

Semantic Search for Quantity Expressions

Bachelor Thesis DRAFT*

EdN:1

Tom Wiesing

Supervisor: Michael Kohlhase

Co-supervisor: Tobias Preusser

Jacobs University, Bremen, Germany

May 8, 2015

Abstract

In this proposal we describe how to introduce units to MathWebSearch. The aim of the project is build an extensible 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. ²³

EdN:2

EdN:3

*EdNOTE: Remove draft status

²EdNOTE: Physics search

³EdNOTE: Re-write abstract

Contents

1	Introduction	3
2	The structure of Mathematics: Theories, Views and Imports	4
2.1	Modeling Mathematics with the help of theories	4
2.2	Extending theories using imports	4
2.3	Views as mappings between theories	5
2.4	Building theory graphs	5
2.5	Using MMT to write down terms and theories	5
3	The structure of Quantity Expressions	7
3.1	Compositional behaviour of Quantity Expressions	7
3.2	Dimensions of Quantity Expressions	7
3.3	Mathematical Theory of Quantity Expressions	8
3.4	Transforming Quantity Expressions from one form into another	9
4	Making Quantity Expressions searchable	11
4.1	An idea for an extensible system	11
4.2	Normalisation of Quantity Expressions	12
4.3	Serialising Quantity Expressions to XML	13
4.4	Finding units inside documents	14
5	Putting it together: The implementation and its limits	15
5.1	Structure	15
5.2	Backend	15
5.3	Frontend	15
5.4	Caveats	15
6	Future work	16
6.1	Extension of the theory graph of units	16
6.2	Integration with MathWebSearch	16
6.3	More documents	16
7	Conclusion	17

1 Introduction

In this thesis 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 semantically 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 L^AT_EX documents¹ for mathematical formulae [HKP14]. As it is semantics-aware it not only searches for formulae in a simple-minded “text search” way but 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. In this way MWS abstracts from variable names.

We want to design a search engine that is similar to the existing MathWebSearch system in that it abstracts from the form of the quantity expression. This system should be very flexible with respect to the units it supports. We want a system in which it is easy to add new units. Eventually we also want to integrate our system into MathWebSearch so that we do not have to re-write existing search algorithms and can also use wildcards as described above.

This thesis is organised as follows: In section 2 we start by giving an introduction to mathematical theory modeling. We proceed in section 3 to talk about how quantity expressions can be formalised and in section 4 we apply these insights in order to start building in our search engine. In section 5 we present the implementation we designed and discuss its advantages and disadvantages. After discussing future work in section 6 we conclude in section 7.

¹Technically, MWS itself can only search XHTML documents. However with the help of L^AT_EXML [Mil] it also handles L^AT_EX documents.

2 The structure of Mathematics: Theories, Views and Imports

Before we start looking at Quantity Expressions and how to build a search engine for them, we want to give an introduction to meta-mathematical structure. We will use this knowledge later to build a better search engine. For this we first need to take a look at the concept of theories.

2.1 Modeling Mathematics with the help of theories

Theories, in this sense, are simply a set of symbols. Each of the symbols optionally can have a type and a definition. Within each theory, we can then use these symbols to write down terms (or expressions) within this theory. Types and definitions of these symbols are terms themselves². As a simple example of this, let us consider the theory of semigroups as seen in 1

Semigroup	
G	: type
\circ	: $G \rightarrow G \rightarrow G$
assoc	: $\text{ded } (\forall x \in G. \forall y \in G. \forall z \in G. (x \circ y) \circ z = x \circ (y \circ z))$

Figure 1: The theory of semigroups.

In this theory, we define 3 symbols: G , \circ and `assoc`. In the first line we define a type G . Next we define a function \circ that takes 2 arguments of type G and returns another term of type G . In the last line, we make the statement that associativity holds⁴.

EdN:4

2.2 Extending theories using imports

Sometimes we want to extend theories without having to define everything again. For example, we want to say that a Monoid is a semi-group along with an identity element. In the semi-group example above, we have also used terms from other theories to define G as a type.

We can model this concept by using imports. An import from one theory into another makes symbols of the imported theory available in the target theory. In figure 2 we can easily define a monoid.

Monoid	
import Semigroup	
e	: G
id	: $\text{ded } (\forall x : G. x \circ e = e \circ x = x)$

Figure 2: The theory of monoids which imports the theory of semigroups as defined in figure 1.

²They are not terms over the same theory however.

⁴EdNOTE: possibly explain / mention Curry–Howard isomorphism

2.3 Views as mappings between theories

However imports are not the only way theories can be related. If we have 2 theories, we sometimes want to have a map between them. In addition to the theory of monoids above we define the theory of non-negative integers in figure 3.

Non-negative integers	
\mathbb{Z}_0^+	: type
0	: \mathbb{Z}_0^+
+	: $\mathbb{Z}_0^+ \rightarrow \mathbb{Z}_0^+ \rightarrow \mathbb{Z}_0^+$
assoc	: $\text{ded } (\forall x : \mathbb{Z}_0^+ . \forall y : \mathbb{Z}_0^+ . \forall z : \mathbb{Z}_0^+ . (x \circ y) \circ z = x \circ (y \circ z))$
id	: $\text{ded } (\forall x \in \mathbb{Z}_0^+ . x + 0 = 0 + x = x)$

Figure 3: The theory of non-negative integers.

A map from the theory of monoids to the theory of positive integers should map all symbols from the theory of monoids to symbols from the theory of positive integers. Furthermore, such a map should be truth preserving, i. e. if I write down a true statement as a term over the theory of monoids and translate this term, it should still be true in the theory of positive integers. Such a mapping is called a *View* from the theory of monoids to the theory of Positive integers. Such a view ϕ could look as follows:

$$\phi = \left\{ \begin{array}{l} G \mapsto \mathbb{Z}_0^+ \\ e \mapsto 0 \\ \circ \mapsto + \\ \text{assoc} \mapsto \text{assoc} \\ \text{id} \mapsto \text{id} \end{array} \right\}$$

If we take a closer look at this view, we notice that we also have to map the imported symbols. This is needed so that we can translate any term or statement in theory of monoids to a term or statment into the theory of non-negative integers.

2.4 Building theory graphs

⁵ We have seen in the examples above that we can model mathematics with the help of theories, views and imports. To make this structure even more obvious, we can represent it in a graph, a so-called theory graph. We consider the theories as vertices of such a graph and the views and imports as edges. An example can be found in figure 4. EdN:5

2.5 Using MMT to write down terms and theories

MMT is a **M**odule system for **M**athematical **T**heories [RK13]. With the help of MMT we can represent theories, views and imports in *.mmt* files. It is easy to write these files and anyone without programming knowledge can easily extend existing ones. The objects defined in these files can then be used via an API to write down terms and transform these using definitions and views.

⁵EdNOTE: Remove section heading or extend this section.

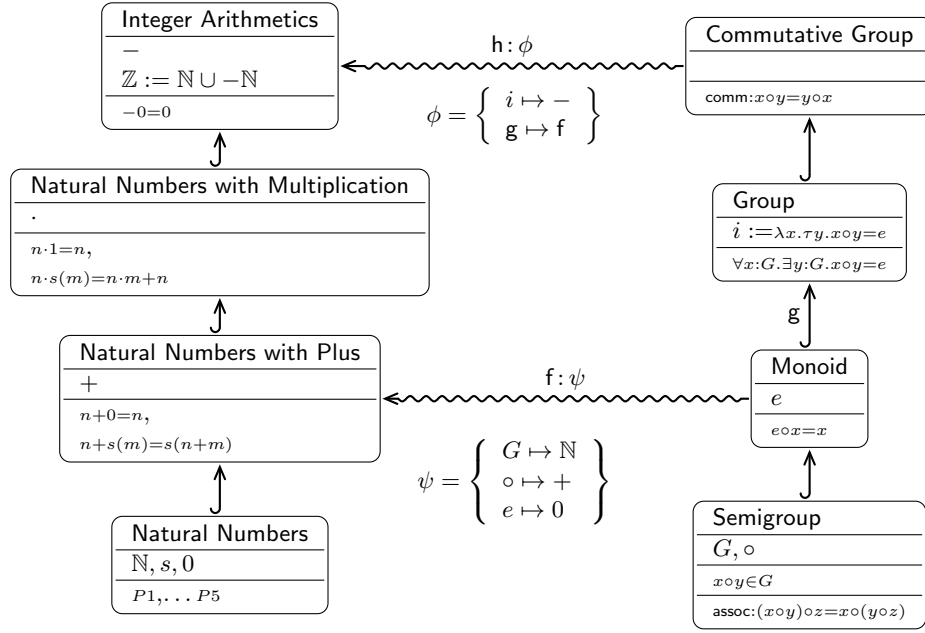


Figure 4: A simple theory graph. Imports are represented as solid edges and views as squigged edges.

This furthermore allows us to easily model an extensible system for units for use within a search engine. We will come back to this later in section 4.1.

3 The structure of Quantity Expressions

The first step in developing a good search engine for Quantity Expressions is to take a closer look at Quantity Expression.

3.1 Compositional behaviour of Quantity Expressions

For this purpose let us take a look at:

$$x = 25 \frac{\text{m}}{\text{s}}$$

We notice that x consists of 2 parts, a scalar (25) and a scalar-free $\frac{\text{m}}{\text{s}}$. Furthermore, the unit consists of two primitive units m and s . Since they are divided by one another, we can conclude that x describes a velocity.

While m and s are certainly primitive units (they can not be decomposed further), it is not easy to define a unit as a composition of simple units. Consider the following example:

$$y = \frac{\text{L}}{100 \text{ km}}$$

y is certainly a unit and also a quantity expression. It consists of 2 sub-expressions, L and 100 km. The first one is a primitive unit and the second one a multiplication of a number and the unit km. It is thus reasonable to define the following 6 types of quantity expressions:

1. A primitive unit, such as m (meter). This is the most obvious one.
2. the Multiplication of a quantity expression with a scalar.
3. the Division of a quantity expression by a scalar.
4. The multiplication \cdot which takes 2 existing quantity expression and generates a new one, for example $\cdot (100, \text{m}) = 100 \text{ m}$
5. The division \backslash which again takes 2 quantity expressions and generates a new one, for example $\backslash (\text{m}, \text{s}) = \frac{\text{m}}{\text{s}}$
6. The addition of two quantity expressions

This allows us to easily generate the quantity expressions x and y from primitive units m, s, L and km.

3.2 Dimensions of Quantity Expressions

Now let us briefly examine the dimensions of Quantity Expressions. The dimension of a quantity expression is the type of quantity it expresses. For example 5 m describes some length. According to the International System of Units⁶ there are seven basic dimensions: EdN:6

1. length
2. mass

⁶EdNOTE: Quote something properly here

3. time
4. electric current
5. temperature
6. luminous intensity
7. amount of substance.

In addition to these dimensions there are 2 more special dimensions: The *count* dimension (used for counting of objects) and the *none* dimension for dimensionless quantities. ⁷

EdN:7

Similar to the compositional behaviour of quantity expressions, dimensions can be multiplied and divided. Unlike quantity expressions however they can not be multiplied with numbers. Furthermore, when multiplying to quantity expressions of dimensions a and b their resulting dimension is $a \cdot b$, the multiplication of the dimensions. The same goes for division. In this regard the dimension of a quantity expression behaves like a type.

3.3 Mathematical Theory of Quantity Expressions

The realisation that dimensions are types of quantity expressions leads us to our first formalisation of quantity expressions. We first define a theory of dimensions in figure 5 and then import it to define a theory of Quantity Expressions in figure 6.

Dimension	
dimension	: type
none	: dimension
count	: dimension
length	: dimension
mass	: dimension
time	: dimension
current	: dimension
temperature	: dimension
luminous	: dimension
amount	: dimension
.	: dimension \rightarrow dimension \rightarrow dimension
\	: dimension \rightarrow dimension \rightarrow dimension

Figure 5: A formalisation of the theory of dimensions.

Figure 5 defines the 9 basic dimensions and then dimension composition via multiplication and division. Then we move on in 6 to define quantity expressions. Each quantity expression has a dimension (via the *QE* constant). This allows us to define basic units (which we will actually do in the next figure.). With the *QENMul* and *QENDiv* symbols we can multiply and divide quantity expressions by numbers (for this case we actually need to import some theory of numbers to allow us actually write this down as a Term). Then

⁷EdNOTE: Explain these more

Quantity Expression	
import Dimension	
import Number	
QE	: dimension \rightarrow type
QENMul	: $\forall x : \text{dimension}. \text{number} \rightarrow \text{QE}(x) \rightarrow \text{QE}(x)$
QENDiv	: $\forall x : \text{dimension}. \text{QE}(x) \rightarrow \text{number} \rightarrow \text{QE}(x)$
QEMul	: $\forall x : \text{dimension}. \forall y : \text{dimension}. \text{QE}(x) \rightarrow \text{QE}(y) \rightarrow \text{QE}(\cdot(x, y))$
QEAdd	: $\forall x : \text{dimension}. \text{QE}(x) \rightarrow \text{QE}(x) \rightarrow \text{QE}(x)$
QEDiv	: $\forall x : \text{dimension}. \forall y : \text{dimension}. \text{QE}(x) \rightarrow \text{QE}(y) \rightarrow \text{QE}(\backslash(x, y))$

Figure 6: The chosen formalisation of the theory of quantity expressions.

we define multiplication and division of quantity expressions in such a way that dimensions multiply and divide appropriately.

Now we need to introduce some basic units. Let us start by just defining meter in figure 7. We can now write a term in this theory that expresses any number of meters.

Meter	
import QuantityExpression	
Meter	: QE(length)

Figure 7: A theory defining the primitive unit Meter.

3.4 Transforming Quantity Expressions from one form into another

This is a very nice start of a unit system⁸. Let us now define a few more units of length. In figure 8 we show a few non-si units. Here we first define thou as a quantity expression of length and then one-by-one define more units in terms of the previous one. EdN:8

Non SI Lengths	
import QuantityExpression	
Thou	: QE(length)
Foot	: QE(length) = QENMul(1000, Thou)
Yard	: QE(length) = QENMul(3, Foot)
Chain	: QE(length) = QENMul(22, Yard)
Furlong	: QE(length) = QENMul(10, Chain)
Mile	: QE(length) = QENMul(8, Furlong)

Figure 8: A theory of some non-SI units of length.

We now want to relate quantity expressions with units from the *Meter* theory to units from the *Non SI Lengths* theory. It is known that $1\text{Thou} = 0.0000254\text{Meter}$. This can be easily expressed with a view ψ between these two theories:

⁸EdNOTE: Better formulation

$$\psi = \{ \text{Thou} \mapsto \text{QENMul}(0.0000254, \text{Meter}) \}$$

Even though the view just maps the symbol *Thou* to some Term in the Meter theory, we can also use it to transform any other term from the Thou Theory. Since all units are defined in terms of the previous one, we can just expand all definitions to get an expression containing only numbers and the unit *Thou*. Then we can use the view as normal to get a Quantity Expression in the *Meter* theory.

4 Making Quantity Expressions searchable

Now that we know what quantity expressions are and how we can convert between equivalent representations, it is time to start about concrete algorithms for our search engine.

When querying this search engine with a quantity expression, we want to be able find occurrences of equivalent quantity expressions. For this we need several components:

1. a component that allow the user to enter quantity expressions,
2. a software to find quantity expressions within documents
3. a way to send quantity expressions from the user to the search engine
4. an algorithm to convert between different forms of a quantity expressions (or at least determine if 2 quantity expressions are equivalent) and finally
5. a unit system that is flexible enough to handle different kinds of units.

4.1 An idea for an extensible system

Let us start with point (5), a unit system that is flexible enough to handle all different of units. In figure 9 you can find a small part of the unit graph we developed. At first, in the gray area, we define a basic version of our system as defscribe in section 3.3. Next we continue in the top left corner to define all the basic SI units. For the following we only look at the theories descrobing length and area units (however there are also some for all of the other dimensions).

We start with defining the non-SI unit *Thou* in the *Imperial Lengths A1*. In order to be able to convert between these units we apply what we learned in section 3.4 and create a view that maps this unit back to standard SI units. We then expand into *Imperial Lengths A2* and define a few more units incrementally based on *Thou* (we use an import to be able to do this).

Next, we define a new unit *Link* in the *Imperial Lengths B1* theory and define the unit *Rod* in *Imperial Lengths B2*. We then also make a view from *Imperial Lengths B1* to *Imperial Lengths A2*. Since *Imperial Lengths B2* imports *Imperial Lengths B1* this allows us to easily transform both of the new units back to SI if needed.

Next we want to define some units with the dimension Area. For this we first collect all known lengths and make them available in a *Imperial Lengths* theory via imports. Then we import this theory into a new *Imperial Area* theory where we can then define *Perch*, *Rod* and *Acre*.

But why do we define so many different theories? Why do we not simply define all the units of length in a single theory? Obviously, this one length theory would be very big. Furthermore, it would require anyone who wants to add more units to the system to add to this one theory, which can cause conflicts if there are 2 units with the same name (such as *Mile* and (*nautical*) *Mile*).

Defining units in the way we do it allows us to incrementally add new units to the system. If we have a view back to already known theories, it is easy to compare them back to the existing units.

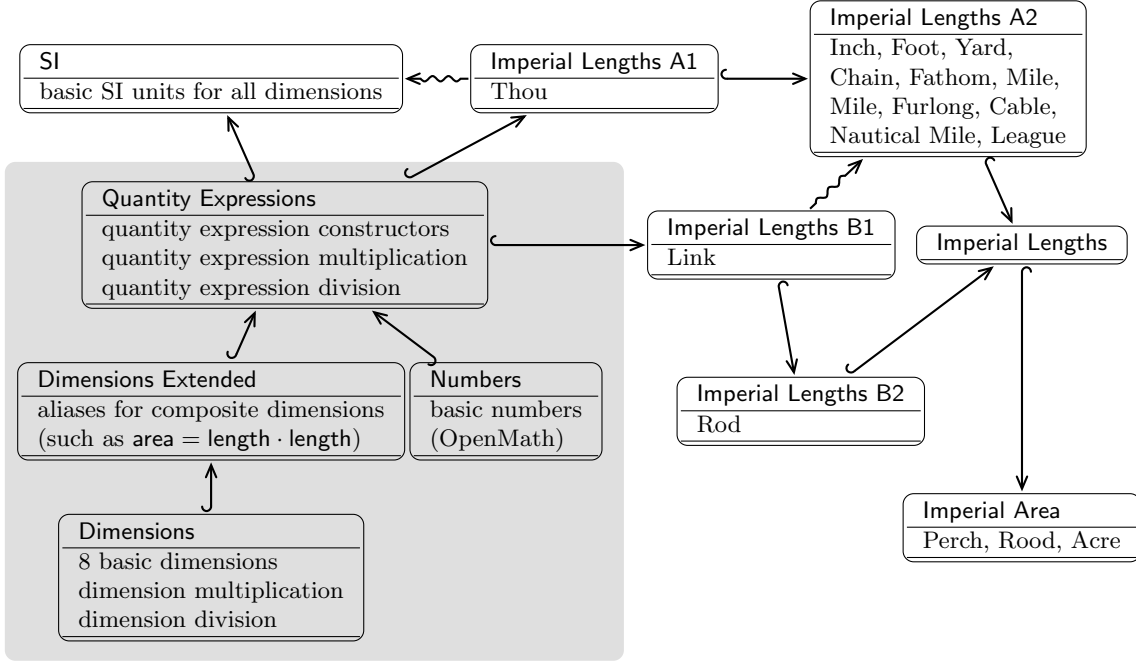


Figure 9: A small part of the graph containing unit theories. Imports are represented as solid edges and views as squigged edges. The gray area contains basic theories that are used to define quantity expressions.

4.2 Normalisation of Quantity Expressions

To build a search engine however, it is inefficient to try to directly compare a search query with known quantity expressions. One way to compare quantity expressions is to normalise them to a suitable normal form. The normal form we came up with works in 2 steps.

We first turn a Term representing a Quantity Expression into a Term over the theory containing only standard units. In general we can choose any theory we want for this. We choose the *SI* theory for this so that the normal forms can be understood easily.

In order to transform a term to standard units, we just expand each unit (represented by constants) into a Term over the standard theory. This is achieved by translating it along a path in the theory graph. A path in the theory graph can consist out of 2 types of edges, Views and Imports. For views we can obviously use the mappings to translate the unit from the source to the target theory of the respective view. For imports however, it is not that simply. We can translate terms along them via definition expansion. This takes us to the imported theory from the theory that imports. We can use these edges as described and use standard graph theory algorithms to find a path for a unit to be recursively translated upon.

In the second step, we separate the scalar from the units. For this, we first compute the scalar value of a quantity expression and then cancel out units which occur both in numerator and denominator of the quantity expression. This step is much easier than the first one, as we no longer have to use the theory graph. For this reason, we compute the normalisation of each unit only once by caching the result. Then we no longer need to look at the graph for our normalisation procedure at all.

Let us illustrate the process of normalisation a bit more. We will try to normalise the quantity expression

$$q = 42 \frac{\text{Furlong}}{\text{Fortnight}}$$

First, we need to normalise *Furlong* (a quantity expression of length) to *Meter*. For this we first do definitional expansion:

$$\text{Furlong} = 10 \text{ Chain} = 10 (22 \text{ Yard}) = \dots = 10 (22 (3 (12 (1000 \text{ Thou}))))$$

Next, we need to use a view to turn this into meters:

$$10 (22 (3 (1000 \text{ Thou}))) = 10 (22 (3 (12 (1000 (0.0000254 \text{ Meter}))))$$

Similary, we use definitional expansion to find that:

$$\text{Fortnight} = (2 (7 (24 (60 (60 \text{ Second}))))))$$

Substituting into q gives us the normal form of the quantity expression after the first step:

$$42 \frac{10 (22 (3 (1000 (0.0000254 \text{ Meter}))))}{2 (7 (24 (60 (60 \text{ Second}))))}$$

Next, we want to extract the scalar component and compute it:

$$42 \frac{10 (22 (3 (12 (1000 (0.0000254))))))}{2 (7 (24 (60 (60))))} = 0.006985$$

We also extract the unit component:

$$\frac{\text{Meter}}{\text{Second}}$$

Finally we compose these again to get the normal form:

$$0.006985 \frac{\text{Meter}}{\text{Second}}$$

4.3 Serialising Quantity Expressions to XML

Now that we have a normal form of units, we want to be able to exchange quantity expressions between the user, the search engine and the harvesting component. The W3C has written a note on how to handle Units in MathML [Gro]. MathML is an XML serialisation of Mathematical expression. We orient ourselves on the format. For this serialisation we start with a term over some unit theory. The serialisation of the term

$$0.006985 \frac{\text{Meter}}{\text{Second}}$$

can be found in figure 10. Let us take a closer look to understand this serialisation. We translate each component of the Term individually and then assemble them together.

We translate the unit (constant) *Meter* to `<symbol cd='SIBase'>Meter</symbol>`. This XML Tag expresses that we are talking about the symbol *Meter*. The cd attribute stands for content directory. In our case this is the name of the theory where the symbol was declared.

```

<apply>
  <times />
  <cn type='float'>0.006985</cn>
  <apply>
    <divide />
    <csymbol cd='SIBase'>Meter</csymbol>
    <csymbol cd='SIBase'>Second</csymbol>
  </apply>
</apply>

```

Figure 10: Serialisation of a simple Quantity Expression.

We translate multiplication and division by using `<apply><times />...</apply>` and `<apply><divide />...</apply>` respectively. We can use quantity expressions directly here. If we want to multiply with or divide by numbers, we use `<cn type='float'>...</cn>` instead of the `<csymbol>` tags.

This format allows us to easily exchange quantity expressions. It is also easy to send quantity expressions from the user to the backend (the core component of the search engine.)

4.4 Finding units inside documents

Apart from the frontend, which we will talk about in section 5.3, there is only one major component left: Finding and matching quantity expressions inside documents. As we have outlined in the previous sections, we can easily compare them with the help of normalisation. If we have a list of quantity expressions and where they occur, we can use the following algorithm to easily find them:

1. First find quantity expressions inside documents and store them in a list along with their origin.
2. Normalise each of them as described before.
3. Store the normalised forms in an efficient index along with their original versions and origins.
4. When a query is sent, normalise the query.
5. Next look in the index for the normalised query.
6. Then deliver the original versions and their origins from the index.

This algorithm allows us to easily and efficiently find quantity expression in documents after we have built up an index. In order to build up this index we need a list of quantity expressions and occurrences inside documents. A software that fullfills this task is called a Spotter. In our case the spotter is developed by Stiv Sherko in a seperate effort [She14].

5 Putting it together: The implementation and its limits

We can now join all these components together to build a search engine. To recap, we need these components:

- MMT, so that we can use Terms and theories. Specifically, the mmt-api.
- Theories and Views in MMT notation, contained in mmt files.

5.1 Structure

5.2 Backend

5.3 Frontend

* Frontend: HTML / JS * Harveste

5.4 Caveats

* type equalities * requirements of the graph * (in)stabilities of MMT * davenport units (relative units)

6 Future work

9

EdN:9

6.1 Extension of the theory graph of units

6.2 Integration with MathWebSearch

6.3 More documents

⁹EdNOTE: TODO: Write me

7 Conclusion

¹⁰

EdN:10

* what did we want to do? * what did we achieve? * were we successfull * end with: there is still a lot to be done

¹⁰EdNOTE: TODO: Write me

References

- [Gro] W3C Working Group. Units in MathML. <http://www.w3.org/TR/2003/NOTE-mathml-units-20031110/>. Note 10 November 2003.
- [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.
- [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.