

Authoring and Publishing Units and Quantities in Semantic Documents

Mihai Cîrlănar¹, Deyan Ginev¹, and Christoph Lange^{1,2}

¹ Computer Science, Jacobs University Bremen, Germany
{m.cirlanaru,d.ginev,ch.lange}@jacobs-university.de

² Universität Bremen, Germany

Abstract. This paper shows how an explicit representation of units and quantities can improve the experience of semantically published documents, and provides a first authoring method in this respect. To exemplify the potential and practical advantages of encoding explicit semantics regarding units w.r.t. user experience, we demonstrate a *unit system preference* service, which enables the user to choose the system of units for the displayed paper. By semantically publishing units, we obtain a basis for a wide range of applications and services such as *unknown unit lookup*, *unit and quantity semantic search* and *unit and quantity manipulation*. Enabling semantic publishing for units is also presented in the context of a large collection of legacy scientific documents (the ARXMLIV corpus), where our approach allows to non-invasively enrich legacy publications.

1 Motivation

Units and quantities¹, although widely spread, lack a formal standard representation for semantic publishing. A multitude of problems [45] arise from the different flavors (country specific unit standards) and formats (abbreviations, special cases of occurrence) of units, making it hard for the untrained reader to fully understand the information provided. Semantic publishing solves most such problems by disambiguating the unit and quantity occurrences, thus enabling a wide range of applications and services to interact with them.

A **unit** is *any determinate quantity, dimension, or magnitude adopted as a basis or standard of measurement for other quantities of the same kind and in terms of which their magnitude is calculated or expressed* [32], but from the top-most level of perception, it simply provides information on a wide range of quantifiable aspects. Concrete examples for the great extent of units and quantities include cooking recipes, medical prescriptions, scientific papers and many other. Semantic publishing can provide the middle layer that would ensure an (automated) way of identifying and understanding these occurrences, which can enable the evolution of useful technologies and services.

¹ We chose the *units and quantities* wording in order to emphasize the semantic dependence between the *unit* and its *quantity* (i.e. amount, magnitude). An in-depth analysis of the two concepts is provided in section 2 of [9].

At the perception level, aside from quantifying properties and relations between objects, units bring the meaning of scale. Moreover, units have allowed scientists to better transmit and exchange knowledge among themselves.

In real life, the misinterpretation of units and their quantities has often caused accidents with harsh/expensive consequences. Consider losing a \$125 million satellite [30] because of the differences between metric and imperial unit systems, or running out of fuel in mid-flight with an aircraft whose fuel sensors were faultily configured in displaying the units [2]. Fields like medicine, commerce, civil engineering have also been marked by such types of errors and pitfalls [45].

Providing semantics to units and their quantities for the publishing industry, either by supplying semantic authoring tools or by semantically enriching their occurrences in legacy documents, has high impact benefits. It will enable transparent exchange of scientific knowledge between different academic communities, typical of technical papers with a high occurrence of units and quantities, and also enhance the reader’s experience, via novel interactive services with day-to-day published material, e.g. cooking recipes or technical manuals.

In the following sections, we introduce preliminaries (section 2), outline our approach to semantic units and quantities, and review relevant state of the art (section 4). That provides a basis for presenting *unit and quantity interaction services* (section 6). We outline immediate strategies (section 5) for extending the benefits of semantic units to legacy documents (section 7) and conclude with a summary of our mid-term outlook of future work (section 8).

2 Preliminaries

The core of semantic publishing resides in open and standardized markup languages used to encapsulate semantics. OPENMATH and *Content* MATHML are the most widely used semantic markup (also called “content markup”) languages for mathematical expressions, which are ubiquitous in science and engineering.

2.1 OpenMath and Content MathML

OPENMATH [7] and the semantically equivalent Content MATHML [4] are standards for the representing the semantics of mathematical expressions [28] – as annotations to visual renderings, or for the purpose of communication between computational services. Our investigations focus on these two languages.²

Structurally, both OPENMATH and MATHML provide a valuable basis for machine processing of mathematical expressions; that makes them suitable for

² The prevalence of XML-based semantic markup languages for representing mathematical expressions – as opposed to RDF – has historical reasons but is also due to the complex n -ary and ordered structures of mathematical expressions, which are hard to break down into RDF triples, be it standoff RDF markup or embedded RDFa annotations. Both representations have in common that the vocabulary terms (here: functions, operators, sets, constants) are identified by URIs. We refer to [29] for an in-depth treatment.

semantic publishing of units and quantities. The expressivity of MATHML, provided by its vocabulary having close to 100 XML elements for mathematical functions and operators [28] and multiple *unit and quantity* representation possibilities [13], and the modularity and extensibility of OPENMATH’s vocabulary by way of modular ontologies (“Content Dictionaries”, abbreviated as CDs), enable the development of applications and services (some of which are discussed in section 6.2) that build upon the semantic publishing of units and quantities.

2.2 The Semantic Publishing Pipeline

Semantic Publishing, conceived as a process, consists of at least three components, namely *authoring*, *publishing* and *interaction*. Usually these processes imply three different groups of contributors – authors, publishers and readers. Incorporating the full publishing lifecycle into a single system, striving for integration and collaboration between the different participants, brings great benefits. In this paper, we take the benefits of the social web for well-established and accepted³ and focus on the more novel semantic aspects of the publishing realm. To this extent, we develop our work in the context of the Planetary eMath3.0 system [34,27], which provides, on top of a stable, well-established Web 2.0 framework, an architecture for semantic services that interact with semantically annotated mathematical and technical documents.

In our work on units and quantities, we have concentrated on setting the necessary technological foundation, hence building on the languages introduced in section 2.1 to select and enhance the authoring and interaction aspects.

3 Semantic Units – Idea Outline

In order to understand how a semantic representation of units and quantities will integrate with the publishing flow of our framework of choice, one first needs to pinpoint what they comprise and how they are *represented*.

A computational *semantic entity* is an object with explicit *structure*, representable in a machine-understandable form, and denoting a corresponding real-world entity. The denotation is usually encoded via a machine-readable ontology. This definition is directly applicable to semantic units and quantities, which are exactly the machine-readable representations of their physical counterparts.

For the *representation* we choose OPENMATH, since it encompasses units through modular ontologies, called Content Dictionaries (CDs) [9]. OPENMATH CDs enable extensibility through the creation of new such ontologies that can add new symbols, or simply through the extension of the existing unit ontologies/CDs.

³ For mathematics, including the mathematical foundations of science and engineering, see, e.g., the PlanetMath free encyclopedia [35] and the Polymath wiki/blog-based collaboration effort [5].

As a running example, we consider a semantic representation of the physical *quantity* 100 km/h; one possible OPENMATH representation is⁴:

```

<OMA>
  <OMS cd="arith1" name="times" />
  <OMI>100</OMI>
  <OMA>
    <OMS cd="arith1" name="divide" />
    <OMA>
      <OMS cd="units_ops1" name="prefix" />
      <OMS cd="units_siprefix1" name="kilo" />
      <OMS cd="units_metric1" name="metre" />
    </OMA>
    <OMS cd="units_time1" name="hour" />
  </OMA>
</OMA>
    
```

Listing 1. OPENMATH representation of 100 km/h

4 State of the Art

We review the relevant prior work involving units and quantities in the context of semantic publishing. Note that we do not cover the publishing dimension itself; we consider it a stand-alone level within a semantic publishing framework, independent of the processed content.

4.1 Representation

The semantic publishing aspect of units in scientific documents has not yet accumulated a sizable body of prior work. Previous research has mainly been concerned with the standardization of unit and quantity representation, which is far from complete (not covering every unit occurrence possibility) or sufficiently machine comprehensible. There is a number of units-related semantic web ontologies: The authors of the Measurement Units Ontology [6] review a number of ways of representing units in RDF. The SWEET (Semantic Web Earth and Environmental Terminology [43,37]) and QUDT (Quantities, Units, Dimensions and Data Types [21]) ontologies are particularly remarkable for linking units to related machine-comprehensible information. SWEET 2.0 describes just 91 units but comprises 150 modules that cover many different sciences as well as common foundations of science; it links units to the SWEET descriptions of the fields of science where they occur. QUDT 1.1 covers 807 units and links all quantities⁵,

⁴ This is one out of several ways of representing units (cf. [13]). For a detailed description of the XML schema see section 3.1.2 of [7].

⁵ Differing from this paper, QUDT uses the term “quantity” for “an observable property of an object [...] that can be measured and quantified numerically”, and uses the term “quantity value” for the numerical value of a quantity [21].

units, and dimensions to their counterparts in the DBpedia dataset [12], where users or automated agents can then explore further relations.

For OPENMATH, a representation of units and quantities has been proposed (cf. [13]), and several CDs covering common units have been provided. The in-depth analysis of the prospective representations of units and their dimensions that [13] proposes (taking into account the pros and cons of each approach) allows for a broader view on the multitude of semantic publishing possibilities. The two most significant sets of OPENMATH unit CDs have been developed by James Davenport and Jonathan Stratford [38] and Joseph Collins [9], respectively. The former are remarkable for their explicit representation of conversion rules (see also Section 4.4). The latter ones provide a standards-compliant implementation of SI⁶ quantities and units, providing strong insight on the concepts of *quantity* and *unit* and on the prospects of capturing more of their semantics in the representation.

4.2 Authoring

In “pre-semantic” environments, such as L^AT_EX, there are first approximations of content-oriented macros that represent units. A prominent example is the L^AT_EX package *SIunits* [20] which covers the full range of base and derived units in the SI system, as well as SI prefixes, a range of widely accepted units external to SI and generic mechanisms for creating custom author-specified unit constructs. The package enables a large set of abbreviative commands, which are internally built up from the compositional application of atomic building blocks. In this sense, the authoring process via *SIunits* is *nearly semantic* on the interface level, but *entirely presentational* on the output side.

Still, all major semantic authoring systems (e.g. the semantic L^AT_EX extensions sT_EX [25], SALT [19], the Ontology Add-in for Microsoft Office Word [14], or the semantic content management system PAUX [33]) have so far neglected the specific use case of units. This can be partially explained by the lack of a widely agreed standard representation, as well as different primary development foci – mathematics for sT_EX, rhetorical structures for SALT, life sciences terminology for the Word ontology add-in, and educational texts from areas unrelated to physics, such as law, for PAUX. Notably, sT_EX could, in principle, support units already, as its wide coverage of the conceptual model of OPENMATH and its generic mechanism for defining new symbols and concepts could easily be utilized for specifying the relevant unit and quantity symbols. Section 5 presents how we have done that in a way that does not disrupt existing L^AT_EX authoring practices. While L^AT_EX is commonly used in mathematics, science, and engineering, our solution is unlikely to appeal to life scientists, where Microsoft Office is more widely used; however, we leave unit support for word processors to future work.

⁶ The International System of Units [39].

4.3 Computation

This section briefly covers computation as a common prerequisite of interacting with units and quantities, which is covered in the following section. To realize why unit conversion requires more powerful computation facilities than just multiplication, consider conversions of dates between different calendars, such as the Gregorian and the Julian calendar with their different notions of months and leap years [47].

OPENMATH has been designed for exchanging mathematical expressions between computer algebra systems and automated theorem provers; any OPENMATH-aware computer algebra system can therefore, in principle, perform unit conversions on OPENMATH expressions (cf. [46] for details). In contrast, RDF and OWL do not allow for defining mathematical operators and functions in a way as straightforward as in OPENMATH. SWEET and QUDT introduce custom OWL properties for describing conversion factors (e.g. *qudt:conversionMultiplier*), for which applications would have to provide hard-coded support – until recently. With SPIN (SPARQL Inferencing Notation [22]), there is an emerging standard for representing rules and constraints on RDF graphs, which has been utilized for converting quantities described using QUDT [24]. For computation, SPIN draws on the basic arithmetic operations supported by the SPARQL RDF query language [36]. More complex functions can, in principle, be provided as SPIN rules; so far, there is, however, just one library that implements SPIN [23].

4.4 Interaction

Applications taking advantage of the semantic publishing of units and their quantities have been experimented with by various authors, albeit the lack of authoring support. The unit conversion service [42,38] by Jonathan Stratford, which users can easily extend by uploading new Content Dictionaries (CDs) with new units and conversion rules, provides a good example of the power of semantically annotated units. Besides having implemented a service, Stratford has also identified the difficulties of unit conversion and the limitations of OPENMATH’s current state with regard to unit representation.

Stratford’s conversion service is interactive in that users can enter quantities into a web form and upload definitions of new units. We have additionally made it interactively accessible from web documents that contain MATHML formulas with OPENMATH annotations, as created by the publishing pipeline explained in section 2.2 (cf. [16]). This interaction with units in publications has, however, remained a proof of concept so far, as *producing* suitably annotated documents required manual authoring of quantity expressions in OPENMATH XML markup – a barrier that we are trying to overcome with the work presented in this paper.

Wolfram|Alpha [48], another interactive (web) service, provides unit conversion capabilities through its API [44] and widgets⁷ [49]. However, as its preferred input representation is natural language, and its output representation does not make the semantic structures explicit, we did not consider it for our research.

⁷ Mini-apps built on top of Wolfram|Alpha queries [50].

5 Semantic Authoring of Units and Quantities

We have revised the available methods and technologies and established that semantic authoring support for units does not formally exist at present. Consequently, we set out to make the first steps towards extending one of the more prepared software solutions, namely sTeX , with a special authoring module for units, by building on the existing pre-semantic toolbox of the *SIunits* $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ package. sTeX [25] is essentially a collection of $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ packages that offer semantic macros. sTeX can be translated into XML markup using $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{XML}$ [31] bindings, thus enabling easier subsequent processing – including semantic web publishing (cf. [11]). Our units extension follows a similar approach⁸.

As described in section 4.2, *SIunits* provides an sTeX -like content authoring interface. For our running example, we are interested in authoring $\boxed{100 \text{ km/h}}$ in order to create the content representation shown in Listing 1. There are many ways to author the representation in $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$, e.g. via $\text{\$}\text{\texttt{textrm}\{100\}\text{\texttt{,km/h}}\text{\$}$. The *SIunits* package makes the process less ad-hoc by focusing on the content and factoring out the presentational quirks, in the form of package options. Hence, one would instead write the more semantic $\text{\texttt{\unit}\{100\}\{\text{\texttt{kilo}}\text{\texttt{metre}}\text{\texttt{per}}\text{\texttt{hour}}\}}$. Moreover, the use of sTeX (unit) modules enables a more appropriate semantic (markup) representation of a quantity-unit pair by eliminating the inadequate times operator (cf. Listing 1) with a generic quantity constructor of the form: $\text{quantityFN} : \text{real} \times \text{unit} \rightarrow \text{quantity}$. Also, individual unit constructors can be defined (e.g. $\text{unitFN} : \text{real} \rightarrow \text{quantity}$) to further simplify the authoring process, e.g. $\text{\texttt{\unit}\{100\}\{\text{\texttt{gramme}}\}}$ would be authored as $\text{\texttt{\gramme}\{100\}}$.

It is interesting to observe that a completely different motivation than ours, namely to provide a convenient and centralized interface to control the *presentation* of the unit entities on a document level, essentially leads to the same result that we desire – a *semantics-oriented* authoring interface.

In our effort to leverage this functionality, we first created a $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{XML}$ binding for the *SIunits* package. It helped us to pinpoint the semantic map between the interface and the OPENMATH representation and provided a non-invasive semantic enrichment for $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ documents based on the package. Next, we use the gained understanding in building a native sTeX module for units, roughly based on the *SIunits* interface. Table 1 shows a small snippet comparing the different stages. One easily notices the abbreviative power of the sTeX approach, which hides the verbose and overly complex binding declaration under its hood, exposing the author to a controlled $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ vocabulary and facilitating reuse.

The Planetary eMath3.0 system, into which we have integrated the components of our semantic publishing pipeline, provides in-browser editing of sTeX documents with semantic syntax highlighting as well as context-aware autocompletion of semantic macros and links [27]. In the same environment, the user can

⁸ The *SIunits* bindings and sTeX extension will be released in the respective bundles (the arXMLiv binding library and the sTeX package on CTAN) with the authors' strong commitment to free software licenses compatible with the originals.

Table 1. Definitions for `\kilo\metre`, typeset as ‘km’

Language	Definition	Semantics
\LaTeX	<pre>\newcommand{\kilo}{\ensuremath{\mathrm{k}}} \newcommand{\metre}{\ensuremath{\mathrm{m}}}</pre>	✗
L ^A T _E X _M L	<pre>DefConstructor('kilo','' <ltx:XMApp> <ltx:XTok meaning="prefix" cd="units_ops1"/> <ltx:XTok meaning="kilo" cd="units_siprefix1"> k </ltx:XTok> #1 </ltx:XMApp>'); DefConstructor('metre','' <ltx:XTok meaning="metre" cd="units_metric1"> m </ltx:XTok>');</pre>	✓
$\S\TeX$	<pre>\symdef[name=kilo,cd=units_siprefix1]{kiloPX}{\mathrm{k}} \symdef[name=metre,cd=units_metric1]{metre}{\mathrm{m}} \symdef[name=prefix,cd=units_siprefix1]{prefixFN}{} \symdef{kilo}[1]{\mixfixii}{\kiloPX}{\prefixFN}{#1}{}</pre>	✓

interact with the published versions of these documents, as we will explain in the following section.

6 Interaction with Units and Quantities

Given the provisions for authoring support, we move to the added-value benefits one could reap from interacting with a published document. This section details relevant use cases and explains the prerequisites that are already available.

6.1 Unit (System) Preference Service

A concrete scenario for a prospective service that would take advantage of semantically published papers, based on the ideas from section 3, can be evolved on top of common published material like *cooking recipes*. These provide a good use case thanks to the high density of units and quantities they contain. Moreover, the physical quantities are restricted to a small subset (quantity/mass related units) including special types of *units* [1] which are not formally defined and might prove to be misleading:

$$\begin{aligned}
 1 \text{ teaspoon (tsp)} &\approx 5 \text{ millilitres (mL)} \\
 1 \text{ tablespoon (tbsp)} &\approx 15 \text{ millilitres (mL)} \\
 1 \text{ cup} &\approx 250 \text{ millilitres (mL)}
 \end{aligned}$$

The idea of the *unit (system) preference* service is to allow the user/reader to choose a preferred system of units (e.g. imperial, metric) or simply preferred

types of units (e.g. “minutes” instead of “hours”, “kilogrammes” instead of “grammes”) for the representation of physical quantities and then seamlessly adapt the document to these preferences. This can only be achieved at the end of the semantic publishing pipeline: The publishing process implemented by the Planetary eMath3.0 system requires a machine-comprehensible representation of knowledge (here: units and quantities) as described in section 3 and generated, e.g., via the authoring support introduced in section 5, and then applies a *semantics-preserving transformation*, resulting in a human-comprehensible published document with user-invisible but machine-readable annotations (here: XHTML with OPENMATH-annotated MATHML formulae) [27]. Into these annotations, the Planetary frontend hooks interactive services, utilizing the JOBAD library (Javascript API for OMDoc-based Active Documents [16]), which provides for communication with web services, manipulation of the user-visible as well as the machine-readable parts of the document, and providing user interface primitives such as a context menu. In our *unit (system) preference* service for Planetary, the computational facilities required for converting quantities are provided by the Universal OPENMATH Machine web service [51], which reasons and computes with OPENMATH objects. Figure 1 visualizes the architecture and data flow.

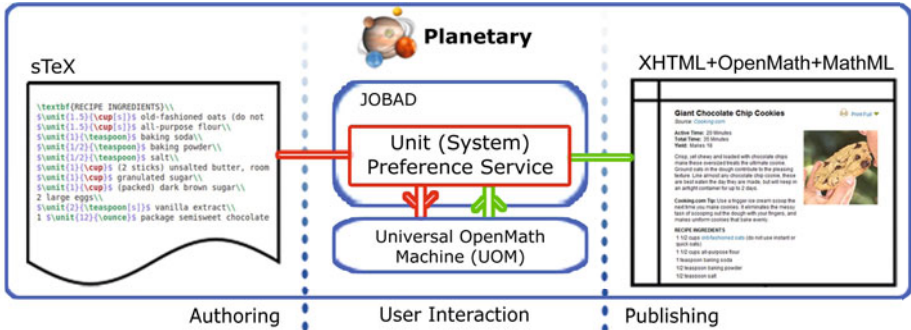


Fig. 1. Architecture of the Unit (System) Preference service, and data flow in the Chocolate Chip Cookies recipe [10] use-case

We chose a *cooking recipe* use-case as a proof of concept for such a generic (web) service taking advantage of semantically published units; this should not be treated as a fully fledged, independent software product. The design choice was to present the user only with the generic, commonly-used, unit systems (e.g. metric, imperial) and not with specific units as options for preference, especially due to the numerous existing types of units (some of which might not even apply to the subset of unit occurrences in the document). Once the user selects a unit preference, either from the *Available Units Preference Settings* bar on top, or from the context menu, each semantically annotated unit in the document passes through the data flow shown in figure 1. It is replaced with the converted quantity

and the new unit (again the choice was made for the most commonly used units, e.g. *tablespoons* would be converted to *milli litres* for the metric preference), rendered in a human-readable way, but with the semantic structure preserved in a machine-readable annotation. Figure 2 demonstrates the interactive user interface of the unit (system) preference service as well as Planetary’s rendering of the sample cooking recipe.⁹ For further technical details about the architecture of the service, we refer to [8].

Giant Chocolate Chip Cookies

Written by: Mihai Cîrlănuș

Description: Cooking Recipe use case for the Unit (System) Preference Service

Available Unit Preference Settings: metric imperial food

Active Time: (20 mins)

Total Time: (35 mins)

Yield: Makes 1 Convert all to metric
imperial

Crisp, yet chewy and loaded with food chips make these oversized treats the ultimate cookie. Ground oats in the dough contribute to the pleasing texture. Like almost any chocolate chip cookie, these are best eaten the day they are made, but will keep in an airtight container for up to (2 days).

RECIPE INGREDIENTS

- (1.5 cups) old-fashioned oats (do not use instant or quick oats)
- (1.5 cups) all-purpose flour
- (1 teaspoon) baking soda
- (1/2 teaspoon) baking powder
- (1/2 teaspoon) salt
- (1 cup) (2 sticks) unsalted butter, room temperature
- (1 cup) granulated sugar
- (1 cup) (packed) dark brown sugar
- 2 large eggs
- (2 teaspoons) vanilla extract
- 1 (12 oz) package semisweet chocolate chips

DIRECTIONS

Preheat oven to (350 °F).

Fig. 2. Screenshot of the Unit (System) Preference service for the cooking recipe

Note that Google’s cooking recipe search [18] offers a related web service. However, their semantic markup for cooking recipes [17] does not extend down to an explicit representation of units and quantities, and thus the interactive conversion capabilities are limited or non-existent.

6.2 Prospective Services Based on Semantically Published Units

Having described in detail one service that enhances the user experience by publishing units semantically, we now list further potential services and applications that the same technology could enable:

⁹ http://trac.mathweb.org/planetary/wiki/Demo_PlanetBox

- **Mapping Natural Sciences Concepts to their Respective Units:** defining Content Dictionaries that would describe the connection of units to general natural sciences concepts like *force* (measured in Newtons: $N = \frac{kgm}{s^2}$ or any variant of the ratio) or *energy* (measured in Joules: $J = Nm = \frac{kgm^2}{s^2} = \dots$) and plenty of other examples. The interconnection of concepts in sciences: $Energy = Force \times displacement$ can further enable scientific formula “spell checking” which might prove to be of great value to physicists, astronomers and many others.
- **Unknown Unit Lookup:** In theoretical scientific papers authors usually use abbreviations for concepts (e.g. N for *Newtons* – the unit for *force*) without mentioning anything about units/dimensions, which might turn out to be difficult for the readers who would be interested to know, for example, the order of measurement (magnitude) for the unknown physical quantities and also a (small) description of the respective concept (e.g. Pa is the unit for *pressure*). Defining a generic way in which semantics can be added to such unknown symbols will enable showing/hiding units for expressions/formulas.
- **Unit and Quantity Semantic Search:** a library-level service that would allow searching for units by their type, name and magnitude and return the relevant results independently of the measuring standard of the occurrences in the paper (e.g. imperial or metric) and also independent of their form (N or $\frac{kgm}{s^2}$).¹⁰
- **Quantity and Unit’s Magnitude Manipulation:** a document interaction service that is able to transform for example $100N \rightarrow 0.1kN$ or 0.1×10^3N or $0.1 \times 10^3 \frac{kgm^2}{s^2}$. This can be useful when it comes to simplifying representations and adapting them consistently to a certain type of magnitude (for example *all occurrences of force expressions should have their unit represented in kN*).

As detailed at the beginning of this paper, having a standard, uniform understanding of units and quantities can prevent hazards and even eliminate entire compatibility check processes in industry. The presented list of prospective enabling technologies shows only a few of the numerous opportunities of interacting with units and quantities in semantically published documents and serves as a strong motivation for future research in this direction.

7 Enabling Semantic Units in Legacy Corpora

The ARXMLIV corpus is the ideal environment for the identification of units and quantities since it contains a collection of more than 600,000 scientific publications. It is based on Cornell University’s ARXIV e-Print archive [3] originally typeset in L^AT_EX, converted to XML in order to achieve easy machine-readability, partial semantics recovery and clear separation of document modalities such as natural language and mathematical expressions [41]. Currently, the project

¹⁰ In contrast, state-of-the-art scientific publication search services, such as Springer’s L^AT_EX search [40], do not support the semantics of units.

has achieved a successful conversion rate of nearly 70% to a semantically enriched XHTML+MATHML representation, natively understandable by modern web browsers [26].

A proof-of-concept check, performed via the ARXMLIV build system (see [41]) revealed roughly 150 ARXIV articles using the *SIunits* package, with an outlook for close to tripling the number when considering sibling packages such as *units* and *SIunitsx*. This gives our work on creating a semantic binding for *SIunits* an even stronger benefit, as we can directly and non-invasively enrich legacy publications, putting them one step further on the path to semantic publishing. An additional, mid-term benefit is the opportunity to build a linguistic *Gold Standard* for units; we created both legacy (to presentational MATHML) and semantic (to OPENMATH) bindings in order to provide a raw, presentational output and its annotated, semantic counterpart. Having both as a basis, unit spotters can then be developed using methods of Computational Linguistics and Machine Learning, further enriching the ARXMLIV corpus.

Such enhancements not only enable the interactive services of semantic publishing on legacy corpora, but also provide a tempting outlook to the development of an ecosystem of linguistic analysis modules, which can draw on the captured semantics of units and quantities, as originally envisioned by the LAMA-PUN project [15].

8 Conclusions and Future Work

Units and quantities are sufficiently wide-spread and important to not be disregarded from the context of semantic documents. Unfortunately, by now, there have been only isolated approaches (see section 4) to exploit the semantic power of units. Moreover, the wide range of existing unit types and representations makes it almost impossible to identify and semantically enrich all of them, especially when we are talking about occurrence contexts as unrelated as cooking recipes, medical prescriptions, technical documents or scientific papers.

We have emphasized the importance of three major components of the semantic publishing process for units – *representation*, *authoring* and *interaction* –, and detailed technologies for improving each of them. Moreover, by providing a cooking recipe interaction use-case as well as a series of further potential services and applications on top of semantically published units, we contribute means of better manipulation and interpretation of *units and quantities* to the Semantic Publishing Industry and to legacy corpora.

Acknowledgments. The authors would like to thank Michael Kohlhase for his extensive support and advice regarding the writing of this paper, Anton Antonov for writing the LATEXML bindings for the *SIunits* L^AT_EX package, and the other Planetary developers¹¹ for providing the context for developing our authoring and interaction services. This extended and revised version of a submission to the ESWC 2011 workshop on Semantic Publishing (SePublica)

¹¹ <http://trac.mathweb.org/planetary/wiki/people>

has benefited from the extensive helpful suggestions provided by the anonymous peer reviewers, and from constructive feedback and further suggestions given by the participants of the workshop.

References

1. Code of Federal Regulations – Food and Drugs, http://edocket.access.gpo.gov/cfr_2004/aprqtr/21cfr101.9.htm
2. Aviation Safety – Air Canada Accident Report, <http://aviation-safety.net/database/record.php?id=19830723-0> (visited on October 25, 2010)
3. arXiv.org e-Print archive, <http://www.arxiv.org>
4. MathML 3.0. Recommendation. W3C (2010), <http://www.w3.org/TR/MathML3>
5. Barany, M.J.: [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: WikiSym (2010)
6. Measurement Units Ontology, http://forge.morfeo-project.org/wiki_en/index.php/Measurement_Units_Ontology (visited on April 16, 2011)
7. Buswell, S., et al.: OpenMath 2.0. Tech. rep. The OpenMath Society (2004), <http://www.openmath.org/standard/om20>
8. Cîrlănu, M.: Authoring, Publishing and Interacting with Units and Quantities in Technical Documents. BSc. Thesis. Jacobs University Bremen (2011)
9. Collins, J.B.: OpenMath Content Dictionaries for SI Quantities and Units. In: Carette, J., Dixon, L., Coen, C.S., Watt, S.M. (eds.) MKM 2009, Held as Part of CICM 2009. LNCS, vol. 5625, pp. 247–262. Springer, Heidelberg (2009)
10. Cooking.com – Giant Chocolate Chip Cookies, <http://www.cooking.com/recipes-and-more/recipes/Giant-Chocolate-Chip-Cookies-recipe-5112.aspx> (visited on March 5, 2011)
11. David, C., et al.: Publishing Math Lecture Notes as Linked Data. In: Aroyo, L., Antoniou, G., Hyvönen, E., ten Teije, A., Stuckenschmidt, H., Cabral, L., Tudorache, T. (eds.) ESWC 2010. LNCS, vol. 6089, pp. 370–375. Springer, Heidelberg (2010)
12. DBpedia, <http://dbpedia.org> (visited on January 23, 2010)
13. Davenport, J.H., Naylor, W.A.: Units and Dimensions in OpenMath (2003), <http://www.openmath.org/documents/Units.pdf>
14. Fink, J.L., et al.: Word add-in for ontology recognition: semantic enrichment of scientific literature. BMC Bioinformatics 11, 103 (2010)
15. Ginev, D., et al.: An Architecture for Linguistic and Semantic Analysis on the arXMLiv Corpus. In: Applications of Semantic Technologies Workshop at Informatik (2009), http://www.kwarc.info/projects/lamapun/pubs/AST09_LaMaPUn+appendix.pdf
16. Giceva, J., Lange, C., Rabe, F.: Integrating Web Services into Active Mathematical Documents. In: Carette, J., Dixon, L., Coen, C.S., Watt, S.M. (eds.) MKM 2009, Held as Part of CICM 2009. LNCS, vol. 5625, pp. 279–293. Springer, Heidelberg (2009)
17. Google Cooking Recipe Publishing Schema, <http://www.google.com/support/webmasters/bin/answer.py?answer=173379> (visited on June 5, 2011)
18. Google Cooking Recipe Search, <http://www.google.com/landing/recipes/> (visited on May 6, 2011)

19. Groza, T., et al.: SALT – Semantically Annotated LATEX for Scientific Publications. In: Franconi, E., Kifer, M., May, W. (eds.) ESWC 2007. LNCS, vol. 4519, pp. 518–532. Springer, Heidelberg (2007)
20. Heldoorn, M.: The SIunits package: Consistent application of SI units, <http://mirror.ctan.org/macros/latex/contrib/SIunits/SIunits.pdf> (visited on March 13, 2011)
21. QUDT – Quantities, Units, Dimensions and Data Types in OWL and XML, <http://www.qudt.org> (visited on July 15, 2011)
22. SPIN – Overview and Motivation. Member Submission. W3C, <http://www.w3.org/Submission/2011/SUBM-spin-overview-20110222/>
23. The TopBraid SPIN API, <http://topbraid.org/spin/api/> (visited on July 15, 2011)
24. Units ontology with SPIN support published, <http://composing-the-semantic-web.blogspot.com/2009/08/units-ontology-with-spin-support.html> (visited on July 15, 2011)
25. Kohlhase, M.: Using LATEX as a Semantic Markup Format. Mathematics in Computer Science, 2.2 (2008)
26. Kohlhase, M., et al.: MathWebSearch 0.4, A Semantic Search Engine for Mathematics (2008), <http://mathweb.org/projects/mws/pubs/mkm08.pdf>
27. Kohlhase, M., et al.: The Planetary System: Web 3.0 & Active Documents for STEM. In: Procedia Computer Science 4 (2011): International Conference on Computational Science (ICCS). Finalist Executable Papers Challenge (2011)
28. Kohlhase, M., Rabe, F.: Semantics of OpenMath and MathML3. In: 22nd OpenMath Workshop (2009)
29. Lange, C.: Ontologies and Languages for Representing Mathematical Knowledge on the Semantic Web. Semantic Web Journal (accepted, 2011) <http://www.semantic-web-journal.net/content/new-submission-ontologies-and-languages-representing-mathematical-knowledge-semantic-web>
30. CNN – NASAs metric confusion caused Mars orbiter loss, http://articles.cnn.com/1999-09-30/tech/9909_30_mars.metric_1_mars-orbiter-climate-orbiter-spacecraft-team?_s=PM:TECH (visited on October 29, 2010)
31. LaTeXXML: A LATEX to XML Converter, <http://dlmf.nist.gov/LaTeXXML/> (visited on March 3, 2011)
32. Oxford English Dictionary. “unit” definition, <http://dictionary.oed.com/entrance.dtl> (visited on October 29, 2010)
33. PAUX Technologies, <http://paux.de> (visited on October 10, 2010)
34. Planetary Developer Forum, <http://trac.mathweb.org/planetary/> (visited on January 20, 2011)
35. PlanetMath.org – Math for the people, by the people, <http://planetmath.org> (visited on January 6, 2011)
36. SPARQL Query Language for RDF. Recommendation. W3C, (2008), <http://www.w3.org/TR/2008/REC-rdf-sparql-query-20080115/>
37. Raskin, R.G., Pan, M.J.: Knowledge representation in the semantic web for Earth environmental terminology (SWEET). Computers & Geosciences 31 (2005)
38. Stratford, J., Davenport, J.H.: Unit Knowledge Management. In: Autexier, S., Campbell, J., Rubio, J., Sorge, V., Suzuki, M., Wiedijk, F. (eds.) AISC 2008, Calculemus 2008, and MKM 2008. LNCS (LNAI), vol. 5144, pp. 382–397. Springer, Heidelberg (2008)
39. The International System of Units (SI) Bureau International des Poids et Mesures 8 edn. (2006), http://www.bipm.org/utls/common/pdf/si_brochure_8_en.pdf

40. Springer, (ed.) LATEX Search, <http://www.latexsearch.com> (visited on April 16, 2011)
41. Stamerjohanns, H., et al.: Transforming large collections of scientific publications to XML. *Mathematics in Computer Science*, 3.3 (2010)
42. Stratford, J.: Creating an extensible Unit Converter using OpenMath as the Representation of the Semantics of the Units. Tech. rep. 2008-02. University of Bath, <http://www.cs.bath.ac.uk/pubdb/download.php?resID=290>
43. Semantic Web for Earth and Environmental Terminology (SWEET). NASA, <http://sweet.jpl.nasa.gov/> (visited on August 22, 2010)
44. Wolfram|Alpha API, <http://www.wolframalpha.com/developers.html> (visited on May 5, 2011)
45. US Metric Association “Unit Mixups” article, <http://lamarcolostate.edu/~hillger/unit-mixups.html> (visited on October 25, 2010)
46. Vrandečić, D., et al.: Semantics of Governmental Statistics Data. In: *Web Science (2010)*, <http://journal.webscience.org/400/>
47. Wikipedia: Hebrew calendar, http://en.wikipedia.org/wiki/Hebrew_calendar
48. Wolfram|Alpha, <http://www.wolframalpha.com> (visited on May 5, 2011)
49. Wolfram|Alpha Units and Measures Widgets, <http://developer.wolframalpha.com/widgets/gallery/category/?cat=units> (visited on May 5, 2011)
50. Wolfram|Alpha Widgets, <http://developer.wolframalpha.com/widgets/> (visited on May 5, 2011)
51. Zamdzhev, V.: Universal OpenMath Machine. BSc. Thesis. Jacobs University Bremen (2011)