

# Semantic Search for Quantity Expressions

## Guided Research Proposal

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

January 25, 2015

### Abstract

In this proposal we describe how to approach Semantic search for Quantity Expressions, in particular how to make the existing MathWebSearch system aware of physical units.<sup>1</sup>

EdN:1

## 1 Introduction

In this paper we want to give an approach to Semantic Search for Quantity Expressions. A quantity expression is a number<sup>2</sup> together with a physical unit, for example  $25 \frac{\text{m}}{\text{s}}$  where  $\frac{\text{m}}{\text{s}}$  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.

EdN:2

MathWebSearch (MWS for short) is an existing semantics-aware system to search  $\text{\LaTeX}$  documents<sup>3</sup> <sup>1</sup> for mathematical formulae [2]. 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.

EdN:3

---

<sup>1</sup>EdNOTE: Continue and finish abstract

<sup>2</sup>EdNOTE: Find a better expression for this.

<sup>3</sup>EdNOTE: Do not quote the damon but something general here.

<sup>1</sup>Technically, MWS itself can only search XHTML documents. However with the help of  $\text{\LaTeX} \text{Xml}$  [1] it also handles  $\text{\LaTeX}$  documents.

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 will be useful in physical papers. 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<sup>4</sup>: In section 2 we shortly describe and discuss the existing MathWebSearch system and then proceed in section 3 to describe in detail the proposed extension. Finally in section 4 we discuss possible problems with this approach and related work. EdN:4

## 2 Background - The existing MathWebSearch system

As mentioned above, MathWebSearch is a search engine for formulae. MWS consists of 3 components as well as a frontend <sup>2</sup> [3].

The backend consists of 3 main components,

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

The crawler, as its name suggests, crawls corpora for formulae. These formulae are then transformed into a readable format. Next, the core system build a search index and gets fed queries from the public API. It then processes these queries and send results back to the client via the restful API.

The frontend, 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  $\text{\LaTeX}$  formula to search for, the  $\text{\LaTeX}$ ml daemon [1] is used to transform this into an XHTML query. 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.

An extension to MWS, Thema-Search, <sup>5</sup> allows to search for Math and EdN:5

---

<sup>4</sup>EdNOTE: Update this if we change the structure

<sup>2</sup>The frontend is not part of MWS directly but rather built on top of the HTTP API, more on this later.

<sup>5</sup>EdNOTE: Quote Thema-search

Text simultaneously, however this done with the help of a separate search engine.

For each corpus MWS uses, a separate crawler has to be implemented. One current crawler exists for the ArXiV system,<sup>6</sup> which has a corpus of approximately ??? documents<sup>7</sup>. The frontend is deployed at ???<sup>8</sup>

- <sup>9</sup>. Disadvantages of the current approach
- has to be re-generated each time a document is added

### 3 The proposed extension

The goal of the guided research is to get a complete system that searches a corpus of documents for units as already indicated above. This system should be nicely accessible to the end user. Furthermore, it should be extendable with respect to

- *the unit system*, Adding new units should be easy and most translations should be deducted by the system automatically.
- *and the searchable corpus*. It should be easy to search a different corpus of documents provided units are properly marked up inside it.

As with the existing system, the frontend should be a web page that works in all modern browsers. It should be mobile friendly. It should communicate directly with the backend via a RESTful API. The frontend should incorporate 4 main elements:

1. a search input for a (real) quantity,
2. a search input for a unit, described further in section 4,
3. an additional search input for text and
4. a result page that displays results and links to the found documents.

The REST-based backend should have 3 major tasks:

1. translate raw unit input into a standardised form (see section 4 for details),

---

<sup>6</sup>EdNOTE: Quote this

<sup>7</sup>EdNOTE: Get an estimate here

<sup>8</sup>EdNOTE: Quote the frontend here, maybe mention that this is thema-search

<sup>9</sup>EdNOTE: Give an overview of the advantages / disadvantages here

2. search for documents using the input from the frontend and
3. make the documents available to the enduser.

The corpus <sup>10</sup>

EdN:10

- should consist of a lot of tex documents
- should have marked up units
- ideally, if a single document is added, only the new corpus should have to be re-scanned (procedural generation)
- should be easily exchangeable

The unit transition system <sup>11</sup>

EdN:11

- should be a graph
- should have few connected components and each of the components should be sparse (i. e. few connections)
- translation are:
  - either a factor towards a single unit
  - or a composition of a factor together with a product or fraction of units <sup>12</sup>
  - Perhaps include prefixes somehow?

EdN:12

## 4 Problems and related Work

There are several problems with this approach. <sup>13</sup>

EdN:13

The units have to be entered in some fashion in the search engine. 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 si-prefixes (such as kilo or mini) should also be recognised which is an additional problem. There are a few alternatives to use as input formats: AsciiMath, L<sup>A</sup>T<sub>E</sub>X and MathML, to name only 3. Independent of the

---

<sup>10</sup>EDNOTE: Write form here onwards

<sup>11</sup>EDNOTE: Write this part.

<sup>12</sup>EDNOTE: Figure out more details about this

<sup>13</sup>EDNOTE: Continue intro paragraph

input format, the end result (delivered to the backend of the search engine) should be MathML containing the unit in some yet-to-be-determined standard form.

The unit equivalences should be in the form of a theory graph<sup>14</sup> The equivalence (translation) itself contains both a translation for the quantity and the unit. The values searched for should be simple real numbers together with the formula. However since all numbers must be represented in a finite fashion, translations formulae can give problems. A natural quantity in one unit can be a repeating decimal in another unit. Also rounding of quantities might depend on the unit used. It is thus preferable not to search for an exact value but instead for a range of values. These ranges should depend on the unit. These ranges have just been implemented in MWS by Radu<sup>15, 16 17</sup>

Finally, there are the problems of finding a corpus and marking up units in this corpus. The second problem, marking up units in a corpus, will be taken on by ???<sup>18</sup> in a separate thesis. The first problem will be postponed for now. In case there is still time, a suitable corpus can be marked up manually and plugged into the system. For first testing, I will create a dummy corpus of lorem-ipsu style documents which will contain a random sampling of units.

<sup>19</sup>

## References

- [1] Deyan Ginev. The L<sup>A</sup>T<sub>E</sub>X<sub>ml</sub> daemon: Editable math for the collaborative web. URL: <http://latexml.mathweb.org>.
- [2] 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. URL: <http://research.nii.ac.jp/ntcir/workshop/OnlineProceedings11/pdf/NTCIR/Math-2/05-NTCIR11-MATH-HambasanR.pdf>.

<sup>14</sup>EDNOTE: Reference this properly.

<sup>15</sup>EDNOTE: Quote needed; reformulate

<sup>16</sup>EDNOTE: Write exactly how to use ranges

<sup>17</sup>EDNOTE: Check if the reference in (2) is clear

<sup>18</sup>EDNOTE: Who is doing unit finding?

<sup>19</sup>EDNOTE: Write final paragraph

- [3] Michael Kohlhase and Corneliu Prodescu. Mathwebsearch manual. Web manual, Jacobs University. URL: <https://svn.mathweb.org/repos/mws/doc/manual/manual.pdf>.