-LOD, състоящи се от трансформирани в RDF и интегрирани в база данни твърдения за биоразнообразието. Твърденията са извлечени от XML на научни статии, публикувани от Пенсофт, XML таксономични дискусии, публикувани от Плаци, и CSV-дъмп на таксономичния гръбнак на GBIF. OpenBiodiv-LOD са достъпни под graph.openbiodiv.net и са описани в глава 3. Този резултат отговаря на Задача 3.

Резултат 4. Софтуерен пакет за средата за програмиране R, RDF4R. RDF4R дава възможност за манипулиране на RDF в R. С негова помощ е осъществено трансформирането на научни публикации от полу-структуриран XML формат в структуриран семантичен RDF. Този резултат е обсъден в глава 4 и Senderov et al., 2016 и изпълнява Задача 4. Пакетът е достъпен онлайн като свободен софтуер под github.com/pensoft/rdf4r. Освен това допълнителен изходен код (не-оптимизиран), описващ XML схемите на Пенсофт и Плаци и работещ в тандем с RDF4R за конвертиране на XML в RDF може да бъде намерен под github.com/pensoft/ropenbio.

Резултат 5. Процеси за работа, които позволяват обогатяването на XML статии, публикувани от Пенсофт, чрез автоматично импортиране на данни от основните международни хранилища за данни за биологичното разнообразие: BOLD, GBIF, iDigBio, както и PlutoF. Освен това, благодарение на резултатите на тази дисертация е възможно автоматично да се създадат ръкописи от метаданни, кодирани в Ecological Metadata Language (EML). Обсъждането на тези автоматизирани работни процеси – автоматично генериране на data papers и автоматичен импорт на occurrence records - се извършва в глава 5. Този резултат изпълнява Задача 5.

Резултат 6. Уеб-сайт openbiodiv.net, съдържащ семантична търсачка и приложения. Уеб-сайтът е обсъден в глава 6 и изпълнява Задача 6.

Дискусия, изводи и бъдещи насоки

OpenBiodiv-LOD са свързани отворени данни за биоразнообразието на основата на научни публикации. OpenBiodiv-O е базата на свързаните отворени данни OpenBiodiv-LOD. Чрез разработването на онтология, насочена към биологичната таксономия, са запълнени празнините между онтолозиите за ресурси от биологичното разнообразие като Darwin-SW и онтолозиите за семантично публикуване като SPAR. Счита се, че е правилно да се моделира самият таксономичен процес, а не някакво конкретно състояние на знанието, тъй като процесът по описването на биоразнообразието на земята далеч не е завършил. На дадения етап, покритието на онтологията е достатъчно, за да може тя да бъде основата за създаването на LOD. В този смисъл онтологията е завършена. От друга страна добавяне на класове и свойства за нови типове данни за биоразнообразието е възможно и желателно.

Свързаните отворени данни, подобно на онтологията, вече представляват солиден ресурс за биолозите, тъй като включват информация от повечето статии, публикувани от Пенсофт и Плаци и наброяват над 600 милиона твърдения. Подобно на онтологията, те следва да бъдат разширявани. Статия, относно генерирането на данни по онтология е подготвена и представя пред научното списание *Cybernetics and Information Technologies*.

Софтуерният пакет RDF4R беше използван успешно за създаването на LOD, и като такъв, той може да се счита за завършен. Както всеки един софтуерен пакет, обаче, той би следвало да бъде поддържан и развиван. Очаква се публикуването му под формата на Software Description Paper.

Уеб-сайтът е все още в бета-версия. Функционалността, която работи отлично е семантичната търсачка. За някои основни типове данни има създадени макети (templates) за визуализация. Въпреки всичко, сайтът търпи разширение и повечето потребители на системата използват езика за SPARQL за търсене.

Важен извод, който може да се направи като обобщение на постигнатите резултати е, че е възможно използването на семантичен граф за интегрирането на голям обем от данни за биологичното разнообразие. Неочаквано ни се отдаде възможност да илюстрираме мощта на графа при анализа на щетите от трагичния пожар в Museu Nacional в Рио де Жанейро. Освен това илюстрирахме, че е възможно да се напишат сравнително прости логически изводи, позволяващи проверка на валидността на дадено таксономично име.

Поради значителния размер на данните е установено, че макар и използването на семантичен граф да е възможно, някои от първоначално избраните технологии са неприложими или трудно приложими. Извършено е наблюдението (вж. глава 3), че практическото приложение на пълния логически модел OWL е трудно, поради проблеми с бързодействието. Вместо него се спряхме на по-маломощната, но по-бърза RDFS. Друго наблюдение, което направихме е, че макар и програмната среда R да дава известни предимства при бързото създаване на прототипа на системата, при увеличаването на сложността на програмния код, необходимо при реално-съществуващата система за покриване на всички частни случаи, език с динамични типове като R създава главоболия при дебъгинг. Същевременно останахме впечатлени от мощния инструментариум на функционалното програмиране, който R предоставя.

Основна трудност се оказва дизамбигуацията на ресурси като имена на автори или таксономични имена. Във функционалния дизайн на пакета RDF4R заложихме модули, които позволяват при търсенето на идентификатор за даден ресурс да вмъкнем списък от функции/ правила за неговата дизамбигуация. Въпреки

това имаме само ограничен успех с дизамбигуацията на базата на правила и по тази причина в производствената система тя е изключена към момента.

Имайки предвид тези и други “поуки” бъдещето развитие на проекта OpenBiodiv може да се очертае по следния не непременно изчерпателен начин:

1. Като непосредствени цели да разширят LOD и онтологията с нови типове данни и нови източници на данни, използвайки съществуващия фреймуърк. Такива данни са напр. геномните данни, данните за место-находища на видове, (био-)географски данни, визуални данни, описателни данни и пр.
2. Да се търси още по-тясна интеграция с други съществуващи хранилища за данни за биоразнообразието освен GBIF. Напр. BioImages, iNaturalist, BOLD, и т.н.
3. Като по-дългосрочна задача да се изследва преминаването от семантичен граф към технология, където машината за извод е разграничена от семантичния граф като WikiData или Neo4j. Освен увеличеното бързодействие, това ще даде допълнителна гъвкавост на проекта като например позволи използването на машини за извод различни от RDF-базирането като напр. Euler.
4. Да бъде продължена разработката на системния софтуер с още по-широко приложение на функционалното програмиране и портирането му към функционален език като напр. Haskell или O’Caml.
5. Да бъде изследван проблемът за дизамбигуацията и сродните проблеми за named entity recognition на интересни ресурси от биоразнообразието, както и различни задачи от разпознаването на образи, от гледната точка на машинното обучение.
6. Разширяване на уеб-сайта на системата с повече темплейти и нови приложения.

Основни научни и научно-приложни приноси

Резултатите, дискутирани в предишните два раздела, обуславят следните научни и научно-приложни приноси:

1. Научен принос: създаване на онтология и формален модел на областта на публикуване на знание за биологичното разнообразие.
2. Научно-приложен принос: анализ на източниците на информация и създаване на OpenBiodiv-LOD.
3. Научно-приложен принос: имплементация на софтуерните модули на OpenBiodiv.

OpenBiodiv-O запълва уникалната ниша между библиографски онтологии като SPAR и онтологии за биоразнообразието като Darwin-SW и като такава без съмнение е от голям научен интерес за общността на информатиката на биоразнообразието. Работата има и сериозен научно-приложен характер, като в допълнение на онтологията предоставя и свързани отворени данни, използващи онтологията, както и програмно обезпечение за потребителите и разработчиците на системата.

Списък с публикации

Публикации в международни научни списания

1. Viktor Senderov and Lyubomir Penev. 2016. "The Open Biodiversity Knowledge Management System in Scholarly Publishing". *Research Ideas and Outcomes* 2, no. e7757 (). ISSN: 2367-7163. doi:10.3897/rio.2.e7757. <http://rio.pensoft.net/articles.php?id=7757>. Намерени 3 уникални цитата от Franz and Sterner, 2018, Ordynets et al., 2017 и Burt and Mengual, 2017.
2. Sarah Faulwetter, Evangelos Pafilis, Lucia Fanini, Nicolas Bailly, Donat Agosti, Christos Arvanitidis, Laura Boicenco, Terry Catapano, Simon Claus, Stefanie Dekeyzer, Teodor Georgiev, Aglaia Legaki, Dimitra Mavraki, Anastasis Oulas, Gabriella Papastefanou, Lyubomir Penev, Guido Sautter, Dmitry Schigel, Viktor Senderov, Adrian Teaca, and Marilena Tsompanou. 2016. "EMODnet Workshop on mechanisms and guidelines to mobilise historical data into biogeographic databases". *Research Ideas and Outcomes* 2, no. e10445 (). ISSN: 2367-7163. doi:10.3897/rio.2.e10445. <http://rio.pensoft.net/articles.php?id=10445>. Намерен уникален цитат от Pyron, 2018.
3. Pedro Cardoso, Pavel Stoev, Teodor Georgiev, Viktor Senderov, and Lyubomir Penev. 2016. "Species Conservation Profiles compliant with the IUCN Red List of Threatened Species". *Biodiversity Data Journal* 4 (): e10356. ISSN: 1314-2828, 1314-2836. doi:10.3897/BDJ.4.e10356. <http://bdj.pensoft.net/articles.php?id=10356>. Индексация в WoS SCOPUS, както и в SJR 0.465. 4 уникални цитата от Bachman et al., 2018, Lin et al., 2017, Li et al., 2017 и Milano et al., 2017.
4. Viktor Senderov, Teodor Georgiev, and Lyubomir Penev. 2016. "Online direct import of specimen records into manuscripts and automatic creation of data papers from biological databases". *Research Ideas and Outcomes* 2 (): e10617. ISSN: 2367-7163. doi:10.3897/rio.2.e10617. <http://rio.pensoft.net/articles.php?id=10617>. Намерен уникален цитат от Ordynets et al., 2017.
5. Lyubomir Penev, Daniel Mitchen, Vishwas Chavan, Gregor Hagedorn, Vincent Smith, David Shotton, Éamonn Ó Tuama, Viktor Senderov, Teodor Georgiev, Pavel Stoev, Quentin Groom, David Remsen, and Scott Edmunds. 2017b. "Strategies and guidelines for scholarly publishing of biodiversity data". *Research Ideas and Outcomes* 3, no. e12431 (). ISSN: 2367-7163. doi:10.3897/rio.3.e12431. <http://riojournal.com/articles.php?id=12431> Намерени 8 уникални цитата от Tennant et al., 2017, Marwick and Birch, 2017, Kissling et al., 2018, Mathieu, 2018, Шашков and Иванова, 2018, Шашков et al., 2017, Filippova et al., 2017 и Филиппова et al., 2017.
6. Lyubomir Penev, Teodor Georgiev, Peter Geshev, Seyhan Demirov, Viktor Senderov, Pliyana Kuzmova, Iva Kostadinova, Slavena Peneva, and Pavel Stoev. 2017a. "ARPHA-BioDiv: A toolbox for scholarly publication and dissemination of biodiversity data based on the ARPHA Publishing Platform". *Research Ideas and Outcomes* 3, no. e13088 (). ISSN: 2367-7163. doi:10.3897/rio.3.e13088. <http://riojournal.com/articles.php?id=13088>
7. Emmanuel Arriaga-Varela, Matthias Seidel, Albert Deler-Hernández, Viktor Senderov, and Martin Fikáček. 2017. "A review of the Cercyon Leach (Coleoptera, Hydrophilidae, Sphaeridiinae) of the Greater Antilles". *ZooKeys* 681 (): 39–

93. ISSN: 1313-2970, 1313-2989. doi:10.3897/zookeys.681.12522. <https://zookeys.pensoft.net/articles.php?id=12522>. Индексация в WoS IF 1.079, Q3 SCOPUS, SJR 0.533.

8. Viktor Senderov, Kiril Simov, Nico Franz, Pavel Stoev, Terry Catapano, Donat Agosti, Guido Sautter, Robert A. Morris, and Lyubomir Penev. 2018. "OpenBiodiv-O: ontology of the OpenBiodiv knowledge management system". *Journal of Biomedical Semantics* 9, no. 5 (). ISSN: 2041-1480. doi:10.1186/s13326-017-0174-5. <https://jbiomedsem.biomedcentral.com/articles/10.1186/s13326-017-0174-5> Индексация в WoS IF 1.6, Q3 SCOPUS, SJR 0.952. Намерени 3 уникални цитата от Michel et al., 2018, Page, 2018b и Page, 2018a.
9. OpenBiodiv-LOD: Populating the OpenBiodiv Ontology. Непубликувана статия, представена в списанието Cybernetics and Information Technologies.

Представен е списък с публикации, свързани с дисертацията. Изброените статии са публикувани без изключение в четири международни научни списания: пет статии в Research Ideas and Outcomes, една статия в ZooKeys (WoS IF 1.079, Q3 SCOPUS, SJR 0.533), една статия в Biodiversity Data Journal (WoS SCOPUS, SJR 0.465) и една статия в Journal of Biomedical Semantics (WoS IF 1.6, Q3 SCOPUS, SJR 0.952). Общият брой намерени цитати, които кандидатът е натрупал, изключвайки само-цитати е 20, като конкретните цитиращи статии са посочени в списъка по горе. Общият брой намерени цитати, включително само-цитати и цитати на работа извън обхвата на дисертацията, е 48 (Google Наука).

[1] е ранна версия на въведението, както и на глава 1 и съдържа труд по Задача 2 (Архитектура). Текстовете на публикациите [2, 3, 5, 6, 7] не са част от текста на дисертацията дословно, но съдържат труд по Задача 5 (Методи за работа). Представените в тези публикации идеи са в голяма степен включени в глава 5, чиито текстови гръбнак се формира от публикация [4]. [8] е най-важната публикация в рамките на тази дисертация и е публикувана във реномираното списание Journal of Biomedical Semantics. [8] съставлява съдържанието на глава 2 и е основната част от работата, изпълняваща Задача 1 (Онтология). Статията беше на заглавната страница на JBS в продължение на няколко месеца (фиг. 6.2). Глава 3 и глава 4, които формират Задачи 3, 4, се подготвят като ръкописи в международни списания. Освен това софтуерната библиотека RDF4R, описана в глава 4, се подготвя да бъде изпратена в хранилището с отворен код rOpenSci¹.

Апробация на резултатите

Доклади пред научен семинар на ПНЗ

1. Доклад пред научен семинар на ИБЕИ на БАН на 26.10.2015 г. ("Публикуване, визуализация и разпространение на първични и геномни данни за биологичното разнообразие на основата на откритата система за управление на информацията").
2. Доклад пред научен семинар в ИИКТ на БАН на 31.03.2016 г. (Open Biodiversity Knowledge Management System)
3. Доклад пред научен семинар на ИИКТ на БАН за 23.03.2018 г. (OpenBiodiv: a knowledge-based system of biodiversity information)

¹"We build software with a community of users and developers, and educate scientists about transparent research practices." <https://ropensci.org/>

11. Доклад на европейската конференция на биосистематиците, BioSyst.eu 2016 на 15.08.2017 г. (The OpenBiodiv Knowledge System: The Future of Access to Biodiversity Knowledge)
12. Доклад на международния симпозиум TDWG 2017 в Отава, Канада от 1. до 6.10.2017 г. (OpenBiodiv Computer Demo: an Implementation of a Semantic System Running on top of the Biodiversity Knowledge Graph)
13. Доклад на международния симпозиум TDWG 2017 в Отава, Канада от 1. до 6.10.2017 г. (OpenBiodiv: an Implementation of a Semantic System Running on top of the Biodiversity Knowledge Graph)
14. Постер на международния симпозиум TDWG 2017 в Отава, Канада от 1. до 6.10.2017 г. (OpenBiodiv: an Implementation of a Semantic System Running on top of the Biodiversity Knowledge Graph)
15. Доклад по време на работната среща на BIG4 в Ла Палма, Испания от 30. окт. до 3 ноем. 2017 г. (Midterm Progress Report)
16. Доклад пред научен семинар на групата по биоинформатика (група Ронкуист) в Кралския природо-научен музей в Стокхолм на 29.11.2017 г.

Bibliography

- Agosti, D., C. Klingenberg, N. Johnson, C. Stephenson, and C. Catapano. 2007. *Why not let the computer save you time by reading the taxonomic papers for you?* Tech. rep. Zenodo. doi:[10.5281/zenodo.15584](https://doi.org/10.5281/zenodo.15584).
- Agosti, Donat. 2006. "Biodiversity data are out of local taxonomists' reach". *Nature* 439, no. 7075 (): 392–392. ISSN: 0028-0836, 1476-4687. doi:[10.1038/439392a](https://doi.org/10.1038/439392a). <http://www.nature.com/articles/439392a>.
- Arriaga-Varela, Emmanuel, Matthias Seidel, Albert Deler-Hernández, Viktor Senderov, and Martin Fikáček. 2017. "A review of the Cercyon Leach (Coleoptera, Hydrophilidae, Sphaeridiinae) of the Greater Antilles". *ZooKeys* 681 (): 39–93. ISSN: 1313-2970, 1313-2989. doi:[10.3897/zookeys.681.12522](https://doi.org/10.3897/zookeys.681.12522). <https://zookeys.pensoft.net/articles.php?id=12522>.
- Auer, Sören, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. "DBpedia: A Nucleus for a Web of Open Data". In *The Semantic Web*, ed. by David Hutchison, Takeo Kanade, Josef Kittler, Jon M. Kleinberg, Friedemann Mattern, John C. Mitchell, Moni Naor, Oscar Nierstrasz, C. Pandu Rangan, Bernhard Steffen, Madhu Sudan, Demetri Terzopoulos, Doug Tygar, Moshe Y. Vardi, Gerhard Weikum, Karl Aberer, Key-Sun Choi, Natasha Noy, Dean Allemang, Kyung-Il Lee, Lyndon Nixon, Jennifer Golbeck, Peter Mika, Diana Maynard, Riichiro Mizoguchi, Guus Schreiber, and Philippe Cudré-Mauroux, 4825:722–735. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN: 978-3-540-76297-3 978-3-540-76298-0. doi:[10.1007/978-3-540-76298-0_52](https://doi.org/10.1007/978-3-540-76298-0_52). http://link.springer.com/10.1007/978-3-540-76298-0_52.
- Bachman, Steven P, Eimear M Nic Lughadha, and Malin C Rivers. 2018. "Quantifying progress toward a conservation assessment for all plants". *Conservation Biology* 32 (3): 516–524.
- Barrasa, Jesús. 2017. *RDF Triple Stores vs. Labeled Property Graphs: What's the Difference?* Published online. Visited on 01/11/2019. <https://neo4j.com/blog/rdf-triple-store-vs-labeled-property-graph-difference/>.
- Baskauf, Steve, and Campbell O. Webb. 2016. "Darwin-SW: Darwin Core-based terms for expressing biodiversity data as RDF". *Semantic Web Journal* 7 (6): 629–643. doi:[10.3233/SW-150203](https://doi.org/10.3233/SW-150203).
- Beck, Kent, Mike Beedle, Arie van Bennekum, Alistair Cockburn, Ward Cunningham, Martin Fowler, James Grenning, Jim Highsmith, Andrew Hunt, Ron Jeffries, Jon Kern, Brian Marick, Robert C. Martin, Steve Mellor, Ken Schwaber, Jeff Sutherland, and Dave Thomas. 2001. *Manifesto for Agile Software Development*. Published online. Visited on 01/11/2019. <http://www.agilemanifesto.org/>.
- Berendsohn, Walter G. 1995. "The concept of "potential taxa" in databases". *Taxon*: 207–212. Visited on 07/22/2017. <http://www.jstor.org/stable/1222443>.

- Berners-Lee, Tim, James Hendler, and Ora Lassila. 2001. "The Semantic Web". *Scientific American* 284, no. 5 (): 34–43. ISSN: 0036-8733, visited on 04/13/2018. doi:10.1038/scientificamerican0501-34. <http://www.nature.com/doifinder/10.1038/scientificamerican0501-34>.
- Blomquist, HL. 1948. *The grasses of North Carolina*. Duke University Press.
- Boettiger, Carl, Scott Chamberlain, Edmund Hart, and Karthik Ram. 2015. "Building software, building community: lessons from the rOpenSci project". *Journal of Open Research Software* 3 (1).
- Bollacker, Kurt, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. "Freebase: a collaboratively created graph database for structuring human knowledge", 1247. ACM Press. ISBN: 978-1-60558-102-6, visited on 07/22/2018. doi:10.1145/1376616.1376746. <http://portal.acm.org/citation.cfm?doid=1376616.1376746>.
- Burt, Trevor, and Ximo Mengual. 2017. "Origin and diversification of hoverflies: a revision of the genera Asarkina and Allobaccha—A BIG4 Consortium PhD project". *Research Ideas and Outcomes* 3:e19860.
- Cardoso, Pedro, Pavel Stoev, Teodor Georgiev, Viktor Senderov, and Lyubomir Penev. 2016. "Species Conservation Profiles compliant with the IUCN Red List of Threatened Species". *Biodiversity Data Journal* 4 (): e10356. ISSN: 1314-2828, 1314-2836. doi:10.3897/BDJ.4.e10356. <http://bdj.pensoft.net/articles.php?id=10356>.
- Catapano, Terence. 2010. "TaxPub: an extension of the NLM/NCBI journal publishing DTD for taxonomic descriptions". Published online. Visited on 07/11/2017. <https://www.ncbi.nlm.nih.gov/books/NBK47081/>.
- Challenges in irreproducible research*. 2010. Published online. Visited on 04/15/2018. <https://www.nature.com/collections/prbfkwmwvz>.
- Chavan, Vishwas, and Lyubomir Penev. 2011. "The data paper: a mechanism to incentivize data publishing in biodiversity science". *BMC Bioinformatics* 12 (Suppl 15): S2. ISSN: 1471-2105, visited on 02/13/2018. doi:10.1186/1471-2105-12-S15-S2. <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-12-S15-S2>.
- Chen, Mingmin, Shizhuo Yu, Nico Franz, Shawn Bowers, and others. 2014. "Euler/X: a toolkit for logic-based taxonomy integration". Published online, *arXiv preprint arXiv:1402.1992*. Visited on 08/11/2017. <https://arxiv.org/abs/1402.1992>.
- Claerbout, Jon F., and Martin Karrenbach. 1992. "Electronic documents give reproducible research a new meaning", 601–604. Society of Exploration Geophysicists. Visited on 04/15/2018. doi:10.1190/1.1822162. <http://library.seg.org/doi/abs/10.1190/1.1822162>.
- Constantin, Alexandru, Silvio Peroni, Steve Pettifer, David Shotton, and Fabio Vitali. 2016. "The Document Components Ontology (DoCO)". Ed. by Oscar Corcho. *Semantic Web* 7, no. 2 (): 167–181. ISSN: 22104968, 15700844, visited on 08/13/2017. doi:10.3233/SW-150177. <http://www.medra.org/servlet/aliasResolver?alias=iospress&doi=10.3233/SW-150177>.
- Deans, Andrew R., Matthew J. Yoder, and James P. Balhoff. 2012. "Time to change how we describe biodiversity". *Trends in Ecology & Evolution* 27, no. 2 (): 78–84. ISSN: 01695347, visited on 07/11/2017. doi:10.1016/j.tree.2011.11.007. <http://linkinghub.elsevier.com/retrieve/pii/S0169534711003302>.

- Dmitriev, D.A., and M. Yoder. 2017. *NOMEN*. Published online. Visited on 07/22/2017. <https://github.com/SpeciesFileGroup/nomen>.
- Egloff, Willi, David Patterson, Donat Agosti, and Gregor Hagedorn. 2014. "Open exchange of scientific knowledge and European copyright: The case of biodiversity information". *ZooKeys* 414 (): 109–135. ISSN: 1313-2970, 1313-2989, visited on 04/13/2018. doi:10.3897/zookeys.414.7717. <http://zookeys.pensoft.net/articles.php?id=3830>.
- Filippova, Nina V., Ilya V. Filippov, Dmitry S. Schigel, Natalia V. Ivanova, and Maxim P. Shashkov. 2017. "Biodiversity informatics: global trends, national perspective and regional progress in Khanty-Mansi Autonomous Okrug". *Environmental Dynamics and Global Climate Change* 8, no. 2 (): 46–56. ISSN: 2541-9307, 2218-4422, visited on 03/15/2018. doi:10.17816/edgcc8246-56. <http://journals.eco-vector.com/EDGCC/article/view/7080>.
- Franz, N.M., and R.K. Peet. 2009. "Perspectives: Towards a language for mapping relationships among taxonomic concepts". *Systematics and Biodiversity* 7, no. 1 (): 5–20. ISSN: 1477-2000, 1478-0933, visited on 07/22/2017. doi:10.1017/S147720000800282X. <http://www.tandfonline.com/doi/abs/10.1017/S147720000800282X>.
- Franz, Nico M, and Beckett W Sterner. 2018. "To increase trust, change the social design behind aggregated biodiversity data". *Database* 2018.
- GBIF Secretariat. 2017. "GBIF Backbone Taxonomy. Checklist dataset". Published online. Visited on 06/30/2018. doi:10.15468/39omei. <https://doi.org/10.15468/39omei>.
- Godtsenhoven, Karen, van, Mikael Karstensen Elbaek, Barbara Sierman, Magchiel Bijsterbosch, Patrick Hochstenbach, Rosemary Russell, and Maurice Vanderfeesten. 2009. *Emerging Standards for Enhanced Publications and Repository Technology : Survey on Technology*. Amsterdam: Amsterdam University Press. ISBN: 978-90-8964-189-2, visited on 04/15/2018. doi:10.5117/9789089641892. <http://dare.uva.nl/aup/nl/record/316870>.
- Gruber, Thomas R. 1993. "A translation approach to portable ontology specifications". *Knowledge Acquisition* 5, no. 2 (): 199–220. ISSN: 10428143, visited on 08/07/2017. doi:10.1006/knac.1993.1008. <http://linkinghub.elsevier.com/retrieve/pii/S1042814383710083>.
- Harris, Larry R., Jeffrey M. Hill, Dayton Marcott, and Timothy F. Rochford. 1993. Knowledge base management system. Pat. US5228116A, **patentfiled** July 1993. <https://patents.google.com/patent/US5228116A/en>.
- Heath, Tom, and Christian Bizer. 2011. *Linked data: evolving the web into a global data space*. 1. ed. Synthesis lectures on the semantic web: theory and technology 1. OCLC: 732238828. San Rafael, Calif.: Morgan & Claypool. ISBN: 978-1-60845-430-3 978-1-60845-431-0.
- International Commission on Zoological Nomenclature. 1999. *International Code of Zoological Nomenclature*. Fourth Edition. London, UK: The International Trust for Zoological Nomenclature. ISBN: 0 85301 006 4.
- . 2017. *The Official Registry of Zoological Nomenclature*. Published online. Visited on 08/11/2017. <http://zoobank.org/>.

- International code of nomenclature for algae, fungi and plants (Melbourne code): adopted by the Eighteenth International Botanical Congress Melbourne, Australia, July 2011.* 2012. Regnum vegetabile v. 154. OCLC: ocn824722354. Königstein, Germany: Koeltz Scientific Books. ISBN: 978-3-87429-425-6.
- Jarke, Matthias, Bernd Neumann, Yannis Vassiliou, and Wolfgang Wahlster. 1989. “KBMS Requirements of Knowledge-Based Systems”. In *Foundations of Knowledge Base Management*, ed. by Michael L. Brodie, John Mylopoulos, Joachim W. Schmidt, Joachim W. Schmidt, and Constantino Thanos, 381–394. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN: 978-3-642-83399-1 978-3-642-83397-7, visited on 07/22/2018. doi:[10.1007/978-3-642-83397-7_17](https://doi.org/10.1007/978-3-642-83397-7_17). http://www.springerlink.com/index/10.1007/978-3-642-83397-7_17.
- Kennedy, Jessie B., Robert Kukla, and Trevor Paterson. 2005. “Scientific Names Are Ambiguous as Identifiers for Biological Taxa: Their Context and Definition Are Required for Accurate Data Integration”. In *Data Integration in the Life Sciences: Second International Workshop, DILS 2005, San Diego, CA, USA, July 20-22, 2005. Proceedings*, ed. by Bertram Ludäscher and Louiqa Raschid, 80–95. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN: 978-3-540-31879-8. doi:[10.1007/11530084_8](https://doi.org/10.1007/11530084_8).
- Kissling, W Daniel, Ramona Walls, Anne Bowser, Matthew O Jones, Jens Kattge, Donat Agosti, Josep Amengual, Alberto Basset, Peter M Van Bodegom, Johannes HC Cornelissen, and others. 2018. “Towards global data products of Essential Biodiversity Variables on species traits”. *Nature ecology & evolution*: 1.
- Koperski, M, M Sauer, W Braun, and S Gradstein. 2000. *Referenzliste der Moose Deutschlands*. Vol. 34. Schriftenreihe Vegetationsk.
- Kraker, Peter, Derick Leony, Wolfgang Reinhardt, N.A. Gü, and nter Beham. 2011. “The case for an open science in technology enhanced learning”. *International Journal of Technology Enhanced Learning* 3 (6): 643. ISSN: 1753-5255, 1753-5263, visited on 04/15/2018. doi:[10.1504/IJTEL.2011.045454](https://doi.org/10.1504/IJTEL.2011.045454). <http://www.inderscience.com/link.php?id=45454>.
- Li, Diyan, Tiandong Che, Binlong Chen, Shilin Tian, Xuming Zhou, Guolong Zhang, Miao Li, Uma Gaur, Yan Li, Majing Luo, and others. 2017. “Genomic data for 78 chickens from 14 populations”. *GigaScience* 6 (6): 1–5.
- Lin, Qiang, Ying Qiu, Ruobo Gu, Meng Xu, Jia Li, Chao Bian, Huixian Zhang, Geng Qin, Yanhong Zhang, Wei Luo, Jieming Chen, Xinxin You, Mingjun Fan, Min Sun, Pao Xu, Byrappa Venkatesh, Junming Xu, Hongtuo Fu, and Qiong Shi. 2017. “Draft genome of the lined seahorse, *Hippocampus erectus*” [inlangen]. *GigaScience* 6, no. 6 (): 1–6. ISSN: 2047-217X, visited on 03/15/2018. doi:[10.1093/gigascience/gix030](https://doi.org/10.1093/gigascience/gix030). <https://academic.oup.com/gigascience/article-lookup/doi/10.1093/gigascience/gix030>.
- Linnaeus, Carl von. 1758. *Systema naturae per regna tria naturae: secundum classes, ordines, genera, species, cum characteribus, differentiis, synonymis, locis. Volume 1*.
- Mallet, James. 2001. “Species, concepts of”. In *Encyclopedia of biodiversity*, 5:427–440. Published online. Visited on 07/11/2017. <http://tarjomefa.com/wp-content/uploads/2016/02/4420-english.pdf>.

- Mansinghka, Vikash, Richard Tibbetts, Jay Baxter, Pat Shafto, and Baxter Eaves. 2015. "BayesDB: A probabilistic programming system for querying the probable implications of data". ArXiv: 1512.05006, *arXiv:1512.05006 [cs]* (). Visited on 07/22/2018. <http://arxiv.org/abs/1512.05006>.
- Marwick, Ben, and Suzanne Birch. 2017. *A Standard for the Scholarly Citation of Archaeological Data as an Incentive to Data Sharing*. Tech. rep. SocArXiv. Visited on 03/15/2018. doi:10.17605/OSF.IO/PY4HZ. <https://osf.io/preprints/socarxiv/py4hz/>.
- Mathieu, Jérôme. 2018. "EGrowth: A global database on intraspecific body growth variability in earthworm". *Soil Biology and Biochemistry* 122:71–80.
- Michel, Franck, Catherine Faron-Zucker, Sandrine Terceirie, and Gargominy Olivier. 2018. "Modelling Biodiversity Linked Data: Pragmatism May Narrow Future Opportunities". *Biodiversity Information Science and Standards* 2:e26235.
- Michener, William K., James W. Brunt, John J. Helly, Thomas B. Kirchner, and Susan G. Stafford. 1997. "NONGEOSPATIAL METADATA FOR THE ECOLOGICAL SCIENCES". *Ecological Applications* 7, no. 1 (): 330–342. ISSN: 1051-0761, visited on 02/14/2018. doi:10.1890/1051-0761(1997)007[0330:NMFTE]2.0.CO;2. [http://doi.wiley.com/10.1890/1051-0761\(1997\)007\[0330:NMFTE\]2.0.CO;2](http://doi.wiley.com/10.1890/1051-0761(1997)007[0330:NMFTE]2.0.CO;2).
- Mietchen, Daniel. 2014. "The Transformative Nature of Transparency in Research Funding". *PLoS Biology* 12, no. 12 (): e1002027. ISSN: 1545-7885, visited on 04/15/2018. doi:10.1371/journal.pbio.1002027. <http://dx.plos.org/10.1371/journal.pbio.1002027>.
- Milano, Filippo, Paolo Pantini, Stefano Mammola, and Marco Isaia. 2017. "LA CONSERVAZIONE DELL'ARANEOFAUNA IN ITALIA E IN EUROPA". *ATTI DELL'ACCADEMIA NAZIONALE ITALIANA DI ENTOMOLOGIA. RENDICONTI* 65:91–103.
- Miller, Jeremy, Torsten Dikow, Donat Agosti, Guido Sautter, Terry Catapano, Lyubomir Penev, Zhi-Qiang Zhang, Dean Pentcheff, Richard Pyle, Stan Blum, Cynthia Parr, Chris Freeland, Tom Garnett, Linda S Ford, Burgert Muller, Leo Smith, Ginger Strader, Teodor Georgiev, and Laurence Bénichou. 2012. "From taxonomic literature to cybertaxonomic content". *BMC Biology* 10 (1): 87. ISSN: 1741-7007, visited on 04/13/2018. doi:10.1186/1741-7007-10-87. <http://bmcbiol.biomedcentral.com/articles/10.1186/1741-7007-10-87>.
- Miller, Jeremy, Donat Agosti, Lyubomir Penev, Guido Sautter, Teodor Georgiev, Terry Catapano, David Patterson, David King, Serrano Pereira, Rutger Vos, and Soraya Sierra. 2015. "Integrating and visualizing primary data from prospective and legacy taxonomic literature". *Biodiversity Data Journal* 3 (): e5063. ISSN: 1314-2828, 1314-2836, visited on 04/13/2018. doi:10.3897/BDJ.3.e5063. <http://bdj.pensoft.net/articles.php?id=5063>.
- Mindell, David P., Brian L. Fisher, Peter Roopnarine, Jonathan Eisen, Georgina M. Mace, Roderic D. M. Page, and Richard L. Pyle. 2011. "Aggregating, Tagging and Integrating Biodiversity Research". Ed. by Sean A. Rands. *PLoS ONE* 6, no. 8 (): e19491. ISSN: 1932-6203, visited on 04/13/2018. doi:10.1371/journal.pone.0019491. <http://dx.plos.org/10.1371/journal.pone.0019491>.

- Momtchev, Vassil, Deyan Peychev, Todor Primov, and Georgi Georgiev. 2009. “Expanding the pathway and interaction knowledge in linked life data”. *Proc. of International Semantic Web Challenge*. Visited on 07/22/2017. <http://challenge.semanticweb.org/documents/Linked%20Life%20Data-LLD%20semantic%20web%20challenge%202009.pdf>.
- Mulsant, E. 1866. “Monographie des Coccinellides”. *Mém. L’Acad. Imp. Lyon, Cl. Sci., LMém. L’Acad. Imp. Lyon*. 15:1–112.
- Neo4J Developers. 2012. *Neo4J Graph NoSQL Database*. Published online. <http://neo4j.com>.
- Obitko, Marek. 2007. “Translations between ontologies in multi-agent systems”. Ph. D. Dissertation, , Czech Technical University, Faculty of Electrical Engineering.
- Ototext. 2018. *GraphDB 8.6*. Published online. Visited on 07/17/2018. <http://graphdb.ontotext.com/>.
- Ordynets, Alexander, Anton Savchenko, Alexander Akulov, Eugene Yurchenko, Vera Malysheva, Urmaz Kõljalg, Josef Vlasák, Karl-Henrik Larsson, and Ewald Langer. 2017. “Aphylophoroid fungi in insular woodlands of eastern Ukraine”. *Biodiversity Data Journal* 5 (): e22426. ISSN: 1314-2828, 1314-2836, visited on 03/15/2018. doi:10.3897/BDJ.5.e22426. <https://bdj.pensoft.net/articles.php?id=22426>.
- Page, R. D. M. 2008. “Biodiversity informatics: the challenge of linking data and the role of shared identifiers”. *Briefings in Bioinformatics* 9, no. 5 (): 345–354. ISSN: 1467-5463, 1477-4054, visited on 04/13/2018. doi:10.1093/bib/bbn022. <https://academic.oup.com/bib/article-lookup/doi/10.1093/bib/bbn022>.
- . 2015. *Putting some bite into the Bouchout Declaration*. Published online. <http://iphylo.blogspot.bg/2015/05/putting-some-bite-into-bouchout.html>.
- . 2014. *The vision thing - it’s all about the links*. Published online. <http://iphylo.blogspot.bg/2014/06/the-vision-thing-it-all-about-links.html>.
- Page, Roderic D. M. 2006. “Taxonomic names, metadata, and the Semantic Web”. *Biodiversity Informatics* 3, no. 0 (). ISSN: 15469735, visited on 04/13/2018. doi:10.17161/bi.v3i0.25. <https://journals.ku.edu/index.php/jbi/article/view/25>.
- Page, Roderic DM. 2018a. “Liberating links between datasets using lightweight data publishing: an example using plant names and the taxonomic literature”. Published online, *bioRxiv*: 343996.
- . 2018b. “Ozymandias: A biodiversity knowledge graph”. Published online, *bioRxiv*: 485854.
- Parr, Cynthia S., Robert Guralnick, Nico Cellinese, and Roderic D.M. Page. 2012. “Evolutionary informatics: unifying knowledge about the diversity of life” [inlangen]. *Trends in Ecology & Evolution* 27, no. 2 (): 94–103. ISSN: 01695347, visited on 04/13/2018. doi:10.1016/j.tree.2011.11.001. <http://linkinghub.elsevier.com/retrieve/pii/S0169534711003247>.
- Patterson, D.J., J. Cooper, P.M. Kirk, R.L. Pyle, and D.P. Remsen. 2010. “Names are key to the big new biology”. *Trends in Ecology & Evolution* 25, no. 12 (): 686–691. ISSN: 01695347, visited on 07/11/2017. doi:10.1016/j.tree.2010.09.004. <http://linkinghub.elsevier.com/retrieve/pii/S0169534710002181>.

- Patterson, David J., David Remsen, William A. Marino, and Cathy Norton. 2006. "Taxonomic Indexing—Extending the Role of Taxonomy". Ed. by Rod Page. *Systematic Biology* 55, no. 3 (): 367–373. ISSN: 1076-836X, 1063-5157, visited on 11/16/2018. doi:10.1080/10635150500541680. <https://academic.oup.com/sysbio/article/55/3/367/1667279>.
- Pellissier Tanon, Thomas, Denny Vrandečić, Sebastian Schaffert, Thomas Steiner, and Lydia Pintscher. 2016. "From Freebase to Wikidata: The Great Migration", 1419–1428. ACM Press. ISBN: 978-1-4503-4143-1, visited on 07/22/2018. doi:10.1145/2872427.2874809. <http://dl.acm.org/citation.cfm?doid=2872427.2874809>.
- Penev, Lyubomir, Teodor Georgiev, Peter Geshev, Seyhan Demirov, Viktor Senderov, Iliyana Kuzmova, Iva Kostadinova, Slavena Peneva, and Pavel Stoev. 2017a. "ARPHA-BioDiv: A toolbox for scholarly publication and dissemination of biodiversity data based on the ARPHA Publishing Platform". *Research Ideas and Outcomes* 3, no. e13088 (). ISSN: 2367-7163. doi:10.3897/rio.3.e13088. <http://riojournal.com/articles.php?id=13088>.
- Penev, Lyubomir, W. John Kress, Sandra Knapp, De-Zhu Li, and Susanne Renner. 2010a. "Fast, linked, and open – the future of taxonomic publishing for plants: launching the journal PhytoKeys". *PhytoKeys* 1, no. 0 (). ISSN: 1314-2003, 1314-2011, visited on 07/22/2017. doi:10.3897/phytokeys.1.642. http://www.pensoft.net/journal_home_page.php?journal_id=3&page=article&type=show&article_id=642&abstract=1.
- Penev, Lyubomir, Terence Catapano, Donat Agosti, Teodor Georgiev, Guido Sautter, and Pavel Stoev. 2012. "Implementation of TaxPub, an NLM DTD extension for domain-specific markup in taxonomy, from the experience of a biodiversity publisher". Published online. Visited on 07/11/2017. <https://www.ncbi.nlm.nih.gov/books/NBK100351/>.
- Penev, Lyubomir, Donat Agosti, Teodor Georgiev, Terry Catapano, Jeremy Miller, Vladimir Blagoderov, David Roberts, Vincent Smith, Irina Brake, Simon Rycroft, Ben Scott, Norman Johnson, Robert Morris, Guido Sautter, Vishwas Chavan, Tim Robertson, David Remsen, Pavel Stoev, Cynthia Parr, Sandra Knapp, W. John Kress, Frederic Thompson, and Terry Erwin. 2010b. "Semantic tagging of and semantic enhancements to systematics papers: ZooKeys working examples". *ZooKeys* 50 (): 1–16. ISSN: 1313-2970, 1313-2989, visited on 04/15/2018. doi:10.3897/zookeys.50.538. <http://zookeys.pensoft.net/articles.php?id=2215>.
- Penev, Lyubomir, Daniel Mietchen, Vishwas Chavan, Gregor Hagedorn, Vincent Smith, David Shotton, Éamonn Ó Tuama, Viktor Senderov, Teodor Georgiev, Pavel Stoev, Quentin Groom, David Remsen, and Scott Edmunds. 2017b. "Strategies and guidelines for scholarly publishing of biodiversity data". *Research Ideas and Outcomes* 3, no. e12431 (). ISSN: 2367-7163. doi:10.3897/rio.3.e12431. <http://riojournal.com/articles.php?id=12431>.
- Penev, Lyubomir, Christopher Lyal, Anna Weitzman, David Morse, David King, Guido Sautter, Teodor Georgiev, Robert Morris, Terry Catapano, and Donat Agosti. 2011. "XML schemas and mark-up practices of taxonomic literature". *ZooKeys* 150 (): 89–116. ISSN: 1313-2970, 1313-2989, visited on 07/04/2018. doi:10.3897/zookeys.150.2213. <http://zookeys.pensoft.net/articles.php?id=3038>.
- Peroni, Silvio. 2014. "The semantic publishing and referencing ontologies". In *Semantic Web Technologies and Legal Scholarly Publishing*, 1st ed., 15:121–193. Springer. Visited on 07/22/2017.

- Peroni, Silvio, and David Shotton. 2012. "FaBiO and CiTO: Ontologies for describing bibliographic resources and citations". *Web Semantics: Science, Services and Agents on the World Wide Web* 17 (): 33–43. ISSN: 15708268, visited on 08/13/2017. doi:10.1016/j.websem.2012.08.001. <http://linkinghub.elsevier.com/retrieve/pii/S1570826812000790>.
- Poorani, J., and Roger Booth. 2016. "Harmonia manillana (Mulsant), a new addition to Indian Coccinellidae, with changes in synonymy". *Biodiversity Data Journal* 4 (): e8030. ISSN: 1314-2828, 1314-2836, visited on 08/13/2017. doi:10.3897/BDJ.4.e8030. <http://bdj.pensoft.net/articles.php?id=8030>.
- Pyle, Richard. 2016. "Towards a Global Names Architecture: The future of indexing scientific names". *ZooKeys* 550 (): 261–281. ISSN: 1313-2970, 1313-2989, visited on 08/12/2017. doi:10.3897/zookeys.550.10009. <http://zookeys.pensoft.net/articles.php?id=6241>.
- Pyron, Robert Alexander. 2018. "A 21st Century Vision for Neotropical Snake Systematics". *Revista Latinoamericana de Herpetología* 1 (1): 58–62.
- R Core Team. 2016. *R: A Language and Environment for Statistical Computing*. Published online, Vienna, Austria. <https://www.R-project.org/>.
- Rebholz-Schuhmann, Dietrich, Harald Kirsch, and Francisco Couto. 2005. "Facts from text—is text mining ready to deliver?" *PLoS biology* 3 (2): e65. Visited on 07/22/2017. <http://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.0030065>.
- Sarah Faulwetter, Evangelos Pafilis, Lucia Fanini, Nicolas Bailly, Donat Agosti, Christos Arvanitidis, Laura Boicenco, Terry Catapano, Simon Claus, Stefanie Dekeyzer, Teodor Georgiev, Aglaia Legaki, Dimitra Mavraki, Anastasis Oulas, Gabriella Papastefanou, Lyubomir Penev, Guido Sautter, Dmitry Schigel, Viktor Senderov, Adrian Teaca, and Marilena Tsompanou. 2016. "EMODnet Workshop on mechanisms and guidelines to mobilise historical data into biogeographic databases". *Research Ideas and Outcomes* 2, no. e10445 (). ISSN: 2367-7163. doi:10.3897/rio.2.e10445. <http://rio.pensoft.net/articles.php?id=10445>.
- Senderov, Viktor, and Lyubomir Penev. 2016. "The Open Biodiversity Knowledge Management System in Scholarly Publishing". *Research Ideas and Outcomes* 2, no. e7757 (). ISSN: 2367-7163. doi:10.3897/rio.2.e7757. <http://rio.pensoft.net/articles.php?id=7757>.
- Senderov, Viktor, Teodor Georgiev, and Lyubomir Penev. 2016. "Online direct import of specimen records into manuscripts and automatic creation of data papers from biological databases". *Research Ideas and Outcomes* 2 (): e10617. ISSN: 2367-7163. doi:10.3897/rio.2.e10617. <http://rio.pensoft.net/articles.php?id=10617>.
- Senderov, Viktor, Kiril Simov, Nico Franz, Pavel Stoev, Terry Catapano, Donat Agosti, Guido Sautter, Robert A. Morris, and Lyubomir Penev. 2018. "OpenBiodiv-O: ontology of the OpenBiodiv knowledge management system". *Journal of Biomedical Semantics* 9, no. 5 (). ISSN: 2041-1480. doi:10.1186/s13326-017-0174-5. <https://jbiomedsem.biomedcentral.com/articles/10.1186/s13326-017-0174-5>.
- Senderov, Viktor, Nico M. Franz, and Kiril Simov. 2017. *OpenBiodiv Ontology and Guide*. Published online. Visited on 08/09/2017. <http://openbiodiv.net/ontology>.

- Shotton, David. 2009. "Semantic publishing: the coming revolution in scientific journal publishing". *Learned Publishing* 22, no. 2 (): 85–94. ISSN: 09531513, visited on 04/15/2018. doi:10.1087/2009202. <http://doi.wiley.com/10.1087/2009202>.
- Singhal, Amit. 2012. *Introducing the knowledge graph: things, not strings*. Published online. Visited on 01/11/2019. <https://googleblog.blogspot.com/2012/05/introducing-knowledge-graph-things-not.html>.
- Staab, Steffen, and Rudi Studer, eds. 2009. *Handbook on Ontologies*. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN: 978-3-540-70999-2 978-3-540-92673-3, visited on 08/07/2017. doi:10.1007/978-3-540-92673-3. <http://link.springer.com/10.1007/978-3-540-92673-3>.
- Sterner, Beckett, and Nico M. Franz. 2017. "Taxonomy for Humans or Computers? Cognitive Pragmatics for Big Data". *Biological Theory* 12, no. 2 (): 99–111. ISSN: 1555-5542, 1555-5550, visited on 07/11/2017. doi:10.1007/s13752-017-0259-5. <http://link.springer.com/10.1007/s13752-017-0259-5>.
- Taxonomic Names and Concepts Interest Group. 2006. *Taxonomic Concept Transfer Schema (TCS), version 1.01*. Published online. Visited on 01/11/2019. <http://www.tdwg.org/standards/117>.
- Tennant, Jonathan P, Jonathan M Dugan, Daniel Graziotin, Damien C Jacques, François Waldner, Daniel Mietchen, Yehia Elkhatib, Lauren B Collister, Christina K Pikas, Tom Crick, and others. 2017. "A multi-disciplinary perspective on emergent and future innovations in peer review". *F1000Research* 6.
- The Bouchout Declaration for Open Biodiversity Knowledge Management*. 2014. Published online. Visited on 01/11/2019. http://www.bouchoutdeclaration.org/fileadmin/Dateien/PDF/Bouchout_Declaration_EN.pdf.
- Tzitzikas, Yannis, Carlo Allocca, Chryssoula Bekiari, Yannis Marketakis, Pavlos Fafalios, Martin Doerr, Nikos Minadakis, Theodore Patkos, and Leonardo Candela. 2013. "Integrating heterogeneous and distributed information about marine species through a top level ontology". In *Research Conference on Metadata and Semantic Research*, 289–301. Springer. Visited on 07/22/2017. http://link.springer.com/chapter/10.1007/978-3-319-03437-9_29.
- University of Copenhagen, University of Turku, Institute for Systematic Zoology and Evolutionary Biology, Zoologisches Forschungsmuseum Alexander Koenig, Naturhistorisches Museum Wien, Pensoft Publishers company, ERA7 Bioinformatics, Swedish Museum of Natural History, and Decuria IT company. 2014. *BIG4 - Biosystematics, Informatics and Genetics of the big 4 insect groups: training tomorrow's researchers and entrepreneurs*. Published online. Visited on 01/11/2019. <http://big4-project.eu>.
- Vrandečić, Denny, and Markus Krötzsch. 2014. "Wikidata: a free collaborative knowledgebase". *Communications of the ACM* 57, no. 10 (): 78–85. ISSN: 00010782, visited on 07/22/2018. doi:10.1145/2629489. <http://dl.acm.org/citation.cfm?doid=2661061.2629489>.
- Was ist Open Science?* Published online. Visited on 01/11/2019. <http://openscienceasap.org/open-science/>.
- What is GBIF?* Published online. Visited on 08/12/2017. <http://www.gbif.org/what-is-gbif>.
- Wickham, Hadley. 2017. *httr: Tools for Working with URLs and HTTP*. Published online. Visited on 07/17/2018. <https://cran.r-project.org/web/packages/httr/>.

- Wickham, Hadley, James Hester, and Jeroen Ooms. 2018. *xml2: Parse XML*. Published online. Visited on 07/17/2018. <https://cran.r-project.org/web/packages/xml2/index.html>.
- Wieczorek, John, David Bloom, Robert Guralnick, Stan Blum, Markus Döring, Renato Giovanni, Tim Robertson, and David Vieglais. 2012. “Darwin Core: An Evolving Community-Developed Biodiversity Data Standard”. Ed. by Indra Neil Sarkar. *PLoS ONE* 7, no. 1 (): e29715. ISSN: 1932-6203, visited on 07/22/2017. doi:10.1371/journal.pone.0029715. <http://dx.plos.org/10.1371/journal.pone.0029715>.
- Williams, Antony J., Lee Harland, Paul Groth, Stephen Pettifer, Christine Chichester, Egon L. Willighagen, Chris T. Evelo, Niklas Blomberg, Gerhard Ecker, Carole Goble, and Barend Mons. 2012. “Open PHACTS: semantic interoperability for drug discovery”. *Drug Discovery Today* 17, numbers 21-22 (): 1188–1198. ISSN: 13596446, visited on 07/22/2017. doi:10.1016/j.drudis.2012.05.016. <http://linkinghub.elsevier.com/retrieve/pii/S1359644612001936>.
- Witteveen, Joeri. 2015. “Naming and contingency: the type method of biological taxonomy”. *Biology & Philosophy* 30, no. 4 (): 569–586. ISSN: 0169-3867, 1572-8404, visited on 08/13/2017. doi:10.1007/s10539-014-9459-6. <http://link.springer.com/10.1007/s10539-014-9459-6>.
- Wolfram|Alpha, *Making the world's knowledge computable*. Published online. Wolfram Alpha LLC. Visited on 06/10/2018. <https://www.wolframalpha.com/>.
- pro-iBiosphere project final report*. 2014. Published online. Visited on 04/13/2018. http://adm.pro-ibiosphere.eu/getatt.php?filename=oo_4751.pdf.
- pro-iBiosphere*. Visited on 08/12/2017. <http://wiki.pro-ibiosphere.eu/>.
- Филиппова, НВ, ИВ Филиппов, ДС Щигель, НВ Иванова, and МП Шашков. 2017. “Информатика биоразнообразия: мировые тенденции, состояние дел в России и развитие направления в Ханты-Мансийском Автономном Округе”. *Динамика окружающей среды и глобальные изменения климата* 8 (2): 46–56.
- Шашков, МП, ИФ Чадин, and НВ Иванова. 2017. “Методические рекомендации по стандартизации данных для публикации через глобальный портал GBIF.ORG и подготовке статьи о данных”. *Труды Кольского научного центра РАН*, no. 6-5 (8).
- Шашков, МП, and НВ Иванова. 2018. “Стандарты и веб-инструменты для публикации данных через глобальные порталы по биоразнообразию”. *Доклады Международной конференции “Математическая биология и биоинформатика”* 7:e98. doi:10.17537/icmbb18.55.

Abstracts of Dissertations

Number 2, 2020

INSTITUTE OF INFORMATION AND COMMUNICATION TECHNOLOGIES
BULGARIAN ACADEMY OF SCIENCES

БЪЛГАРСКА АКАДЕМИЯ НА НАУКИТЕ

ИНСТИТУТ ПО ИНФОРМАЦИОННИ И КОМУНИКАЦИОННИ ТЕХНОЛОГИИ

Брой 2, 2020

Автореферати на дисертации