

# Units for MathWebSearch\*

## Guided Research Proposal

EdN:1

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

January 3, 2015

### Abstract

In this proposal we describe an approach to introduce Units to MathWebSearch<sup>2</sup>

EdN:2

## 1 Introduction

3

EdN:3

MathWebSearch (MWS for short<sup>4</sup>) is a system to search (latex) documents for mathematical formulae. Additionally it can also search for text in the documents<sup>5</sup>. However it not only searches for formulae in a simple-minded string way but also includes simple transformation rules, such as  $a + b = b + a$ . Additionally, it is possible to search with wildcards such as  $\textcolor{red}{x} + \sqrt{\textcolor{red}{x}}$ . In this example MWS delivers results of the given form where  $\textcolor{red}{x}$  is substituted with any sub-formula<sup>6</sup>.

EdN:4

EdN:5

EdN:6

MWS has been shown to be very useful for mathematicians<sup>7</sup>. The transformation system it uses can currently be used only for mathematical formulae which limits its applications. In this paper we propose an extension for physical units<sup>8</sup>. Instead of transforming mathematical formulae, the search

EdN:7

EdN:8

---

\*EdNOTE: Preliminary Title

<sup>2</sup>EdNOTE: Write abstract properly

<sup>3</sup>EdNOTE: Write an introductory sentence / paragraph?

<sup>4</sup>EdNOTE: should I really use abbreviations here?

<sup>5</sup>EdNOTE: Or is this the multi-faceted search that is currently planned? Do we really need this sentence?

<sup>6</sup>EdNOTE: Re-formulate this and link to an example

<sup>7</sup>EdNOTE: Quote needed

<sup>8</sup>EdNOTE: Reformulate this?

engine should transfer physical units. 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>9</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:9

## 2 Background - The existing MathWebSearch system

As mentioned above, MathWebSearch is a search engine for mathematical formulars in documents. It has a corpus of ??? <sup>10</sup> documents and is currently deployed and used by ??? and ??? <sup>11</sup>. EdN:10  
EdN:11

The frontend, running client-side in a web browser, is written in HTML5, CSS and JavaScript. It accesses a REST backend and depends on MathML support to render Mathematics. It simply accessing a REST backend via AJAX. When the client enters a LaTeX forumar to search for, the backend renders MathML which is then sent back to the client. Next, the client renders the MathML (to show the formular the user is searching for) and also sends back the MathML to the server to search for it. Finally the server sends back results to the client which then shows a list of them.

The backend, written in ?? <sup>12</sup>, as can be seen from the procedure during a common search, has 2 independent components. The LaTeX to MathML rendering is not part of MathWebSearch but rather uses LaTeXXML <sup>13</sup>. The searching is handled by MWS directly. EdN:12  
EdN:13

For each corpus a big search index is generated before search <sup>14</sup>. Then during runtime only this index is searched. <sup>15</sup> EdN:14  
EdN:15

<sup>16</sup> EdN:16

Advantages of this approach / how MWS is received <sup>17</sup> EdN:17

Disadvantages of the current approach

- has to be re-generated each time a document is added

---

<sup>9</sup>EdNOTE: Update this possible if we change the structure

<sup>10</sup>EdNOTE: Get an estimate here

<sup>11</sup>EdNOTE: where?

<sup>12</sup>EdNOTE: What is the backend written in?

<sup>13</sup>EdNOTE: Really? Reference needed.

<sup>14</sup>EdNOTE: Explain this better

<sup>15</sup>EdNOTE: Explain theory graphs here

<sup>16</sup>EdNOTE: Explain how MWS uses those here

<sup>17</sup>EdNOTE: Find points here

### 3 The proposed extension

- want a complete system
- searches a corpus of documents for units
- which is presentable to the end user
- should be extendable with respect to
  - the corpus. Plugging in a new corpus should be as easy as running a script somewhere.
  - the units. Adding new units should be simple by just adding a conversion to one already known unit.

The frontend

- a web page
- should work in modern browsers, preferably mobile-friendly
- should only be a frontend for a REST backend
- has an input for a unit
- has an input for a value
- maybe have facetted search on top

The backend

- REST based
- based on the existing system
- has to have a format of units
- has to receive text queries
- has to receive exact values or ranges or automatically generated ranges

The corpus

- 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

- 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>18</sup>
  - Perhaps include prefixes somehow?

EdN:18

19

EdN:19

## 4 Problems and related Work

There are several problems with this approach. <sup>20</sup>

EdN:20

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, LaTeX and MathML, to name only 3. Independent of the 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.

As described above, the unit equivalences should be in the form of a theory graph <sup>21</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 formulal. However since all numbers must

EdN:21

<sup>18</sup>EdNOTE: Figure out more details about this

<sup>19</sup>EdNOTE: Write this

<sup>20</sup>EdNOTE: Continue intro paragraph

<sup>21</sup>EdNOTE: Reference this properly.

be represented in a finite fashion, translations formulas 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>22</sup>. <sup>23</sup>

EdN:22

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>24</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.

EdN:23

EdN:24

25

EdN:25

---

<sup>22</sup>EdNOTE: Quote needed; reformulate

<sup>23</sup>EdNOTE: Write exactly how to use ranges

<sup>24</sup>EdNOTE: Who is doing unit finding?

<sup>25</sup>EdNOTE: Write final paragraph