

Ontologies and Languages for Representing Mathematical Knowledge on the Semantic Web

Editor(s): Aldo Gangemi, ISTC-CNR Rome, Italy

Solicited review(s): Claudio Sacerdoti Coen, University of Bologna, Italy; Alexandre Passant, DERI, National University of Galway, Ireland; Aldo Gangemi, ISTC-CNR Rome, Italy

Christoph Lange

*FB 3 (Mathematics and Computer Science),
University of Bremen, Germany
Computer Science, Jacobs University Bremen,
Germany
E-mail: ch.lange@jacobs-university.de*

Abstract. Mathematics is a ubiquitous foundation of science, technology, and engineering. Specific areas of mathematics, such as numeric and symbolic computation or logics, enjoy considerable software support. Working mathematicians have recently started to adopt Web 2.0 environments, such as blogs and wikis, but these systems lack machine support for knowledge organization and reuse, and they are disconnected from tools such as computer algebra systems or interactive proof assistants. We argue that such scenarios will benefit from Semantic Web technology.

Conversely, mathematics is still underrepresented on the Web of [Linked] Data. There are mathematics-related Linked Data, for example statistical government data or scientific publication databases, but their *mathematical* semantics has not yet been modeled. We argue that the services for the Web of Data will benefit from a deeper representation of mathematical knowledge.

Mathematical knowledge comprises structures given in a logical language – formulae, statements (e.g. axioms), and theories –, a mixture of rigorous natural language and symbolic notation in documents, application-specific metadata, and discussions about conceptualizations, formalizations, proofs, and (counter-)examples. Our review of vocabularies for representing these structures covers ontologies for mathematical problems, proofs, interlinked scientific publications, scientific discourse, as well as mathematical meta-

data vocabularies and domain knowledge from pure and applied mathematics.

Many fields of mathematics have not yet been implemented as proper Semantic Web ontologies; however, we show that MathML and OpenMath, the standard XML-based exchange languages for mathematical knowledge, can be fully integrated with RDF representations in order to contribute existing mathematical knowledge to the Web of Data.

We conclude with a roadmap for getting the mathematical Web of Data started: what datasets to publish, how to interlink them, and how to take advantage of these new connections.

Keywords: mathematics, mathematical knowledge management, ontologies, knowledge representation, formalization, linked data, XML

1. Introduction: Mathematics on the Web – State of the Art and Challenges

A review of the state of the art of mathematics on the Web has to acknowledge traditional web sites that are in day to day use: review and abstract services such as Zentralblatt MATH [34] and MathSciNet [41], the arXiv pre-print server [46], libraries of formalized and machine-verified mathematical content such as the Mizar Mathematical Library (MML [9]), and reference works such as the Digital Library of Mathematical Functions (DLMF [3]) or Wolfram MathWorld [8].

These sites have facilitated the *access* to mathematical knowledge. However, (i) they offer a limited degree of interaction and do not facilitate collaboration, and (ii) the means of automatically retrieving, using, and adaptively presenting knowledge through automated agents are restricted. Concerning the Web in general, problem (i) has been addressed by Web 2.0 applica-

tions, and problem (ii) by the Semantic Web. This section reviews to what extent these developments have been adopted for mathematical applications and suggests a new combination of Web 2.0 and Semantic Web to overcome the remaining problems.

1.1. How Working Mathematicians have Embraced Web 2.0 Technology

An increasing number of working mathematicians has recently started to use Web 2.0 technology for collaboratively developing new ideas, but also as a new publication channel for established knowledge. Research blogs and wiki encyclopedias are typical representatives.

1.1.1. Research Blogs and Wiki Encyclopedias

As a publication channel for established knowledge, blogs reflect the traditional mathematical practice of publishing short reviews of previously published material, as done by the review and abstract services mentioned initially. Researchers have also found blogs useful to gather early feedback about preliminary findings, whereas traditional publications, even in the state submitted for peer review, would rather present mature ideas while hiding alternatives that had once been considered and then discarded – a working method that effectively hinders collaboration outside of small research groups. Successful collaborations among mathematicians not knowing each other before have started in blogs and converged into conventional articles [55]. GOWERS, an active blogger, initiated the successful Polymath series, where blogs are the exclusive communication medium for proving theorems in a massive collaborative effort [14,121,58], including the recent collaborative review of a claimed but wrong proof of $P \neq NP$ [15]. Compared to research blogs, the MathOverflow forum [7], where users can post their problems and solutions to others' problems, offers more instant help with smaller problems. By its reputation mechanism, it acts as an agile simulation of the traditional scientific publication and peer review processes.

For evolving ideas emerged from a blog discussion, or for creating permanent, short, interlinked descriptions of topics, wikis have been found more appropriate. The nLab wiki [32], a companion to the n-Category Café blog [31], is a prominent example for that, and also for the emerging practice of *Open Notebook Science*, i.e. “making the entire primary record of a research project public”, including “failed, less significant, and otherwise unpublished experi-

ments” [228]. The Polymath maintainers have also set up a companion wiki for “collect[ing] pertinent background information which was no longer part of the active ‘foreground’ of exchanges on the [...] blog entries” [58]. Finally, where MathOverflow focuses on concrete problems and solution, the Tricky [25], also initiated by GOWERS, is a wiki repository of general mathematical techniques – reminiscent of a Web 2.0 remake of PÓLYA’s classic “How to Solve It” [194].

Wikis that collect *existing* mathematical knowledge, for educational and general purposes, are more widely known. PlanetMath [192], counting more than 8,000 entries at the time of this writing, is a mathematical encyclopedia. The general-purpose Wikipedia with 15 million articles in over 250 languages also covers mathematics [227]. Targeting a general audience, it omits most formal proofs but embeds the pure mathematical knowledge into a wider context, including, e.g., the history of mathematics, biographies of mathematicians, and information about application areas. The lack of proofs is partly compensated by linking to the technically similar ProofWiki [17], containing over 2,500 proofs, or to PlanetMath. Finally, Connexions [86], technically driven by a more traditional content management system, is an open web repository specialized on courseware. Connexions promotes the contribution of small, reusable course modules – more than 17,000, about 4,000 from mathematics and statistics, and about 6,000 from science and technology – to its *content commons*, so that the original author, but also others can flexibly combine them into collections, such as the notes for a particular course.

The wikis mentioned so far have been set up from scratch, hardly reusing content from existing knowledge bases, but maintainers of established knowledge bases are also starting to employ Web 2.0 frontends – for example the recently developed prototypical wiki frontend for the Mizar Mathematical Library (MML) [39,216], a large library of formalized and machine-verified mathematical content. The wiki intends to support common workflows in enhancing and maintaining the MML and thus to disburden the human library committee.

1.1.2. Critique – Little Reuse, Lack of Services

Web 2.0 sites facilitate collaboration but still require a massive investment of manpower for compiling a knowledge collection. Machine-supported intelligent knowledge reuse, e.g. from other knowledge collections on the Web, does not take place. Different knowledge bases are technically separated from each other

by using document formats that are merely suitable for knowledge presentation but not for representation, such as XHTML with \LaTeX formulae. The only way of referring to other knowledge bases is an untyped hyperlink. The proof techniques collected in the Tricky cannot be automatically applied to a problem developed in a research blog, as neither of them are sufficiently formalized. Conversely, the Polymath community does not have any automated verification tools at hand but exclusively relies on the crowdsourcing principle that “given enough eyeballs, all bugs [here: errors in a proof] are shallow” [226].¹

Intelligent information retrieval, a prerequisite for finding knowledge to reuse and to apply, is poorly supported. For example, Wikipedia states the Pythagorean theorem as $a^2 + b^2 = c^2$ and files it into the categories “Articles containing proofs” and “Mathematical theorems” [229]. The \LaTeX representation of the formulae does not support search by functional structure. Putting aside the fact that Wikipedia cannot search formulae at all, a search for equivalent expressions such as $x^2 + y^2 = z^2$ or $c = \sqrt{a^2 + b^2}$ would not yield the theorem, unless they explicitly occur in the article.² From the categorization it is neither clear for a machine (albeit very likely for a human) whether the article contains a proof of *that* theorem, nor whether it is correct. Similarly, the Polymath collaborators had to search previous publications of refutations of $P \neq NP$ “proofs” by keyword.

Formalized repositories such as the MML use specialized search engines [56]. While they support *internal* knowledge reuse by formalizing new mathematical concepts of existing ones and proving new theorems by applying ones that have already been proven, they do not support links to external repositories. Thus, the maintainers of each knowledge collection, informal or formalized, hope to receive a critical mass of contributions that makes it sufficiently self-contained for the desired application.

Finally, the integration of mathematical Web 2.0 sites with automated reasoning and computation services is scarce. Interactive computation is available in mathematical e-learning systems, such as ActiveMath [36] or MathDox [174] – where document authors have sufficiently formalized the underlying mathematics in separate editing tools before publishing –, but less so in general-purpose digital libraries and collaboration environments. Mashups, which have otherwise been a driving force of Web 2.0 development, scarcely exist for mathematical tasks.³

1.2. Early Adoption of the Semantic Web in Mathematical Knowledge Management on the Semantic Web

In the early 2000s, when XML was increasingly used for mathematics, particularly for formulae (cf. section 4.1.1), the first building blocks of the Semantic Web vision, such as RDFS, approached standardization. These developments sparked interest in the emerging interdisciplinary mathematical knowledge management (MKM) community, consisting of computer scientists, computer-savvy mathematicians, and digital library researchers, whose objective is “to develop new and better ways of managing mathematical knowledge using sophisticated software tools” [114]⁴, or, more specifically, “to serve (i) mathematicians, scientists, and engineers who produce and use mathematical knowledge; (ii) educators and students who teach and learn mathematics; (iii) publishers who offer mathematical textbooks and disseminate new mathematical results; and (iv) librarians and mathematicians who catalog and organize mathematical knowledge” [114]⁵. They hoped that Semantic Web technologies would help to address their challenges. This seemed technically feasible, particularly as both communities made use of XML as a serialization format and URIs for identifying things [170].

The two main lines of applying Semantic Web technologies to MKM focused on *digital libraries* – improving information retrieval and giving readers access to automated reasoning and computation services –, and *web services* – providing self-describing interfaces to automated reasoning and computation on the Web, so that they could solve problems sent to them by humans or other agents.

1.2.1. Digital Libraries – MathNet, HELM, and their Spin-Offs

Mathematical institutes participating in MathNet [6], an early effort to build “a distributed, efficient and user-driven information and communication system for mathematics” [93], were advised to put up uniformly structured homepages, to publish preprints, and to annotate both with RDF. Some of the 180 MathNet homepages that existed in 2002 [208] are still online; however, the central services, including a preprint search engine⁶ and a browser for MathNet pages, have been out of order since 2007.

Independently, HELM, the Hypertextual Electronic Library of Mathematics [5,49], aimed at “integrat[ing] the current tools for the automation of formal rea-

soning and the mechanization of mathematics [...] with the most recent technologies for the development of web applications and electronic publishing” [5]. In contrast to MathNet and other traditional digital libraries, HELM intended to explicitly represent the fine-grained structures of mathematical expressions to expose them, e.g., to automated reasoners, but also to enrich their publication on the Web. For example, mathematical formulae were rendered in MathML in such a way that actions could be performed on them, e.g. simplifying a selected (sub)expression using an automated reasoning backend attached to the library. HELM completely relied on XML and RDF not only for publishing, but also for its internal knowledge representation. Formalizations of mathematical statements and proofs were encoded in one XML dialect per underlying logical system; primarily, the library of the Coq higher-order proof assistant (cf. section 4.2) was used in HELM. Relevant structural properties, interrelations, and metadata were represented in RDF (cf. section 4.3).

The HELM developers had to carry out a lot of foundational research and development, as suitable reusable implementations were not available for many of the planned features. As none of the prototypical RDF query engines available in 2003 satisfied the HELM requirements⁷, a new one was developed [131, 129]. As browsers did not sufficiently support MathML, a MathML rendering widget suitable for embedding into desktop applications was developed [189].

1.2.2. *Web Services – MONET and Related Architectures*

The MONET project pioneered an architecture for mathematical web services built on Semantic Web technologies [179,81]. MONET services give access to numeric and symbolic computation systems; access to proof assistants or digital libraries was envisaged but not pursued. MONET services come with a machine-comprehensible description and can be registered with a central broker. Mathematical expressions in queries or computation requests to the broker were represented by their functional structure using OpenMath (cf. section 4.3). MONET also required foundational work to be done. OWL and the RACER reasoner were already found suitable for the internal description of services and problems and computing matches. However, the frontend XML languages for service descriptions and queries (which the broker then translated to OWL) had to be designed from scratch. Furthermore, the OWL reasoners of that time could not

efficiently deal with a large number of instances (here: concrete problems instantiating problem descriptions), which required a specific database/reasoner hybrid to be developed, but then, again, the separate treatment of classes and instances constrained the design of the MONET ontologies in that they had to model every object as a class [82]. Part of MONET’s query language is still used in the MathDox e-learning system [174,88]. More importantly, MONET and the competing MathBroker architecture for symbolic computation web services [57] influenced each other. The latter, however, made less use of Semantic Web service technologies. The MathServe architecture, influenced by both of the former but focusing on automated reasoning, made extensive use of more recent Semantic Web service technologies, such as OWL-S service profiles [232].

1.2.3. *Critique – Frustration and Discontinuation*

Semantic Web approaches to MKM have so far failed to fulfill the hopes set in them, the aftermath of HELM and MONET being an instructive example. Both groups of researchers were initially enthusiastic about the possibilities of the emerging Semantic Web, but then it turned out that few stable and reusable implementations existed, and hence a considerable amount of resources had to be invested into developing fundamental building blocks.⁸ The application of Semantic Web technology to MKM was stopped even before solutions had matured to an extent that would have allowed for, e.g., a wide deployment to working mathematicians and usability studies with them.

After 2004, when the HELM and MONET activities had ended, the MKM community has given up using Semantic Web technologies on a large scale, and the Semantic Web community has focused on different application areas. It has been suggested that, after the pioneering phase, it was hard to obtain research funding for applications of Semantic Web technologies to MKM.⁹

From the MKM perspective, the further development was as follows: Parts of the HELM technology have survived in an interactive desktop proof assistant [50], whereas the web frontend and the RDF-based components have been discontinued. Contemporary mathematical web services work without Semantic Web technology. While large parts of the influential OpenMath community had been involved into MONET, which heavily relied on Semantic Web technologies, the current driving force of research symbolic computation web services, the SCIENCE project (Sym-

bolic Computation Infrastructure for Europe [203]), does not use “standard” Semantic Web service technologies at all: SCSCP (Symbolic Computation Software Composability Protocol [133]) is a lightweight XML protocol using TCP sockets, or alternatively SOAP, whose communication semantics heavily relies on a custom OpenMath vocabulary.

1.3. Mathematics on the Semantic Web – Why Retry Now?

Web 2.0 applications are attracting an increasing number of working mathematicians. The usage of Semantic Web technologies to improve MKM has been investigated, albeit without becoming mainstream yet. This section argues why a new combination of Web 2.0 and Semantic Web technologies is needed to address them, and why such a solution is now feasible.

1.3.1. Combining Sem. Web and Web 2.0 for MKM

The combination of Web 2.0 and Semantic Web technology has already proven successful in some fields, including semantic wikis and Linked Data mashups [43]; however, it has hardly been applied to MKM yet. From the point of view that mathematicians are already using the Web 2.0, it seems feasible to incrementally enrich such existing applications with Semantic Web technology without scaring users away – an approach that has succeeded with general-purpose systems such as Semantic MediaWiki [22] or the RDF-enabled Drupal 7 [91], which are now mainstream. A similar rewrite of the system powering the PlanetMath encyclopedia (cf. section 1.1.1) is currently in progress [152]. From the point of view that applying Semantic Web to MKM had been tried without success before the Web 2.0 era, ZACCHIROLI gave two reasons why a hypothetical retry of HELM (cf. section 1.2.3) would benefit from Web 2.0 technology [231]: Mathematical content would become interactively editable directly on the Web, and projects like PlanetMath have proven that there is “a community of people interested in collaboratively authoring rigorous mathematics on the web” [231].¹⁰ Furthermore, HELM or MONET would now benefit from a much wider availability of stable libraries and tools. With SPARQL, for example, there is now a standardized and widely supported query language for RDF.

1.3.2. What MKM can Contribute to the Sem. Web

Conversely, there are now also opportunities for MKM to give back to the Semantic Web. Mathemat-

ical semantics is needed to improve, or even *enable*, certain applications of Linked (Open) Data.

Mathematics is a ubiquitous foundation of science, technology, and engineering. Some of these application areas are already well represented on the Web of Data, but their mathematical foundations are not. Having them represented as well would enable a whole range of new applications:

General-purpose Mathematical Knowledge: The inadequate representation of mathematical knowledge in Wikipedia has been criticized in section 1.1.2. DBpedia, the linked open dataset obtained from Wikipedia [99], inherits these limitations. Such limitations – in DBpedia and elsewhere – forced the Polymath collaborators mentioned in section 1.1.1 to search for previous publications of refutations of $P \neq NP$ “proofs” by keyword.

Statistics: Public sector information, increasingly being published as Linked Data by the US, UK, and other governments [205,103], has been used to provide, e.g., localized information retrieval about political representatives, crime statistics, and hospital waiting list statistics [186]. Statistical datasets contain values derived from ground values, or from other derived values using mathematical functions. The derivation can be as simple as counting; note that counts not only exist in proper *statistical* datasets, but also in statistics about *any* kind of datasets: VoID (Vocabulary of Interlinked Datasets) defines properties for expressing statistical characteristics of datasets, such as the number of distinct subjects [40, section 4.6]. Planning data collection from statistical datasets and interpreting collected data requires a notion of mathematical provenance of their data points. (Section 5.2 outlines a possible solution.)

Publication Databases: The RKB Explorer ACM linked dataset [18] classifies the scientific publications of the ACM according to their Computing Classification System (cf. section 4.3.5). Still, it is impossible for a Linked Data agent to understand that a publication merely classified as “F.1.3 Complexity Measures and Classes” actually deals with the P and NP complexity classes, and how they are defined. Thus, a mathematician who has developed a theory is unable to find out whether or how other mathematicians have built on it: Contemporary publication datasets only show what publications cited the original one, but

not if they reused the mathematical *concept* in question.

Enterprise Applications: Linked Data do not have to be open; the architecture, as defined in [65], also works in enterprise intranets. Renault has used them for retrieving information about spare car parts [204]. Now consider decisions to be made when designing whole cars: They ultimately require mathematical understanding. An engineer looking for an efficient engine for a projected city car might feed inputs such as the weight of the car, the average length and duration of a trip, the most widely available type of fuel and the average environment temperature when starting the engine into a mathematical model of the engine in order to predict its fuel consumption under these constraints.

e-Science: The above use case is actually about reproducing an experiment – one of the key principles of e-science [60]. Publishing descriptions of scientific experiments as Linked Data not only makes the provenance of their result data explicit [177] but also makes whole experiments more easily accessible and thus reproducible. Fine-grained reproducibility once more demands a representation of the mathematical models. Some e-science datasets include them, e.g. the SysMO SEEK “‘assets catalogue’ describing data, models, . . . , workflows and experiment[s]” [60] from systems biology of microorganisms [23], whose publication as Linked Data is in progress (cf. [60]). Currently, the mathematical models are given as Content MathML formulae (cf. section 4.1.1) deeply nested into XML files and thus not directly accessible via URIs. A mathematician who has developed the mathematical model for a scientific experiment cannot see how the model is applied. (The OntoMODEL tool, reviewed in sections 2.4 and 4.3.6, is a step into that direction.)

Thus, in order to enhance current applications of Linked Data towards mathematics, dataset publishers need a mathematical vocabulary. The quality of Linked Data vocabularies – often designed in an ad hoc mapping of existing database structures to RDF – and hence of the linked datasets is often low (see, e.g., [142]). ZIMMERMANN has observed the following reasons for vocabularies being of a low quality [233]:

1. ontologies defining the domain of interest do not exist;
2. they exist but are difficult to find because developed by small groups for experimentation, lacking advertisement;
3. they exist and can be found but they are of poor quality, not complying with standards or best practices;
4. they exist and can be found but there are too many, of mixed quality, and it is difficult to assess which ones are appropriate for a specific use case.

High-quality machine-readable vocabularies for mathematics do exist: The official OpenMath 2.0 Content Dictionaries (CDs), for example, defining 260 mathematical symbols – operators, functions, sets, constants –, have undergone a strict human-driven review process (cf. section 4.1.2), and there are large machine-verified libraries of formalized mathematics (cf. section 4.2). Large parts of the MKM community accept them as standard vocabularies for representing mathematical expressions, but for the rest of the world – including the publishers and ultimately the consumers of linked data – ZIMMERMANN’s criterion (2) applies. Besides a technical mismatch – they are not available as RDF¹¹ – there is a *cultural* mismatch. Mathematics, due to its practice of rigorously reasoning about abstract concepts in a self-contained way using a symbolic notation, is generally perceived as hard and inaccessible (see, e.g., [112]). The average computer scientist, whose work builds on a very restricted area of applied discrete mathematics, is not immune to such stereotypes. By integrating mathematics into the Web of Data, using the techniques explained in this article, we can take it out of the Ivory Tower.

1.3.3. Structure of this Article

This article reviews vocabularies – ontologies and languages – that are suitable for contributing mathematics to the Semantic Web, particularly the Web of [Linked] Data. The remaining sections are structured as follows: Section 2 provides an abstract overview of the structures of mathematical knowledge and thus the background knowledge about the domain that is needed to assess the aptitude of existing ontologies and languages for representing mathematical knowledge adequately to the applications described above. Section 3 defines the scope of this survey and establishes requirements for ontologies and languages. Section 4 provides a comprehensive review and concludes with recommendations on what ontologies and lan-

guages should be used on the Web of Data. Section 5 explains techniques for integrating non-RDF representations, which are still ubiquitous in the MKM domain, into the Web of Data, using the ontologies reviewed. Section 6 tries to predict the benefits of that and points out further research directions.

2. Background: Structures of Mathematical Knowledge

Before we can represent mathematical knowledge on the Semantic Web, we have to understand its structures. Realistic MKM applications, in domains where mathematics is applied but also in pure mathematics, do not only operate on logical and functional structures but also require information about the (non-mathematical) application *context*, about project organization and management (such as “What theorems are still lacking a proof”), about discussions that authors and users hold *about* the mathematical knowledge (such as “I don’t understand what we need this definition for”), etc. There is little literature about these structures. Working mathematicians often use them without reflecting on them. Computer scientists and knowledge engineers have to reflect on them but often do so from the point of view of a system specialized for a particular task – e.g. checking first-order logic proofs – and its particular conceptual model and representation language. Thus, the review given here is influenced by literature on concrete systems, models, languages, and ontologies, but tries to abstract from that.

2.1. Logical and Functional Structures

Mathematical knowledge has a three-layered *logical* structure of objects – composed of symbols –, statements, and theories¹². Symbols comprise operators, functions, sets, and constants. New mathematical concepts (i.e. symbols) can be defined, possibly based on concepts defined previously. A mathematical *object* can be a single symbol, or a compound, such as a complex number, an application of a function to arguments, or a derivative. Here, we call their structure “functional”, as they are built recursively from applying function or constructor symbols to other objects. Some of the properties of mathematical symbols or objects are specified as axioms. Axioms are expressed as formulae in a logical language, e.g. first-order logic (FOL). By applying rules of that logic, other proper-

ties of the mathematical concepts can be inferred. In a usual mathematical document, such properties are first asserted and then proven – or refuted. Often, the choice of what properties of a concept to model as axioms is arbitrary and merely follows established conventions. All kinds of properties of concepts are sometimes subsumed under the term *statements*. This is the case in the OMDoc representation language (cf. section 4.1.3), which distinguishes symbol declarations and axioms, definitions¹³, assertions (theorems, lemmas, corollaries, etc.), proofs (which prove assertions by applying inference rules to axioms and previously proven theorems), and examples. Not all assertions in a realistic mathematical knowledge base have to be true: There can be conjectures whose truth is not yet known, as well as wrong assertions that have been refuted by counter-examples but are kept for instructive purposes. Groups of closely related symbols and their properties form *theories*. When reusing mathematical symbols, their names are often qualified by their theory for disambiguation, i.e. theories act as namespaces for symbols; this is also reflected by speaking of the “home theory” of a statement. For example, both the theory of real numbers and the theory of functions on real numbers have an “addition” operator. The latter can be defined pointwise in terms of the former, but both remain different; for example, one cannot use either of them to add a number to a function.

In the context of theories, statements can be distinguished more precisely into *constitutive* statements (axioms and definitions), which determine the meaning of a theory, and *non-constitutive* ones, such as theorems or proofs, which “only illustrate the mathematical objects in the theory by explicitly stating the properties that are implicitly determined by the constitutive statements” [147, chapter 15.1]. Moreover, the logical language used to express the statements in a theory can itself be modeled as a theory, then called *meta-theory*. For example, the theory of commutative groups can be formalized with FOL as a meta-theory. FOL provides the universal quantifier that is needed for stating the group axiom of commutativity as $\forall a, b \in G. a \circ b = b \circ a$.

For knowledge management tasks, which involve, e.g., reuse of theories or management of theory changes, it has been found useful to build theories on a minimal set of axioms, and to model a whole field of mathematics as a strongly interconnected graph of “little theories” reusing each other (cf. [113]). The connections are called *theory morphisms* or *views*. Some of these views are given by definition – then called *imports* –,

others are postulated and then have to be proven. Theory graphs allow for modeling hierarchies from abstract to concrete concepts; for example, the theory of real numbers would, via some intermediate steps, import the theory of groups, and the morphisms along these imports would map the general operator of a group via multiple inheritance to the specific addition and multiplication operators of the real numbers. This use of theories extends beyond mere namespacing and has particularly been adopted for the structured specification and verification of software (see, e.g., [53]).

Logical/functional structures can be expressed at different *degrees of formality*: Often, an author starts a document by sketching a few formulae and some textual notes. Later, the content is elaborated both into the formal and into the informal direction: A sloppy formula is written more rigorously, rigorous text is formalized in a certain mathematical foundation (meta-theory)¹⁴, taking previously formalized knowledge into account, and natural language explanations are added to formalized knowledge (see, e.g., [147, chapter 4]). Both directions can, in principle, be automated: *Natural language processing* techniques can aid formalization (see, e.g., [117]), whereas proof explanation helps to generate natural language from formalized knowledge (see, e.g., [115]). These solutions, however, can not yet cope with the full complexity of mathematical knowledge as it occurs in practice. Particularly the automated disambiguation of symbolic notation (see below) is hard, as the surrounding text often has to be taken into account for disambiguation [117].

One aspect that is not restricted to logical/functional structures but has been investigated most deeply for them is *dependency*. For the purpose of managing mathematical knowledge, dependency can be defined such that B depends on A in the way d_p iff a change to A may have an impact on the property p of B . To make this definition precise, one has to fix the property p . Different conceptual models and representation languages have done that in different ways. The formal language MMT, which constitutes a subset of the above-mentioned OMDoc, considers dependency w.r.t. logical well-formedness [195, chapter 8.4]; examples for that are given in 4.3.2. The MathLang language (cf. section 4.1.4) considers dependency w.r.t. the reader's ability to understand a knowledge item [199,143].

2.2. Rigorous Language and Symbolic Notation

Mathematical textbooks and other publications predominantly consist of natural language (in its very own style; cf. figure 1 and [234,215]) intermixed with formulae. In such sources of mathematical knowledge, the higher-level logical structures (e.g. theories) are often less obvious, which suggests making the knowledge accessible from its discourse structure. Rhetorical Structure Theory (RST [169]) is a general-purpose model for that, which divides a text into spans, often down to the level of subordinate clauses. RST has a rich vocabulary of relations between *nuclei* (essential text spans) and their *satellites* (spans that provide additional information); for example, a satellite can give evidence to a nucleus, provide background information to facilitate understanding, or define the context in which the nucleus is to be interpreted. Figure 2 models a phrase from figure 1 according to RST.

Note, however, that statement- and object-level structures are usually obviously identifiable in mathematical text as shown in figure 1. Here, a sufficiently fine-grained model of these structures – covering, e.g., inline definitions and proof steps – may lead to a mathematics-specific refinement of RST.

For sections and chapters on the upper levels of a document, several models of discourse in scientific publications have introduced more convenient coarse-grained blocks that correspond to the usual sections of a publication, and reserve RST for markup of an intermediate granularity [128,77]. A typical document in one of these models starts with an abstract and a motivation and ends with a conclusion and a list of references, and has some sections in between that provide background knowledge, explain the actual contribution of the paper, demonstrate practical applications, summarize the results of experiments or evaluations, review the state of the art and related work, etc.

Mathematical formulae are communicated to human readers in a two-dimensional notation, whose complexity is owed to the possibility to define new semantic symbols at will.¹⁵ Choosing an intuitive notation for the concepts dealt with is of great importance to understanding and communication (see, e.g., [194]). The notation of a symbol is usually introduced with its first declaration, typical phrases being “We will denote by Z the set ...”, “The notation aRb means that ...”, etc. [215]. Notation can be conceived as a many-to-many mapping of structures of mathematical knowledge – primarily logical/functional structures – to an arrangement of glyphs on paper. The notation chosen

$$\text{Let } M \text{ be } \dots \left| \begin{array}{l} \text{Suppose that } \dots \\ \text{Assume that } \dots \\ \text{Write } \dots \end{array} \right. \text{Then } \dots, \left| \begin{array}{l} \text{provided } m \neq 1. \\ \text{unless } m = 1. \\ \text{with } g \text{ a constant satisfying } \dots \end{array} \right.$$

Fig. 1. Typical phrase patterns for theorems [215]

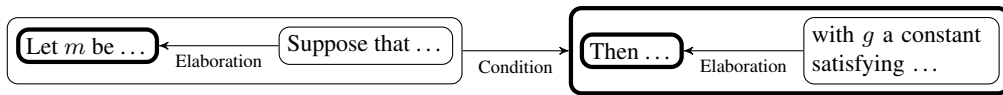


Fig. 2. RST markup of a theorem (nuclei with thick outline)

for a particular object in a particular document is determined by a number of presentation context dimensions (examples taken from [176] and [182]):

language and culture: the French/Russian notation of the binomial coefficient C_n^k vs. the German/English notation $\binom{n}{k}$; see [167] for details

level of expertise: the explicit notation of multiplication as $a \cdot b$, which is common in primary school, vs. the more advanced omission of the operator symbol in the notation ab

area of application: The square root of -1 is written as i in most fields, whereas electrical engineers write it as j to distinguish it from the current I .

community of practice: People with a set theory background tend to include 0 in the set of natural numbers \mathbb{N} , whereas those with a number theory background tend to start with 1 .¹⁶

individual preference: Some mathematicians, who prefer completely idiosyncratic notations when working on their own, translate other articles into their own notation and translate their own articles to a more conventional notation before publication [138, pp. 166–167].

The greatest notational variety has been observed for mathematical symbols. From the level of statements upwards, notation – such as the font chosen for keywords like “Definition” – is more standardized and therefore usually not a subject of research.¹⁷

2.3. Mathematics-Specific Metadata

In an expressive representation of mathematical knowledge, it is hard to draw a line between data and metadata. This article, in line with prior work on metadata in MKM [119], considers all of the previously mentioned structures *data* of primary interest, whereas the remaining, mainly administrative and application-specific information are considered *metadata*.

Metadata can be embedded into the data they describe, or point to the data from outside (“standoff markup”). This section focuses on metadata that are so closely related to the data that embedding them makes most sense in the interest of uniform management workflows. In mathematical practice, subjects annotated with metadata can be very fine-grained, as the following blackboard-style example shows:

$$b^{-1}(\boxed{a^{-1}a})b = b^{-1}(eb) = \dots$$

← We learned that last week

Administrative metadata describe the lifecycle and revision history of a resource, the data format and the usage requirements, copyright information, as well as other general-purpose information. The most widely used metadata vocabulary is Dublin Core [26], which covers general bibliographical information, but also elementary licensing and versioning information. Its semantics is rather weak, but it is widely supported, e.g. by search engines. The Dublin Core Metadata Element Set [108] provides a basic vocabulary, which the more modern DCMI Metadata Terms vocabulary extends in a backwards-compatible way [100]. Dublin Core also paves the path to a more comprehensive domain-specific description of resources, in that it is designed to be complemented by domain-specific classification schemes, whose entries are usually alphanumeric codes. For mathematical publications, the MSC (Mathematics Subject Classification [35]) prevails; this article would be classified as 68T30, where 68 is computer science, 68T is artificial intelligence, and 68T30 is “knowledge representation”. GAMS, the Guide to Available Mathematical Software [4], classifies more fine-grained things, namely mathematical problems, for example, H2a1 = “one-dimensional finite interval quadrature”.

Further metadata vocabularies cover the settings in which mathematical knowledge is applied. For example, the Learning Object Metadata (LOM [140]) describe educational properties of resources, such as their level of difficulty or interactivity, their coverage of topics (e.g. in terms of classification systems), and their intended audience. A vocabulary inspired by LOM has been used in the mathematical e-learning system ActiveMath [176].

2.4. *The Application Environment*

Again, the application environment can be modeled as data rather than metadata, given an appropriate conceptual model of the application domain. We mention three notable examples:

SWEET (Semantic Web Earth and Environmental Terminology [21,198]) is an OWL ontology that describes 4600 concepts in 150 modules from fields related to mathematics, such as physics, chemistry, biology, geology, and astronomy. These modules build on a foundation of general concepts of mathematics (e.g. functions), natural science, and space (e.g. coordinates). SWEET's model of mathematics does not intend to be as elaborate as the structural ontologies reviewed here, but SWEET provides a good showcase of how to integrate knowledge about mathematics with knowledge about its scientific application domains. One example of how SWEET integrates mathematics and science is the concept of a gravity field, defined as a vector field whose force is gravity. A vector field is a subconcept of a function whose result is a vector, and a vector is defined as an array of scalar elements.

GeoSkills is an OWL ontology describing topics, competencies and educational contexts related to interactive geometry [166], albeit without a connection to a structural model.

OntoMODEL is a tool for managing mathematical models in pharmaceutical product development [211]. Its underlying OWL ontology of mathematical models captures notions such as assumptions, parameters, and dependent variables in a way independent from the specific application domain.

2.5. *Discussions in Mathematical Collaboration*

Previous research has produced ontologies for two kinds of scientific discourse: One kind is embedded

into scientific publications, which, e.g., make claims and argue about claims made in other, cited publications. This has often been studied in combination with rhetorical structures; see [128] for an overview.

The other kind of scientific discourse is held externally of the representations of its subjects, e.g. in discussion forums. This perspective has been studied in the context of collaborative problem solving; the most common models employed in knowledge engineering have been derived from IBIS (Issue-Based Information System [153]). The DILIGENT argumentation model has been developed in the context of the namesake collaborative ontology engineering methodology with the design goal of making arguments more focused than in plain IBIS in order to make design decisions more traceable [213,214]. A DILIGENT argumentative thread starts with raising an issue, e.g. verbalizing a requirement for the ontology to be designed or pointing out a problem with its current state. An issue can be resolved by implementing a proposed and approved solution idea. About issues and ideas, the participants can state objective arguments or their subjective position.

The generic notion of issues and ideas can be refined to capture domain-specific problem and solution types. The most common problem with items of mathematical knowledge in knowledge collections, as reported by the 25 participants of a survey that we have conducted among domain experts, is that they are wrong, followed by being incomprehensible, their truth being uncertain, being underspecified, or redundant [164]. Further cases include knowledge items of which it was not clear whether they were useful, and knowledge items expressed in an uncommon style. Problems are most commonly solved by directly improving the affected knowledge item, by splitting it into more than one, or by deleting it altogether. Knowledge items, issues, and ideas cannot be combined arbitrarily. For example assertions, proofs, and examples can be wrong, whereas a notation can rather be inappropriate, misleading, or hard to read and write. Then, if some knowledge item is wrong, it could be deleted, or fixed in place, or kept as an instructive bad example, whereas splitting it into two parts would not solve that problem.

Combining both perspectives on discourse remains to be done for mathematics but would allow for capturing further important mathematical practices. In his work on "Proofs and Refutations" [154], LAKATOS has studied how discussions about mathematical knowledge items materialize into new mathematical knowl-

edge. Consider a discussion thread in which a problem with a proof is pointed out, e.g. that it only covers a specific case and should be generalized. This discussion provides the rationale for a later, generalized re-statement of the respective theorem and its new proof and therefore could be integrated into the text that encloses the theorem and its proof.

3. Problem Statement and Requirements for Representing Mathematical Knowledge on the Semantic Web

This article reviews existing ontologies and languages for representing mathematical knowledge w.r.t. their potential to enable better MKM on the Web, i.e. to advance the state of the art discussed in section 1. In accordance with the wide definition of “MKM” cited in section 1.2, we are interested in what potential the *representational capabilities* of these ontologies and languages have (i) for supporting mathematicians, scientists, and engineers in producing and using mathematical knowledge, (ii) for supporting educators and students in teaching and learning mathematics, (iii) for publishing mathematical textbooks and disseminating new mathematical results, and (iv) for supporting librarians and mathematicians in cataloging and organizing mathematical knowledge.

3.1. Limiting the Scope of this Survey

We consider the following two aspects out of scope of this article:

Tools and Services: This article does not review existing MKM *tools* and *services* based on the ontologies and languages reviewed. Instead we require ontologies and languages to provide for a machine-comprehensible representation of mathematical knowledge (see requirement **C.A** below), and conjecture that for such languages it should in principle be possible to develop any desired tool support; we refer to [160, chapter 6] for an in-depth treatment of state-of-the art MKM tools.

Reasoning: Another aspect that we consider out of scope is mathematical *reasoning* in the narrow sense, i.e. within the dimension of logical and functional structures. While Semantic Web technology is concerned with reasoning, it is, in the interest of scalability, largely limited to decid-

able subsets of first-order logic, such as Description Logic (DL) and Horn Rules (cf. [111] for an overview). In contrast to that, most areas of mathematics require first-order or even higher-order logic to be captured faithfully. Tools for verifying or even constructing proofs in these logics (section 4.2 mentions some of them) exist, within the inherent limitations of calculi for these logics. They usually do not scale across the Web but instead require a full representation of a problem and all of its prerequisites and foundations in main memory. A proof of concept for how well a subset of first-order logic can *approximate* the formalization of mathematical concepts has been presented by BRÖCHELER (cf. [76], which includes references to earlier related approaches). However, this investigation, as well as related ones, had not been motivated by the desire to provide an alternative to existing proof assistants, but by information retrieval use cases, which we briefly discuss in section 4.3.6. In summary, we treat mathematical reasoning as a special task to be supported by specialized tools. We expect knowledge representations in terms of the ontologies and languages that we review to *guide* agents and their users in *finding* the right specialized tool for a mathematical reasoning task – in analogy to the position that the MONET architecture reviewed in section 1.2.2 took on computation. For that reason, we require the ontologies and languages reviewed to capture the logical and functional structures of mathematical knowledge as comprehensively as needed (see requirement **S.L** below), so that it can be passed on to specialized tools – possibly after a translation into their native language, but preferably without losing information.

3.2. Requirements

From the review of the state of the art of mathematics on the Web, we can infer as design goals for Semantic Web applications for MKM the ability to reuse knowledge across knowledge bases, information retrieval adequate to the structures of the knowledge, and integration with mathematical services, such as automated reasoning and computation, without compromising comprehensibility for human end-users. Having reviewed the structures of mathematical knowledge in section 2 and having defined the scope of this

survey in section 3.1, we are now ready to specify more precise requirements for ontologies and languages¹⁸:

S: All of the previously reviewed structures of mathematical knowledge SHOULD be supported; where this is impossible, missing dimensions MUST be compensated for by language extensions along the criteria L.E and L.→ below. We subdivide this criterion as follows:

S.L.{O,S,T}: logical/functional structures: mathematical objects, statements, theories

S.R: rigorous language or rhetorical structures

S.N: notation

S.M: metadata

S.D: discussions

F: Mathematical knowledge occurs in different degrees of formality; applications targeting human users and automated agents require both formal and informal representations. Therefore,

F.R: the language SHOULD be able to represent knowledge in a wide range from informally to fully formalized, and

F.C: many degrees of formality, including formalizations in multiple foundations, SHOULD be able to coexist in one document, interlinked with each other.

L: In real-world applications, mathematical knowledge is combined with multiple dimensions of non-mathematical knowledge. Therefore, a language SHOULD support interlinking of these dimensions by rich annotation facilities, but also give authors the freedom to represent some knowledge by external means and link it to representations in the given language. In detail,

L.A: the language MUST allow for attaching non-mathematical metadata and annotations to mathematical knowledge items, regardless of their granularity,

L.→¹: it MUST allow for linking mathematical knowledge items to external mathematical or non-mathematical resources, and

L.←²: it MUST be possible to address all mathematical knowledge items expressed in the given language from outside, in order to link external representations to them, for example standoff markup or existing representations in different languages.

C: Knowledge represented in a language SHOULD be comprehensible

C.A: to arbitrary automated agents – therefore, the knowledge SHOULD be self-describing in a machine-comprehensible way.

C.H: to human users – therefore, published human-comprehensible documents generated from representations in the given language SHOULD retain semantic annotations, so that assistive services can retrace the original knowledge and make it available to the user on request, e.g. integrated into a user interface.

Tables 1 and 2 at the end of this section summarize how well the languages and ontologies reviewed satisfy these requirements.

4. Review of Languages and Ontologies

This section reviews existing languages and ontologies for representing mathematical knowledge according to the requirements established in section 3.2. Besides Semantic Web ontologies in the narrow sense, other machine-comprehensible representation languages are taken into account. This is because they are widely used in science, technology, engineering, and mathematics, even preferred over ontologies in certain settings, and most existing machine-comprehensible mathematical knowledge is available in these languages rather than RDF. We do, however, pay attention to the possibility to translate such representations to RDF.

4.1. XML-based Semantic Markup Languages

XML languages share the URI foundation with RDF. All semantic XML languages allow for assigning IDs for the inner nodes of their tree-structured representation (i.e. for the *elements*), via XML ID [171]. Together with the URI of the XML document, that allows for global identification and linking and thus satisfies requirement L.←. XPointer [122] allows for identifying more complex subsets of an XML representation, such as node or text ranges, but is rarely supported. Note, however, that additional work is required to integrate semantic XML markup with RDF-based Linked Data; this is discussed in section 5.1.

4.1.1. Content MathML 3 and OpenMath 2 Objects

MathML (Mathematical Markup Language [52]) is an XML language that was originally conceived for embedding mathematical formulae into HTML web pages – which is still its main purpose. It fea-

¹read “L out[going]”

²read “L in[coming]”

tures a presentation-oriented sublanguage (Presentation MathML) but also a semantics-oriented one (Content MathML); the latter is covered here. MathML allows for a fine-grained mix of semantic and presentation markup (“parallel markup”) – thus addressing requirements **F.*** but also **C.H**, considering it as an enabling technology for interacting with formulae in published documents (cf. [116]). It is possible to embed or link non-mathematical information (requirements **L.A** and **L.→**). The task of defining notation, i.e. mapping semantic structures to a human-readable presentation, is left to other, non-mathematical languages such as XSLT.

The related OpenMath [80] language has originally been invented in the mid-1990s to facilitate data exchange between computer algebra systems (CAS) but has since been aligned closely with Content MathML, leaving only syntactic differences in the latest versions of both languages [97].

Both languages are limited to representing the functional tree structure of mathematical objects (requirement **S.L.O**) but are frequently integrated into host languages that cover further structures. Many XML languages already embed MathML or OpenMath officially (see below), whereas others allow for extending their vocabularies accordingly. Listing 1 shows a document that contains MathML and OpenMath objects.

Mathematical objects consist of numbers, variables, symbols, and applications of objects to other objects. Content MathML comes with a default supply of symbols that cover high school and introductory university education. Mathematical objects reference these symbols by URI. Their semantics is defined in external vocabularies called *OpenMath Content Dictionaries* (CDs); authors can create and use additional CDs as needed. The machine-comprehensibility of MathML/OpenMath representations depends on the degree of formality of the CDs. Suggested alternative RDF representations of MathML are discussed in section 4.3.2.

4.1.2. Extensible Mathematical Vocabularies:

OpenMath 2 Content Dictionaries

An OpenMath CD is a collection of (usually closely related) definitions of symbols. The abstract model of a CD covers mathematical objects – used to formally represent properties of a symbol, as opposed to plain text descriptions –, a weak variant of axioms/definitions and examples on the statement level, and a basic notion of theories, plus a limited metadata vocabulary [80, section 4] (requirements **S.L.*** and **S.M**). The

reference encoding of that model is a lightweight XML language, roughly comparable to RDFS in expressivity; the OMDoc language (cf. section 4.1.3) offers a more expressive but compatible alternative. Alternatively, the model has been implemented by ontologies (cf. section 4.3.2).

The *official* CD collection reviewed by the OpenMath Society (cf. [80, section 4.5]) defines 260 symbols from arithmetics, set theory, FOL, algebra, calculus, as well as transcendental and statistical functions [27]. These CDs do not provide a full formalization; instead, developers of, e.g., CAS are supposed to use the CDs as specification manuals when implementing *phrasebooks*, which translate OpenMath objects into the native languages of such systems.

4.1.3. More Expressive CDs and Documents:

OMDoc 1.3

OMDoc (Open Mathematical Documents [12,147]) is an XML language for representing mathematical knowledge that has a particularly rich vocabulary for logical/functional structures (requirements **S.L.***). OMDoc supports MathML and OpenMath objects, and its theories are compatible to OpenMath CDs. Beyond that, OMDoc adds vocabulary for formal and informal statements, modular theories, as well as narratively ordered documents, rhetorical structures (requirement **S.R**; cf. SALT in section 4.3.3), notation definitions (requirement **S.N**). By integrating RDFa [37] for arbitrary metadata and links [160, chapter 5], conforming to the specification of an RDFa host language, OMDoc furthermore satisfies requirements **S.M**, **L.A** and **L.→**, and makes representations more comprehensible to agents (if they are aware of RDFa and if the vocabularies used for annotation are published as Linked Data; requirement **C.A**). On each structural level, OMDoc supports a wide range of degrees of formality, from unstructured text to a full formalization – thus eliminating the need for phrasebooks at least on object level –, which can coexist in an interspersed, literate programming style to serve the needs of human- and machine-oriented applications. OMDoc particularly allows for statement-level parallel markup that interweaves textbook-style natural language with a purely formal representation of the same statements, down to linking subclauses in the text to the corresponding subterms in formulae, and linking words to symbols or variables.¹⁹

Listing 1 gives an example of the syntax of OMDoc 1.3, whose specification is currently being final-

Listing 1: An OMDoc theory with a declaration (type given in OpenMath) and implicit definition of the exponential function (given in Content MathML)

```

<theory xml:id="transc">
  <imports from="sts#sts"/>          <!-- Small Type System for OpenMath [95] -->
  <!-- alternatively, stronger type systems may be used -->
  <imports from="setname1#setname1"/>  <!-- numbers and other basic sets -->
  <symbol name="exp">
    <meta property="dc:description">the exponential function</meta>
    <type><!--  $\mathbb{R} \rightarrow \mathbb{R}$  -->
      <om:OMOBJ>
        <om:OMA>          <!-- OMA applies a constructor or function to arguments -->
          <om:OMS cd="sts" name="mapsto"/>
          <om:OMS cd="setname1" name="R"/>
          <om:OMS cd="setname1" name="R"/>
        </om:OMA>
      </om:OMOBJ>
    </type>
  </symbol>
  <definition xml:id="exp-def" for="exp" type="implicit">
    <CMP>
      <phrase verbalizes="#equal-deriv">
        <term cd="transc" name="exp">The exponential function</term>
        equals its derivative</phrase>
        and evaluates to 1 for an argument of 0.</CMP>
      <FMP><!--  $\exp' = \exp \wedge \exp(0) = 1$  -->
        <m:math> <!-- here, we use the symbol vocabulary built into Content MathML -->
          <m:apply> <!-- as an alternative to explicitly referring to imported CDs -->
            <m:and/>
            <m:apply id="equal-deriv"> <!-- <m:csymbol cd="logic1">and</m:csymbol> -->
              <m:eq/>
              <m:apply>
                <m:diff/>
                <m:csymbol cd="transc">exp</m:csymbol>
              </m:apply>
              <m:csymbol cd="transc">exp</m:csymbol>
            </m:apply>
            <m:apply>
              <m:eq/>
              <m:apply>
                <m:csymbol cd="transc">exp</m:csymbol>
                <m:cn type="integer">0</m:cn>
              </m:apply>
              <m:cn type="integer">1</m:cn>
            </m:apply>
          </m:math>
        </FMP>
    </definition>
  </theory>

```

ized.²⁰ An OWL ontology covering a large subset of OMDoc is reviewed in section 4.3.2.

OMDoc has been used for exchanging knowledge between systems doing structured specification, automated verification, and interactive theorem proving, for documenting Semantic Web ontologies [161], for publishing human-readable documents for interactive browsing [98] and adapted to different audiences [182], and for e-learning in the ActiveMath system mentioned before.

4.1.4. Narrative Documents: MathLang

MathLang [144] is similar to OMDoc but puts an even higher emphasis on formalization of informal, but highly conventionalized mathematical text. From its structural annotations, “proof skeletons” can be generated, i.e. templates in languages for formalized mathematics [144]. MathLang is on a par with OMDoc in its coverage of object- and statement-level logical/functional structures (requirements **S.L.**{**O,S**}), narrative document structures, and literate-programming-style combinations of text and formalizations (requirements **F.**{**R,C**}). However, there are no theory level and no rhetorical structures, the metadata vocabulary is restricted, and links to non-mathematical knowledge are not supported. MathLang has an XML encoding that is used for most processing tasks except authoring and presentation.

MathLang’s “Document Rhetorical aspect” (DRa), which does not cover rhetorical structures in the sense of RST but rather statement-level logical structures and narrative document structures, has also been implemented as an OWL ontology (cf. section 4.3.2).

4.1.5. Languages for Books and Manuals

A formal view on technical specifications reveals structural similarities to mathematical theories. While fully formalized specifications can be written in the same languages as general formalized mathematics (cf. section 4.2) and then be verified automatically, there are different XML languages targeting human audiences, such as engineers implementing a specification, or developers using an API; DocBook, TEI, and DITA are reviewed here. Similar languages exist for e-books and courseware; we review EPUB/DTBook and CNXML/CollXML.

None of these languages is directly suitable for representing mathematical knowledge other than objects; therefore, this review covers them rather briefly. All of them satisfy requirement **S.L.O** by supporting MathML either natively or as an extension²¹, some have limited built-in support for statement-level logical

structures, none supports theories. We first review the individual languages, then discuss further extension possibilities. Generally, these languages do not have a machine-comprehensible semantics, and the publication tools available for them do not support generating semantically annotated human-comprehensible documents.

DocBook 5: DocBook, the most widely used XML language for technical manuals [220,221], focuses on representing structures pertinent to its main application area of software documentation. Other than MathML objects, titled equations, and examples, DocBook does not support mathematical structures, and it has a fixed idiosyncratic metadata vocabulary with a coverage similar to Dublin Core (requirement **S.M**). It hardly offers native markup for representing knowledge in different degrees of formality and interlinking such representations. Embedding arbitrary literal-valued metadata into a document is not supported.

A notable application of DocBook in MKM is the MathDox e-learning system [174,92], whose compound document format combines DocBook with OpenMath objects and further XML vocabularies for programming constructs, requesting user input, queries to MONET services (cf. section 1.2.2), and exercises.

TEI P5: TEI (Text Encoding Initiative [78]) is, due to its focus on digitalization and edition of paper-born documents from the humanities, social sciences, and linguistics, obviously not suited for representing *mathematical* knowledge. Nevertheless, it is a prime example of an expressive semantic markup language. The TEI guidelines recommend using any available representation language for mathematics, according to the requirements, and explicitly mention MathML, OpenMath, and OMDoc [78, chapter 14.2]; a combined TEI+MathML XML schema is available. In its own domain of literature, TEI can express knowledge in a wide range of degrees of formality (requirement **F.R**). It supports fine-grained interlinking of different representations of the same knowledge (requirement **F.C**). TEI has an elaborate but finite metadata vocabulary for representing the provenance of documents and even smallest fragments of text. Arbitrary additional information can be embedded into a document (requirement **L.A**), or provided as standoff markup pointing into the original document (requirement **L.←**), whereas linking to external resources (requirement **L.→**) is restricted in that no link types are supported. Certain sub-vocabularies of TEI have been given a formal se-

mantics by mapping them to relevant domain ontologies [24,188].

DITA 1.1: DITA (Darwin Information Typing Architecture [104]) does not support any particular application scenario by default but is rather intended to offer a framework for developing languages for topic-based technical documentation that are specialized to a particular domain of application.²² The topic paradigm is in contrast to DocBook’s focus on contiguous, narratively ordered manuals. DITA can be extended by MathML and supports [definitions of] concepts and examples. DITA’s built-in metadata vocabulary primarily focuses on the context in which an object can be (re)used, such as the intended audience or keywords. DITA performs as badly as DocBook w.r.t. the **S.*** and **F.*** requirements, except that topics can be interlinked. However, DITA offers stronger support for adding arbitrary metadata and links (requirements **L.***).

EPUB 2.0.1 and DTBook 3: EPUB is a standard for general-purpose e-books, not primarily technical manuals. A complete e-book is a bundle of content files with a Dublin Core metadata record [11,13]. Besides XHTML – which could carry RDFa –, DTBook (DAISY²³ Digital Talking Book [1]) is the recommended format for them. DTBook is a semantically structured format inspired by DocBook but simpler and with less support for customization. RDFa support is planned for the next versions of EPUB and DTBook.

CNXML 0.7 and CollXML: CNXML, the language of the course modules of Connexions (cf. section 1.1.1) [29], is comparable to a subset of DocBook in expressivity. CNXML recommends using Content MathML for mathematical objects and supports definitions, “rules” – comprising, e.g., axioms and theorems – and examples (requirement **S.L.S**); thanks to its educational focus, it also supports exercises [197].

Course modules written in CNXML are combined into collections represented in the CollXML container format [29]. CollXML was preceded by a partial representation of a collection’s structure in RDF [19]. A CollXML document models dependencies between modules – so-called “featured links”, which can be of type “prerequisite”, “supplemental”, or “example”, in three degrees of strength. There is an idiosyncratic metadata vocabulary with a coverage similar to Dublin Core.

Extensibility Towards Further Mathematical Structures: There are two principal approaches to introducing further mathematical structures: (i) literally

reusing elements of sufficiently expressive mathematical markup languages, such as OMDoc, or, (ii) reusing an appropriate ontology for mathematical structures (cf. section 4.3) – provided that the host language supports referencing arbitrary metadata vocabularies on any relevant structural level without first introducing new container elements for them via approach (i), i.e. if there is an RDFa-like infrastructure (cf. section 4.3.1).

Approach (i) works in DocBook, TEI, and DITA; via that extension path, DUCHARME has proposed integrating RDFa into DocBook and DITA [110]. DocBook, TEI, and DITA offer varying degrees of support for approach (ii). There is a workaround for adding RDF-compatible annotations to DocBook: Any DocBook element can carry XLink attributes, which can have a role (= predicate), and from which RDF can be harvested [94]. TEI documents can reference external objects by XPointers, but without any possibility to specify a predicate type; thus, it does not allow for harvesting RDF. DITA provides the *othermeta* element for arbitrary key–value pairs, for which URIs could be used to emulate RDF, or, even more appropriately, the *data* element, which allows for constructing nested data structures and supports RDF’s distinction of URI- and literal-typed objects, as well as datatypes. Similarly, DITA supports links with arbitrary roles from topics to related topics.

4.2. Languages for Formalized Mathematics

Languages for formalized mathematics, such as those of the proof assistants Mizar [10], Isabelle [222], or Coq [30], are of interest here insofar as they also support informal content. They obviously support logical/functional structures of mathematical knowledge on the object, statement, and theory levels (requirements **S.L.***), with different approaches to modular theories (cf. [196] for an overview). Symbols and statements have identifiers, which are not compatible with URIs; however, for exchange purposes, the systems often offer an XML export (cf. [195]). Except for notation definitions (requirement **S.N**), there is little support for other structures; the “lowest common denominator” is to put such information into comment lines, which are post-processed by a specialized tool. Isabelle and Coq formalizations can be interspersed with informal text (partly addressing requirements **F.***), and certain parts of the formalized content can be marked as hidden for human-readable output. In Isabelle, informal text can contain formalized expressions as an-

tiquotations, which the proof assistant evaluates when exporting the document [222, chapter 4].

These languages do not support links out of or into formalizations (requirements $\mathbf{L}.\rightarrow$ and $\mathbf{L}.\leftarrow$). Each language comes with its own set of services that understand formalizations in the respective language, which have a strong model- or proof-theoretic semantics for logical and functional structures. These languages are usually committed to a particular first- or higher-order logical foundation; different assumptions made by different foundations, as well as the fact that part of the knowledge is not explicitly formalized but implied by the underlying foundation, generally make it hard to translate from one language into another one for reuse. Existing translations have usually been hard-coded for pairs of two specific languages (cf. [195, chapter 1.1.3.3], [196] and the Hets system [180,181]). On entailments implied by the choice of logical foundation, recall that Linked Datasets commonly make them explicit, as to circumvent scalability issues of reasoners (cf. [137,105]). With the RDFS or lightweight OWL vocabularies that most Linked Datasets employ, this explication is usually practically feasible; however, such vocabularies would only be capable of partially capturing formalized mathematics.

Each of the libraries that ship with the above-mentioned systems comprises several hundreds of theory files (cf. [223] for exact figures) of a very high quality: Firstly, it took a considerable effort to produce them – for a number of different proof assistants and sources of mathematical content it has been consistently experienced that fully formalizing one rigorous textbook page may take an author 1 to 1.5 weeks [47]; secondly, these formalizations have been machine-verified. These libraries usually have a good coverage of discrete mathematics; for example, Isabelle’s library covers elementary number theory, algebra, set theory, but also analysis. In the Flyspeck project for developing a machine-verified proof of the Kepler Conjecture²⁴, which employed several proof assistants, most of the required formalizations of trigonometry, geometry, topology, measure theory, etc., first had to be developed by the members of the project [132]. After HELM, no serious effort has been undertaken to fully integrate such libraries into the Semantic Web.

4.3. Structural Ontologies for Representing Mathematical Knowledge in RDF

While the languages reviewed so far are machine-comprehensible in their own ways, they do not in-

tegrate into the Semantic Web without translation to RDF (cf. section 5). Representing mathematical knowledge in RDF not only makes it accessible to Semantic Web agents, but also offers powerful ways of interlinking mathematical and non-mathematical knowledge (requirements $\mathbf{L}.*$), formal and informal representations (requirements $\mathbf{F}.*$), etc. However, a fine-grained interlinking of formal and informal representations as in literate programming or MathML’s parallel markup requires considerable effort, due to the absence of a native notion of order. RDF, when published in compliance with the Linked Data principles [65], is always machine-comprehensible in the sense that a machine can simply retrieve information about resources by dereferencing URIs (requirement $\mathbf{C.A}$). However, the informative value of the latter information, and the power of RDF in general, stands and falls by the availability of appropriate vocabularies, i.e. ontologies.

This section reviews ontologies that allow for representing mathematical knowledge natively in RDF, or that offer themselves as translation targets for knowledge originally represented in non-RDF languages. The review includes obsolete ontologies insofar as aspects of their design are still instructive today.

4.3.1. Different Approaches to Representing Mathematical Knowledge in RDF

Usually, when representing knowledge in RDF, one finds or develops an appropriate vocabulary and chooses a suitable RDF serialization, e.g. RDF/XML or XHTML+RDFa. In the presence of mathematical objects, this decision becomes harder due to their inherent complexity. Therefore, we briefly discuss possible representations before reviewing concrete ontologies.

Complete RDF Representations: Due to their n -ary ordered tree structure, mathematical objects are not amenable to a straightforward representation as RDF triples. With the narrative order of (not only) mathematical text, e.g. in textbooks, one faces a similar challenge. The use of linked lists or ordered sets, either the collections or sequences built into RDF [61] or custom remakes, is unavoidable. However, such data structures are not generally supported by RDF software, and they do not go well along with DL reasoning²⁵ and querying²⁶. The N3 Vocabularies reviewed below demonstrate this approach for mathematical objects, the SALT ontology for rhetorical structures.

Mathematical Objects as XML Literals: Compared to RDF triples, XML offers a much more intuitive representation of n -ary ordered trees. With Content MathML and OpenMath, there are standardized semantic XML representations of mathematical objects, which are widely understood by mathematical software (e.g. CAS phrasebooks). Therefore, reusing them as XML literals of *rdf:XMLLiteral* datatype while representing other structures of mathematical knowledge as RDF triples suggests itself (cf. the OpenMath CD ontology in section 4.3.2 and the OntoMODEL ontology in section 4.3.6 for examples). From a Semantic Web perspective, this has, however, the drawback that XML literals are largely opaque to contemporary RDF tools. The Virtuoso triple store [187] allows for filtering XML literals matched by a SPARQL graph pattern by XPath node tests [62]. The Corese RDF engine can additionally reuse variables from the proper SPARQL part of a query in XPath expressions [90]. None of these extensions has made it into the SPARQL standard yet.

Embedding RDFa into XML: RDFa is a set of XML attributes for embedding RDF graphs into X[HT]ML documents [37]. That allows for focusing on those structures that can easily be represented in RDF, while leaving the representation of n -ary and ordered structures to XML. However, queries that need both kinds of information have to be implemented separately. The upcoming RDFa 1.1 API [209], which remains to be implemented by browsers, will at least give in-browser scripts similar means of accessing embedded RDF as the Document Object Model (DOM) offers for X[HT]ML. The XSPARQL [38] query language combines SPARQL and XQuery; however, such a query would still rely on a separate service that makes the RDFa annotations available as queryable RDF.

The first official RDFa host languages were the presentation-oriented XHTML and SVG languages [42], which allow human-comprehensible documents to carry as much semantic annotation as needed by agents, such as assistive services. RDFa can also be embedded into semantic markup languages; that has been done for OMDoc (cf. [161], and [160, chapter 5] for full details). MathML has supported fine-grained annotation of presentational or semantic markup long before RDFa, with a similar expressivity (e.g. `<annotation definitionURL="link-type" src="link-target"/>`). OpenMath has a similar annotation syntax, albeit without URI support. Neither the MathML nor the OpenMath developers are cur-

rently planning to support the RDFa syntax. When using RDFa in semantic markup, one has to take care that the RDFa annotations do not interfere with the native semantics of the host language.²⁷

Standoff Markup: Finally, one can maintain parallel representations of the same concepts both in RDF and in one of the specialized languages reviewed above. In such a setting, the RDF graph acts as standoff markup pointing to fragments of the other representation and adding information to them, such as additional metadata, links, or semantic abstractions not supported by the original language. Conversely, information about n -ary structures and order would only be represented in the latter language. Most of the knowledge is usually represented redundantly in RDF and the other language – one of them possibly generated by automatic translation from the other one – to provide a maximum amount of information to agents that only understand one representation. This has so far been the most common approach in MKM (cf. section 4.3.2).

4.3.2. Logical and Functional Structures

Few approaches to completely representing logical/functional structures of mathematical knowledge in RDF have been made so far. A larger number of ontologies exists for representing mathematical statements, whereas the theory level has rarely been covered so far. The ontologies reviewed in this section have most commonly been used in standoff markup for XML representations.

N3 Vocabularies and RDF Encodings of Content MathML: The cwm [63] and Euler [101] reasoners natively use the N3 [64] superset of RDF. The standard N3 vocabularies cover a limited subset of object- and statement level structures, constrained to FOL as a meta-theory. Beyond domain knowledge, i.e. a library of basic mathematical functions (cf. section 4.3.6 for details), the N3 “math” vocabulary provides weak formalizations of general structural concepts such as the concept of a function. When a concrete function f is used as the predicate of an RDF triple, whose subject is a collection $(x_1 \dots x_n)$ holding the arguments, the reasoner infers $f(x_1, \dots, x_n)$ as the value of the object. When the object is identified by a URI or blank node ID, it can be reused in the subject of another mathematical expression. Listing 2 shows a sample set of facts and rules yielding `:ABC :side3 5`. Few RDF processors support the full N3 syntax. When an N3-aware processor is not available, the n -ary ordered tree structure or mathematical formulae has to be broken

Listing 2: The Pythagorean Theorem in N3

```
:ABC :side1 3 ; :side2 4 .

{?triangle :side1 ?a ; :side2 ?b .
 ?c is math:exponentiation of
  (((?a 2)!math:exponentiation
   (?b 2)!math:exponentiation)
   !math:sum 0.5) . }
=> { ?triangle :side3 ?c } .
```

down into explicit RDF triples. Combining RDF reification and N3’s “reason” vocabulary, which models the structure of proofs, allows for partially capturing the statement level. The coverage of the N3 vocabularies is determined by the needs of a FOL reasoner and thus not suitable for representing *arbitrary* mathematical knowledge. The semantics of mathematical functions is not fully specified in N3; cwm and Euler merely have built-in support for evaluating them.

Two RDF encodings of Content MathML have been suggested independently from N3. These representations look similar to N3, except that the application of a function is usually modeled with the [reified] application being the subject, and the function symbol and the arguments being the object(s). That makes nested expressions easier to write without the additional syntactic sugar of N3. An encoding proposed by MARCHIORI [170]²⁸ has obvious design flaws – such as introducing, for no obvious reason, two different ways of referencing symbols in CDs and applying them to arguments –, which another, similar representation independently developed by ROBBINS avoids [200]. Both suggestions have neither been implemented nor taken up by the MKM community.²⁹

As an advantage of representing formulae in RDF, MARCHIORI points out that it allows for making references to bound variables more explicit: Indeed, a bound variable is always represented as a unique RDF resource, be it on declaration or on usage. Content MathML, however, optionally supports a similar explication by making occurrences of the bound variable refer to the place where it is declared via *@xref* and *@id* attributes. MARCHIORI developed an ad hoc vocabulary from the Content MathML element and attribute names, which has little value from a Linked Data perspective. ROBBINS only uses a special vocabulary for the object constructors of Content MathML but the canonical OpenMath CD URIs (e.g. <http://www.openmath.org/cd/arith1#plus>) for

symbols. The latter are compatible with Linked Data, as explained in section 5.2.

Ontologies for OpenMath CDs: Two ontologies implement the data model of OpenMath CDs; due to the simplicity of that model they have a limited coverage of logical and functional structures and can therefore be treated briefly.

In an early phase of the above-mentioned MONET project, an RDFS vocabulary for representing OpenMath Content Dictionaries (CDs) was developed [79]. The RDFS vocabulary covered OpenMath’s logical/functional structures on the theory and statement levels, as well as metadata, by classes and properties, and represented mathematical objects as XML literals.

More recently, we have developed a more expressive OWL ontology [156], which covers more of the theory and statement levels (but still within the limits of the OpenMath CD model and therefore not comparable to more expressive ontologies). However, its representation of mathematical objects only covers flat occurrences of symbols, following the approach of the OMDoc ontology explained below.

MONET Problem Ontology: Rather than original structures of mathematical knowledge, the MONET OWL ontologies (cf. section 1.2.2 for MONET) describe mathematical problems and the software used to solve them. It is, however, instructive to study how the MONET problem ontology represents mathematical objects. It focuses on the operator or constructor symbol at the root of the functional tree representation of a mathematical object. Suppose the MONET broker knows a web service for computing definite integrals constructed with the *oms:calculus1#defint* symbol [81]. The type of problem that that service solves can be modeled as follows:

$$p:\textit{definite_integration} \sqsubseteq$$

$$p:\textit{Problem} \sqcap g:\textit{GamsH2a}$$

$$\sqcap = 1p:\textit{openmath_head.oms:calculus1\#defint}$$

The deeper structure is only represented in OpenMath; it is not used for service matching, but sent to a matching service for computation.

HELM: The HELM system (cf. section 1.2.1) generates from an original formalized representation in a non-XML language both a full XML representation and a standoff RDF graph containing a structural outline of properties relevant for searching.

HELM’s RDFS ontologies distinguish terms (corresponding to mathematical objects in our terminology), objects (roughly corresponding to statements), and theories. There is a notion of dependency, such as a corollary being a consequence of a theorem (*hth:isConsequenceOf*), or a lemma being a prerequisite of a theorem (*hth:isPremiseOf*). Terms can have occurrences of other HELM objects, i.e. symbols. Such an occurrence is reified as a resource, which has an *h:position* and an integer *h:depth* counting the number of premises, including universal quantifiers. Among the positions that have been found relevant for answering queries, e.g. for finding applicable theorems for proving something (cf. [202,130]), there are the following, explained using the example of the theorem $\forall a : \mathbb{N}. \forall b : \mathbb{N}. \forall c : \mathbb{N}. a \leq b \wedge b \leq c \Rightarrow a \leq c$:

h:MainHypothesis: the head symbol of a hypothesis; here: \wedge (depth 0³⁰)

h:InHypothesis: any other symbol anywhere else in a hypothesis; here, either of the two \leq

h:MainConclusion: the head symbol of the conclusion; here: \leq (depth 4)

h:InConclusion: any other symbol anywhere else in a conclusion (none in this example)

From the point of view of representing logical and functional structures in general w.r.t. requirement **S.L** (i.e. not just in the specific setting of the HELM system), the HELM ontologies do not sufficiently abstract from the native knowledge representation of the Coq library (cf. section 4.2). Concrete examples for that are the implicit relation between theories and their statements and the relatively idiosyncratic mechanism (judged from a modern Linked Data perspective) for identifying theory items – once by pointers into the XML representation³¹, and secondly by identifiers that are similar to blank node IDs, the latter being used for modeling dependencies.³² Both circumstances would make it hard to represent, e.g., the logical and functional structure of Mizar articles in terms of the HELM ontologies.

MoWGLI: An RDFS ontology with a wide coverage of logical/functional structures, including informal representations, educational content, and a rich set of metadata, was developed in the MoWGLI project [120]. MoWGLI reused vocabulary from the HELM ontologies, existing general and educational metadata ontologies, the XML schema of the OMDoc markup language (cf. section 4.1.3), and the metadata vocabularies of the ActiveMath e-learning system.

For the latter two, an RDFS model was newly developed. The MoWGLI ontology (merely called “meta-data model” due to its standoff usage) does not make further assumptions about the format in which the full knowledge is represented.

Summarizing, MoWGLI serves as an instructive example of a comprehensive integrated mathematical ontology. However, various shortcomings³³ make it technically unusable. It is not clear whether it has ever been applied; except for its specification, no trace in the form of annotated documents is left.

OMDoc 1.3: The OMDoc OWL ontology [158, 160] has been modeled after the conceptual model and XML schema of the OMDoc language (cf. section 4.1.3). While not yet as comprehensive as the OMDoc language³⁴, the ontology has a richer statement- and theory-level vocabulary and more notions of dependency than the other ontologies reviewed, which justifies a slightly longer treatment.

Figure 3 shows the core classes and properties. Some of the depicted classes have subclasses. Definitions can, e.g., be pattern-based, implicit, or recursive, and types can be declared or asserted. Assertions comprise theorems, lemmas, and corollaries, and they can have different truth values. Moreover, the ontology covers sub-statement structures such as proof steps. Different degrees of formality are distinguished by a property. The definition in listing 1 is formal but not fully computerized to a degree an automated theorem prover would understand; actually, it consists of a formal and an informal part.³⁵

There are three orthogonal properties that relate mathematical knowledge items to each other, each with a hierarchy of subproperties. Whole-part properties link, e.g., theories to their statements and proofs to their steps. The two parts of the definition in listing 1 are related by a verbalizes/formalizes relation; similar relations can occur on all structural levels. Thirdly, there is dependency w.r.t. logical well-formedness, validity, and presentation. If, for example, one symbol is defined in terms of other symbols, such as the exponential function in terms of differentiation, its well-formedness depends on them. This is reflected by the following property hierarchy:

o:hasDefinition \circ *o:usesSymbol*
 \sqsubseteq *o:hasOccurrenceOfInDefinition*
 \sqsubseteq *o:wellFormednessDependsOn*
 \sqsubseteq *o:dependsOn*

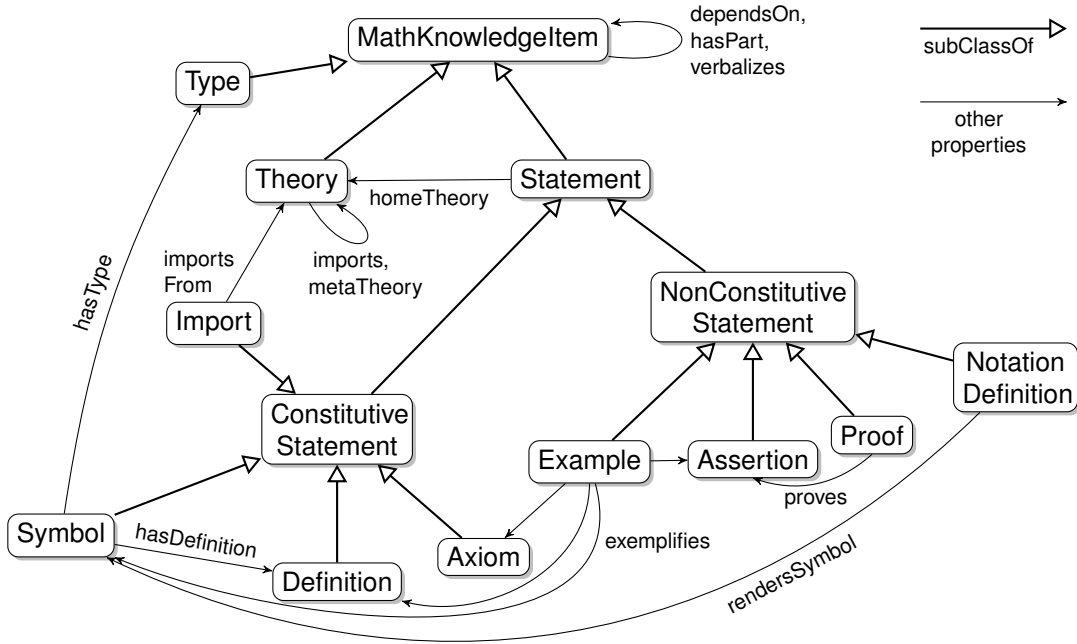


Fig. 3. The core of the OMDoc ontology (slightly simplified) [160]

The *o:usesSymbol* property flattens the functional structure of a mathematical object by treating all occurrences of symbols equally, regardless of the depth of the expression tree in which they occur. For dependency w.r.t. validity, there is so far merely a proof of concept, namely the dependency of a proof on an inference rule or any other axiom or proven assertion used to justify a proof step:

$$\begin{aligned}
 & o:hasStep \circ o:stepJustifiedBy \\
 & \sqsubseteq o:validityDependsOn \sqsubseteq o:dependsOn
 \end{aligned}$$

Note that, with a structural ontology alone, we cannot check *whether* an expression is well-formed, or whether a proof is valid, as we are outside of a particular foundation; the validity itself has to be determined by other means (cf. the discussion in section 3.1).

The OMDoc ontology also covers notation definitions for symbols; for example, the exp symbol from our example could be defined to render as e^x . When an OMDoc document is published, the presentation of any formula using that symbol *possibly* depends on

that notation definition:

$$\begin{aligned}
 & o:usesSymbol \circ o:hasNotationDefinition \\
 & \sqsubseteq o:possiblyUsesNotationDefinition \\
 & \sqsubseteq o:presentationDependsOn \sqsubseteq o:dependsOn
 \end{aligned}$$

By means of the ontology, this cannot be decided definitely, as the knowledge base might have alternative notations for different presentation contexts. Static context matching exceeds the expressivity of a DL ontology, and in a dynamic setting, where the presentation context depends on the profile of the user viewing the published document, its notational dependencies can only be determined at runtime.

MathLang DRa: The “Document Rhetorical aspect” (DRa) of the MathLang representation language covers larger chunks of mathematical text – document sections as well as mathematical statements – and their interrelations, such as a proof justifying a theorem [199,144]. A generic dependency relation has been defined, which is used for validating whether the narrative order of a document respects the logical dependencies. Conceptually, this is similar to the statement level of the OMDoc ontology. The OWL imple-

mentation of the DRa vocabulary merely serves as a formal specification of the DRa semantics, whereas the validator processes an XML representation of the DRa [199]. A drawback of the DRa ontology is that it cannot easily be extended by, e.g., additional statement types and additional dependency relations.

PML: PML (Proof Markup Language), an “interlingua for sharing explanations generated by various automated systems such as hybrid web-based question answering systems, text analytics, theorem proving, task processing, web services execution, rule engines, and machine learning components” [175], has been implemented as an OWL ontology consisting of modules for provenance, information manipulation or justification, and trust. PML assumes that facts and proofs have been written in some other language and merely adds standoff markup. Resources annotated that way can be referenced by URI or, in the case of text-based languages such as KIF, by byte offset. The justification module supports unproven conclusions or goals, assumptions, direct assertions, and antecedent→consequent justifications backed by inference rules. The provenance module has a vocabulary for describing inference rules – again, not down to the object level. So far, this is similar to the OMDoc ontology. Finally, the trust module allows for expressing degrees of belief in informations and trust in agents.

4.3.3. *Scientific Documents*

While logical/functional structures of mathematical knowledge may occur on their own, e.g. in formalized knowledge bases, rhetorical structures are usually studied in the context of documents written in, e.g., \LaTeX or an XML language. Two very similar families of ontologies suitable for modeling rhetorical structures in mathematical documents are SALT [126,127] and OntoReST [183]; further related models and ontologies have been reviewed in [128,77]. SALT and OntoReST are relevant for the following reasons:

- Both have a good coverage of RST-style rhetorical structures.
- Either use case is related to mathematical collaboration: SALT focuses on annotating and linking scientific publications on the Web. OntoReST focuses on consistency checking in concurrent collaborative writing.
- Both allow for an arbitrarily fine-grained annotation of phrases. SALT additionally focuses on cross-document links for justifying statements by citing the claims made [and justified] in external publications [127].

- Both are, in principle, open for integration with arbitrary domain knowledge – which would be mathematical knowledge in our case.

Both approaches comprise three ontologies; here, we explain the model of SALT:

The Document Ontology models the outline of the document – sections, paragraphs, sentences, and text chunks (in OntoReST: “spans”) below sentence level [124]. The latter remain in the original representation of the document; SALT provides standoff markup via start and end pointers to their positions in the full text. Additionally, one can represent the linear order of document units by numbering them.

The Annotation Ontology connects instances of the document ontology with annotations of their rhetorical structure and with background domain knowledge, such as the topic of a section [123]. While rhetorical structures are the primary focus of SALT, the mechanism is sufficiently general to also permit annotation of other structural dimensions.

The Rhetorical Ontology covers RST-style rhetorical relations [125]. Their nuclei and satellites, subsumed as “rhetorical elements”, are linked to text spans in the document via the annotation ontology. Coarse-grained rhetorical blocks that can be applied on top level of a document are offered as an alternative. (This part of SALT has now evolved into the Ontology of Rhetorical Blocks (ORB [84]); the Document Components Ontology (DoCo [206]) provides another recent, more comprehensive alternative.) In previous research, we have investigated the particular suitability of SALT for modeling rhetorical structures in mathematical textbooks by aligning the rhetorical markup of OMDoc (cf. section 4.1.3) to it [160, chapter 3.3]. OntoReST provides a stronger OWL formalization of RST that supports consistency checking [183].

About the above-mentioned DoCo, note furthermore that it is just one member of a family of Semantic Publishing and Referencing Ontologies (SPAR [20]), which additionally cover citations, bibliographies, as well as publishing workflows and people involved.

4.3.4. *Scientific Discourse Ontologies*

Ontologies formalizing discourse inside scientific publications are closely related to the above-mentioned ontologies for rhetorical structures; in fact, SALT sup-

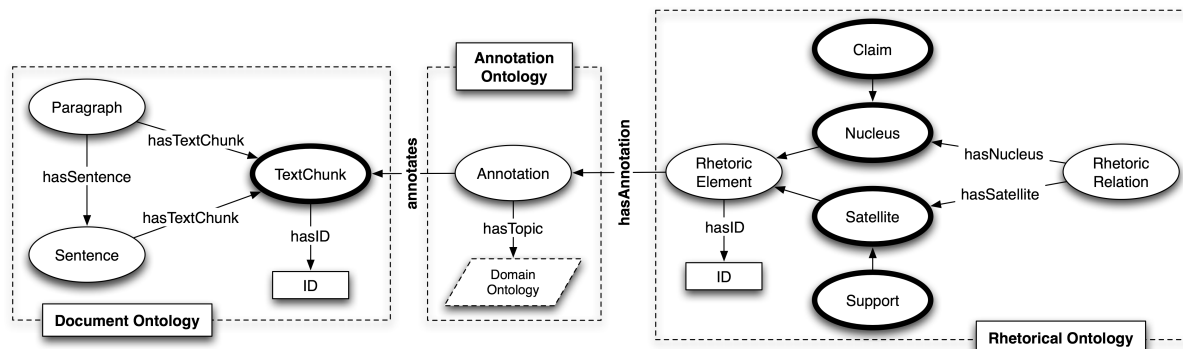


Fig. 4. The three-layered architecture of the SALT ontologies (simplified) [127]

ports both. GROZA et al. provide an overview of further related document formats and ontologies [128].

The DILIGENT argumentation model introduced in section 2.5 has been implemented in several variants [213,214,102]. The SIOC (Semantically Interlinked Online Communities [71,68,70]) ontology, which models user-generated content on the Web and is widely supported by Web 2.0 applications, has a DILIGENT-inspired argumentation module that allows for a slightly more flexible thread structure than the original DILIGENT implementations, which makes it applicable in a wider range of Web 2.0 settings [163]. An ontology of mathematical problems and solutions has been provided as an extension of the SIOC argumentation module [164] – which is the only known occurrence of an ontology specifically capturing (a subset of) *mathematical* discourse.

Combining both models of scientific discourse in one ontology has been pioneered by the alignment of SWAN (Semantic Web Applications for Neuromedicine [128]) with SIOC [191]. SWAN models scientific discourse, not exactly in publications, but in a distributed knowledge base by pointers to bibliographic records and entities from domain ontologies. Its primary target is neuromedicine, but, as SALT, it supports arbitrary domain ontologies in principle.

4.3.5. Mathematical Metadata Vocabularies

Most of the metadata vocabularies mentioned in section 2.3 have been implemented as ontologies. This is the case with Dublin Core [184] and LOM [141, 185]. The mathematics-specific metadata vocabulary of OpenMath CDs (cf. section 4.1.2) has been implemented as an extension to Dublin Core [156]. Their review status is documented by metadata fields such as status (official, experimental, private, or obsolete), version, and the date of the next review.

When a classification scheme is not available as an ontology, one can use the identifiers of its categories as literal values of metadata fields such as *dc:subject*. A proper ontology implementation, where each category is a resource of its own, has further advantages: The hierarchy of categories can be made explicit, and URIs can be used more flexibly in queries. In the MONET project, a simple class hierarchy of the GAMS problems has been implemented [178]. DOLOG et al. have turned the ACM CCS, a computer science classification scheme, into an ontology [106], drawing on the “classification” vocabulary of LOM; the ACM themselves are working on an official Linked Dataset³⁶. Similarly, the MSC 2010 has been implemented as a Linked Dataset, using the SKOS ontology [33] (cf. [44,162]), an official release being expected in early 2012.

4.3.6. Pure and Applied Mathematical Domain Ontologies

Any instance document of the XML languages and formalized languages and any dataset expressed in terms of the structural ontologies reviewed in this section can be considered a mathematical *domain ontology*. In particular, the following knowledge collections have been designed for reuse, reviewed or machine-verified to ensure a high quality, and published with stable identifiers – however, usually neither satisfying the Linked Data criteria nor linked to other collections: The official OpenMath CDs have been mentioned in section 4.1.2. Reusable OMDoc implementations of a large number of logics and translations between them have been published in a “Logic Atlas”, including various variants of first-order, higher-order, modal, and description logic [151,87].³⁷ The collection of mathematical, scientific and technological courseware in

the Connexions repository has been mentioned in section 1.1.1, libraries of formalized mathematics in section 4.2; however, the latter employ custom identification mechanisms that would first have to be translated into URIs. The N3 “math” vocabulary (cf. section 4.3.2) declares a fixed set of basic mathematical functions, roughly corresponding to the *arith1*, *relation1*, and *trans1* OpenMath CDs, without the purpose of formalizing them.

Formalization of mathematical domain knowledge in the decidable first-order logic subsets employed on the Semantic Web has been attempted, but the possibilities are limited. For example, the facts that a differentiable function is a function that has a derivative function and that differentiable functions are continuous functions can be represented in DL, as exemplified by BRÖCHELER [76] – but the fact that a differentiable function satisfies a certain ε/δ criterion cannot, as it would require higher order logic. BRÖCHELER nevertheless points out use cases for querying a DL formalization of domain knowledge in alternation with a representation of logical *structures* in terms of a DL ontology: finding (via the structural ontology) examples for, e.g., groups (instances of a mathematical concept, determined via the domain ontology), or finding applicable theorems or definitions of mathematical concepts and (again via the domain ontology) all related concepts.

In addition to ontologies representing knowledge from the mathematical domain, there are domain ontologies from fields related to mathematics: GAMS (cf. section 4.3.2) features a directory of software that solves mathematical problems [4], but only a subset has been made available within MONET. The SWEET domain ontology for science and GeoSkills for interactive geometry have already been reviewed in section 2.4. The mathematical model ontology of the OntoMODEL tool, also introduced in section 2.4, is particularly notable for its usage of Content MathML (embedded into RDF as XML literals) for equations and variable dependencies [211].

4.3.7. Upper Level Ontologies

Upper level ontologies, also called foundational ontologies, describe general concepts shared by many domains; see [172] for an overview. They often aim at capturing common sense and providing ontological background knowledge to natural language processing applications [172]. As upper level ontologies provide a shared foundation to which designers of domain-specific ontologies can link the latter, they may serve

as a tool for aligning domain-specific ontologies – such as the various ontologies that model structures of mathematical knowledge as well as mathematical domain knowledge. In a scenario where multiple agents perform different actions on a heterogeneous collection of mathematical knowledge, whose representation makes use of multiple domain-specific ontologies, such an alignment helps to reduce misunderstandings among different services (cf. [173]).

As an example, we point out relations between the DOLCE upper level ontology [173] and domain-specific ontologies reviewed before; however, this merely serves as a pointer towards possible future research, as no such alignment has been performed so far for any of the ontologies reviewed. From a **structural** point of view, DOLCE covers, for example, several notions of parthood, which are sufficiently generic to comprise both, e.g., a symbol being part of a theory and a section being part of a document, and thus may serve an agent that explores a knowledge collection along these different structural dimensions. From a mathematical **domain knowledge** point of view, some of DOLCE’s concepts can be considered abstractions of mathematical concepts; for example, a *quality* such as “the value of the sin function for $x = 0$ ” corresponds to a mathematical property of a mathematical object, and its *quale* (here: the value 0), which is a (possibly point-sized) *region* in a *quality space*, corresponds to a member or subset of some set (here: $0 \in \mathbb{R}$).

4.4. Conclusion

Table 1 shows at a first glance that no single semantic markup language satisfies all of our requirements for representing mathematical knowledge on the Semantic Web, but that expressive XML languages and RDF complement each other. OMDoc’s good results have, in fact, been achieved only recently, by integrating RDFa into the language [160, chapter 6]. The XML languages, headed by OMDoc (which includes MathML or OpenMath for mathematical objects), lead the way w.r.t. coverage of different structures of mathematical knowledge, as well as combining formal and informal representations. Moreover, they are reasonably well accepted by the MKM community, as opposed to RDF, and most of today’s mathematical domain knowledge is available in these XML languages, or in formalized languages that have XML translations.

Table 1: How the languages reviewed satisfy the knowledge representation requirements

St. ^a	Requirement	Structures						Formality			Linking			Compr.	
		S.L.*		S.R	S.N	S.M	S.D	F.R	F.C	L.A	L.→	L.←	C.A	C.H	
		O	S	T											
W	MathML 3	++	-	-	-	-	++	++ ³⁹	+	+	+	+	+	+	
W	OpenMath 2 Objects	++	-	-	-	-	+	○	○	-	-	+	+	○	
U	OpenMath 2 CDs	++ ^c	○	○	-	-	○	○	-	-	-	-	-	○	
U	OMDoc 1.2/1.3	++ ^c	++	+	+ ^{+/++^{ef}}	-	++	+	○/++ ^f	- ^{+/++^f}	- ^{+/++^f}	- ^{+/++^f}	- ^{+/++^f}	+	
C	MathLang	++	++	-	-	-	++	+	-	-	-	○	○	-	
~W	DocBook 5	++ ^c	-	-	-	-	-	+ ^g	-	-	+	+	-	-	
~W	TEI P5	++ ^d	-	-	-	-	++	○	+	-	+	+	-	-	
~W	DITA 1.1	++ ^d	-	-	-	-	-	+ ^g	+	+	+	+	-	-	
~W	EPUB 2.0.1/ DTBook 3	++ ^d	-	-	-	-	-	+	-	-	-	+	+	-	
W	CNXML 0.7/ CollXML/mdml	++ ^c	+	-	-	-	-	○	-	-	-	+	-	-	
W	Formalized languages	++	++	+ ^{+/++}	-	-	-	○	○	-	-	-	-	-	
~W	RDF(a) ^h 1.0/1.1	(depends on vocabulary, see table 2)					○	+	++	++	++	++	○	+	

^a Current status of development/deployment: W = widely used; ~W = widely used, but rarely for mathematical knowledge; U = used in a small, open community; C = used in a small, closed community; O = obsolete, no longer in use

^b Legend (summarizing section 3.2): S.* = coverage of knowledge structures (S.L., {O,S,T}) = logical/functional structures: mathematical objects, statements, theories; S, {R,N,M,D} = rigorous language or rhetorical structures, notation, metadata, discussions), F.* = degrees of formality supported (FR = range of degrees, FC = coexistence of different degrees), L.* = ability to link to non-mathematical knowledge (L-A = annotation, L.→ = outgoing links, L.← = incoming links), C, {A,H} = comprehensibility to agents/humans

Symbols: ++ requirement satisfied excellently (and setting an example for other languages), + very well (but leaving room for improvement), ○ sufficiently well for basic modeling tasks, - insufficiently (includes the case that a feature is missing by design)

^c built-in support for MathML/OpenMath objects

^d via MathML extension

^e Dublin Core and similar vocabularies built in, others available via built-in RDFa

^f an improvement of OMDoc 1.3 over the previous version 1.2 that has been achieved by integrating RDFa

^g Dublin Core and similar vocabularies built in, others available via non-RDFa extensions

^h While most features summarized here apply to RDF in general, we specifically consider the serialization of RDF as RDFa embedded into XML here, as it facilitates the annotation of XML markup (requirements L-A and L.→).

Table 2
Structural coverage of vocabularies/ontologies

St. ^a	Structures	Logical/functional			Rhetorical Notation	Metadata	Discussion
		Objects	Statements	Theories			
~W	N3 Vocabularies	+	○	–	–	–	–
O/P ^b	OpenMath CD	○	○	○	–	○	–
O	HELM	+	+	○	–	– ^c	–
O	MoWGLI	+	++	○	–	– ^c	+
P	OMDoc	○	++	+	– ^d	+	–
C	MathLang DRa	–	+	–	–	–	–
P	PML	–	++ ^e	–	–	–	–
P	SALT	–	–	–	++	–	+ ^f
P	OntoReST	–	–	–	++	–	–
P	DILIGENT	–	–	–	–	–	+ ^f
P	SIOC Argumenta- tion	–	–	–	–	–	++
W	Dublin Core	–	–	–	–	–	++ ^g

^a see table 1 for a symbol legend

^b This line merges both OpenMath CD ontologies reviewed; the older one (cf. [79]) is no longer in use.

^c While the XML markup languages employed by HELM and MoWGLI allow for describing notations, their RDF vocabularies do not.

^d intentionally delegated to SALT

^e proofs only

^f does not cover mathematical discourse, but is extensible to specific domains

^g can be combined with mathematical classification schemes

RDF, above all, has superior linking capabilities. Table 2 shows that, given a combination of suitable structural ontologies⁴⁰, RDF is capable of covering all structures of mathematical knowledge. Finally, RDF's linking capabilities and the vast supply of RDF vocabularies for non-mathematical knowledge – such as vocabularies describing application domains of mathematics (cf. the discussion in section 1.3.2), or user profiles (e.g., FOAF [74]), or projects (e.g., DOAP [109]) – provide a unique possibility that is not offered by any other representation language reviewed: modeling and processing mathematical knowledge in the context where it is applied or reused, in a machine-comprehensible way. However, some structures of mathematical knowledge are not yet covered by ontologies as deeply as by markup languages (most prominently theories), and others have hardly been covered at all; the latter affects mathematical discourse (both rhetorical structures and discussions) and mathematical notation. The circumstance that each structural ontology only covers at most two structural dimensions is not a problem, as the RDF data model allows for combining different ontologies by design.⁴¹

In MKM practice, RDF standoff markup pointing to full XML representations has so far proven most useful. Particularly in the case of mathematical objects, only selected information relevant for, e.g., information retrieval, is represented in RDF, as opposed to representing full n -ary ordered trees using RDF collections. RDFa embedded into XML has so far only been used in a few OMDoc documents but also seems a promising way to satisfy the given representation requirements.

Besides these considerations, the availability of tools also advises a division of responsibilities between XML and RDF.⁴² Editors and publishing tools for semantic representations of mathematical knowledge are almost exclusively available for XML or formalized languages. For most XML languages, translations to the native languages of computer algebra systems or proof assistants have been implemented. Conversely, RDF is preferable for information retrieval, except on the object level. Both representations are good to have for a thorough validation; browsers also exist for both.

5. Representing Mathematical Knowledge on the Semantic Web

The previous section has concluded with the finding that both XML and RDF representations of mathematical knowledge are needed on the Semantic Web. This section explains techniques for integrating both. While the concrete examples are taken from representations for mathematical knowledge, the results apply to any domain where semantic XML markup is in use, such as the humanities with TEI.

5.1. Bridging Mathematical XML Markup and Ontologies by Translation

This section explains how to translate XML representations of mathematical knowledge to RDF.

Why not just Combined XML and RDF Queries? When an overall knowledge collection is represented partly in XML and partly in RDF (for example the logical/functional and document structures in XML, and links to discussions and applications in RDF), there is the possibility of employing a query language that supports both representations, such as XSPARQL (cf. section 4.3.1). However, such languages are not yet supported by either XML or RDF databases out of the box. XSPARQL has so far been implemented by rewriting the SPARQL part of a query into XQuery and therefore requires a special execution environment that may not always be easy to provide. In contrast to that, support for *either* XQuery *or* SPARQL queries is wide available, and allows query developers to concentrate on one data model. Therefore, we do not devote further attention to combined XML and RDF queries.

Rationale for Translating XML to RDF: Specifically, we consider XML→RDF translations that enable knowledge collections to be queried as RDF. The reverse translation has the drawback that XML-based querying approaches have less built-in support for abstraction and traversing links in both directions, whereas most triple stores offer a basic level of abstraction (via RDFS entailment) and link traversal in queries (via SPARQL) for free. Assuming, for example, the two OMDoc+RDFa fragments ...

```
<theory about="#t">
  <imports about="#i" from="#u"/>
```

and

```
<proof about="#p">
  <derive about="#step">
    <FMP>¬T = ⊥</FMP>
    <method><!-- proof by known axiom -->
    <premise xref="#axiom1"/>
```

... it would require a considerable effort of declaring, e.g., XML Schema datatypes and implementing XQuery functions to determine that both *#i* depends on *#u* and *#p* depends on *#axiom1*, whereas an OMDoc→RDF translation would generate the triples ...

```
<#t> o:hasImport <#i> .
<#i> o:importsFrom <#u> .
<#p> o:hasStep <#step> .
<#step>
  o:stepExternallyJustifiedBy <#axiom1> .
```

... from which a triple store with DL entailment support would infer ...

```
<#t> o:dependsOn <#u> .
<#p> o:dependsOn <#axiom1> .
```

... and SPARQL would allow for querying these links in both directions.

Requirements for Translating Mathematical XML Markup to RDF: In our previous work on translating OMDoc documents and OpenMath CDs to RDF, we have identified the following general requirements for translating semantic XML markup to RDF, independently from the XML language and the ontology [160, chapter 3.7]:

All structural entities that correspond to concepts covered by the given ontologies **MUST** be given an *identifier* by applying the first of the following rules that matches:

1. If the XML language is an RDFa host language [37], the identifier – URI or blank node ID – **MUST** be determined according to the RDFa processing rules for identifying a *new subject* [37, section 7.5].
2. If the XML language specifies how to generate a URI for an entity represented by an XML fragment, that URI **MUST** be used.
3. If the XML language specifies how to generate an ID for an entity, e.g. via XML ID [171], that ID **MUST** be used as a fragment ID if possible w.r.t. the syntax of URIs [66]; appending it to the document's URI yields the URI.

4. If the XML language specifies how to generate an ID for an entity, which does not qualify as a fragment ID, the translator **MUST** generate a fragment ID, which **SHOULD** reflect the original ID.
5. In any case, the translator **MUST** generate a resource. It **SHOULD** be identified by a minted URI, but it **MAY** also be a blank node with an ID, or an anonymous blank node. Minted URIs **MUST NOT** conflict with URIs generated for other entities in the XML document.

In practice, most semantic XML markup languages support IDs on all elements, but authors only use them when an element is a target of an explicit link in the markup. Many RDF properties, such as whole-part relations, are, however, not represented by explicit XML links but by a parent-child relation, but triples using these properties require identifiable subjects and objects. Also note that manually maintained IDs may not survive refactorings, a common situation in libraries of formalized mathematics. So far there is no ready-to-use solution for this problem, but URBAN has pointed out the problem and suggested automatic generation of identifiers for mathematical objects and statements based on their *content* [217]. This is non-trivial due to the n -ary ordered tree structure of the content but can be made practically manageable by applying a cryptographic hash function to XML representations of such content [217].

For authors and developers, reusing the identifiers from the XML markup in the RDF representation emphasizes the correspondence of both representations. For agents, it improves retrievability, e.g., of RDF standoff markup for an XML representation: If a structural entity always has the same identifier, regardless of the representation format – semantic markup, RDF, or even a human-comprehensible presentation –, and if its different representations are published according to the “cool URI” best practices [201], all of them can be made available under the same URI. A client – agent or browser – would select the desired representation via HTTP content negotiation.

The complexity of semantic XML markup languages for mathematical knowledge entails a number of challenges to the declaration and implementation of an XML→RDF mapping, for example⁴³:

URI Format Differences: OpenMath specifies a canonical URI syntax for symbols. The “namespace base URI”, called *CDBase*, may – and, in practice, usually is – omitted and defaults to `http://www.openmath.org/`. However, when Open-

Math objects occur inside OMDoc theories, the default base URI of a symbol is determined from the theory from which the symbol has been imported.

Mapping Elements to Classes: Generally, OMDoc XML elements correspond to classes from the OMDoc ontology. However, the ontology has been designed with its utility for RDF-based applications in mind, not necessarily to represent the OMDoc XML markup literally. Therefore, instances of some subclasses are represented by the same element, only differing in the value of a certain attribute, or even by elements with different names.

Markup Choices for Representing Relations: Relations between two entities can be represented as a parent-child relation in XML markup, as a sibling relation, or by URI- or ID-valued attributes. As stated for classes above, the exact type of a relation is sometimes influenced by additional attributes on the same element.

Markup Choices for Representing Literal-valued Properties: Literal-valued properties can be represented by text-valued immediate child elements, by descendant elements nested more deeply, or by attributes.

Implicit Structures: The target ontologies reify certain concepts that do not have an explicit representation in the semantic markup. This is, e.g., the case with informal/formal property pairs in OpenMath CDs, and with document units, annotations, and rhetorical relations in the mapping of OMDoc’s rhetorical markup to SALT.

Alternative Representations of Classification Schemes: Where a classification scheme has been implemented as an ontology, its categories are represented as classes or individuals. Otherwise, they are represented as literals. Similarly, metadata with a finite value space can be represented as RDF literals or as instances of an enumerated class.

In our previous work, we have found existing declarative XML→RDF mappings too restricted and instead chose to implement a library of XSLT convenience functions and templates, which facilitates the implementation of frequently occurring translation patterns but gives access to the full power of XSLT if necessary [165,155,159].

5.2. Contributing Mathematics to the Web of Data

Benefits ... of publishing knowledge as Linked Data include easier development of interactive mashups (see, e.g., [135,212]) and the possibility to detect previously unknown links (see, e.g., [136]). Given that mathematical knowledge is likely to be available partly in XML and partly in RDF, as explained in section 4.4, data providers should publish both representations – which is possible, as outlined in section 5.1.

... for agents ... In [219,157], we describe a scenario where an agent accesses both RDF datasets and OpenMath CDs by dereferencing URIs: The rules for computing derived values in statistical datasets are represented as RDF annotations pointing to a function – a symbol from an OpenMath CD – and other values from the dataset that should be passed as arguments to the function, using the SCOVOLink extension of the SCOVO vocabulary [219]. When an agent wants to verify the derived value, it has to construct an OpenMath object from this RDF representation and send it to an OpenMath-aware computation service (cf. sections 1.2.2 and 1.2.3). When the function is not called using positional arguments or an argument list, but using named arguments, the agent has to consult the XML representation of the CD to get their order right. Additionally, the agent can utilize CDs to look up the definitions of functions for which it does not have built-in support.

... and humans: The semantic representations for mathematical knowledge reviewed in this article allow for preserving the full semantics in documents published for human readers, so that, e.g., assistive services can utilize it. For anything except mathematical objects, i.e. formulae, XHTML+RDFa is a suitable publication format. Assistive services for formulae have previously been realized for Presentation MathML with Content MathML or OpenMath annotations [116,190]. With HTML5 becoming mainstream [139], which includes MathML without requiring the strict XML conformance that authors and widget toolkits often fail to achieve, more browsers can soon be expected to support MathML.⁴⁴ We have implemented a library that publishes OMDoc as XHTML+MathML+RDFa [98]. Based on that, we have realized interactive declaration lookup for symbols in formulae by dereferencing their canonical URIs (pointing to a symbol declaration in a CD) from the formula's annotation and transforming the OpenMath

declarations thus obtained to human-readable Presentation MathML.

Challenges: In our previous work on publishing and consuming OMDoc documents and OpenMath CDs as Linked Data [98,157], we have identified three challenges to publishing mathematical knowledge as Linked Data: specifying MIME types for XML languages, bad practices of authors, and restrictions in the URI formats of XML languages.

The HTTP Content Negotiation mechanism outlined in section 5.1 distinguishes representation formats by MIME type. MathML 3, for example, has officially registered MIME types [52], OMDoc specifies an unofficial one [147], whereas MIME types for OpenMath objects and CDs have merely been proposed so far [157].

Authoring practices that are bad from a Linked Data point of view⁴⁵ result from the fact that, where semantic XML languages for representing mathematical knowledge support URIs, authors use them wrongly or not at all. For example, hardly any OpenMath CD that has been contributed to `openmath.org` specifies a CDBase URI or references symbols by full URIs, which indicates a lack of awareness. The fallback value `http://www.openmath.org/` is not suitable for non-official CDs mainly used by one research group, even independently from Linked Data considerations, as they do not control the `openmath.org` domain. Finally, if authors are aware of the fact that CDs and symbols have URIs, they usually merely consider it a globally unique *name*, but not a means of retrieving information about these resources [157].

Thirdly, while RDF publishers can freely choose URIs (cf., e.g., [67,137,105]), the URI formats of non-RDF languages often impose restrictions that complicate Linked Data publishing and have to be worked around. For example, OpenMath's schema of `cdbase/cd#name` symbol URIs, which has also been adopted by Strict Content MathML, complies well with linked data practices – unless CDs grow large. As resolving fragments after the # (hash) in a URI is up to the client, the consequent use of hash URIs for OpenMath symbols forces clients to always download a complete CD from the server, in which it could then locate the symbol with the desired name. Publishers of large CDs would have to set up a redirect, where an initial request for a hash URI would result in an RDF graph that merely redirects, via *rdfs:seeAlso* links, hash URIs to slash URIs, from which the client would be able to retrieve the desired fine-grained information.

The URIs of the upcoming OMDoc 1.6 have a slash-like format [196], but, again, without alternatives. Another problem of OMDoc 1.3 and its hash URI format is that symbols and theories have to be declared as fragments of the same document, which is not compatible with OpenMath’s *cdbase/theory#symbolname* schema, even though OMDoc uses OpenMath objects. Combined with the possibility of having multiple theories in a document, redirect workarounds may not be possible.⁴⁶ A final problem with old languages such as the OpenMath CD language, is that certain entities – including mathematical properties – cannot be given IDs. An XML→RDF translator might generate some, but an agent interested in retrieving XML representations would also need them. As a use case that would require such fine-grained links, consider the DLMF [3]. It contains a large number of equations describing or defining mathematical functions, which could be linked to the corresponding mathematical properties in OpenMath CDs.

6. Research Directions Towards a Mathematical Web of Data

Large collections of mathematical knowledge exist in non-RDF representations – including some that are already machine-comprehensible (cf. section 4.3.6), but much larger ones that are not yet. The ontologies reviewed in section 4.3 and the translation techniques outlined in section 5 now enable us to contribute them to the Web of Data and fill a gap that existing Linked Datasets about, e.g., statistical government data or scientific publications, have left. The roadmap outlined in this section is primarily presented from this Linked Data point of view but also covers other aspects of where there is potential for advancing the state of the art.

6.1. Bootstrapping the Mathematical Web of Data by Publishing the MSC and the OpenMath CDs

Two of the most foundational datasets that have to be published as Linked Data in order to get mathematics on the Web of Data started are the MSC (due to its wide use in digital libraries) and the official OpenMath CDs (which define the semantics of Content MathML 3). Initial Linked Data implementations proved feasible with the technology available and have been finished (cf. [44,162] for the MSC and for the OpenMath CDs [159]). This is, however, only the first step

in making the knowledge contained in these datasets accessible.

The second step is mutually interlinking them, and linking to them from mathematics-related existing datasets, so that services for these existing datasets can be extended by mathematical functionality. The use case sketched in section 5.2 points out how statistical datasets can benefit from links to OpenMath CDs. The inevitable DBpedia [99] is a further candidate for linking, with two expected benefits: (i) DBpedia→OpenMath links would give DBpedia’s large audience a more formal perspective on mathematics, whereas (ii) OpenMath→DBpedia links would help to relate mathematical formulae to non-mathematical background knowledge, such as the history of the respective area of mathematics, or its applications in industry. Note that most of these links across mathematical datasets will have to be established from scratch; even links across existing pre-linked-data collections of formalized mathematics hardly exist to date [39].

Mathematical knowledge collections that are already available on the Web, but not currently in a semantic representation, should also be semantically annotated – not necessarily as deeply as, e.g., OMDoc documents, but at least with mathematical metadata and links to relevant OpenMath CDs. For example, the DLMF [3] could benefit from access to computation services via OpenMath, whereas the benefit for PlanetMath, which is currently being overhauled to make it more interactive and more semantic [152], would be similar as for DBpedia. Note, however, that the possibility of fine-grained links across mathematical resources entails the unsolved challenge of how to *present* the target of a link to a human reader. When the target resource is to be displayed in the context of the link source, should the presentation context be determined for the source document, or for the target document?

A possible gateway into annotating the mathematical semantics of scientific publications is the arXiv [46]. Its documents are mostly available as presentation-oriented \LaTeX ; however, a long-term effort to automatically annotate their mathematical structure using linguistic techniques is in progress [117], the translation of 300,000 of the 500,000 publications to XHTML+MathML, which has at least more semantic structure than \LaTeX , being a first success [210]. Publishing a basic metadata record for each arXiv publication as Linked Data is feasible, as the metadata are available as XML, and the publications have sta-

ble URIs. Next, much harder steps would be interlinking with publication databases already existing as Linked Data, such as DBLP [2], identifying mathematical symbols that could be linked to the OpenMath CDs, and further on automatically identifying statement- and theory-level logical/functional structures. While these steps benefit from the availability ontologies for these structural aspects, scientific discourse in mathematics as well as mathematical notation have hardly been covered by ontologies so far.

On interlinking mathematical datasets, recall once more their differing degrees of formality. One can expect datasets of formalized mathematics to represent, e.g., theory morphisms across logics as faithfully as possible in RDF, whereas links between informal datasets (e.g. between PlanetMath and DBpedia) will hardly have a stronger semantics than the catch-all *rdfs:seeAlso*. Links across degrees of formality, e.g. employing the verbalizes/formalizes properties of OMDoc (cf. section 4.3.2), entail a specific challenge: In the past, they have only been employed with representations originating from the same source, but in the Linked Open Data cloud, one might also want to express that a DBpedia article “roughly” verbalizes a Mizar article – with some differences in notation and terminology. Such scenarios may require a more differentiated linking vocabulary.

6.2. Possibilities for Mathematical Computation and Reasoning

The availability of true mathematical knowledge as Linked Data would also allow for taking a serious view on the April fool’s joke “Linked Open Numbers”, a huge dataset describing billions of natural numbers [218]. It provided descriptions as trivial as the name of each number in natural language, its predecessor and its successor. But how about a dataset of non-trivial properties of numbers? Accessing, for example, prime factor decompositions of large numbers – an information relevant for cryptography – in a linked dataset, could be much faster than computing it once more, provided a supercomputer has already done the computation once and published the results. Another source of non-trivial knowledge about numbers, which deserves being published as Linked Data, is the Online Encyclopedia of Integer Sequences [207].

An issue related to the combination of information retrieval and computation is the development of suitable query processors and reasoners, which operate on large-scale RDF graphs linking mathematical and non-

mathematical knowledge but are also capable of mathematical reasoning and computation. N3 reasoners already support a limited set of mathematical functions (cf. section 4.3.2). The upcoming SPARQL 1.1 supports basic arithmetics. Additionally, many query processors allow for defining extension functions; a path for supplying arbitrary functions to query processors via OpenMath should be investigated. Taking a formal semantics and computational complexity into account, such an extension could even be specified as an entailment regime [118]⁴⁷, which makes a query return a well-defined set of additional, *entailed* results beyond the information that is explicitly encoded in the RDF graph being queried; at the same time, the basic properties of such an extension could be *described* using the SPARQL service description vocabulary [230] or an OpenMath-aware extension thereof. Combining RDF queries and mathematical reasoning, such as proof checking, would complement the inferring potential of structural ontologies outlined in section 4.3.2. This challenge can possibly also be addressed by making specialized tools for mathematical reasoning accessible from the RDF world via entailment regimes.

6.3. Conclusion

Even without these (non-trivial) steps to bridge Semantic Web querying and reasoning and the much stronger (and much less scalable) first-order and higher-order calculi required to verify logical/functional structures of mathematics such as theorems or proofs, contemporary Semantic Web and Linked Data technology already provides a solution to managing mathematical knowledge not just on its own but in the wider context of its reuse and application. As, however, experts in the mathematical domain tend to lack technical expertise in Linked Data publishing, the process of publishing existing collections of mathematical knowledge needs to be supported by automatic translation from the languages that these experts are more fluent in: languages for formalized mathematics (reasonably easy via XML; cf. section 4.2), and \LaTeX (hard for plain \LaTeX , as it requires linguistic techniques, cf. section 6.1; reasonably easy with \LaTeX extensions that provide semantic markup, such as $\mathcal{S}\TeX$ [148]). Secondly, further applications need to be developed to convince domain experts of the promises of a mathematical Web of Data; however, Linked Data technology can in turn facilitate the development of integrated service platforms for science, technology, engineering,

and mathematics, which attract domain experts with a Web 2.0 interface that they are already familiar with (cf. [152] for a proof of concept).

A mathematical Web of Data would not only provide mathematicians with better information retrieval support for the next collaborative Web-based review of a $P \neq NP$ proof or with better social interaction support in the next collaborative effort to formalize a proof of a theorem like the Kepler Conjecture – consider particularly newcomers who are not yet familiar with the structure of a formalized library and would therefore appreciate guidance by simplified annotations and links. A mathematical Web of Data would also boost already successful Web of Data applications to statistics, e-science, business, etc., by taking into account their mathematical foundations.

Acknowledgments: The author would like to thank FLORIAN RABE and MICHAEL KOHLHASE for their help with defining the scope of this survey between knowledge management and automated reasoning, furthermore MICHAEL KOHLHASE for his explanations on rhetorical structures of mathematical knowledge, as well as the reviewers (CLAUDIO SACERDOTI COEN, ALEXANDRE PASSANT, and ALDO GANGEMI) for their insightful and constructive remarks.

Notes

¹We are not aware of any mathematical Web 2.0 site that integrates formal verification.

²Part of this problem has been solved by semantic formula search engines based on MathML or OpenMath, such as MathWebSearch [150], which employs substitution tree indexing. A complete solution would additionally require the term rewriting capabilities of a computer algebra system. Integrating MathWebSearch into a web-based publishing environment is currently in progress (cf. [152]); less powerful formula search engines, which employ full text indexing and therefore lose more of the semantic structure, are already in use in publication environments – for example in the ActiveMath e-learning system [168].

³ProgrammableWeb [16], a directory of mashups, lists 3 mashups tagged with “math”, out of more than 6,000 mashups overall. The recently released “widgets” for the Wolfram Alpha “computational knowledge engine” [28] are a first step towards more mashups, albeit limited to acting as frontends to Wolfram Alpha.

⁴This notion of the term “knowledge management” is wider than that of its traditional definition as “a range of practices used in an organisation to identify, create, represent, distribute and enable adoption of insights and experiences. Such insights and experiences comprise knowledge, either embodied in individuals or embedded in organisational processes or practice.” [225]

⁵enumeration added by the author

⁶MathNet which actually featured the first working implementation of Dublin Core [193]!

⁷The HELM requirements were independence of a concrete RDF syntax (such as RDF/XML), disjunction, data source identification, and a well-defined formal semantics [129].

⁸The HELM developers made no secret out of their frustration: “It is a pity that [...] most of the expectations about XML technologies [including RDF, which the authors treat as an ‘XML-based technology’] have not been fulfilled due to intrinsic deficiencies in their design and implementation. MathML failed to be adopted by major browsers; [...] and RDF never really went beyond the project phase.” [51] Personal communication with ASPERTI on 2010-07-09 confirmed that that statement referred to the immaturity of these technologies at the time of developing HELM.

⁹personal communication with MICHAEL KOHLHASE on 2010-11-12

¹⁰Similarly, BAEZ suggests that the release of a \TeX formula editor plugin for the popular WordPress blog engine was a major incentive for mathematicians to start blogging [55].

¹¹That can be accommodated for, as explained in section 5.1.

¹²reusing the terminology introduced by KOHLHASE in [147, chapters 2.3 and 3.2] and refined in [145]

¹³Definitions typically occur in textbook-style mathematics. From a logical point of view they are merely a variant of axioms.

¹⁴Note that formalization is not necessarily a linear process. While each step of writing down a sloppy formula more rigorously, and then formalizing it in a certain mathematical foundation produces a “more formal” result in that it restricts the space of possible interpretations, one cannot say that a formalization in one foundation is more formal than a formalization in another foundation.

¹⁵In fact, mathematical notation has an infinitely extensible vocabulary as well as grammar. In contrast, consider musical notation, which is also a two-dimensional notation, but has a largely fixed vocabulary and grammar.

¹⁶This can also be considered a difference w.r.t. the area of application. For example, in theoretical computer science it is advantageous to include 0, as many of the required induction proofs start at 0, whereas negative integers are rarely needed.

¹⁷We note that on these levels the above-mentioned context dimensions also have an influence – but not on *how* a mathematical knowledge item is presented to the users, but *what* knowledge items are chosen: for example *which* definition of the same concept, which proof for the same theorem, or which example for the same thing (see, e.g., [146] on the context-sensitivity of examples, and [182] on generating documents from snippets using contextual information).

¹⁸Capitalized keywords are used in accordance with RFC 2119 [73].

¹⁹Listing 1 gives a partial example (notice the usage of the *term* and *phrase* elements); however, for reasons of space we refer to [160, listing 4.1] for a fully elaborated example and to [147, chapter 14.4] for documentation.

²⁰The possibility of giving implicit definitions actually depends on the foundation having a definite description operator, as discussed in [149].

²¹Due to the design of the MathML schema, supporting MathML always comprises both Presentation and Content MathML. The TEI P5 guidelines explicitly mention both (cf. [78, chapter 14.2]); CNXML explicitly recommends Content MathML [89].

²²That is where the reference to CHARLES DARWIN comes from.

²³DAISY = Digital Accessible Information Systems

²⁴This conjecture, posed in 1611, states that the density of a packing of unit spheres in 3 dimensions is at most $\pi/(3\sqrt{2})$. This reflects the intuitive observation that the way, in which, e.g., oranges

in a market booth are stacked, is optimal. However, it turned out exceedingly complex to prove.

²⁵In an OWL setting, one has to avoid RDF collections, as the RDF encoding of OWL uses them internally for representing n -ary DL expressions. Instead, one has to create one's own linked lists [107].

²⁶At least support for querying RDF collections, which some query processors already support by non-standard extensions, will be standardized in the upcoming SPARQL 1.1 [134].

²⁷See [160, chapter 5] for a discussion of concrete examples from OMDoc.

²⁸His encoding differs from the N3 encoding in that order is represented using RDF's built-in container membership properties *rdf:_n* ($n = 1, 2, \dots$) instead of RDF collections, but that is a secondary issue.

²⁹A possible explanation in MARCHIORI's case is that his proposal did not originate out of the MKM community but that he was an external (Semantic Web) expert invited to give a keynote, which consisted of a rather ad hoc sketch of possible applications of Semantic Web technology to MKM.

³⁰The depth of a symbol in *h:MainHypothesis* position is computed relatively to the hypothesis (which may have its own universal quantifiers).

³¹HELM uses relative XPath's pointing into an XML document that is assumed to contain one theory; cf. [202, example 6.1].

³²Each *hth:TheoryItem* has an *hth:ident* property pointing to an *hth:HelmlD* resource that is identified relatively to the current RDF graph, like a blank node ID. Dependencies (*hth:dependence*) are expressed indirectly as links between theory items and the HELM IDs of dependent theory items.

³³ambiguities and errors in its own modeling, tampering with the semantics of reused vocabularies (such as DCMES), limited documentation, and use of bad RDFS modeling practices (cf. [160, chapter 2.4.10.2] for detailed examples)

³⁴The ontology does not cover complex theory morphisms, abstract datatypes, and presentation contexts.

³⁵The names *CMP* = Commented Mathematical Property and *FMP* = Formal Mathematical Property are for historical reasons and OpenMath compatibility.

³⁶personal communication with BERNARD ROUS, ACM Director of Publications, 2011-06-08

³⁷Another large collection of knowledge represented in OMDoc – proven in use but neither reviewed nor validated – is formed by KOHLHASE's computer science lecture notes [148]; in contrast to the Logic Atlas, they are in textbook style. A proof-of-concept Linked Data version of them has been developed [98], which is, however, not sufficiently stable for reuse, as both its URI format and its underlying implementation are experimental.

³⁸While the official OpenMath CDs, whose symbols are commonly used in parallel MathML markup, do not fully specify the formal semantics of their symbols, the parallel markup *mechanism* is open to arbitrary CDs, including CDs with a stronger semantics, implemented e.g. in OMDoc.

³⁹While the official OpenMath CDs, whose symbols are commonly used in parallel MathML markup, do not fully specify the formal semantics of their symbols, the parallel markup *mechanism* is open to arbitrary CDs, including CDs with a stronger semantics, implemented e.g. in OMDoc.

⁴⁰This table excludes some of the ontologies reviewed in section 4.3: The MONET Problem Ontology does not model the mathematical structures that are of interest here but serves as an example of

a division of responsibilities between an RDF and OpenMath XML representation, classification schemes can only be applied to mathematical documents via a bibliographical metadata ontology such as Dublin Core, and domain ontologies do not cover mathematical structures.

⁴¹For example, a combination of ontologies for mathematical objects (e.g. N3), statements (e.g. the OMDoc ontology), and rhetorical structures (e.g. SALT) would allow for reproducing the literate programming style of OMDoc or MathLang in RDF, but with a considerably larger effort, once more due to the fact that order and n -ary structures are implicitly supported by XML but have to be modeled explicitly in RDF.

⁴²[160, chapter 6] provides a comprehensive overview.

⁴³see [160, chapter 3.7] for details

⁴⁴At the moment, only Mozilla/Firefox supports MathML well enough to allow for interactive manipulation via scripts.

⁴⁵For an overview of Linked Data best practices, see [137,105].

⁴⁶See [160, chapter 6.4.1.3] for a detailed discussion.

⁴⁷This possibility has been pointed out by DENNY VRANDEČIĆ (personal communication, 2010-06-02).

References

- [1] Specifications for the Digital Talking Book. Technical report, DAISY Consortium, 2005. URL <http://www.daisy.org/z3986/2005/Z3986-2005.html>. ANSI/NISO Z39.86-2005. ISSN 1041-5653.
- [2] D2R server publishing the DBLP bibliography database. URL <http://dblp.l3s.de/d2r/>.
- [3] Digital Library of Mathematical Functions. National Institute of Standards and Technology (NIST), 2010. URL <http://dlmf.nist.gov>.
- [4] GAMS: Guide to Available Mathematical Software. National Institute of Standards and Technology (NIST). URL <http://gams.nist.gov>.
- [5] HELM. Hypertextual Electronic Library of Mathematics. URL <http://helm.cs.unibo.it>.
- [6] Math-Net, an International Information and Communication System. International Mathematical Union (IMU). URL <http://www.math-net.org>.
- [7] MathOverflow. URL <http://mathoverflow.net>.
- [8] Eric W. Weisstein, editor. Wolfram MathWorld, the web's most extensive mathematics resource. URL <http://mathworld.wolfram.com>.
- [9] Mizar mathematical library. URL <http://www.mizar.org/library>.
- [10] Mizar system. URL <http://mizar.org/system/>.
- [11] OEBPS Container Format (OCF). Recommended specification, International Digital Publishing Forum, 2006. URL <http://www.idpf.org/ocf/ocf1.0/download/ocf10.htm>.
- [12] OMDoc. URL <http://omdoc.org>.
- [13] Open Packaging Format (OPF), version 2.0 v1.0. Recommended specification, International Digital Publishing Forum, 2007. URL http://www.idpf.org/2007/opf/OPF_2.0_final_spec.html.
- [14] The polymath blog. URL <http://polymathprojects.org/>.
- [15] Deolalikar P vs NP paper. 2010. URL <http://michaelnielsen.org/polymath1/index.php?>

Table 3
 Namespace prefix→URI bindings used in this article

Prefix	URI	Language/Ontology
<i>dc</i>	http://purl.org/dc/elements/1.1/	Dublin Core Metadata Element Set
<i>g</i>	http://gams.nist.gov#	GAMS
<i>h</i>	http://www.cs.unibo.it/~schena/schema-h.rdf#	HELM objects
<i>hth</i>	http://www.cs.unibo.it/~schena/schema-hth.rdf#	HELM theories
<i>math</i>	http://www.w3.org/2000/10/swap/math#	N3 math functions
<i>o</i>	http://omdoc.org/ontology#	OMDoc ontology
<i>oms</i>	http://www.openmath.org/cd/	OpenMath symbols
<i>p</i>	http://monet.nag.co.uk/problems/	MONET problems

- title=Deolalikar_P_vs_NP_paper&oldid=3654.
- [16] ProgrammableWeb. Mashups, APIs, and the Web as Platform. URL <http://www.programmableweb.com>.
- [17] ProofWiki. URL <http://www.proofwiki.org>.
- [18] acm.rkbexplorer.com. Advanced Knowledge Technologies (AKT). URL <http://acm.rkbexplorer.com>.
- [19] Rhaptos Trac: Collection structure redesign / inception. 2009. URL <https://trac.rhaptos.org/trac/rhaptos/wiki/CollectionStructureRedesign/Inception?version=54>.
- [20] SPAR – Semantic Publishing and Referencing. 2011. URL <http://purl.org/spar>.
- [21] Semantic Web for Earth and Environmental Terminology (SWEET). NASA, 2011. URL <http://sweet.jpl.nasa.gov/>.
- [22] Semantic MediaWiki. URL <http://semantic-mediawiki.org>.
- [23] SysMO-DB SEEK. URL <http://www.sysmo-db.org/seek/>.
- [24] TEI – ontologies SIG. URL <http://www.tei-c.org/Activities/SIG/Ontologies/>.
- [25] Tricki. A repository of mathematical know-how. URL <http://www.tricki.org>.
- [26] Dublin Core Metadata Initiative. URL <http://www.dublincore.org>.
- [27] OPENMATH content dictionaries. URL <http://www.openmath.org/cd/>.
- [28] WolframAlpha widgets. URL <http://developer.wolframalpha.com/widgets/>.
- [29] Connexions – XML languages. URL <http://cnx.org/help/authoring/xml>.
- [30] The Coq proof assistant. URL <http://coq.inria.fr/>.
- [31] The n-Category Café. A group blog on math, physics and philosophy. URL <http://golem.ph.utexas.edu/category/>.
- [32] nLab. URL <http://ncatlab.org/>.
- [33] SKOS Simple Knowledge Organization System. World Wide Web Consortium (W3C). URL <http://www.w3.org/2004/02/skos/>.
- [34] Zentralblatt MATH. URL <http://www.zentralblatt-math.org>.
- [35] Mathematics Subject Classification MSC2010, 2010. URL <http://msc2010.org>.
- [36] ActiveMath. ACTIVEMATH. URL <http://www.activemath.org>.
- [37] Ben Adida, Mark Birbeck, Shane McCarron, and Ivan Herman. RDFa Core 1.1. Syntax and processing rules for embedding RDF through attributes. W3C Working Draft, World Wide Web Consortium (W3C), March 2011. URL <http://www.w3.org/TR/2011/WD-rdfa-core-20110331/>.
- [38] Waseem Akhtar, Jacek Kopecký, Thomas Krennwallner, and Axel Polleres. XSPARQL: Traveling between the XML and RDF worlds – and avoiding the XSLT pilgrimage. In Bechhofer et al. [59], pages 432–447.
- [39] Jesse Alama, Kasper Brink, Lionel Mamane, and Josef Urban. Large Formal Wikis: Issues and Solutions. In Davenport et al. [96], pages 133–148.
- [40] Keith Alexander, Richard Cyganiak, Michael Hausenblas, and Jun Zhao. Describing Linked Datasets with the VoID Vocabulary. W3C Interest Group Note, World Wide Web Consortium (W3C), March 2011. URL <http://www.w3.org/TR/2011/NOTE-void-20110303/>.
- [41] American Mathematical Society. MathSciNet Mathematical Reviews on the Net. URL <http://www.ams.org/mathscinet/>.
- [42] Ola Andersson, Robin Berjon, Erik Dahlström, Andrew Emons, Jon Ferraiolo, Anthony Grasso, Vincent Hardy, Scott Hayman, Dean Jackson, Chris Lilley, Cameron McCormack, Andreas Neumann, Craig Northway, Antoine Quint, Nandini Ramani, Doug Schepers, and Andrew Shellshear. Scalable Vector Graphics (SVG) Tiny 1.2 specification. W3C Recommendation, World Wide Web Consortium (W3C), December 2008. URL <http://www.w3.org/TR/2008/REC-SVGTiny12-20081222/>.
- [43] Anupriya Ankolekar, Markus Krötzsch, Thanh Tran, and Denny Vrandečić. The two cultures: Mashing up Web 2.0 and the Semantic Web. *Web Semantics*, 6(1):70–75, 2008.
- [44] Ioannis Antoniou, Charalampos Bratsas, Anastasia Dimou, Patrick Ion, Christoph Lange, and Wolfram Sperber. Mathematics Subject Classification Linked Wiki. URL <http://sci-class.math.auth.gr/MSCLW/>.
- [45] Lora Aroyo, Grigoris Antoniou, Eero Hyvönen, Annette ten Teije, Heiner Stuckenschmidt, Liliana Cabral, and Tania Tudorache, editors. *The Semantic Web: Research and Applications (Part I)*, number 6088 in Lecture Notes in Computer Science, 2010. Springer Verlag, Berlin/Heidelberg.

- [46] ArXiv. [arxiv.org](http://www.arxiv.org) e-Print archive. URL <http://www.arxiv.org>.
- [47] Andrea Asperti and Claudio Sacerdoti Coen. Some considerations on the usability of interactive provers. In Autexier et al. [54], pages 147–156.
- [48] Andrea Asperti, Bruno Buchberger, and James Harold Davenport, editors. *Mathematical Knowledge Management, MKM'03*, number 2594 in Lecture Notes in Computer Science, 2003. Springer Verlag, Berlin/Heidelberg.
- [49] Andrea Asperti, Luca Padovani, Claudio Sacerdoti Coen, Ferruccio Guidi, and Irene Schena. Mathematical knowledge management in HELM. *Annals of Mathematics and Artificial Intelligence, Special Issue on Mathematical Knowledge Management, Kluwer Academic Publishers*, 38(1–3):27–46, May 2003.
- [50] Andrea Asperti, Claudio Sacerdoti Coen, Enrico Tassi, and Stefano Zacchiroli. User interaction with the Matita proof assistant. *Journal of Automated Reasoning*, 39(2):109–139, 2007.
- [51] Andrea Asperti, Herman Geuvers, and Raja Natarajan. Social processes, program verification and all that. *Mathematical Structures in Computer Science*, 19(5):877–896, October 2009.
- [52] Ron Ausbrooks, Stephen Buswell, David Carlisle, Giorgi Chavchanidze, Stéphane Dalmas, Stan Devitt, Angel Diaz, Sam Dooley, Roger Hunter, Patrick Ion, Michael Kohlhase, Azzeddine Lazrek, Paul Libbrecht, Bruce Miller, Robert Miner, Murray Sargent, Bruce Smith, Neil Soiffer, Robert Sutor, and Stephen Watt. Mathematical Markup Language (MathML) version 3.0. W3C Recommendation, World Wide Web Consortium (W3C), 2010. URL <http://www.w3.org/TR/MathML3>.
- [53] Serge Autexier, Dieter Hutter, Till Mossakowski, and Axel Schairer. Maya: Maintaining structured documents. In OMDOC – *An open markup format for mathematical documents [Version 1.2]* Kohlhase [147], chapter 26.12. URL <http://omdoc.org/pubs/omdoc1.2.pdf>.
- [54] Serge Autexier, Jacques Calmet, David Delahaye, Patrick D. F. Ion, Laurence Rideau, Renaud Rioboo, and Alan P. Sexton, editors. *Intelligent Computer Mathematics*, number 6167 in Lecture Notes in Artificial Intelligence, 2010. Springer Verlag, Berlin/Heidelberg. ISBN 3642141277.
- [55] John Baez. Math blogs. *Notices of the AMS*, page 333, March 2010. URL <http://www.ams.org/notices/201003/rtx1003003333p.pdf>.
- [56] Grzegorz Bancerek. Information retrieval and rendering with MML Query. In Borwein and Farmer [72], pages 266–279.
- [57] Rebhi S. Baraka. *A Framework for Publishing and Discovering Mathematical Web Services*. PhD thesis, Johannes Kepler Universität Linz, 2006. URL http://www.risc.uni-linz.ac.at/publications/download/risc_2945/06-04.pdf.
- [58] Michael J. Barany. ‘[b]ut this is blog maths and we’re free to make up conventions as we go along’: Polymath1 and the modalities of ‘massively collaborative mathematics’. In Phoebe Ayers and Felipe Ortega, editors, *Proceedings of the 6th International Symposium on Wikis and Open Collaboration (WikiSym)*, ACM Press, 2010. URL <http://www.wikisym.org/ws2010/Proceedings/>.
- [59] Sean Bechhofer, Manfred Hauswirth, Jörg Hoffmann, and Manolis Koubarakis, editors. *The Semantic Web: Research and Applications*, number 5021 in Lecture Notes in Computer Science, 2008. Springer Verlag, Berlin/Heidelberg.
- [60] Sean Bechhofer, John Ainsworth, Jiten Bhagat, Iain Buchan, Philip Couch, Don Cruickshank, David De Roure, Mark Delderfield, Ian Dunlop, Matthew Gamble, Carole Goble, Darius Michaelides, Paolo Missier, Stuart Owen, David Newman, and Shoab Sufi. Why linked data is not enough for scientists. In *6th IEEE e-Science conference*, 2010, pages 300–307.
- [61] Dave Beckett. RDF/XML syntax specification (revised). W3C Recommendation, World Wide Web Consortium (W3C), February 2004. URL <http://www.w3.org/TR/2004/REC-rdf-syntax-grammar-20040210/>.
- [62] Anders Berglund, Scott Boag, Don Chamberlin, Mary F. Fernández, Michael Kay, Jonathan Robie, and Jérôme Siméon. XML Path Language (XPath) 2.0. W3C Recommendation, World Wide Web Consortium (W3C), 2007. URL <http://www.w3.org/TR/2007/REC-xpath20-20070123/>.
- [63] Tim Berners-Lee. cwm. 2009. URL <http://www.w3.org/2000/10/swap/doc/cwm.html>.
- [64] Tim Berners-Lee. Notation 3 (N3) – a readable RDF syntax. Technical report, World Wide Web Consortium (W3C), 2006. URL <http://www.w3.org/DesignIssues/Notation3.html>.
- [65] Tim Berners-Lee. Design issues: Linked data. Technical report, World Wide Web Consortium (W3C), 2006. URL <http://www.w3.org/DesignIssues/LinkedData.html>.
- [66] Tim Berners-Lee, Roy T. Fielding, and Larry Masinter. Uniform resource identifier (URI): Generic syntax. RFC 3986, Internet Engineering Task Force (IETF), 2005. URL <http://www.ietf.org/rfc/rfc3986.txt>.
- [67] Diego Berrueta and Jon Phipps. Best practice recipes for publishing RDF vocabularies. W3C Working Group Note, World Wide Web Consortium (W3C), August 2008. URL <http://www.w3.org/TR/2008/NOTE-swbp-vocab-pub-20080828/>.
- [68] Diego Berrueta, Dan Brickley, Stefan Decker, Sergio Fernández, Christoph Görn, Andreas Harth, Tom Heath, Kingsley Idehen, Kjetil Kjernsmo, Alistair Miles, Alexandre Passant, Axel Polleres, and Luis Polo. *SIOC Core Ontology Specification*, March 2010. URL <http://rdfs.org/sioc/spec/>.
- [69] Chris Bizer, Sören Auer, and Gunnar Aastrand Grimnes, editors. *Scripting and Development for the Semantic Web (SFSW)*, number 449 in CEUR Workshop Proceedings, Aachen, May 2009. URL <http://CEUR-WS.org/Vol-449/>.
- [70] Uldis Bojars and John G. Breslin. SIOC core ontology specification. W3C Member Submission, World Wide Web Consortium (W3C), June 2007. URL <http://www.w3.org/Submission/2007/SUBM-sioc-spec-20070612/>.
- [71] Uldis Bojars, John G. Breslin, Vassilios Peristeras, Giovanni Tummarello, and Stefan Decker. Interlinking the Social Web with Semantics. *IEEE Intelligent Systems*, 23(3), pages 29–40, 2008.
- [72] Jon Borwein and William M. Farmer, editors. *Mathematical Knowledge Management (MKM)*, number 4108 in Lecture Notes in Artificial Intelligence, 2006. Springer Verlag, Berlin/Heidelberg.

- [73] Scott Bradner. Key words for use in RFCs to indicate requirement levels. RFC 2119, Internet Engineering Task Force (IETF), 1997. URL <http://www.ietf.org/rfc/rfc2119.txt>.
- [74] Dan Brickley and Libby Miller. FOAF vocabulary specification 0.98. Technical report, ILRT Bristol. URL <http://xmlns.com/foaf/spec/20100809.html>.
- [75] Dan Brickley, Vinay K. Chaudhri, Harry Halpin, and Deborah L. McGuinness, editors. *Proceedings of the AAAI Spring Symposium on Linked Data Meets Artificial Intelligence*, 2010.
- [76] Matthias Bröcheler. A mathematical semantic web. Bachelor's thesis, Computer Science, Jacobs University, Bremen, 2007. URL <http://www.eecs.jacobs-university.de/archive/bsc-2007/broecheler.pdf>.
- [77] Simon Buckingham Shum, Tim Clark, Anita de Waard, Tudor Groza, Siegfried Handschuh, and Ágnes Sándor. Scientific discourse on the semantic web: A survey of models and enabling technologies. accepted for publication in the *Semantic Web Journal*, 2011. URL <http://www.semantic-web-journal.net/content/scientific-discourse-semantic-web-survey-models-and-enabling-technologies>.
- [78] Lou Burnard and Syd Bauman. TEI P5: Guidelines for electronic text encoding and interchange. Technical report, TEI Consortium. URL <http://www.tei-c.org/release/doc/tei-p5-doc/en/html/>.
- [79] Stephen Buswell. RDF/XML test cases for RDF logic, web ontology and maths content. Deliverable 5.3b, Semantic Web Advanced Development for Europe (SWAD-Europe), 2001. URL http://www.w3.org/2001/sw/Europe/reports/xml_test_cases/wp53.html.
- [80] Stephen Buswell, Olga Caprotti, David P. Carlisle, Michael C. Dewar, Marc Gaëtano, and Michael Kohlhase. The Open Math standard, version 2.0. Technical report, The Open Math Society, 2004. URL <http://www.openmath.org/standard/om20>.
- [81] Olga Caprotti, Mike Dewar, and Daniele Turi. Mathematical service matching using description logic and OWL. In Andrea Asperti, Grzegorz Bancerek, and Andrej Trybulec, editors, *Mathematical Knowledge Management, MKM'04*, number 3119 in Lecture Notes in Artificial Intelligence, pages 73–87. Springer Verlag, 2004, Berlin/Heidelberg.
- [82] Olga Caprotti, Mike Dewar, and Daniele Turi. Mathematical service matching using description logic and OWL. Technical report, The MONET Consortium, 2004. URL http://monet.nag.co.uk/monet/publicdocs/monet_onts.pdf.
- [83] Jacques Carette, Lucas Dixon, Claudio Sacerdoti Coen, and Stephen M. Watt, editors. *MKM/Calculus Proceedings*, number 5625 in Lecture Notes in Artificial Intelligence, July 2009. Springer Verlag, Berlin/Heidelberg. ISBN 9783642026133.
- [84] Paolo Ciccarese and Tudor Groza. Ontology of Rhetorical Blocks (ORB). W3C Interest Group Note, World Wide Web Consortium (W3C), October 2011. URL <http://www.w3.org/TR/2011/NOTE-hcls-orb-20111020/>.
- [85] Tim Clark, Joanne S. Luciano, M. Scott Marshall, Eric Prud'hommeaux, and Susie Stephens, editors. *Semantic Web Applications in Scientific Discourse (SWASD)*, number 523 in CEUR Workshop Proceedings, Aachen, 2009. URL <http://CEUR-WS.org/Vol-523/>.
- [86] CNX. Connexions. URL <http://cnx.org>.
- [87] Mihai Codescu, Fulya Horozal, Michael Kohlhase, Till Mossakowski, and Florian Rabe. Project Abstract: Logic Atlas and Integrator (LATIN). In Davenport et al. [96], pages 289–291.
- [88] Arjeh M. Cohen, Hans Cuypers, and Rikko Verrijzer. Mathematical context in interactive documents. *Mathematics in Computer Science*, 3(3):331–347, 2010.
- [89] Sarah Coppins and Brent Hendricks. Content MathML (Connexions web site). URL <http://cnx.org/content/m9008/2.15/>.
- [90] Olivier Corby, Leila Kefi-Khelif, Hacène Cherfi, Fabien Gandon, and Khaled Khelif. Querying the semantic web of data using SPARQL, RDF and XML. Technical Report 6847, INRIA Sophia Antipolis, February 2009. URL <http://hal.inria.fr/docs/00/36/23/81/PDF/RR-6847.pdf>.
- [91] Stéphane Corlosquet, Renaud Delbru, Tim Clark, Axel Polleres, and Stefan Decker. Produce and Consume Linked Data with Drupal! In Abraham Bernstein, David R. Karger, Tom Heath, Lee Feigenbaum, Diana Maynard, Enrico Motta, and Krishnaprasad Thirunarayan, editors, *The Semantic Web – ISWC 2009*, number 5823 in Lecture Notes in Computer Science, pages 763–778. Springer Verlag, Berlin/Heidelberg, October 2009.
- [92] Hans Cuypers, Arjeh M. Cohen, Jan Willem Knopper, Rikko Verrijzer, and Mark Spanbroek. MathDox, a system for interactive mathematics. In *Proceedings of the World Conference on Educational Multimedia, Hypermedia & Telecommunications 2008 (ED-MEDIA'08)*, pages 5177–5182. AACE, June 2008. URL <http://go.editlib.org/p/29092>.
- [93] Wolfgang Dalitz, Wolfram Sperber, and Winfried Neun. Math-Net, a model for information and communication systems in sciences. In Jan von Knop, Peter Schirmbacher, and Viljan Mahnic, editors, *EUNIS*, number 13 in Lecture Notes in Informatics, pages 140–145. GI, 2001. ISBN 3-88579-339-3.
- [94] Ron Daniel Jr. Harvesting RDF statements from XLinks. W3C Note, World Wide Web Consortium (W3C), September 2000. URL <http://www.w3.org/TR/2000/NOTE-xlink2rdf-20000929/>.
- [95] James Davenport. A small OpenMath type system. *Bulletin of the ACM Special Interest Group on Symbolic and Automated Mathematics (SIGSAM)*, 34(2):16–21, 2000.
- [96] James Davenport, William Farmer, Florian Rabe, and Josef Urban, editors. *Intelligent Computer Mathematics*, number 6824 in Lecture Notes in Artificial Intelligence, 2011. Springer Verlag, Berlin/Heidelberg.
- [97] James H. Davenport and Michael Kohlhase. Unifying Math Ontologies: A tale of two standards. In Carette et al. [83], pages 263–278. ISBN 9783642026133. URL <http://kwarc.info/kohlhase/papers/mkm09-MMLOM3.pdf>.
- [98] Catalin David, Michael Kohlhase, Christoph Lange, Florian Rabe, Nikita Zhiltsov, and Vyacheslav Zholudev. Publishing math lecture notes as linked data. In Lora Aroyo, Grigoris Antoniou, Eero Hyvönen, Annette ten Teije, Heiner Stuckenschmidt, Liliana Cabral, and Tania Tudorache, editors, *The Semantic Web: Research and Applications (Part II)*, number

- 6089 in *Lecture Notes in Computer Science*, pages 370–375. Springer Verlag, Berlin/Heidelberg, 2010.
- [99] DBpedia. URL <http://dbpedia.org>.
- [100] DCMI Usage Board. DCMI metadata terms. DCMI recommendation, Dublin Core Metadata Initiative, 2008. URL <http://dublincore.org/documents/2008/01/14/dcmi-terms/>.
- [101] Jos De Roo. EulerSharp. URL <http://eulersharp.sourceforge.net/>.
- [102] Klaas Dellschaft, Hendrik Engelbrecht, José Monte Barreto, Sascha Rutenbeck, and Steffen Staab. Cicero: Tracking design rationale in collaborative ontology engineering. In Bechhofer et al. [59], pages 782–786.
- [103] Li Ding, Dominic DiFranzo, Alvaro Graves, James R. Michaelis, Xian Li, Deborah L. McGuinness, and Jim Hendler. Data-gov wiki: Towards linking government data. In Brickley et al. [75].
- [104] DITA. OASIS Darwin Information Typing Architecture (DITA). URL <http://www.oasis-open.org/committees/dita/>.
- [105] Leigh Dodds and Ian Davis. *Linked Data Patterns*. A pattern catalogue for modelling, publishing, and consuming Linked Data. August 2011. URL <http://patterns.dataincubator.org/book/>.
- [106] Peter Dolog, Rita Gavrioloaie, Wolfgang Nejdl, and Jan Brase. Integrating adaptive hypermedia techniques and open RDF-based environments. In *Proceedings of the 12th WWW conference*. ACM Press, 2003.
- [107] Nicholas Drummond, Alan Rector, Robert Stevens, Georgina Moulton, Matthew Horridge, Hai Wang, and Julian Sedenberg. Putting owl in order: Patterns for sequences in OWL. In Bernardo Cuenca Grau, Pascal Hitzler, Connor Shankey, and Evan Wallace, editors, *OWL: Experiences and Directions (OWLED)*, November 2006.
- [108] Dublin Core Metadata Element Set. Dublin Core metadata element set. DCMI recommendation, Dublin Core Metadata Initiative, 2008. URL <http://dublincore.org/documents/2008/01/04/dces/>.
- [109] Edd Dubmill. DOAP – description of a project. URL <http://trac.usefulinc.com/doap>.
- [110] Bob DuCharme. Using RDFa with DITA and DocBook. 2009. URL <http://www.devx.com/semantic/Article/42543/>.
- [111] Thomas Eiter, Giovambattista Ianni, Thomas Krennwallner, and Axel Polleres. Rules and ontologies for the semantic web. In Cristina Baroglio, Piero A. Bonatti, Jan Małuszynski, Massimo Marchiori, Axel Polleres, and Sebastian Schaffert, editors, *Reasoning Web*, number 5224 in *Lecture Notes in Computer Science*, pages 1–53. Springer, 2008. URL <http://axel.deri.ie/publications/eite-et-al-2008.pdf>.
- [112] Hans Magnus Enzensberger. *Drawbridge up*. A K PETERS, LTD., 1999. German original: Zugbrücke außer Betrieb. Die Mathematik im Jenseits der Kultur. English translation by Tom Artin.
- [113] William Farmer, Josuah Guttman, and Xavier Thayer. Little theories. In D. Kapur, editor, *Proceedings of the 11th Conference on Automated Deduction*, number 607 in *Lecture Notes in Computer Science*, pages 467–581, Saratoga Springs, NY, USA, 1992. Springer Verlag, Berlin/Heidelberg.
- [114] William M. Farmer. MKM: A new interdisciplinary field of research. *Bulletin of the ACM Special Interest Group on Symbolic and Automated Mathematics (SIGSAM)*, 38(2):47–52, 2004.
- [115] Armin Fiedler. *Prex*: An interactive proof explainer. In Raveev Goré, Alexander Leitsch, and Tobias Nipkow, editors, *Automated Reasoning — 1st International Joint Conference, IJCAR 2001*, number 2083 in *Lecture Notes in Artificial Intelligence*, pages 416–420, Siena, Italy, 2001. Springer Verlag, Berlin/Heidelberg.
- [116] Jana Giceva, Christoph Lange, and Florian Rabe. Integrating web services into active mathematical documents. In Carette et al. [83], pages 279–293. ISBN 9783642026133. URL <https://svn.omdoc.org/repos/jomdoc/doc/pubs/mkm09/jobad/jobad-server.pdf>.
- [117] Deyan Ginev, Constantin Jucovschi, Stefan Anca, Mihai Grigore, Catalin David, and Michael Kohlhasse. An architecture for linguistic and semantic analysis on the arXMLiv corpus. In *Applications of Semantic Technologies (AST) Workshop at Informatik 2009*, 2009. URL http://www.kwarc.info/projects/lamapun/pubs/AST09_LaMaPUn+appendix.pdf.
- [118] Birte Glimm and Chimezie Ogbuji. SPARQL 1.1 entailment regimes. W3C Working Draft, World Wide Web Consortium (W3C), October 2010. URL <http://www.w3.org/TR/2010/WD-sparql11-entailment-20101014/>.
- [119] George Goguadze. Metadata for mathematical libraries. Deliverable D3.a, MoWGLI, 2003. URL http://mowgli.cs.unibo.it/misc/deliverables/metadata/D3a_metadata_for_math/math_metadata.pdf.
- [120] George Goguadze. Metadata model. Deliverable D3.b, MoWGLI, 2003. URL http://mowgli.cs.unibo.it/misc/deliverables/metadata/D3bmetadata_model/metadata_model.pdf.
- [121] Timothy Gowers and Michael Nielsen. Massively collaborative mathematics. *Nature*, 461(15):879–881, 2009.
- [122] Paul Grosso, Eve Maler, Jonathan Marsh, and Norman Walsh. W3C XPointer framework. W3C Recommendation, World Wide Web Consortium (W3C), March 2003. URL <http://www.w3.org/TR/2003/REC-xptr-framework-20030325/>.
- [123] Tudor Groza and Siegfried Handschuh. SALT annotation ontology. Technical report, Digital Enterprise Research Institute (DERI), 2009. URL <http://salt.semanticsauthoring.org/ontologies/sao>.
- [124] Tudor Groza and Siegfried Handschuh. SALT document ontology. Technical report, Digital Enterprise Research Institute (DERI), 2009. URL <http://salt.semanticsauthoring.org/ontologies/sdo>.
- [125] Tudor Groza and Siegfried Handschuh. SALT rhetorical ontology. Technical report, Digital Enterprise Research Institute (DERI), 2009. URL <http://salt.semanticsauthoring.org/ontologies/sro>.
- [126] Tudor Groza, Siegfried Handschuh, Knud Möller, and Stefan Decker. SALT – semantically annotated L^AT_EX for scientific publications. In Enrico Franconi, Michael Kifer, and Wolfgang May, editors, *The Semantic Web: Research and Applications*, number 4519 in *Lecture Notes in Computer Science*, pages 518–532. Springer Verlag, Berlin/Heidelberg, 2007. ISBN 978-3-540-72666-1.

- [127] Tudor Groza, Knud Möller, Siegfried Handschuh, Diana Trif, and Stefan Decker. SALT: Weaving the claim web. In Karl Aberer, Key-Sun Choi, Natasha Fridman Noy, Dean Allemang, Kyung-Il Lee, Lyndon J. B. Nixon, Jennifer Golbeck, Peter Mika, Diana Maynard, Riichiro Mizoguchi, Guus Schreiber, and Philippe Cudré-Mauroux, editors, *ISWC/ASWC*, number 4825 in Lecture Notes in Computer Science, pages 197–210. Springer Verlag, Berlin/Heidelberg, 2007. ISBN 978-3-540-76297-3.
- [128] Tudor Groza, Siegfried Handschuh, Tim Clark, Simon Buckingham Shum, and Anita de Waard. A short survey of discourse representation models. In Clark et al. [85]. URL <http://CEUR-WS.org/Vol-523/>.
- [129] Ferruccio Guidi. *Searching and Retrieving in Content-based Repositories of Formal Mathematical Knowledge*. PhD thesis, Università di Bologna, 2003.
- [130] Ferruccio Guidi and Claudio Sacerdoti Coen. Querying distributed digital libraries of mathematics. In Thérèse Hardin and Renaud Rioboo, editors, *Proceedings of the 11th Symposium on the Integration of Symbolic Computation and Mechanized Reasoning (Calculus 2003)*, pages 17–30, Rome, Italy, September 2003. URL <http://www.calculumus.net/meetings/rome03/Proceedings/final.pdf>.
- [131] Ferruccio Guidi and Irene Schena. A query language for a metadata framework about mathematical resources. In Asperti et al. [48], pages 105–118.
- [132] Thomas C. Hales, John Harrison, Sean McLaughlin, Tobias Nipkow, Steven Obua, and Roland Zumkeller. A revision of the proof of the Kepler conjecture. *Discrete and Computational Geometry*, pages 1–34, 2010.
- [133] Kevin Hammond, Peter Horn, Alexander Kononov, Steve Linton, Dan Roozmond, Abdallah Al Zain, and Phil Trinder. Easy composition of symbolic computation software: A new lingua franca for symbolic computation. In *Proceedings of the 2010 International Symposium on Symbolic and Algebraic Computation (ISSAC)*, pages 339–346. ACM Press, 2010.
- [134] Steve Harris and Andy Seaborne. SPARQL 1.1 query language. W3C Working Draft, World Wide Web Consortium (W3C), October 2010. URL <http://www.w3.org/TR/2010/WD-sparql11-query-20101014/>.
- [135] Olaf Hartig, Hannes Mühleisen, and Johann-Christoph Freytag. Linked data for building a map of researchers. In Bizer et al. [69]. URL <http://CEUR-WS.org/Vol-449/>.
- [136] Philipp Heim, Steffen Lohmann, and Timo Stegemann. Interactive relationship discovery via the semantic web. In Aroyo et al. [45], pages 303–317.
- [137] Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space*. Morgan & Claypool, San Rafael, CA, 2011. URL <http://linkeddatabook.com>.
- [138] Bettina Heintz. *Die Innenwelt der Mathematik. Zur Kultur und Praxis einer beweisenden Disziplin*. Springer Verlag, Wien, 2000.
- [139] Ian Hickson. HTML5. A vocabulary and associated APIs for HTML and XHTML. W3C Working Draft, World Wide Web Consortium (W3C), May 2011. URL <http://www.w3.org/TR/2011/WD-html5-20110525/>.
- [140] IEEE Learning Technology Standards Committee. Standard for Learning Object Metadata. Technical Report 1484.12.1, IEEE, 2002.
- [141] IEEE Learning Technology Standards Committee. Standard for Resource Description Framework (RDF) binding for Learning Object Metadata data model. Technical Report 1484.12.4, IEEE, 2002.
- [142] Prateek Jain, Pascal Hitzler, Peter Z. Yeh, Kunal Verma, and Amit P. Sheh. Linked data is merely more data. In Brickley et al. [75].
- [143] Fairouz Kamareddine, Manuel Maarek, Krzysztof Retel, and J. B. Wells. Narrative structure of mathematical texts. In Manuel Kauers, Manfred Kerber, Robert Miner, and Wolfgang Windsteiger, editors, *Towards Mechanized Mathematical Assistants. MKM/Calculus*, number 4573 in Lecture Notes in Artificial Intelligence, pages 296–312. Springer Verlag, Berlin/Heidelberg, 2007. ISBN 978-3-540-73083-5.
- [144] Fairouz Kamareddine, J. B. Wells, and Christoph Zengler. Computerising mathematical text with MathLang. *Electron. Notes Theor. Comput. Sci.*, 205:5–30, 2008. ISSN 1571-0661. URL <http://www.cedar-forest.org/forest/papers/drafts/mathlang-coq-short.pdf>.
- [145] Andrea Kohlase and Michael Kohlase. Semantic knowledge management for education. *Proceedings of the IEEE; Special Issue on Educational Technology*, 96(6):970–989, June 2008. URL <http://kwarc.info/kohlase/papers/semkm4ed.pdf>.
- [146] Andrea Kohlase and Michael Kohlase. Spreadsheet interaction with frames: Exploring a mathematical practice. In Carette et al. [83], pages 341–356. ISBN 9783642026133. URL <http://kwarc.info/kohlase/papers/mkm09-framing.pdf>.
- [147] Michael Kohlase. OMDoc – An open markup format for mathematical documents [Version 1.2]. Number 4180 in Lecture Notes in Artificial Intelligence. Springer Verlag, Berlin/Heidelberg, August 2006. URL <http://omdoc.org/pubs/omdoc1.2.pdf>.
- [148] Michael Kohlase. Using L^AT_EX as a Semantic Markup Format. *Mathematics in Computer Science*, 2(2):279–304, 2008. URL <https://svn.kwarc.info/repos/stex/doc/mcs08/stex.pdf>.
- [149] Michael Kohlase. Towards bootstrapping the pragmatic to strict mapping in OMDoc. August 2009. URL <https://svn.omdoc.org/repos/omdoc/trunk/doc/blue/p2s-bootstrap/note.pdf>.
- [150] Michael Kohlase and Ioan Şucan. A search engine for mathematical formulae. In Tetsuo Ida, Jacques Calmet, and Dongming Wang, editors, *Proceedings of Artificial Intelligence and Symbolic Computation, AISC'2006*, number 4120 in Lecture Notes in Artificial Intelligence, pages 241–253. Springer Verlag, Berlin/Heidelberg, 2006. URL <http://kwarc.info/kohlase/papers/aisc06.pdf>.
- [151] Michael Kohlase, Till Mossakowski, and Florian Rabe. Latin: Logic atlas and integrator. URL <http://latin.omdoc.org>.
- [152] Michael Kohlase, Joe Corneli, Catalin David, Deyan Ginev, Constantin Jucovschi, Andrea Kohlase, Christoph Lange, Bogdan Matican, Stefan Mirea, and Vyacheslav Zholudev. The Planetary system: Web 3.0 & active documents for STEM. *Procedia Computer Science*, 4:598–607, 2011. URL <https://svn.mathweb.org/repos/planetary/doc/epc11/paper.pdf>. Finalist at the Executable Papers Challenge.

- [153] Werner Kunz and Horst W. J. Rittel. Issues as elements of information systems. Working paper 131, Institute of Urban and Regional Development, University of California, Berkeley, July 1970.
- [154] Imre Lakatos. *Proofs and Refutations*. Cambridge University Press, 1976.
- [155] Christoph Lange. Krextor – an extensible XML→RDF extraction framework. In Bizer et al. [69]. URL <http://ceur-ws.org/Vol-449/ShortPaper2.pdf>.
- [156] Christoph Lange. wiki.openmath.org – how it works, how you can participate. In James H. Davenport, editor, *22nd OpenMath Workshop*, July 2009. URL <http://staff.bath.ac.uk/masjhd/OM2009.html>.
- [157] Christoph Lange. Towards OpenMath content dictionaries as linked data. In Michael Kohlhase and Christoph Lange, editors, *23rd OpenMath Workshop*, July 2010. URL <http://cicm2010.cnam.fr/om/>.
- [158] Christoph Lange. The OMDoc ontology, 2010. URL <http://kwarc.info/projects/docOnto/omdoc.html>.
- [159] Christoph Lange. Krextor – an extensible framework for contributing content math to the web of data. In Davenport et al. [96], pages 304–306. URL <http://kwarc.info/clange/pubs/krextor-system.pdf>.
- [160] Christoph Lange. *Enabling Collaboration on Semiformal Mathematical Knowledge by Semantic Web Integration*. PhD thesis, Jacobs University Bremen, 2011, number 11 in Studies on the Semantic Web. AKA Verlag and IOS Press, Heidelberg and Amsterdam, 2011.
- [161] Christoph Lange and Michael Kohlhase. A mathematical approach to ontology authoring and documentation. In Carette et al. [83], pages 389–404. ISBN 9783642026133. URL <http://kwarc.info/kohlhase/papers/mkm09-omdoc4onto.pdf>.
- [162] Christoph Lange, Patrick Ion, Anastasia Dimou, Charalampos Bratsas, Wolfram Sperber, Michael Kohlhase, and Ioannis Antoniou. Getting Mathematics Towards the Web of Data: the Case of the Mathematics Subject Classification. Submitted to Extended Semantic Web Conference 2012. URL <http://kwarc.info/clange/pubs/eswc2012-msc-skos.pdf>.
- [163] Christoph Lange, Uldis Bojārs, Tudor Groza, John Breslin, and Siegfried Handschuh. Expressing argumentative discussions in social media sites. In John Breslin, Uldis Bojārs, Alexandre Passant, and Sergio Fernández, editors, *Social Data on the Web (SDoW), Workshop at the 7th International Semantic Web Conference*, number 405 in CEUR Workshop Proceedings, Aachen, 2008. URL <http://ceur-ws.org/Vol-405/paper4.pdf>.
- [164] Christoph Lange, Tuukka Hastrup, and Stéphane Corlosquet. Arguing on issues with mathematical knowledge items in a semantic wiki. In Joachim Baumeister and Martin Atzmüller, editors, *Wissens- und Erfahrungsmanagement LWA (Lernen, Wissensentdeckung und Adaptivität) Conference Proceedings*, volume 448, October 2008.
- [165] Christoph Lange et al. Krextor – the KWARC RDF extractor. URL <http://kwarc.info/projects/krextor/>.
- [166] Paul Libbrecht. Cross curriculum search through the geoskills ontology. In D. Massart, J.-N. Colin, F. Van Asche, and M. Wolpers, editors, *Proceedings of the 2nd International Workshop on Search and Exchange of e-le@rning Material (SE@M)*, located at EC-TEL-2008, CEUR Workshop Proceedings, Aachen, 2008. CEUR-WS.org. URL <http://ceur-ws.org/Vol-385/>.
- [167] Paul Libbrecht. Notations around the world: Census and exploitation. In Autexier et al. [54], pages 398–410. ISBN 3642141277.
- [168] Paul Libbrecht and Erica Melis. Methods for Access and Retrieval of Mathematical Content in ActiveMath. In N. Takayama and A. Iglesias, editors, *Proceedings of ICMS-2006*, number 4151 in Lecture Notes in Artificial Intelligence, pages 331–342. Springer Verlag, Berlin/Heidelberg, 2006. URL <http://www.activemath.org/publications/Libbrecht-Melis-Access-and-Retrieval-ActiveMath-ICMS-2006.pdf>.
- [169] William C. Mann and Maite Taboada. Rhetorical structure theory. URL <http://www.sfu.ca/rst/>.
- [170] Massimo Marchiori. The mathematical semantic web. In Asperti et al. [48], pages 216–223. Keynote.
- [171] Jonathan Marsh, Daniel Veillard, and Norman Walsh. `xml:id` version 1.0. W3C Recommendation, World Wide Web Consortium (W3C), September 2005. URL <http://www.w3.org/TR/2005/REC-xml-id-20050909/>.
- [172] Viviana Mascardi, Valentina Cordi, and Paolo Rosso. A comparison of upper ontologies. In Matteo Baldoni, Antonio Bocalatte, Flavio De Paoli, Maurizio Martelli, and Viviana Mascardi, editors, *WOA*, pages 55–64. Seneca Edizioni Torino, 2007. ISBN 978-88-6122-061-4.
- [173] Claudio Masolo, Stefano Borgo, Aldo Gangemi, Nicola Guarino, Alessandro Oltramari, and Luc Schneider. The wonderweb library of foundational ontologies. WonderWeb Deliverable 17, Laboratory for Applied Ontology – ISTC-CNR, May 2003. URL <http://wonderweb.semanticweb.org/deliverables/documents/D17.pdf>.
- [174] MathDox. MathDox – interactive mathematics. URL <http://www.mathdox.org>.
- [175] Deborah L. McGuinness, Li Ding, Paulo Pinheiro da Silva, and Cynthia Chang. PML 2: A modular explanation interlingua. In Thomas Roth-Berghofer, Stefan Schulz, and David B. Leake, editors, *Proceedings of the AAAI Workshop on Explanation-Aware Computing (ExaCt)*, 2007. URL <http://www.aaai.org/Papers/Workshops/2007/WS-07-06/WS07-06-008.pdf>.
- [176] Erica Melis, Giorgi Gogvadze, Paul Libbrecht, and Carsten Ullrich. Culturally adapted mathematics education with ActiveMath. *AI & Society*, 24(3):251–265, 2009.
- [177] Paolo Missier, Satya S. Sahoo, Jun Zhao, Carole Goble, and Amit Sheth. Janus: from workflows to semantic provenance and linked open data. In *Proceedings of the 3rd Provenance and Annotation Workshop*, 2010.
- [178] MKMNET. MKMNET (Mathematical Knowledge Management Network). URL <http://monet.nag.co.uk/mkm/>.
- [179] Monet. MONET – Mathematics on the net. web page at <http://monet.nag.co.uk>. URL <http://monet.nag.co.uk>.
- [180] Till Mossakowski. Hets: the Heterogeneous Tool Set. URL <http://www.dfki.de/sks/hets>.
- [181] Till Mossakowski, Christian Maeder, and Klaus Lüttich. The Heterogeneous Tool Set. In Orna Grumberg and Michael Huth, editors, *Proceedings of the 13th International Conference on Tools and Algorithms for the Construction and Anal-*

- ysis of Systems TACAS-2007, number 4424 in Lecture Notes in Computer Science, pages 519–522, Berlin, Germany, 2007. Springer Verlag.
- [182] Christine Müller. *Adaptation of Mathematical Documents*. PhD thesis, Jacobs University Bremen, 2010. URL <http://kwarc.info/cmuller/papers/thesis.pdf>.
- [183] Hala Naja-Jazzar, Nishadi de Silva, Hala Skaf-Molli, Charbel Rahhal, and Pascal Molli. OntoReST: A RST-based ontology for enhancing documents content quality in collaborative writing. *INFOCOMP Journal of Computer Science*, 8(3): 1–10, September 2009.
- [184] Mikael Nilsson, Andy Powell, Pete Johnston, and Ambjörn Naeve. Expressing Dublin Core metadata using the Resource Description Framework (RDF). DCMi recommendation, Dublin Core Metadata Initiative. URL <http://dublincore.org/documents/2008/01/14/dc-rdf/>.
- [185] Mikael Nilsson, Matthias Palmér, and Jan Brase. The LOM RDF binding – principles and implementation. In *3rd Annual Ariadne Conference*, 2003.
- [186] Tope Omitola, Christos L. Koumenides, Igor O. Popov, Yang Yang, Manuel Salvadores, Martin Szomszor, Tim Berners-Lee, Nicholas Gibbins, Wendy Hall, mc schraefel, and Nigel Shadbolt. Put in your postcode, out comes the data: A case study. In Aroyo et al. [45], pages 318–332.
- [187] OpenLink Software. OpenLink universal integration middleware – Virtuoso product family. URL <http://virtuoso.openlinksw.com>.
- [188] Christian-Emil Ore and Øyvind Eide. TEI and cultural heritage ontologies: Exchange of information? *Literary and Linguistic Computing*, 24(2):161–172, 2009.
- [189] Luca Padovani. GtMathView. URL <http://helm.cs.unibo.it/mml-widget/>.
- [190] Luca Padovani and Stefano Zacchiroli. From notation to semantics: There and back again. In Borwein and Farmer [72], pages 194–207.
- [191] Alexandre Passant, Paolo Ciccarese, John G. Breslin, and Tim Clark. SWAN/SIOC: Aligning scientific discourse representation and social semantics. In Clark et al. [85].
- [192] PlanetMath.org. PlanetMath.org – math for the people, by the people. URL <http://planetmath.org>.
- [193] Judith Plümer. Erinnerungen an Prof. Dr. Schwänzl, 2004. URL <http://d-mathnet.preprints.org/research/wissinfo/NachrufRS.html>.
- [194] George Pólya. *How to Solve it*. Princeton University Press, 1973.
- [195] Florian Rabe. *Representing Logics and Logic Translations*. PhD thesis, Jacobs University Bremen, 2008. URL <http://kwarc.info/frabe/Research/phdthesis.pdf>.
- [196] Florian Rabe and Michael Kohlhase. A scalable module system. Manuscript, submitted to Information & Computation, 2011. URL <http://kwarc.info/frabe/Research/mmt.pdf>.
- [197] Ricardo Radaelli-Sanchez and Connexions. The intermediate CNXML (Connexions web site). URL <http://cnx.org/content/m9006/2.22/>.
- [198] Robert G. Raskin and Michael J. Pan. Knowledge representation in the semantic web for Earth environmental terminology (SWEET). *Computers & Geosciences*, 31:1119–1125, 2005.
- [199] Krzysztof Retel. *Gradual Computerisation and Verification of Mathematics*. PhD thesis, Heriot-Watt University, Edinburgh, April 2009.
- [200] Andrew Robbins. Semantic MathML. 2009. URL <http://straymindcough.blogspot.com/2009/06/semantic-mathml.html>.
- [201] Leo Sauermann and Richard Cyganiak. Cool URIs for the semantic web. W3C Interest Group Note, World Wide Web Consortium (W3C), December 2008. URL <http://www.w3.org/TR/2008/NOTE-cooluris-20081203/>.
- [202] Irene Schena. *Towards a Semantic Web for Formal Mathematics*. PhD thesis, University of Bologna, March 2002. Technical Report UBLCS–2002–6.
- [203] SCIENCE. The SCIENCE project – Symbolic Computation Infrastructure for Europe. URL <http://www.symbolic-computation.org/>.
- [204] François-Paul Servant. Linking enterprise data. In Christian Bizer, Tom Heath, Kingsley Idehen, and Tim Berners-Lee, editors, *Linked Data on the Web (LDOW)*, number 369 in CEUR Workshop Proceedings, Aachen, April 2008. URL <http://CEUR-WS.org/Vol-369/>.
- [205] John Sheridan and Jeni Tennison. Linking UK government data. In Christian Bizer, Tom Heath, Tim Berners-Lee, and Michael Hausenblas, editors, *Linked Data on the Web (LDOW)*, number 628 in CEUR Workshop Proceedings, Aachen, April 2010. URL <http://CEUR-WS.org/Vol-628/>.
- [206] David Shotton and Silvio Peroni. DoCo, the Document Components Ontology. May 2011. URL <http://purl.org/spar/doco>.
- [207] Neil J. A. Sloane. The on-line encyclopedia of integer sequences. *Notices of the AMS*, 50(8):912, 2003.
- [208] Wolfram Sperber. Math-Net international and the Math-Net page. In Fengshan Bai and Bernd Wegner, editors, *Electronic Information and Communication in Mathematics*, number 2730 in Lecture Notes in Computer Science, pages 169–177. Springer, 2003.
- [209] Manu Sporny, Benjamin Adrian, Nathan Rixham, Mark Birbeck, and Ivan Herman. RDFa API. W3C Working Draft, World Wide Web Consortium (W3C), April 2011. URL <http://www.w3.org/TR/2011/WD-rdfa-api-20110419/>.
- [210] Heinrich Stamerjohanns, Michael Kohlhase, Deyan Ginev, Catalin David, and Bruce Miller. Transforming large collections of scientific publications to XML. *Mathematics in Computer Science*, 3(3):299–307, 2010. URL <http://kwarc.info/kohlhase/papers/mcs10.pdf>.
- [211] Pradeep Suresh, Shuo-Huan Hsu, Pavan Akkisetty, Gintaras V. Reklaitis, and Venkat Venkatasubramanian. OntoMODEL: Ontological Mathematical Modeling Knowledge Management in Pharmaceutical Product Development, 1: Conceptual Framework. *Industrial & Engineering Chemistry Research*, 49(17):7758–7767, 2010.
- [212] Carlos Tejo-Alonso, Sergio Fernández, Diego Berrueta, Luis Polo, María Jesús Fernández, and Víctor Morlán. eZaragoza, a tourist promotional mashup. In Adrian Giurca, Brigitte Endres-Niggemeyer, Christoph Lange, Lutz Maicher, and Pascal Hitzler, editors, *AI Mashup Challenge at ESWC*, June 2010. URL <http://sites.google.com/a/fh-hannover.de/aimashup/home/ezaragoza>.
- [213] Christoph Tempich, H. Sofia Pinto, York Sure, and Steffen Staab. An argumentation ontology for DIstributed, Loosely-controlled and evolvInG Engineering processes of oNTology

- gies (DILIGENT). In Asunción Gómez-Pérez and Jérôme Euzenat, editors, *The Semantic Web: Research and Applications*, number 3532 in Lecture Notes in Computer Science, pages 241–256. Springer Verlag, Berlin/Heidelberg, 2005. ISBN 3-540-26124-9.
- [214] Christoph Tempich, Elena Simperl, Markus Luczak, Rudi Studer, and H. Sofia Pinto. Argumentation-based ontology engineering. *IEEE Intelligent Systems*, 22(6):52–59, 2007. ISSN 1541-1672.
- [215] Jerzy Trzeciak. *Writing Mathematical Papers in English*. Gdańskie Wydawnictwo Oświatowe, 1995.
- [216] Josef Urban, Jesse Alama, Piotr Rudnicki, and Herman Geuvers. A wiki for Mizar: Motivation, considerations, and initial prototype. In Autexier et al. [54], pages 455–469. ISBN 3642141277.
- [217] Josef Urban. Content-based encoding of mathematical and code libraries. In Christoph Lange and Josef Urban, editors, *ITP Workshop on Mathematical Wikis (MathWikis)*, number 767 in CEUR Workshop Proceedings, Aachen, 2011, pages 49–53. URL <http://ceur-ws.org/Vol-767/paper-09.pdf>.
- [218] Denny Vrandečić, Markus Krötzsch, Sebastian Rudolph, and Uta Lösch. Leveraging non-lexical knowledge for the Linked Open Data Web. In Rodolphe Héliot and Antoine Zimmermann, editors, *The Fifth RAFT, the yearly bilingual publication on nonchalant research*, 2010. URL http://km.aifb.kit.edu/projects/numbers/linked_open_numbers.pdf.
- [219] Denny Vrandečić, Christoph Lange, Michael Hausenblas, Jie Bao, and Li Ding. Semantics of governmental statistics data. In *Proceedings of WebSci'10: Extending the Frontiers of Society On-Line*. Web Science Trust, 2010. URL <http://journal.webscience.org/400/>.
- [220] Norman Walsh. The DocBook schema. Committee specification, OASIS, August 2008. URL <http://www.docbook.org/specs/docbook-5.0-spec-cs-01.html>.
- [221] Norman Walsh and Leonard Muellner. *DocBook 5.0: The Definitive Guide*. O'Reilly, 2008.
- [222] Makarius Wenzel, Clemens Ballarin, Stefan Berghofer, Timothy Bourke, Lucas Dixon, Florian Haftmann, Gerwin Klein, Alexander Krauss, Tobias Nipkow, David von Oheimb, Larry Paulson, and Sebastian Skalberg. *The Isabelle/Isar Reference Manual*. URL <http://isabelle.in.tum.de/doc/isar-ref.pdf>.
- [223] Freek Wiedijk. Statistics on digital libraries of mathematics. *Studies in Logic, Grammar and Rhetoric*, 18(31), 2009.
- [224] Wikipedia. Wikipedia, the free encyclopedia. URL <http://www.wikipedia.org>.
- [225] Wikipedia: Knowledge management. Knowledge management, December 2009. URL http://en.wikipedia.org/w/index.php?title=Knowledge_management&oldid=329227520.
- [226] Wikipedia: Linus' Law. Linus' Law, March 2011. URL http://en.wikipedia.org/w/index.php?title=Linus%27_Law&oldid=421629750.
- [227] Wikipedia: Mathematics. Portal: Mathematics, December 2009. URL <http://en.wikipedia.org/w/index.php?title=Portal:Mathematics&oldid=329137789>.
- [228] Wikipedia: Open Notebook Science. Open Notebook Science, July 2010. URL http://en.wikipedia.org/w/index.php?title=Open_Notebook_Science&oldid=372235042.
- [229] Wikipedia: Pythagorean theorem. Pythagorean theorem, November 2009. URL http://en.wikipedia.org/w/index.php?title=Pythagorean_theorem&oldid=328597679.
- [230] Gregory Todd Williams. SPARQL 1.1 Service Description. W3C Working Draft, World Wide Web Consortium (W3C), May 2011. URL <http://www.w3.org/TR/2010/WD-sparql11-service-description-20110512/>.
- [231] Stefano Zacchiroli. *User Interaction Widgets for Interactive Theorem Proving*. PhD thesis, Università di Bologna, 2007.
- [232] Jürgen Zimmer. *MathServe – A Framework for Semantic Reasoning Services*. PhD thesis, Universität des Saarlandes, July 2008.
- [233] Antoine Zimmermann. Ontology recommendations for the data publishers. In Mathieu d'Aquin, Alexander García Castro, Christoph Lange, and Kim Viljanen, editors, *1st Workshop on Ontology Repositories and Editors*, number 596 in CEUR Workshop Proceedings, Aachen, 2010. URL <http://ceur-ws.org/Vol-596/>.
- [234] Claus Zinn. *Understanding Informal Mathematical Discourse*. PhD thesis, Technischen Fakultät der Universität Erlangen-Nürnberg, 2004.