

Semantic Search for Quantity Expressions

Guided Research Proposal

Tom Wiesing
Supervisor: Michael Kohlhase
Jacobs University, Bremen, Germany

March 1, 2015

1

EdN:1

Abstract

In this proposal we describe how to introduce units to MathWebSearch. The aim of the project is build a complete semantics-aware system that searches a corpus of documents for quantity expressions. The project will be based on the existing MathWebSearch system and related technologies. ²

EdN:2

1 Introduction

In this paper we want to give an approach to Semantic Search for Quantity Expressions. A quantity expression is a scalar together with a physical unit, for example $25 \frac{\text{m}}{\text{s}}$ where 25 is the scalar and $\frac{\text{m}}{\text{s}}$ the unit. This quantity expression is equivalent to $90 \frac{\text{km}}{\text{h}}$ (with equivalent in this sense meaning it expresses the same quantity) although the units are different. In a semantic search for quantity expressions we want to be able to search for a certain quantity expression and find any equivalent ones as well.

MathWebSearch (MWS for short) is an existing semantics-aware system to search \LaTeX documents ¹ for mathematical formulae [HKP14]. As it is semantics-aware it not only searches for formulae in a simple-minded “text search” way but also includes simple transformation rules, such as $a + b = b + a$. Additionally it can deal with wildcards such as $x + \sqrt{x}$. For this query MWS would deliver formulae as above where x has been substituted with any sub-formula.

So far the transformation system has been used by MWS exclusively for mathematical formulae. In this paper we propose an extension for quantity expressions which should prove useful for physicists. The end user will search, for example, 100°C and also get results which show 212°F or 373.15K .

This proposal is organised as follows: In section 2 we shortly describe the existing MathWebSearch system. Then we proceed in section ?? to describe the work to be done and finally discuss further aspects of the system in 4.

¹EdNOTE: Find a Math Co-supervisor

²EdNOTE: Spin

¹Technically, MWS itself can only search XHTML documents. However with the help of \LaTeX ml [Mil] it also handles \LaTeX documents.

2 An Overview of the existing MathWebSearch system

2.1 The Components of MathWebSearch

As mentioned above, MathWebSearch is a search engine for formulae. MWS consists of 3 main components as well as a frontend ² [KP].

The backend consists of 3 main components,

1. a crawler,
2. a core system and
3. a public rest API.

The crawler, as its name suggests, crawls corpora for formulae. For each corpus MWS uses, a separate crawler has to be implemented. The crawled formulae are passed to the core system and indexed in an MWS Harvest. The core system is also responsible for parsing queries and sending results back to the REST API. This is done by searching the harvest only. In order to make a semantic search for quantity expressions we will adapt this crawler to find and harvest quantity expression instead.

Because MWS is semantics-aware, the harvest can not only contain the exact formulae that are found in the original corpus but also all versions that are equivalent to it. These are generated with the help of MMT and theory graphs, see sections 2.2 and 2.3

The frontend for MathWebSearch, which is not part of MWS itself but running client-side in a web browser, is written in HTML5, CSS and JavaScript. It accesses the REST backend and depends on MathML support to render Mathematics. When the client enters a \LaTeX formula to search for, the \LaTeX ml daemon [Gin] is used to transform the query in content MathML. Next, the client renders the MathML (to show the formula the user is searching for) and then sends the query off to the MWS API. Upon receiving results, the client renders them and links to the original documents.

There are several implementations of frontends and crawlers as well as extensions of MathWebSearch. One particular implementation is capable of crawling the arXiv corpus, which contains approximately 750.000 documents. A list of demos can be found at [Mat].

2.2 Mathematical Theory Modeling with MMT

MMT is a Module system for Mathematical Theories [RK13]. In MMT we can represent mathematical theories, views and imports.

A mathematical theory in MMT is a set of typed symbols and statements. As a simple example, the theory of semigroups can be written down as following:

Semigroup		
G		
\circ	:	$G \rightarrow G \rightarrow G$
assoc	:	$(x \circ y) \circ z = x \circ (y \circ z)$

In the example above ³, we first define the a symbol G (the set of the Semigroup), as well as a (binary) operation \circ and then the law of associativity for the operation \circ . EdN:3

²The frontend is not part of MWS directly but rather built on top of the REST API, more on this later.

³EDNOTE: Check with miancu if the example is correct this way

Theories in MMT can be related in two ways, with imports and views. If theory A imports theory B then theory A contains all symbols and statements from theory B. In particular, every satisfiable statement that can be made in theory A is also a satisfiable statement in theory B.

As an example of an import, the theory of a *commutative groups* imports the theory of *groups*. The *commutative group* theory additionally has the axiom of commutativity. We can consider an import as a type of inheritance relation.

In contrast to this a view is a type of interpretation relation. For example, there is a view from the theory of *Monoids* to the theory of *Natural numbers*. Again, a satisfiable statement in the viewed theory can be translated into a satisfiable statement in the target theory.

Thus a theorem in one theory can be translated along views as well as along imports. For example, given a theorem in the theory of groups, we automatically get a theorem in the theory of commutative groups.

2.3 Using Theory Graphs

The theories and their relations can be represented with the help of a graph. This graph is then called theory graph. In figure 1 we can see a simple example of a theory graph.

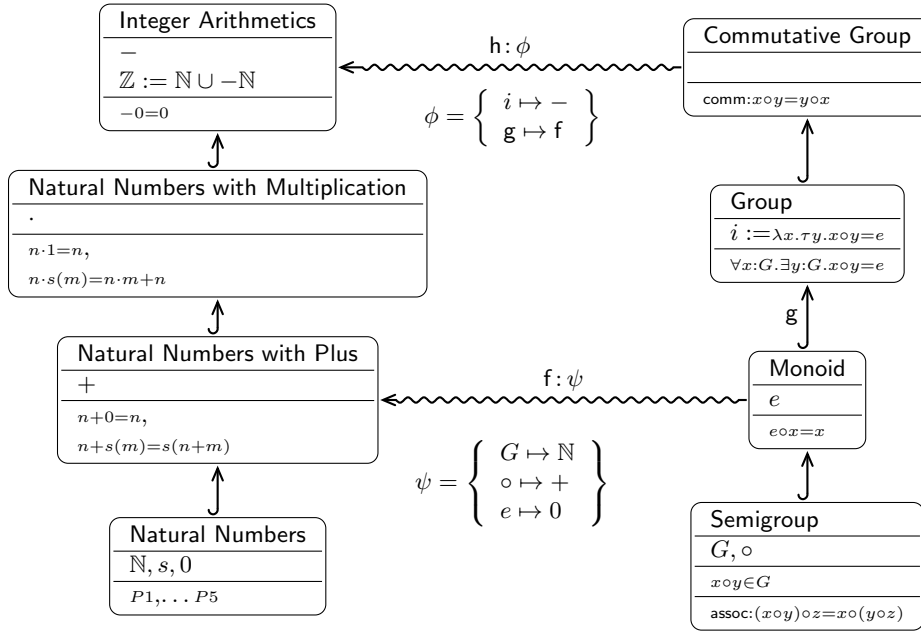


Figure 1: A simple theory graph. Imports are represented as solid edges and views as dashed edges.

Every theory in MMT has a (globally unique) URL, called the MMT URL ⁴. Furthermore, every theorem in MMT can also be represented via a URI. As explained above, if 2 theories are related via a view or import, theorems can be translated along this relation. EdN:4

Sometimes it is useful not to write down a theorem explicitly, but only give an existing theorem and translate this along a view or import. The MMT URL of this induced theorem

⁴EdNOTE: Example URL here

can then be given using the view⁵. This URL contains enough information for MMT to generate the theorem in explicit form [IKP14]. EdN:5

In an MWS Harvest we exploit this property of induced theorems. We have several represented theorems, the formulae in the corpus, and many more induced theorems, other equivalent representations of these formulae. With the help of MMT, the MMT Harvest is built by representing all theses theories (formulae) explicitly.

Applying this principle to quantity expressions, we can consider different units as different theories and different quantity expressions different theorems belonging to their respective unit theories. The unit conversions can then be represented via views. The information on how to translate from one unit to another will be contained within the view. If we know all formulae in a corpus we can then generate an MWS Harvest that contains all the representations.⁶ EdN:6

3 Schedule and Work Packages

We want to build a complete system based on MathWebSearch that can perform a semantic search of quantity expressions on a specific corpus. As such it needs 4 components: (1) *A crawler* that crawls documents for quantity expressions, (2) *a core system* that transforms these quantity expressions and harvests them as described above as well as provides search functionality, (3) *a rest API* that exposes this search functionality and (4) *a frontend* that servers as the endusers⁷ interface to the system. EdN:7

8

EdN:8

We will build the frontend using HTML, CSS and JavaScript. Because it serves as the main interface to the system

Furthermore the system will be flexible with respect to

- *the quantity expressions it understands*, Adding new units should be a simple process so that the system can be easily extended to also capture rare units and
- *the corpus* which should be exchangable easily.

The crawler will be based on existing MWS crawlers. As it is very difficult to find and mark up all units in a corpus automatically, we will restrict ourself to an artificial corpus containing only lorem ipsum⁹ documents with randomly inserted units. This does not solve the problem of writing a crawler, however it will enable us to build a proof-of-concept system first. EdN:9

The core system will almost completely be inherited from existing implementations. As it does not search formulae, it has to be made aware of unit translations. This will be done with the help of theory graphs. It is not yet clear how exactly these will look, a speculation can be found below in ??.

The rest-based API will have to take the user input from the client and pass the entire quantity expression on to the core system. Once it receives a result, it will have to translate these back into human-readable form and then pass these on to the frontend.

⁵EdNOTE: Example URI here

⁶EdNOTE: Should we move this paragraph to the next section?

⁷EdNOTE: Spelling?

⁸EdNOTE: Describe the components again

⁹EdNOTE: What is my contribution, the exact details are not important.

As with the existing system, the frontend will be a web page that incorporates 3 main elements, (1) a search input for a value (2) a search input for a unit and (3) a result page that displays results and links to the found documents.

4 A Semantic Search For Quantity Expressions Based On MathWebSearch

10

EdN:10

A bigger challenge with the approach will be to find a standardised representation for quantity expressions, mainly for units³ ¹¹.

EdN:11

Because we want to translate between these representations however, it is insufficient to just have such a representation. We additionally need to find a standardised way of translating between units need to exist as well. Because MWS uses Theory Graphs and MMT to translate between equivalent formulae, we will have to write down all units in the form of a theory graph. Translations between units will be represented as views.

To translate between units we could use simple translation formulae such as $xK = x + 273.15^\circ C$. However because we want to support composite units (such as meters per second $\frac{m}{s}$) as well, this can cause problems. In machine representation ¹², every quantity expression has to have a finite precision representation. When translating between 2 composite units, this can easily cause rounding errors. In particular, when an author gives an approximation of a quantity in a document, they might round differently depending on the units used.

EdN:12

Hence it will not be sufficient to use exact translations. Thus we should search for a range of values instead. This has recently been implemented as an extension to MathWebSearch [Ham]. It remains to be seen how exactly these ranges should be used.

Furthermore the units will have to be entered in some fashion ¹³ in the frontend. While it is trivial to design an interface where a single unit can be entered, it is non-trivial when we want to recognise composite units as well. Furthermore units with (si-)prefixes (such as kilo or mega) can either be recognised separately or as part of the unit (multiplying directly into the value of the quantity expression). There are several input formats that can be used: AsciiMath, L^AT_EX and MathML, to name only a few. Independent of the input format, the end result (delivered to the backend of the search engine) will have to be represented in one way.

EdN:13

Finally, there is the problem of choosing a suitable corpus and implementing a crawler which will find and mark up quantity expressions. The latter will be taken on by Stiv Sherko in a separate effort [She14]. The task of finding a suitable corpus is postponed for now. In case there is sufficient time a ycorpus can be marked up manually and crawled afterwards.

References

[Gin] Deyan Ginev. TheL^AT_EXml daemon: Editable math for the collaborative web.

¹⁰EDNOTE: Is the section title correct?

³Here, standardised means that all quantity expressions are represented in a standardised, machine-readable fashion, not equivalent quantity expressions being represented in exactly one way.

¹¹EDNOTE: MathML note OM CDs

¹²EDNOTE: Separate problem

¹³EDNOTE: The same?

- [Ham] Radu Hambasan. Ranges in MathWebSearch. private communication.
- [HKP14] Radu Hambasan, Michael Kohlhase, and Corneliu Prodescu. MathWebSearch at NTCIR-11. In Noriko Kando and Hideo Joho and Kazuaki Kishida, editors, *NTCIR 11 Conference*, pages 114–119, Tokyo, Japan, 2014. NII, Tokyo.
- [IKP14] Mihnea Iancu, Michael Kohlhase, and Corneliu-Claudiu Prodescu. Representing, archiving, and searching the space of mathematical knowledge. In Hoon Hong and Chee Yap, editors, *Mathematical Software - ICMS 2014 - 4th International Congress*, volume 8592 of *Lecture Notes in Computer Science*, pages 26–30. Springer, 2014.
- [KP] Michael Kohlhase and Corneliu Prodescu. MathWebSearch manual. Web manual, Jacobs University.
- [Mat] Math WebSearch a semantic search engine. web page at <http://search.mathweb.org>. seen September 2008.
- [Mil] Bruce Miller. **LaTeX**ML: A \LaTeX to XML converter.
- [RK13] Florian Rabe and Michael Kohlhase. A scalable module system. *Information & Computation*, 0(230):1–54, 2013.
- [She14] Stiv Sherko. Extracting quantity and unit semantics from technical documents. Bachelor thesis proposal, Computer Science, Jacobs University, Bremen, 2014.