



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

School of Computer and Commu-
nication Sciences

Computer Science Section

Distributed Information Systems Lab-
oratory (LSIR)

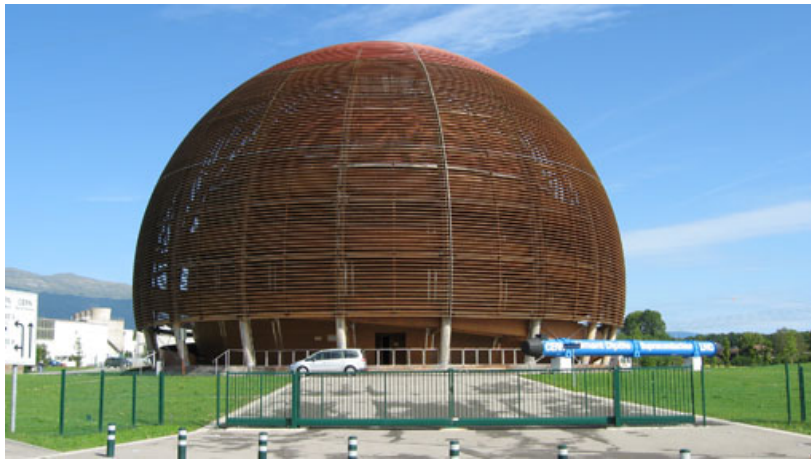


Information and Technology
Department

Digital Library Service

$5e^{x+y}$: A Math Aware Search Engine (for CDS)

MASTER THESIS PROJECT



Author:

Arthur OVIEDO

Supervisors:

CERN: Nikolaos KASIOUMIS

EPFL: Karl ABERER

March 13, 2014



Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 7 |
| 1.1 | Context | 7 |
| 1.1.1 | CERN | 7 |
| 1.1.2 | CDS & Invenio | 7 |
| 1.2 | Project overview | 8 |
| 1.2.1 | Motivation | 8 |
| 1.2.2 | Goals | 10 |
| 2 | Theoretical Background | 11 |
| 2.1 | Boolean Model | 11 |
| 2.2 | Vector Space Model | 12 |
| 2.3 | Probabilistic Relevance Model | 13 |
| 2.4 | Okapi BM25 | 13 |
| 2.5 | Tree Edit Distance | 14 |
| 3 | Technical Considerations | 17 |
| 3.1 | Digital Representations of Mathematics | 17 |
| 3.2 | Math Extraction Tools | 19 |
| 3.3 | Computer Algebra Systems | 20 |
| 4 | Related Work | 23 |
| 4.1 | Status of Mathematical Information Retrieval | 23 |
| 4.2 | Mathematical based search projects | 24 |
| 4.2.1 | MIaS | 24 |
| 4.2.2 | EgoMath | 24 |
| 4.2.3 | DLMFSearch | 25 |
| 4.2.4 | MathWebSearch | 25 |
| 4.2.5 | LaTeXSearch | 26 |
| 4.2.6 | WikiMirs | 26 |

| | | |
|----------|--|-----------|
| 5 | $5e^{x+y}$(CDS) | 27 |
| 5.1 | Features Extraction | 27 |
| 5.1.1 | Notational Features | 27 |
| 5.1.2 | Structural Features | 32 |
| 5.2 | Development with Lucene/Solr | 33 |
| 5.3 | Integration with Mathematica | 35 |
| 5.3.1 | Pattern Extraction | 35 |
| 5.3.2 | Cleaning of variables | 35 |
| 5.3.3 | Expression simplification | 35 |
| 5.3.4 | Recursive interpretation | 38 |
| 5.3.5 | Error handling | 38 |
| 5.3.6 | Timeout Handling | 38 |
| 5.4 | Integration with CDS | 39 |
| 6 | Clustering | 41 |
| 6.1 | Technical details | 41 |
| 6.2 | Estimating number of clusters | 42 |
| 7 | Evaluation | 45 |
| 7.1 | Dataset | 45 |
| 7.2 | Setup | 45 |
| 7.3 | Quality results | 46 |
| 7.4 | Indexing Performance | 47 |
| 7.5 | Querying Performance | 48 |
| 7.6 | Alternative Metrics | 49 |
| 8 | Conclusions | 51 |
| 9 | Future work | 53 |
| 9.1 | Features extraction | 53 |
| 9.2 | Data Analysis | 53 |
| 9.3 | End to end system | 54 |
| | Bibliography | 55 |
| A | Configuration files | 59 |
| A.1 | Solr schema | 59 |
| A.2 | Solr configuration file | 60 |
| B | Test equations | 61 |
| C | Detailed system architecture diagrams | 65 |

List of Figures

| | | |
|-----|---|----|
| 1.1 | CDS search result interface | 8 |
| 1.2 | Invenio global architecture | 9 |
| 5.1 | List of different types of integral related characters in Unicode | 31 |
| 5.2 | List of different types of greater/less than related characters in Unicode | 31 |
| 5.3 | High level deployment architecture | 39 |
| 5.4 | Web interface of $5e^{x+y}$: Query input screen | 40 |
| 5.5 | Web interface of $5e^{x+y}$: Results screen | 40 |
| C.1 | Overview of <code>cern.ch.mathexplorer.lucene.analysis.analyzers</code> package | 65 |
| C.2 | Overview of <code>cern.ch.mathexplorer.lucene.analysis.tokenizers</code> package | 65 |
| C.3 | Overview of <code>cern.ch.mathexplorer.lucene.analysis.filters</code> package | 66 |
| C.4 | Overview of <code>cern.ch.mathexplorer.lucene.query</code> package | 66 |
| C.5 | Overview of the relations between the main packages | 67 |

List of Tables

| | | |
|-----|---|----|
| 3.1 | Comparison of the different languages for expressing mathematical content | 18 |
| 3.2 | Comparison of available Mathematical extraction tools | 20 |
| 3.3 | Comparison of main CAS | 21 |
| 5.1 | Characters with similar glyphs and/or semantics | 30 |
| 5.2 | Comparison of different simplification modes | 37 |
| 6.1 | Cluster measurements for different values of k | 42 |
| 6.2 | Top 20 terms for two different identified clusters | 44 |
| 7.1 | Software used during evaluations | 45 |
| 7.2 | Loose Precision | 46 |
| 7.3 | Strict Precision | 47 |
| 7.4 | Discounted Cumulative Gain | 47 |
| 7.5 | Indexing performance | 48 |
| 7.6 | Query performance in seconds | 48 |
| 7.7 | Distance computation using two different metrics | 49 |
| B.1 | Sample equations used during testing | 63 |

Chapter 1

Introduction

1.1 Context

1.1.1 CERN

CERN (Conseil Européen pour la Recherche Nucléaire or European Organization for Nuclear Research)[1], founded in 1952, is the biggest research center in the area of particle physics. It is located in the French-Swiss border, near Geneva. Even though CERN directly employs around 2400 people, more than 10000 scientist from around 113 countries have visited CERN to deepen their research. As an example of its contributions to science, recently, in July of 2012, CERN announced that two different experiments (ATLAS and CMS) confirmed the existence of the particle named Higgs Boson, which lead to the award of the Nobel Prize to Peter Higgs and François Englert.

Even though the main focus of CERN relies on the study of particle physics, the strong research environment has lead to important advances in several different areas, and one of the main inventions that CERN has contributed to the world is the Web. Tim Berners-Lee, in 1989, proposed a solution[2] to the increasing problem of keeping track of all the information related to the different experiments held at CERN, through a hypertext system.

1.1.2 CDS & Invenio

CDS[3] (CERN Document Server) is the institutional digital library system developed and used at CERN[4]. To this date, it contains more than one million records and more than 400000 full-text documents. It provides tools to manage the complete workflow of a document taking care of the submission, annotation, editing, storage, searching, retrieval and displaying, among

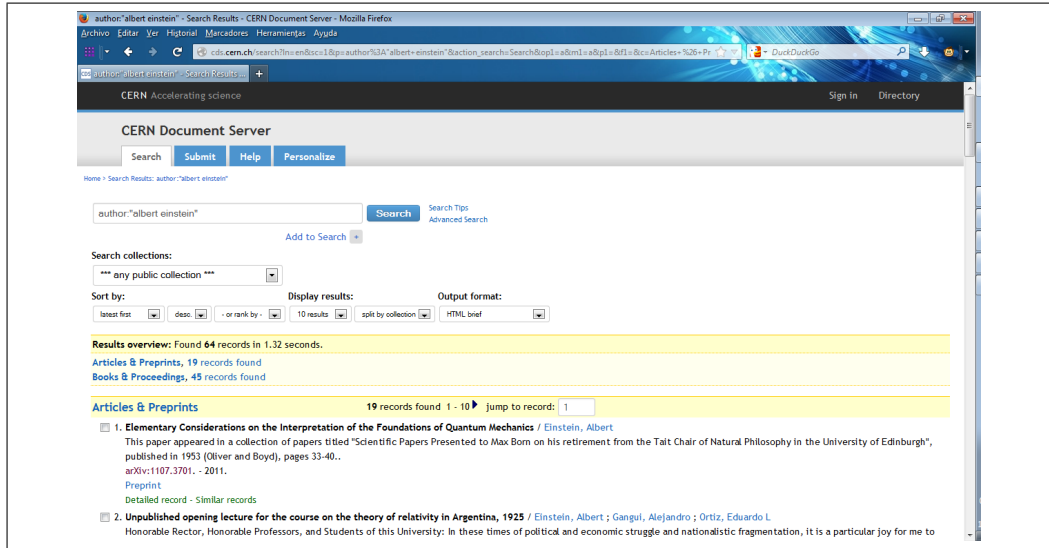


Figure 1.1: CDS search result interface

other phases. Figure 1.1.2 presents the search results for a simple query in CDS.

Invenio[5] is the open-source digital library software platform running behind CDS and it is a project developed in parallel to CDS at CERN. Besides CDS, Invenio supports around thirty scientific institutions worldwide including INSPIRE (a collaboration between Fermilab, CERN, DESY and SLAC) EPFL and ILO. Invenio is composed of several modules, each one of which can be mapped to a specific step in the workflow of a record. Figure 1.2 presents the global architecture of Invenio. A detailed explanation of each module can be accessed in [6]

1.2 Project overview

1.2.1 Motivation

The initial statement for this project was more generic and was encompassed as "New ways to programmatically, accurately and efficiently extract data and metadata from digital files". During a first exploration around this topic, we focused our attention to the mathematical content that is stored in these files. After more discussions we identified that the extraction, storage, indexing and finally searching for mathematical content would be a very useful project for CDS and an interesting research direction. CDS groups records in different categories or Collections such as Published Articles, Preprint, Photos, Books, General Talks among others. Some of these collections contain

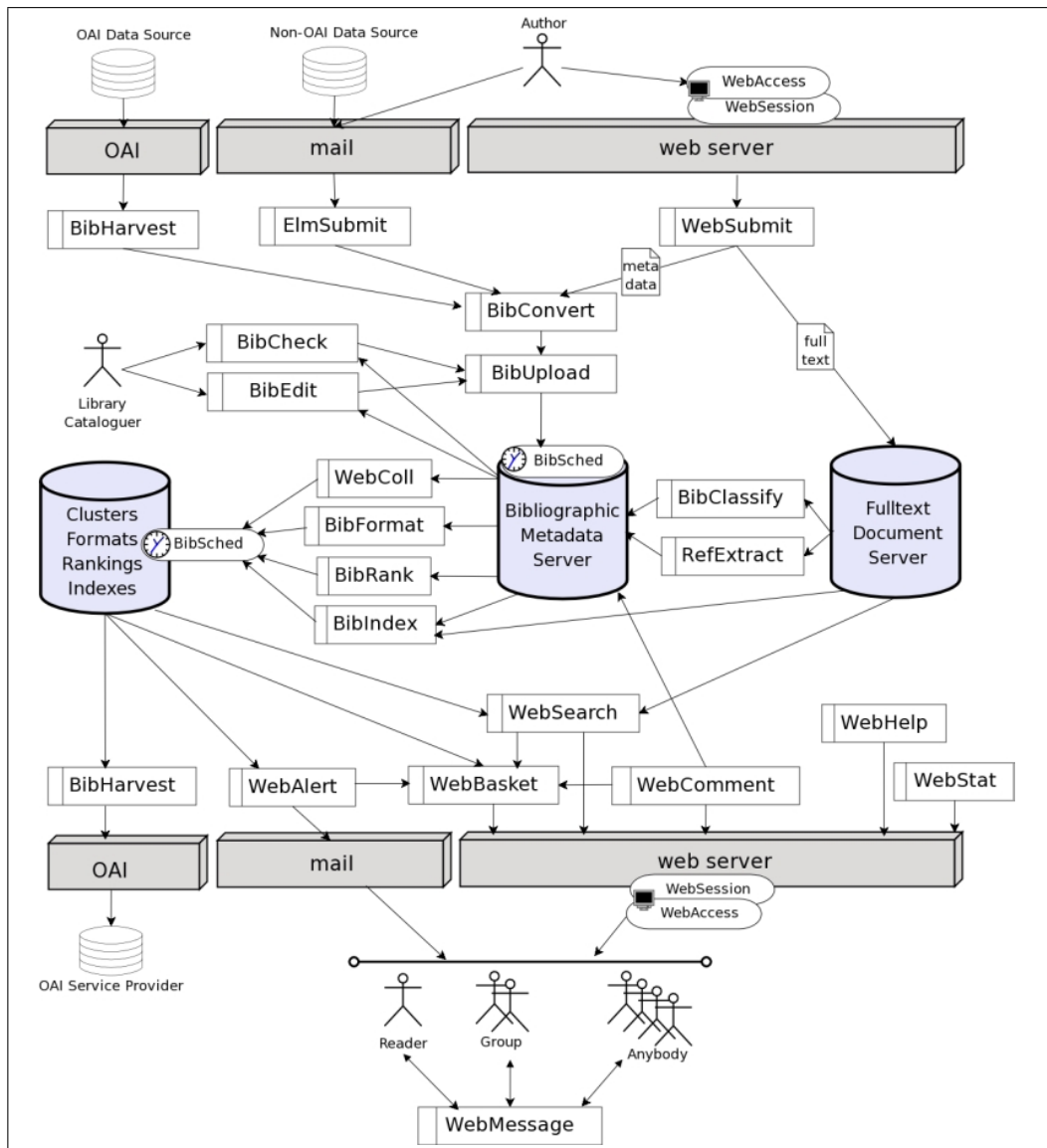


Figure 1.2: Invenio global architecture

scientific documents in all of the different research areas where CERN is involved, and therefore a big amount of mathematical expressions are contained there. Only the Preprints collection contains 698581 records to date, harvesting documents from services like ArXiv[7] where most of the documents are in the areas of physics, mathematics and statistics which are very rich in mathematical content. Currently Invenio provides different ways of searching for records by specifying metadata fields like author:"Albert Einstein" or by keywords. However, there is no way to search for a given mathematical expression. The current workarounds would be to try to search for the name of the given expression if it is common enough to have been named like *Schrödinger equation*. This approach most of the times is not enough since most of the equations are not named and even named ones can be rewritten in several ways and each one may have a different importance to the user. This combination of factors, motivates the development of a complete system that allows users to search for relevant documents, based on mathematical expressions.

1.2.2 Goals

The specific goals of this project can be summarized as:

- Explore ways to automatically extract the mathematical content from a document collection
- Investigate different approaches and formats to store mathematical expressions
- Identify relevant features in a mathematical expression
- Explore ways to efficiently store and retrieve the set of features
- Implement the previous steps into a complete system (a search engine) and integrate it into the Invenio software
- Evaluate different approaches and the quality of the provided results
- Identify deficiencies in our work and propose solutions and further directions

Chapter 2

Theoretical Background

Information Retrieval (IR) consists of the activity of identifying relevant documents from a collection based on some input parameters. Different models have been developed through the evolution of the field. In the subsequent sections, the most relevant models to this work are presented.

2.1 Boolean Model

The Boolean model of Information Retrieval (BIR) was the first one to be developed and is currently one of the most used in the current implementation of IR systems because of its simplicity. It is based on the mathematical concepts of Boolean Logic and Set Theory. A formalization of this framework can be stated as follows:

Given a set of terms $T = t_1, t_2, \dots, t_n$ which are the terms that will be indexed (In a standard search engine this can be the set of words in a language after some normalization steps), a set of documents $D = d_1, d_2, \dots, d_n$ where $d_i \in \mathcal{P}(T)$ that the user will search for. A query element q is a Boolean expression over the index terms and the operators *AND*, *OR* and *NOT*. We define the predicate $Relevant(q, d_i)$ if a document d_i satisfies the query expression q . Finally the result for query over a collection of documents can now be written as $Result(q, D) = \{d | Relevant(q, d)\}$. The benefits of this model are its simplicity, both in its formalism and in its implementation. The main disadvantage of BIR are the fact that all terms are given the same importance which makes it difficult to rank results. Since all results have same importance, it is difficult also to limit the number of retrieved documents and therefore the output set can be difficult to process for the final user.

One of the key data structures for an efficient implementation of this

model is the inverted index where for a given term, the index stores which documents contain it. This structure is used to efficiently solve queries by only scanning document that contain at least one of the terms in the query. Scanning a document can be done by using a standard (Also called forward) index, where for a given document, the index stores the terms associated with it. Both indexes can be implemented using hash tables.

2.2 Vector Space Model

The vector space model[8] which was developed between the 1960s and 1970s alongside the SMART Information Retrieval System. In this model, documents (from a collection D) and queries are represented as vectors d and q in a vector space where each index term corresponds to a dimension. If a term t_i occurs in a given query or document, its i -th coordinate will have a non zero value. Index terms can be keywords or phrases in the context of textual search or any other type of feature that can be extracted from the elements in a given specific context. The scoring of a document for a given query is computed as the cosine of the angle between both vectors which can be expressed as:

$$score(d, q) = \cos\theta(d, q) = \frac{d \cdot q}{\|d\| \|q\|} = \frac{\sum_{i=1}^n weight(d_i, d) * weight(q_i, q)}{\sqrt{\sum_{i=1}^n weight^2(d_i, d)} \sqrt{\sum_{i=1}^n weight^2(q_i, q)}}$$

Different weighting schemes can be employed, and the most standard one is called tf-idf (Term frequency - Inverse document frequency) where terms that appear less frequently in the corpus, are assigned higher values (Intuitively the term "the" will be less useful for identifying a relevant document than the term "Shakespeare" in the context of an English word based search engine). Therefore:

$$weight(t, d) = tf(t, d)idf(t)$$

There are several variations for each of the two components. Some approaches for tf include the raw frequency; Boolean frequency: $tf(t, d) = 1$ if t occurs in d or 0 if not; logarithmically scaled: $tf(t, d) = \log(frequency(t, d) + 1)$ and the augmented frequency: $tf(t, d) = 0.5 + \frac{0.5 * frequency(t, d)}{\max\{frequency(t', d) | t' \in d\}}$. For computing the inverse document frequency the standard approach follows $idf = \log \frac{|D|}{1 + \{d \in D : t \in d\}}$

The main drawbacks of this approach reside in the assumption that the terms are independent.

2.3 Probabilistic Relevance Model

The probabilistic relevance model was developed in 1976 by Robertson and Jones[9]. The similarity of a document d to a query q is given by the probability of d being relevant to q . It assumes there is a set R that is preferred to be the answer set for q . The model assumes that the documents in R are relevant to q and the documents in \bar{R} are not relevant. R then is the set that maximizes:

$$sim(d, q) = \frac{P(R|d)}{P(\bar{R}|d)}$$

The importance of this model is that it serves as a solid base for state of the art information retrieval models.

2.4 Okapi BM25

The Okapi BM25 ranking function is one of the current state-of-the-art models in information retrieval[10]. The actual scoring function is called BM25 and Okapi refers to the first system implementing this function in the 1980s. Given a query q as a vector of keywords q_1, q_2, \dots, q_n then the score for a document d is:

$$score(d, q) = \sum_{i=1}^n IDF(q_i) * \frac{freq(q_i, D)}{freq(q_i, D) + k_1 * ((1 - b) + b * \frac{|D|}{avgdl})}$$

and

$$IDF(q_i) = \log \frac{N - df(q_i) + 0.5}{df(q_i) + 0.5}$$

where $freq(q_i, D)$ is the number of occurrences of keyword q_i in D , $avgdl$ is the average length (In keywords) of a document, and k_1 and b are free parameters. Usually $k_1 = 2$ and $b \in [0, 1]$ controls how important is the length normalization, being set normally to 0.75. N is the size of the collection, $df(q_i)$ is the number of documents that contain term q_i

A further variation that takes into account structure in the form of a document containing different fields and each of them having its own length and importance (boost) can be formulated as follows, known as BM25F can be expressed as follows:

$$score(d, q) = \sum_{i=1}^n \frac{weight(q_i, d)}{k + weight(q_i, d)} * IDF(q_i)$$

and

$$weight(q_i, d) = \sum_{j=1}^m \frac{freq(q_i, D, field_j) * boost_j}{((1 - b_j) + b_j * \frac{|l_j|}{avgfl_j})}$$

where m is the number of different fields, $boost_j$ is the relative importance for each field, b_j is analogous to the b constant in the BM25 group of equations, $|l_j|$ is the length of the $field_j$ and $avgfl_j$ is the average length of the $field_j$.

2.5 Tree Edit Distance

As it will be discussed, one of the most common formats to represent mathematical equations, MathML, is a XML based format, and whose structure represents a rooted tree. Because of this, it makes sense to consider some models to compare and analyse mathematical expressions taking into account their natural tree structure. The most natural way of addressing this, is through the well studied abstract problem of Tree Edit Distance problem, where the idea is to compute a distance between a pair of trees. It should be noted that the development of this problem is very similar to the String Edit Distance problem. The Tree Edit Distance problem can be formulated as follows: Let Σ be a finite alphabet and let $\lambda \notin \Sigma$ be a special empty symbol. Finally let $\Sigma_\lambda = \Sigma \cup \lambda$. Given two trees T_1 and T_2 that are labelled (That is, each node N in the tree has assigned a label $l(N) \in \Sigma$) and ordered (That is, we are provided an order between sibling nodes), we define the following operations ($l(N_1) \rightarrow l(N_2)$) where $(l(N_1), l(N_2)) \in ((\Sigma \cup \lambda) \times (\Sigma \cup \lambda) \setminus (\lambda, \lambda))$:

Relabelling: If $l(N_1) \neq \lambda \wedge l(N_2) \neq \lambda$. We change the label of a node in the tree T

Deletion: If $l(N_2) = \lambda$. The children of N_1 become the children of N_1 's parent (N_1 is not the root of the tree) and occupy its position in the sibling ordering.

Insertion: If $l(N_1) = \lambda$. This is the complement operation of deletion.

For each of this operations $\omega = (l(N_1) \rightarrow l(N_2))$, we assign a cost function $cost(\omega)$ and for a given sequence of this transformations $\Omega = (\omega_1, \omega_2, \dots, \omega_n)$ we assign a cost function $cost(\Omega) = \sum_{i=1}^n cost(\omega_i)$. We also assume that $cost(\omega)$ is a distance metric. With the previous definition, the distance between T_1 and T_2 can be expressed now as $dist(T_1, T_2) = \min \{cost(\Omega) | \Omega \text{ is a sequence of operations that transform } T_1 \text{ into } T_2\}$

A recursive formulation of the problem can be defined in the following way:

$$\begin{aligned}
dist(T, T) &= 0 \\
dist(F_1, T) &= dist(F_1 - v, T) + cost(v \rightarrow \lambda) \\
dist(T, F_2) &= dist(T, F_2 - w) + cost(\lambda \rightarrow w) \\
dist(F_1, F_2) &= \begin{cases} dist(F_1 - v, F_2) + cost(v \rightarrow \lambda) \\ dist(F_1, F_2 - w) + cost(\lambda \rightarrow w) \\ dist(F_1(v), F_2(w)) + dist(F_1 - T_1(v), F_2 - T_2(w)) + cost(v \rightarrow w) \end{cases}
\end{aligned}$$

Where F is a forest, v and w nodes from a free, $F - v$ the forest F with the deletion of the node v , $T(v)$ the tree rooted at v and $F - T(v)$ the forest obtained by removing all the nodes in $T(v)$ from F

The presented set of equations give a simple way to compute the edit distance between a pair of trees, and the most known algorithms like Zhang and Shasha's[11] whose worst time case time bound is $O(|T_1|^2|T_2|^2)$, Klein's[12] that achieves a lower bound of $O(|T_1|^2|T_2|\log|T_2|)$ or Demaine's[13], take into advantage the fact that the recursion relies in the solution of smaller sub-problems, so by carefully choosing which of this sub-problems to solve.

Chapter 3

Technical Considerations

In the previous chapter we presented different abstract models which can be suitable in the building of a math based search system. In this section, we present different aspects that need to be considered when mapping from these theoretical models to a concrete implementation.

3.1 Digital Representations of Mathematics

When building any kind of information retrieval system, selecting the right representations of the data can affect the type of algorithms that one can use, the amount of storage, the quality of the data among other factors. Mathematical content is not an exception, and there are available different formats to encode a mathematical expression. One of the most common ways to represent mathematics in scientific documents, which suits CDS's domain, is \LaTeX .

Another important standard for the representation of mathematics is MathML, which was first proposed by the World Wide Web Consortium in 1998 as the recommended language to be used in web documents. One of the advantages of MathML over \LaTeX , is the fact that it is an application of XML, which makes it very easy process algorithmically. There are two different variations of MathML, presentation and content. The first one is oriented to the visual representation of the expression and the second one is more oriented to the semantics of the expression. Its main component is the `<apply>` tag which represents the function application.

OpenMath is yet another markup language to represent mathematics and can be used to complement presentation MathML to add more semantic information. It is mainly used now for defining mathematical concepts into standard content dictionaries which can be imported. MathML is compatible

with these dictionaries.

Table 3.1 presents a comparison on how different languages can express the quadratic equation $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

| L^AT_EX | Presentation MathML | Content MathML + Presentation Markups | OpenMath |
|--|---|---|---|
| $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ | <pre> <math display="block"> <mrow> <mi>x</mi> <mo>=</mo> <mfrac> <mrow> <mo>-</mo> <mi>b</mi> <mo>± </mo> <msqrt> <msup> <mi>b</mi> <mn>2</mn> </msup> <mo>-</mo> <mn>4</mn> <mo></mo> <mi>a</mi> <mo></mo> <mi>c</mi> </msqrt> </mrow> <mrow> <mn>2</mn> <mo></mo> <mi>a</mi> </mrow> </mfrac> </mrow> </math> </pre> | <pre> <math> <apply> <eq/> <ci>x</ci> <csymbol> <mfrac> <apply> <minus/> <mi>b</mi> ± <mroot> <apply> <minus/> <csymbol> <msup> <mi>b</mi> <mn>2</mn> </msup> </csymbol> <csymbol> <mn>4</mn> <mo>&times;</mo> <mi>a</mi> <mo>&times;</mo> <mi>c</mi> </csymbol> </apply> <mrow></mrow> </mroot> </apply> <apply> <times/> <cn>2</cn> <ci>a</ci> </apply> </mfrac> </csymbol> </apply> </math> </pre> | <pre> <OMOBJ> <OMA> <OMS cd = "relation1" name="eq"/> <OMV name="x"/> <OMA> <OMS cd="arith1" name="divide"/> <OMA> <OMS cd = "arith1" name="minus"/><OMA> <OMS cd = "arith1" name="times"/> <OMA> <OMS cd = "arith1" name="times"/> <OMV name="b"/> <OMV name="+ALE-"/> </OMA> <OMA> <OMS cd="arith1" name="root"/> <OMA> <OMS cd = "arith1" name="minus"/> <OMA> <OMS cd="arith1" name="power"/> <OMV name="b"/> <OMI>2</OMI> </OMA> <OMA> <OMS cd = "arith1" name="times"/> <OMA> <OMS cd = "arith1" name="times"/> <OMI>4</OMI> <OMV name="a"/> </OMA> <OMV name="c"/> </OMA> </OMA> <OMA> <OMI> 2 </OMI> </OMA> </OMA> </OMA> <OMA> <OMS cd = "arith1" name="times"/> <OMI>2</OMI> <OMV name="a"/> </OMA> </OMA> </OMA> </OMOBJ> </pre> |

Table 3.1: Comparison of the different languages for expressing mathematical content

L^AT_EX allows the most compact representations of the equations and it is the closest language to the community. However, its processing is not straightforward because a given equation can include multiple commands from different packages. Also there are several ways to represent the same set of symbols, either using its Unicode code-point or through embedded commands. Another difficulty are custom commands that users can define in the preamble of the document. In overall, the interpretation and processing of a mathematical expression in this language would require writing a complete compiler which is not the purpose of this work. MathML offers a better compromise between the length of the representation and the easiness to be

processed.

3.2 Math Extraction Tools

One of the main steps in any information retrieval task is the extraction of the searchable content from the collection of documents. At CDS the most common storage format for publications is PDF. A big part of the harvested documents in CDS comes from the ArXiv pre-prints service and from here, the source code (\LaTeX) of the submissions is available for download. Smaller portions of other documents are stored in other formats like Microsoft Word, OpenOffice Documents. In practice, the main options from where to do the extraction from, are PDF and \LaTeX files. Also taking into account the previous discussion on the available format for representing mathematics, we focused our attention in tools that output MathML.

A conscious exploration of available software for extracting mathematical content from this two types of documents allowed us to identify different tools with different approaches. \LaTeX files are easier to process and all of the tools implement their own compiler that XML + MathML outputs instead of the typical. Since the original math equation is provided in the document, there is relatively little ambiguity in the equation, and the major issues that the systems struggle with is in the coverage of all the possible commands and packages that are available.

For extracting equations from PDF files, the process turns into an Optical Character Recognition (OCR). The reconstruction of the equation has to infer the mathematical structure based on physical positioning and size of the symbols, the length of the lines and so on.

Table 3.2 presents the main features of the reviewed software.

| Software | Type | Tested version | Input Format(s) | Output Format(s) | Licence Type | Latest Update |
|----------------|-----------------------|----------------|-----------------------------------|------------------|--------------|---------------|
| LaTeXML[14] | Standalone Executable | 0.7.9alpha | \LaTeX | XML + MathML | Open Source | 13 Feb 2014 |
| TeX4ht[15] | Standalone Executable | 2009-01-31-07 | \LaTeX | XML + MathML | Open Source | 11 Jun 2009 |
| SnuggleTex[16] | Java Library | 1.2.2 | \LaTeX (Small fragments) | MathML | Open Source | 24 May 2010 |

| | | | | | | |
|---------------------|------------------------------------|-----------------|---------------------------------|---|------------------|-------------|
| LaTeXmathML[17] | Javascript Library | 30-October-2007 | L ^A T _E X | XML + MathML | Open Source | 30 Oct 2007 |
| Maxtract[18][19] | Standalone Executable | v.1752 | PDF | L ^A T _E X/Annotated PDF | Free to download | 15 Nov 2012 |
| InftyReader[20][21] | Standalone Executable (On Windows) | 2.9.6.2 | PDF, TIFF, BMP, GIF, PNG | L ^A T _E X, XML + MathML | Commercial | 22 Dec 2013 |

Table 3.2: Comparison of available Mathematical extraction tools

For extraction from L^AT_EX files we concluded that the best option is LaTeXML. The tool is still in current development and has lot of support from the academic community. The performance on a small dataset also showed that the results are consistently better than from the other tools. A more thorough test[22] also confirms that LaTeXML provides the best results. For PDF files, the trial version of InftyReader worked very well giving better results and processing successfully more files than Maxtract. As said previously, a big portion of the documents stored in CDS, come from pre-printing services like ArXiv were fortunately the source code of an article is available and because of this reason, our base set for this work will consist on extracting equations from L^AT_EX files. For licensing reasons, for this project we will not integrate our system with InftyReader for the moment and will proceed with LaTeXML as our extraction tool.

3.3 Computer Algebra Systems

Computer Algebra Systems (CASs) are software packages that allow performing symbolic transformations on mathematical objects. They allow to manipulate, simplify and analyse mathematical equations and offer interesting functionalities for a mathematical based search engine. The most important functionalities that we are looking for in such system are simplification and normalization of mathematical expressions and pattern matching. As discussed previously, the language for representing expressions will be MathML, so we also require that such system is able to import this format.

Finally, we are also interested in integrating these functionalities in a bigger system, so integration through APIs was also considered.

| CAS | Imports MathML | Expression normalization | Pattern Matching | License Type | Supported APIs |
|--------------------------------|-----------------------|---------------------------------|-------------------------|---------------------|------------------------------|
| Maple | Yes | Yes | Yes | Commercial | C, Java, VB |
| MatLab (Symbolic Math Toolbox) | Experimental | Yes | Yes | Commercial | C/C++, Fortran |
| Mathematica | Yes | Yes | Yes | Commercial | C, Java, .NET |
| Maxima | Experimental | Yes | Yes | Open Source | C++, Java (External Project) |
| SymPy | No | Yes | Yes | Open Source | Python |

Table 3.3: Comparison of main CAS

Table 3.3 presents a comparison of the desired features in the main CASs. The systems that better fit our requirements are Maple and Mathematica. For licensing reasons, we will continue our project using Mathematica.

Chapter 4

Related Work

4.1 Status of Mathematical Information Retrieval

Even though Information Retrieval has been a research field from around the 1950s where the first metrics for such systems were developed and in the 1960s where SMART (System for Mechanical Analysis and Retrieval of Text), the first computer based system was developed; only in the last couple years focus have been put into adapting the same techniques to mathematical content.

The Mathematics Information Retrieval Workshop[23] held in the context of the Conference of Intelligent Computer Mathematics in July 2012, was the first official event in the area of MIR. Here, a small competition was held to evaluate different systems and approaches.

The NII[24] (National Institute of Informatics) is a Japanese research institute that hosts the project named NTCIR (NII Testbeds and Community for Information access Research) which aims to build evaluation frameworks for specific sub-fields of information retrieval. On 2013, in the context of the 10th NTCIR conference, it hosted the NTCIR-10 Math Task [25] which was the first collaborative effort to evaluate different systems. The Math task was focused on two different subtasks: Math retrieval (Identifying relevant document based on given mathematical expressions and/or keywords) and Math understanding where the idea was to extract natural language information of a mathematical expression in a document. Currently submissions are open for the NTCIR-11 Math Task 2 which will be held on December 2014.

4.2 Mathematical based search projects

4.2.1 MIaS

MIaS[26] (Math Indexer and Searcher), currently, is one of the flagship projects in the indexing and searching of mathematical content. As most of the other projects, it processes documents in MathML format. It allows the user to include in the queries mathematical expressions and textual content. During indexing time, equations are transformed following a set of heuristics that include:

- Ordering of elements: Taking into account commutativity of certain operators (Addition, multiplication), elements are ordered such that $3 + a$ would be converted into $a + 3$ (Since the `<mi>` tag comes first than the `<mn>` tag)
- Unification of variables: This process takes into the account the structure of the equations in despite of the naming of the variables. Expressions like $a + b^a$ and $x + y^x$ would be converted into an expressions of the form $id_1 + id_2^{id_1}$ which would match.
- Unification of constants: This step consists of replacing all occurrences of constants (`<mn>` tags) by a const symbol.

The extracted tokens from a given equation consist of all its valid sub-expressions. The original expression is given a weight value of 1, and following sub-expressions are given smaller weights depending on how general or specific they are. The system, also as most of the current projects, is developed by using the Apache Lucene framework. The system adapts the default scoring equation such that the given weight is taken into account. Publicly available instances of the system can be found at: <http://aura.fi.muni.cz:8085/webmias/ps?n=-1> running on the MREC[27] dataset and <http://aura.fi.muni.cz:8085/webmias-ntcir/> (Running on the NTCIR dataset)

4.2.2 EgoMath

EgoMath[28] and its new version EgoMath²[29], is a system oriented to index the mathematical content from the Wikipedia.org database. The processing steps before indexing are similar to the ones in MIaS (Rearranging of symbols, unification of constants). The extraction process started from a complete dump of the Wikipedia database and filtering only the math articles, which are identified by the "`<math>`" tag. This consists on around thirty thousand

articles, from which 240.000 equations were identified. The processing step includes translating the equation from LaTeX into MathML format and then indexing both representation. Some additional effort is done into splitting large mathematical blocks like tables into single mathematical expressions. At the start of this work, EgoMath was publicly available at <http://egomath.projekty.ms.mff.cuni.cz/>, however at this moment the system seems to be unavailable.

4.2.3 DLMFSearch

The Digital Library for Mathematical Functions [30] is a project launched by the National Institute of Standards and Technology, in 2010, as an online version of the Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables[31]. DLMF's main goal is to compile the mathematical knowledge in the form of equations, functions, tables and make this information useful for researchers and public in general. The system contains a search functionality that allows to input a combination of full text and \LaTeX snippets. The system tries to perform an exact match, and if no results are found, it relaxes the query. DLMF also performs some normalization of the equation and some cleaning of some characters, but no further details are provided. The content indexed by this project is highly curated, which differentiates it from the other projects. In the report published in 2013 [32], it is reported to have indexed around 38000 equations. This project is publicly available at <http://dlmf.nist.gov/>

4.2.4 MathWebSearch

MathWebSearch[33] (Currently on its 0.5 version) is a open-source search engine for mathematical equations. While most of the documentation relates on the architecture of the system and how they address scalability; the indexing technique is also very interesting. The system implements an idea proposed by Peter Graf in [34] called substitution tree indexing. This idea can be viewed as a generalization of the variable and constant unification. It represents the equations as a tree and recursively generalizes each sub-expression. A public available instance of the system is available at <http://search.mathweb.org/tema/> which runs on the Zentralblatt Math database[35].

4.2.5 LaTeXSearch

LaTeXSearch <http://latexsearch.com/> is a system developed by the scientific publisher Springer that allows to search over a database of around eight million latex snippets extracted from their publications. The system, unfortunately, is proprietary and no further details are provided.

4.2.6 WikiMirs

WikiMirs[36] is a system that, as EgoMath, allows to find relevant Wikipedia entries by providing mathematical expression in \LaTeX . It combines the information about the semantic tree with the layout tree of a given expression and uses a similar approach to index different generalization levels. The system is available at <http://www.icst.pku.edu.cn/cpdp/wikimirs/>

Other systems that provide similar search functionalities are Uniquation(<http://uniquation.com>) and Symbolab(<http://www.symbolab.com/>), however detailed information about the systems are not provided. Finally Wolfram Alpha[37] is a very powerful tool that allows to find information about lots of mathematical concepts and a wide variety of other objects (Even Pokemon).

Chapter 5

$5e^{x+y}(\text{CDS})$

In this chapter we present our solution called $5e^{x+y}(\text{CDS})$. We present the key points of our system, and how they are related to previous built systems.

5.1 Features Extraction

Our main goal is to develop a system that compliments the already available search capabilities in Invenio and in particular the CDS instance. Our system was developed on top of the the Apache Lucene/Solr framework. CDS already uses Solr as a search engine layer, so developing $5e^{x+y}(\text{CDS})$ using this technology supposed a reasonable choice. One of the main steps while building an IR tool, is understanding the data and what kind of features can be extracted from it: If the system is based on text, then tasks such as stemming, handling of synonyms, misspelling and so on, should be taken into consideration; if the system being developed is based on images, then examples of features can be histograms of colors, geometrical features and so on. During the literature review, we discovered that most of the projects handled the notational aspects of the expressions, that is, the real naming of the terms in a given equation and each project developed its own set of heuristics to handle the structural component.

5.1.1 Notational Features

Notational features correspond to the elements that relate to the actual naming of the variables, constants and numeric quantities found in an equation. Here we extract tokens that represent single leaf nodes (E.g. $\langle \text{mi} \rangle$, $\langle \text{mn} \rangle$, $\langle \text{mo} \rangle$ tags), and simple constructs like subscript and superscript elements, fractions, roots. We also extract bi-grams and tri-grams as an analogy to

what is done for text indexing in current systems.

After examining a set of documents from our dataset, we identified some heuristics that could be applied to improve the recall of our system:

Unicode Normalizing

The first issue we discovered is that because of the wide range of different characters to represent the possible mathematical symbols, there is need for an automatic way to match different possible variations of a single character. An example of this situation is the specification of an natural number, sometimes denoted by the letter N. In this case, at least three different common representations were identified: $N = 1$ (Using character with unicode code point 0x004e), $\mathcal{N} = 1$ (0xd835 0xdc41 or in latex expressed as `\mathcal{N}`) and $\mathbb{N} = 1$ (0x2115 or in latex expressed as `\mathbb{N}`).

To address this situation, we rely on the Unicode equivalence framework. In the unicode specification, sometimes the same character has two or more different code-points because of backwards compatibility with other encoding standards or because the same character has different essential meanings. Another group of equal characters with different code-points are pre-combined characters such as the Latin letter Ñ which has Unicode point 0x00D1 and also is the sequence of code-points 0x004E (Latin letter N) and 0x0303 (the combining character tilde). Taking this sample into account, Unicode defines the character composition and decomposition operations as replacing the decomposed representation with the pre-composed one and vice-versa. Finally, Unicode defines two types of equivalences between characters: Canonical and compatible. Canonical equivalent characters are assumed to have the same appearance and meaning. Compatible ones are allowed to have slightly different graphical representations but the same meaning in some contexts. One example of compatible equivalent characters is the Roman number 12 **XII** with code-point 0x216b and the sequence of characters XII with codepoints 0x0058 0x0049 0x0049. Taking into account the composition or decomposition of characters and the canonical or compatible equivalences, Unicode specifies 4 different forms of a given character:

- NFD: Normalization Form Canonical Decomposition
- NFC: Normalization Form Canonical Composition
- NFKD: Normalization Form Compatibility Decomposition
- NFKC: Normalization Form Compatibility Composition

Table 5.1 presents a small list of troubling characters classes that were identified in our datasets.

| Character Name | Visual Rendering | NFC | NFD | NFKC | NFKD |
|--|------------------|----------------------|----------------------|---------------------------------|------------------------|
| LATIN CAPITAL LETTER A | A | 0x41 | 0x41 | 0x41 | 0x41 |
| LATIN CAPITAL LETTER A WITH MACRON | Ā | 0x100 | 0x41 0x304 | 0x100 | 0x41 0x304 |
| LATIN CAPITAL LETTER A WITH RING ABOVE | Å | 0xc5 | 0x41 0x30a | 0xc5 | 0x41 0x30a |
| ANGSTROM SIGN (0x212b) | Å | 0xc5 | 0x41 0x30a | 0xc5 | 0x41 0x30a |
| LATIN CAPITAL LIGATURE IJ | IJ | 0x132 | 0x132 | 0x49(I) 0x4a(J) | 0x49 0x4a |
| CYRILLIC CAPITAL LIGATURE A IE | Æ | 0x4d4 | 0x4d4 | 0x4d4 | 0x4d4 |
| LATIN CAPITAL LETTER AE | Æ | 0xc6 | 0xc6 | 0xc6 | 0xc6 |
| DOUBLE-STRUCK CAPITAL N | N | 0x2115 | 0x2115 | 0x2115 | 0x4e |
| LATIN CAPITAL LETTER N | N | 0x4e | 0x4e | 0x4e | 0x4e |
| MATHEMATICAL ITALIC CAPITAL N | <i>N</i> | 0xd835 0xdc41 | 0xd835 0xdc41 | 0x4e | 0x4e |
| ROMAN NUMERAL TWELVE | XII | 0x216b | 0x216b | 0x58 0x49 0x49 | 0x58 0x49 0x49 |
| LATIN CAPITAL LETTER X - LATIN CAPITAL LETTER I - LATIN CAPITAL LETTER I | XII | 0x58 0x49 0x49 | 0x58 0x49 0x49 | 0x58(X) 0x49(I) 0x49(I) | 0x58 0x49 0x49 |
| VULGAR FRACTION ONE QUARTER | ¼ | 0xbc | 0xbc | 0x31(1) 0x2044(/) 0x34(4) | 0x31 0x2044 0x34 |
| DIGIT ONE - SOLIDUS - DIGIT 4 | 1/4 | 0x31 0x2f 0x34 | 0x31 0x2f 0x34 | 0x31 0x2f 0x34 | 0x31 0x2f 0x34 |

| | | | | | |
|--|----------------|------------------------|------------------------|------------------------|------------------------|
| DIGIT ONE - SOLIDUS - DIGIT 4 | $\frac{1}{4}$ | 0x31 0x2f 0x34 | 0x31 0x2f 0x34 | 0x31 0x2f 0x34 | 0x31 0x2f 0x34 |
| DIGIT ONE - FRAC- TION SLASH - DIGIT 4 | $\frac{1}{4}$ | 0x31 0x2044 0x34 | 0x31 0x2044 0x34 | 0x31 0x2044 0x34 | 0x31 0x2044 0x34 |
| THERE EXISTS | \exists | 0x2203 | 0x2203 | 0x2203 | 0x2203 |
| THERE DOES NOT EXIST | \nexists | 0x2204 | 0x2203 0x338 | 0x2204 | 0x2203 0x338 |
| LATIN CAPITAL LETTER OPEN E | \mathfrak{E} | 0x190 | 0x190 | 0x190 | 0x190 |
| EULER CONSTANT | \mathfrak{E} | 0x2107 | 0x2107 | 0x190 | 0x190 |
| PLANCK CON- STANT | h | 0x210e | 0x210e | 0x68 | 0x68 |
| LATIN SMALL LET- TER H | h | 0x68 | 0x68 | 0x68 | 0x68 |
| PLANCK CON- STANT OVER TWO PI | \hbar | 0x210f | 0x210f | 0x127 | 0x127 |
| LATIN SMALL LETTER H WITH STROKE | h̄ | 0x127 | 0x127 | 0x127 | 0x127 |
| GREEK CAPITAL LETTER OMEGA | Ω | 0x3a9 | 0x3a9 | 0x3a9 | 0x3a9 |
| OHM SIGN | Ω | 0x2126 | 0x2126 | 0x3a9 | 0x3a9 |
| INTEGRAL | \int | 0x222b | 0x222b | 0x222b | 0x222b |
| DOUBLE INTE- GRAL | \iint | 0x222c | 0x222c | 0x222b 0x222b | 0x222b 0x222b |
| CONTOUR INTE- GRAL | \oint | 0x222e | 0x222e | 0x222e | 0x222e |
| NABLA | ∇ | 0x2207 | 0x2207 | 0x2207 | 0x2207 |
| WHITE DOWN- POINTING TRIAN- GLE | ∇ | 0x25bd | 0x25bd | 0x25bd | 0x25bd |

Table 5.1: Characters with similar glyphs and/or semantics

Since our goal is to match as many similar characters as possible, our system employs the NFKD representation of a character. For a

given token, we test whether if it is already equal to its NFKD. If not, for all the characters that are not combining nor control characters, we add the token into the same position as the original one. For example, given the token $\langle \text{mi} \rangle \text{\AA} \langle \text{mi} \rangle$, this will lead to the sequence $\langle \text{mi} \rangle \text{\AA} \langle \text{mi} \rangle \langle \text{mi} \rangle \text{A} \langle \text{mi} \rangle$ and the token $\langle \text{mi} \rangle \text{XII} \langle \text{mi} \rangle$ will produce $\langle \text{mi} \rangle \text{XII} \langle \text{mi} \rangle \langle \text{mi} \rangle \text{X} \langle \text{mi} \rangle \langle \text{mi} \rangle \text{I} \langle \text{mi} \rangle \langle \text{mi} \rangle \text{I} \langle \text{mi} \rangle$.

Operator grouping

Even though the Unicode normalization step works for a big part of the cases presented in table 5.1, there are still some groups of characters which would not be matched (even partially), such as some types of integrals. To address these groups, a precomputed table was created for some general groups of operators which have some similar meaning. For each token that belongs to some of this predefined groups, a second token is created identifying the group. For example the token $\langle \text{mo} \rangle \oint \langle \text{mo} \rangle$ is expanded into $\langle \text{mo} \rangle \oint \langle \text{mo} \rangle \langle \text{mo} \rangle \text{INTEGRALS} \langle \text{mo} \rangle$, similarly the token $\langle \text{mo} \rangle \int \langle \text{mo} \rangle$ is expanded into $\langle \text{mo} \rangle \int \langle \text{mo} \rangle \langle \text{mo} \rangle \text{INTEGRALS} \langle \text{mo} \rangle$ and then both will have a partial match in the INTEGRALS token. The Unicode specification does not define a strict rule to group similar characters based on their semantics. For general operators, we used the grouping provided by Xah Lee in his website[38].

Figures 5.1 and 5.2 present lists related of characters for two different categories.

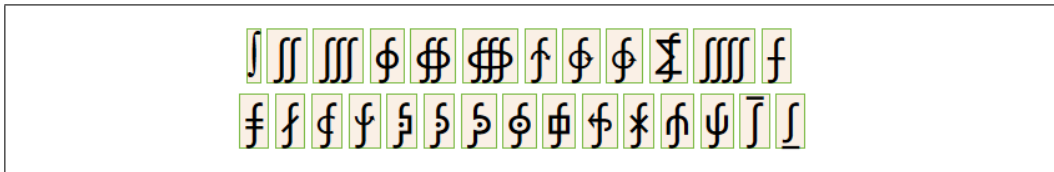


Figure 5.1: List of different types of integral related characters in Unicode

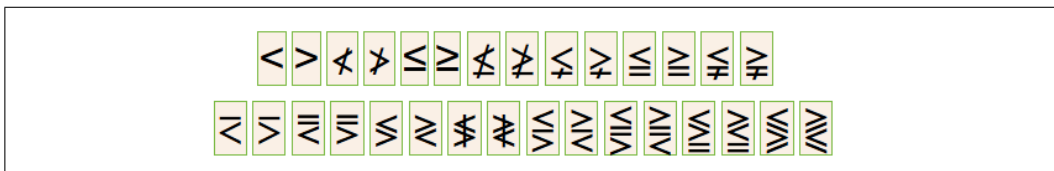


Figure 5.2: List of different types of greater/less than related characters in Unicode

Numeric Approximation

Some equations relate to specific numeric quantities. While most of the previously explored systems handled `<mn>` tags by only representing them as a constant identifier, this approach is not suitable. A lot of articles in CDS relate to specific experiments under certain conditions, and require a high level of precision in the measurements. However, different experiments of the same phenomenon, can result in different measurements and we wanted to include this variability in our system. To handle this, when a token with a certain numeric quantity is discovered, we apply iterative rounding and index all the different approximations in the same position as the original element. For example the token `<mn>2.0135</mn>` will produce the following sequence of tokens: `<mn>2.0135</mn><mn>2.014</mn><mn>2.01</mn><mn>2.0</mn>`.

5.1.2 Structural Features

The structure of a mathematical expression compliments the information that can be obtained over generic substructures where the naming of the variables is irrelevant.

Algebraic structures

The first set of structures consists of purely algebraic relationships between the expressions. Inspired by predefined integral tables for different kinds of expressions, we decided to extract, from a given equation, whether it has an occurrence of these kinds of patterns or not. An example of this type of structures would be: $X_ * (X_ + A_)$ where $X_$ and $A_$ can be any type of expression. Sample equations matching this pattern can be $C(C+1)$ or $\frac{(\sin \theta \cos \theta)^2}{2a} \left(\frac{(\sin \theta \cos \theta)^2}{2a} + \tan \frac{\theta + \phi}{2} \right)$. This structural analysis of equations is more powerful than the unification of variables and constants proposed in previously discussed projects because these patterns, as shown in the second example, can match an arbitrary complex expression. The drawback of our solution is the need to manually select representative and useful patterns. This selection was done based on standard based tables of integrals similar to the ones that CASs employ to compute them.

Subscript/Superscript variations

This set of structures represent common operations that can appear in a given expression, but its identification is based on notation factors. A simple

example of this can be the matrix decomposition $A = VDV^{-1}$ where the actual name of the variables A and V is not meaningful but the real importance is the fact that a given variable appears alongside its inverse in a matrix multiplication form. We would like to assign a partial match to any other kind of this decomposition say $B = MRM^{-1}$. These types of structures live in-between the notational and structural categories but since we identify them using Mathematica, we encompass them as structural features. For this type of structures, we identify simple operations between an element and some possible variations of the same element with some subscript/superscript.

Specific domain structures

Additionally, there are common patterns to specific domains which can improve the quality of the retrieved results. Some specific patterns we identified include: $originalParticle \rightarrow subparticle1, subparticle2, subparticle3$ (for representing particle decaying processes) $S = \int_{t1}^{t2} \mathcal{L} dt$ (for an action in terms of its Lagrangian) and $SU(n)$ (for describing special unitary groups). The main drawback of this approach is the need to manually detect, implement and maintain long lists of patterns and for this reason we did not try to assemble an exhaustive list.

5.2 Development with Lucene/Solr

This section explains how the described components were implemented in the Apache Lucene/Solr framework. One focus of our work was to develop a modular system that allowed easily to deploy new heuristics as they were identified.

The Lucene framework defines some key concepts which ultimately are mapped into concrete classes. A Document is the minimum retrieval unit. For $5e^{x+y}$ we could either use a complete article as a document or a single equation. We choose the latter because it makes it easier to focus on extracting features from a single equation. A document contains a set of fields which represent different types of information that it can store. Each field has different options such as whether it would be indexed and/or stored. A field can also indicate whether it will be a Term Vector. Our fields are as follow:

- `math_notational_field`: This field stores the notational features as a Term Vector.
- `math_structural_field`: This field stores the structural features as a Term Vector.

- `math_normalized_notational_field`: This field stores the notational features of the normalized expression as a Term Vector.
- `filename`: This is an auxiliary field which stores the CDS record / article name where the equation was found.
- `number_occurrences`: Another auxiliary field that is used to score final results.

Each field is associated to a particular Analyzer. An analyzer is a specific class that handles an element from a field and groups all the operations that are to be performed. For this purpose we implemented `SolrNotationalAnalyzer`, `SolrStructuralAnalyzer` and `SolrNormalizedNotationalAnalyzer`. The `Solr` prefix is added to clarify the fact that the class extends a `SolrAnalyzer` which extends a normal `LuceneAnalyzer`.

An analyzer consists of a sequence of `TokenStream` elements. Normally an analyzer contains one `Tokenizer` element (which extends a `TokenStream`) that is in charge of creating the Terms that will be stored in the field, and one or more `TokenFilter` elements which perform operations on the single tokens that were created previously. For the specific cases of `SolrNotationalAnalyzer` and `SolrNormalizedNotationalAnalyzer` we created the following classes:

- `MultiplePatternTokenizer`: This class is in charge of applying a set of regular expressions over the given expression in order to extract the single elements that were presented. Lucene incorporates the class `PatternTokenizer` but it only accepts one regular expression.
- `UnicodeNormalizingFilter`: This class receives single tokens and verifies if the token is in Unicode NFKD. In negative case, it emits an additional token in the normalized form.
- `SynonymExpandingFilter`: This class verifies if the received token belongs to a predefined category of elements and emits an additional token with it.
- `NumericalRoundingFilter`: This class processes `<mn>` tokens. When a number contains a decimal part, it applies the standard rounding mode iteratively and for each one, it emits a token for each step.

`SolrStructuralAnalyzer` is simpler in its workflow: It consists of a `StructuralPatternsTokenizer` that tries to find occurrences of the predefined patterns in the expression. For each match, it emits a token with the pattern that was found.

5.3 Integration with Mathematica

This section describes some aspects that had to be considered while implementing the integration with Mathematica. Our system uses the JLink library to create a connection with the underlying process `MathKernel`. The main class for communicating java code and Mathematica is `KernelLink` and the main used method was `evaluateToOutputForm(String expr, int pageWidth)` where `expr` is the expression to be evaluated and `pageWidth` affects the formatting of the result. Once all the boilerplate has been written, the steps followed to define and find patterns in the expression are as follow: We import the string in MathML format through the command `ImportString["<mathmlExpression>", "MathML"]`. This creates the expression in the box form, but no further analysis is possible from this form. We then try to interpret this form into a fully Mathematica compatible expression by using `MakeExpression`

5.3.1 Pattern Extraction

Once the expression is created, finding a pattern is easy by using the command `Position[<variable holding the interpreted expression>, Pattern]`. A simple pattern for a quadratic expression is: `a_.x^2 + b_.x + c_.`

5.3.2 Cleaning of variables

After each operation, all of the variables should be cleaned or it may lead to deadlocks, wrong results or other issues since the `MakeExpression` command also interprets the expression so if the expression is, for example `a = b + c` Mathematica remembers the assignment of the new value of `a` and uses it in further evaluations. The cleaning is done through the command `Remove["Global`*"]`

5.3.3 Expression simplification

The second feature that we are interested in is automatic expression simplification. Mathematica includes different transformations of expressions. These include: `Simplify` which tries to transform a given equation in an equivalent one but in a reduced form, `FullSimplify` which is similar to the previous one but tries a higher number of heuristics to apply that can simplify expressions further at the expense of more processing time, `Factor` to factorize polynomials, `Refine` which simplifies expressions given certain

assumptions about the variables, TrigReduce that tries to apply trigonometric identities to simplify the expression, FunctionExpand that tries to expand special functions in terms of sums and products of other elementary functions.

The two functions that are more relevant for our work are Simplify and FullySimplify since they apply a combination of the other methods and do not require additional information.

Table 5.2 presents some sample equations found in some the documents in CDS and the different forms of simplification.

| Original Form | Simplify | FullSimplify |
|---|---|--|
| $\frac{\Phi'}{2(1+\Phi)} + \frac{1}{r} = 0$ | $\frac{1}{r} + \frac{\Phi'}{2\Phi+2} = 0$ | $\frac{2}{r} + \frac{\Phi'}{\Phi+1} = 0$ |

Table 5.2: Comparison of different simplification modes

5.3.4 Recursive interpretation

Unfortunately, Mathematica cannot always interpret correctly an input expression. In this case, it can, however interpret some sub-expressions from the original one. Our system takes advantage of this, and the first attempt is to try to interpret the original string. If this step is not successful (the result variable equals the string "{ }"), we identify the occurrences of RowBox which identify all the subexpressions in the tree. For each of this elements we call the MakeExpression and if it was successful we search for structural patterns in it. A first thought would tell us that if sub-tree was interpreted successfully, there is no need to try to go further down in that sub-tree, however, even if a sub-expression is interpreted correctly, the result sometimes is not complete and processing further down gives better results. For this reason, we do not discard any of the subexpressions and try to interpret all of them. For this same reason, we do not store the number of occurrences of a given pattern, but only whether it contains is or not.

5.3.5 Error handling

Sometimes, after certain time, the connection between KernelLink and Mathematica is lost and it launches an exception and all following invocations fail as well. Also, re-instantiating a KernelLink would yield an Exception. For this scenario, killing the associated processes to Mathematica was necessary with `killall -s 9 Mathematica; killall -s 9 MathKernel` or `taskkill /im Mathematica.exe; taskkill /im MathKernel.exe` depending on the OS the system is running on.

With this, the next creation of a new KernelLink works properly and the subsequent invocations work as expected.

5.3.6 Timeout Handling

Unfortunately, checking the quality of a LaTeX document is a difficult job and users have a lot of freedom in the way they write their document since the only visible result is the rendered version. One example we found during our work is the following TeX snippet:

```
"\ \varphi\ are\ shown,\ as\ well\ as\ the\ weight\ and\ tension"
```

A human can easily identify that this corresponds to text and that should be inserted into a proper "mbox" or "text" tag. In this case, the snippet produced the following MathML code:

```
<math><mrow><mi>\varphi</mi><mo></mo><mrow><mi>a</mi><mo></mo>
```

```

<mi>r</mi><mo></mo><mi>e</mi></mrow><mo></mo><mrow><mi>s
</mi><mo></mo><mi>h</mi><mo></mo><mi>o</mi><mo></mo><mi>
w</mi><mo></mo><mi>n</mi></mrow><mo>,</mo><mrow><mi>a</m
i><mo></mo><mi>s</mi></mrow><mo></mo><mrow><mi>w</mi><mo>
</mo><mi>e</mi><mo></mo><mi>l</mi><mo></mo><mi>l</mi></
mrow><mo></mo><mrow><mi>a</mi><mo></mo><mi>s</mi></mrow>
<mo></mo><mrow><mi>t</mi><mo></mo><mi>h</mi><mo></mo><mi>
e</mi></mrow><mo></mo><mrow><mi>w</mi><mo></mo><mi>e</m
i><mo></mo><mi>i</mi><mo></mo><mi>g</mi><mo></mo><mi>h</
mi><mo></mo><mi>t</mi></mrow><mo></mo><mrow><mi>a</mi><m
o></mo><mi>n</mi><mo></mo><mi>d</mi></mrow><mo></mo><mro
w><mi>t</mi><mo></mo><mi>e</mi><mo></mo><mi>n</mi><mo></
mo><mi>s</mi><mo></mo><mi>i</mi><mo></mo><mi>o</mi><mo></
mo><mi>n</mi></mrow></mrow></math>

```

Mathematica took around 10 minutes to interpret it. In the ideal case we would be able to identify the malformed expressions and pass to Mathematica only well-formed equations to be processed but as this simple example shows, that is a very challenging task and may require more work. As a remedial situation, we used a fixed timeout to stop the processing of such expressions.

5.4 Integration with CDS

Currently the Invenio code-base is being under a deep refactoring and rewriting. CDS currently runs in the stable branch master while the new developments are preferred to be implemented in the new branch pu. Our system is deployed as a standalone jar which is added to a Solr instance. Also, some extra configurations should be added. Details on this can be found in appendix A. The communication between Solr and Invenio was through the library `solrpy`. Figure 5.3 presents the main components of our deployment.

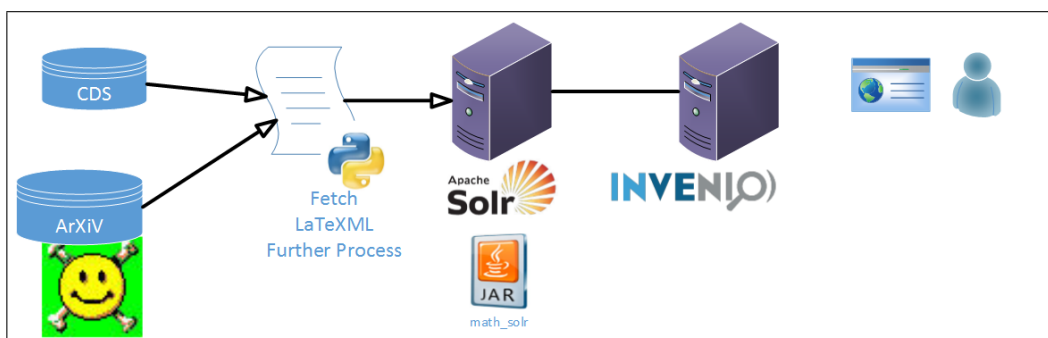
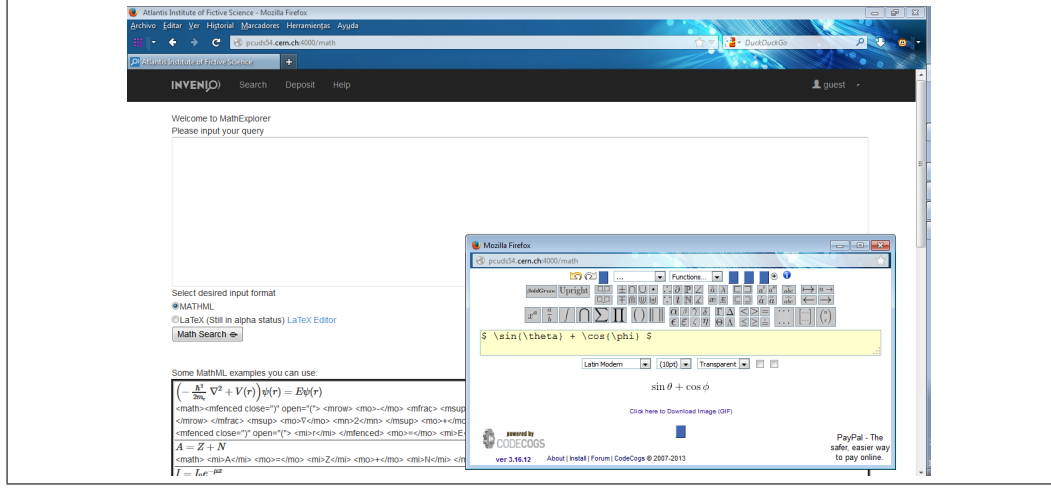
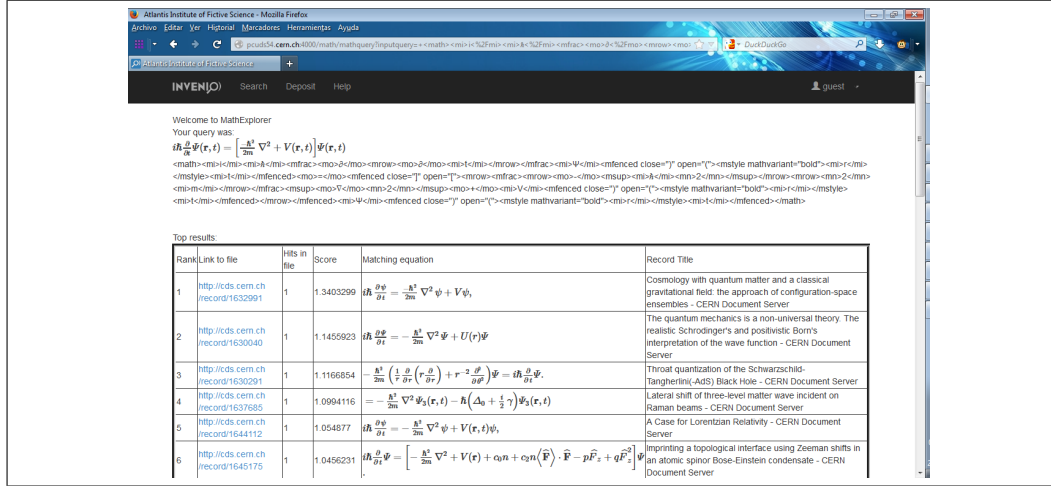


Figure 5.3: High level deployment architecture

A demo interface over the pu branch was developed to showcase the system, however to have a complete integration with CDS, more work is required. The

Figure 5.4: Web interface of $5e^{x+y}$: Query input screenFigure 5.5: Web interface of $5e^{x+y}$: Results screen

demo setup employs the same dataset as used for experimentation. We include an online \LaTeX editor[39] that renders the input expression as the user types it.

Chapter 6

Clustering

In parallel to the development of $5e^{x+y}$ we perform some analysis over our datasets to have a better understanding of the data. Even though this series of analysis did were not directly introduced into the codebase of the system, they provide valuable insights and allowed us to identify future directions of work.

6.1 Technical details

We used Apache Mahout v0.9[40] for performing the clustering task. The following set of commands were used to create and analyse our dataset:

To convert the solr index to a mahout compatible format (Mahout vector):

```
./mahout lucene.vector -d <path to solr index> -o ../re-
sults/vectors -t ../results/dict -f math_notational_field -st
../results/dict_seq -x 20 -md 2 -w TFIDF -n 2
```

To run the canopy clustering on the data:

```
./mahout canopy -t1 0.15 -t2 0.1 -i ../re-
sults/vectors -o ../results/clustering/canopy -dm
org.apache.mahout.common.distance.CosineDistanceMeasure
```

To run k-means algorithm on the resulting canopies:

```
./mahout kmeans -i ../results/vectors -o ../re-
sults/clustering_kmeans -xm sequential -c ../re-
sults/clustering/canopy -x 50 -cd 0.1 -cl -dm
org.apache.mahout.common.distance.CosineDistanceMeasure
```

To get the resulting points with their assigned cluster in a human readable format: `./mahout seqdumper -i ../results/clustering_kmeans/clusteredPoints/part-m-0`
> `../results/clustering_kmeans/clusteredPoints/cluster-points.txt`

6.2 Estimating number of clusters

Since the estimation of a good number of clusters k is a challenging task, we also tried to estimate it by evaluating intra and inter cluster distances for different values of k . Table 6.1 shows that around $k=500$ the inter-cluster distance begins to remain stable and by increasing k to 1024, its value is affected only by less than 1% indicating that a good k number is around 500.

| K | Average Intra-cluster distance | Average Inter-cluster distance |
|----------|---------------------------------------|---------------------------------------|
| 2 | 0.916 | 0.376 |
| 4 | 0.907 | 0.557 |
| 8 | 0.881 | 0.538 |
| 16 | 0.859 | 0.628 |
| 32 | 0.818 | 0.781 |
| 64 | 0.772 | 0.833 |
| 128 | 0.732 | 0.910 |
| 256 | 0.694 | 0.938 |
| 512 | 0.655 | 0.955 |
| 1024 | 0.610 | 0.963 |

Table 6.1: Cluster measurements for different values of k

The resulting clusters can help us to identify common structures and further improve our results by either adding frequent structures that are common to a cluster, or by reducing the number of elements to perform comparisons against and allow to speed-up the querying process. Table 6.2 shows the top terms in 2 different clusters that were identified. By simple inspection, it is possible to identify the first one as a cluster grouping trigonometric expressions (Sines and Cosines) and the second one grouping limits where the variable goes to infinity.

| Term | Relative frequency | Term | Relative frequency |
|-------------|---------------------------|---------------|---------------------------|
| \cos | 0.163661208 | ∞ | 0.233151718 |
| \cos | 0.157439621 | \rightarrow | 0.108867782 |
| \cos | 0.121893736 | \rightarrow | 0.108017583 |
| θ | 0.097241797 | ARROWS | 0.105854015 |
| θ | 0.095168449 | $-$ | 0.105012368 |
| \cos | 0.088535288 | ∞ | 0.102411785 |
| θ | | | |

| | | | |
|--|-------------|--|-------------|
| \cos \cos \cos | 0.088535288 | \rightarrow \rightarrow ∞ \rightarrow | 0.055650946 |
| \sin \sin | 0.068049679 | \rightarrow ∞ | 0.055490528 |
| \sin | 0.066531563 | \rightarrow \rightarrow \rightarrow ∞ \rightarrow | 0.043459138 |
| \rightarrow \sin \sin \cos \rightarrow | 0.032338911 | $\rightarrow,$ ∞ | 0.021523067 |
| \sin \sin θ | 0.032338911 | \rightarrow \rightarrow \rightarrow \rightarrow ∞ \rightarrow | 0.020583379 |
| ϕ | 0.031248212 | t \rightarrow | 0.016400033 |
| Φ | 0.019359824 | \rightarrow \rightarrow \rightarrow \rightarrow ∞ \rightarrow | 0.01478247 |
| \cos | 0.018438684 | t | 0.014620022 |
| \sin Φ | 0.017617149 | \rightarrow \pm ∞ \rightarrow | 0.014477746 |
| \rightarrow \rightarrow \cos \sin \cos \rightarrow | 0.016889963 | \pm ∞ | 0.014348642 |
| \cos \sin Φ | 0.016157276 | $ $ | 0.012372129 |
| \rightarrow \cos \sin Φ \rightarrow | 0.016157276 | \lim | 0.012214223 |

| | | | |
|--------------------------|-------------|----------------------------|-------------|
| $\langle \phi \rangle$ | 0.015344056 | $\langle C \rangle$ | 0.010421859 |
| $\langle \alpha \rangle$ | 0.014598769 | $\langle \epsilon \rangle$ | 0.009871146 |

Table 6.2: Top 20 terms for two different identified clusters

Chapter 7

Evaluation

When building an IR system there are two main aspects that should be considered. One is the quality of the presented results and the second one is the performance of the system since web users currently tolerate some milliseconds of latency for retrieving results.

7.1 Dataset

Our dataset consists on 12234 latex documents fetched from ArXiv. This documents correspond to the latest records harvested for CDS. The documents were processed using LateXML as described in the extraction process. From the given dataset, the simpler equations were filtered out (The ones containing one or two elements). The total number of indexed equations is 1412297.

7.2 Setup

The machine used for experiments has a Intel Core i7 processor running at 3.4Ghz with 4 cores. The available RAM is 8Gb, and the Java Virtual Machine used for running Solr was run with the `-Xms2G -Xmx5G` parameters. The complete specifications of the used software is presented in table 7.1

| Software | Version |
|----------------------|---------------------------------------|
| Operative System | Linux Mint 16 Petra (Ubuntu based) |
| Linux Kernel | 3.11.0-15 |
| Java Virtual Machine | 1.7.0-45 |
| Apache Solr | 4.6 |
| Python | 2.7.5+ |

Table 7.1: Software used during evaluations

7.3 Quality results

In this section we are interested in evaluating the quality of the results provided by our system. The type of queries that $5e^{x+y}$ was designed for, are typical equations that can be found in a physics scientific publication. Our main user case is that a user of CDS is interested in finding documents inside the collection whose equations are related to the one that is provided.

The evaluated queries were provided by physics experts working at CERN from their own research topics. The queries were provided in \LaTeX format as the users used them in their own publications. Table B.1 presents the set of queries that were evaluated, and a simple explanation about each one. Each of the users evaluated the results only for the equations he/she provided. For each of the presented result, the user ranked it as "Strictly Relevant", "Partially Relevant" and "Completely Irrelevant". As all evaluations including human interaction, this scoring is at some level subjective.

For each metric we compared three different scenarios:

- N : Only notational features over the original expression were used
- N+S : Notational and structural features over the original expression
- N+S+NN : Notational and structural features over the original expression and notational features over the normalized string using Mathematica's Simplify function
- N+S+NN (BM25): Same scenario as before but using BM25 similarity function instead of TF-IDF. This scoring function is implemented in Lucene through the class `BM25Similarity`

We also were interested in also using the Fully Simplify function in Mathematica but we found no difference between both forms and the results were the same.

We present results for three metrics: Partial precision (Document either partially relevant or completely relevant) strict precision and discounted cumulative gain ($\text{DCG}_p = \sum_{i=1}^p \frac{2^{rel_i} - 1}{\log_2(i+1)}$) which measures, in addition, the order of the results..

| Scenario | Top 5 | Top 10 | Top 20 |
|---------------|-------|--------|--------|
| N | 0.86 | 0.79 | 0.78 |
| N+S | 0.9 | 0.85 | 0.77 |
| N+S+NN | 0.92 | 0.94 | 0.88 |
| N+S+NN (BM25) | 0.92 | 0.94 | 0.87 |

Table 7.2: Loose Precision

| Scenario | Top 5 | Top 10 | Top 20 |
|---------------|-------|--------|--------|
| N | 0.76 | 0.61 | 0.50 |
| N+S | 0.84 | 0.70 | 0.56 |
| N+S+NN | 0.84 | 0.83 | 0.65 |
| N+S+NN (BM25) | 0.84 | 0.83 | 0.67 |

Table 7.3: Strict Precision

| Scenario | Top 5 | Top 10 | Top 20 |
|---------------|-------|--------|--------|
| N | 7.17 | 9.81 | 12.00 |
| N+S | 7.61 | 10.7 | 14.59 |
| N+S+NN | 7.66 | 11.81 | 16.32 |
| N+S+NN (BM25) | 7.66 | 11.82 | 16.42 |

Table 7.4: Discounted Cumulative Gain

Results show that for the selected queries, the system performs reasonably well achieving a precision of 0.84 in the top 5 results in the best scenario and 0.65 in the top 20 results. Even more important than the absolute values, it is more interesting to see how the addition of structural features and the addition of normalization to the expressions, allows the system to improve its performance reasonably. In the top 20 results there was an increase of 15% from the N scenario to the N+S+NN one. Also we see that in the case of the top 5 results going from N+S to N+S+NN resulted in the same strict precision of 0.84, it still represented and increase in the DCG metric which means still results were improved. A final remark is that there was not a significant difference between the two different similarity functions that were evaluated, and in at least one case, BM25 performed worse than standard tf-idf.

7.4 Indexing Performance

For this set of evaluations, we were only interested in observing how the performance of the indexing stage is affected by the different types of features that are employed. We selected the first 100 records from our dataset and measured the total time it required to index all the elements in Solr using a python script and the solrpy library to do the communication. The default similarity in Solr was used. We also measured the index size at the end of each experiment. Our four scenarios were organized as follows: We tested the same three scenarios as for the quality results and we added a fourth one N+S+FN where we used the fully normalized form of the expressions as opposed to the standard simplified form in the N+S+NN scenario.

Table 7.5 summarizes the results of this set of experiments. The first result is the significant increase in processing time when using Mathematica. Going from the N scenario to N+S represent an increase of almost 15x in the indexing time. Adding standard normalization increases the indexing time by a factor near to 1.44x and full normalization by a factor of 1.55x.

With respect to storage size, we see a small footprint (Less than 5%) in the total size by going from N to N+S. Adding normalization represents an increase by a factor of 1.78x and the size of the fully normalized index represents a smaller increase (1.76x)

We can observe that there is a small trade-off between time and storage size depending on which normalization mode is employed.

| Feature Set | Total time (sec) | Index Size (kb) | Avg. time per equation | Avg. size of equation (kb) |
|-------------|------------------|-----------------|------------------------|----------------------------|
| N | 58 | 13926 | 0.00524 | 1.237 |
| N+S | 884 | 14438 | 0.0785 | 1.282 |
| N+S+NN | 1273 | 25780 | 0.113 | 2.345 |
| N+S+FN | 1377 | 25548 | 0.122 | 2.269 |

Table 7.5: Indexing performance

7.5 Querying Performance

Similarly, we were also interested in evaluating how the query performance of the queries was affected with the different scenarios. For this set of queries we used the complete dataset, and evaluated the N and the N+S+NN scenarios for queries of different length.

| Expression Length (Number of characters in MathML) | N | N+S+NN |
|--|-------|--------|
| 70 | 0.178 | 0.327 |
| 198 | 0.179 | 0.810 |
| 452 | 0.323 | 1.692 |
| 771 | 0.393 | 2.066 |

Table 7.6: Query performance in seconds

Table 7.6 shows that for smaller queries, the querying time increases from an average of 0.178 seconds to 0.327 which represents an increase of 83%. For the largest evaluated set of queries, the querying time passed from 0.393 to 2.066 sec-

onds which represents an increase of around 425%. The absolute values, however, show that in the worst case the slower query took a little more than 2 seconds in average which is still acceptable for a specialized search engine.

7.6 Alternative Metrics

According to some experimental evaluations presented in [41], Tree Edit Distance provides an interesting alternative metric to the standard vector space model using angular distance. In this work, authors present an improvement in the quality of retrieved results compared to the classic approach of indexing the single leaves of the tree and a subset of the possible sub-trees. For evaluating the feasibility of this approach we compared a set of equations using both metrics: Tree edit distance (Using Zheng and Shasha’s algorithm implemented in python [42].) and a standard angular distance that was not heavily optimized. We compared elapsed times for expressions of different length.

| Expression Length (Number of leaf nodes) | Angular Distance Time (msec) | Tree Edit Distance Time(msec) |
|---|---|--|
| 3 | 0.01195 | 2.65 |
| 8 | 0.0142 | 9.97 |
| 7 | 0.0129 | 9.165 |
| 12 | 0.0194 | 15.15 |
| 18 | 0.0201 | 29.5 |
| 20 | 0.0245 | 27.75 |
| 22 | 0.0264 | 32.4 |
| 24 | 0.02502 | 32.1 |

Table 7.7: Distance computation using two different metrics

This result show that for shorter expressions, using tree edit distance represent an increase in time in more than 200x, and a factor of around 1200x for the longer expressions, making it impractical for implementing a system that needs to be responsive. Even though there are algorithms which provide a better asymptotic guarantees and better actual running times like RTED[43] with respect to Zheng and Shasha’s algorithm, the best improvements represent a decrease in processing time of 50% to 80% in trees with around 1000 nodes, while the standard tested equations which are representative of our dataset, contain less than 100 nodes.

Chapter 8

Conclusions

The area of Mathematical Information Retrieval (MIR) is still very young and different systems, algorithms and datasets are being developed. Our work, $5e^{x+y}$, proposes a complete MIR system for a concrete digital library, such as the CERN Document Server based on the Invenio digital library software platform. In terms of our proposed goals we can conclude that:

- We explored and identified different tools and approaches to automatically detect and extract mathematical content from digital files. We concluded that LaTeXML and InftyReader are currently the most suitable set of tools to process Latex and PDF files respectively. We automatized the extraction process and incorporated it into the Invenio workflow.
- Our exploration of the extraction tools, lead us to identify the MathML format as the best format to store and process mathematical content. It is the recommended language for web documents, and the best extraction tools agree on it.
- Once a mathematical expression is extracted and parsed to a suitable format the next step is to identify suitable features to index. We separated our features into notational and structural ones. For the first category, we discovered some issues that have not yet been addressed in previous work and developed heuristics that allow to match more cases by normalizing unicode characters, grouping similar operators and expanding different numeric quantities. For the structural features, we explored the integration of a Computer Algebra System to perform pattern matching and simplification of the given expression.
- We explored different IR frameworks and models to store and retrieve the indexed expressions and found that the standard vector space model is suitable for our needs.

- We implemented our system in top of the Lucene/Solr framework. We took advantage of the modular design of Lucene to structure our work in independent classes where each one has a defined role. This structure also allows to easily implement and incorporate further heuristics to handle more cases. Our system is deployed as a Solr plugin making it very easy to be used from an external application. As a final component in our system, we developed the glue code to integrate our search engine with Invenio and developed a simple web interface.
- We performed different sets of evaluation on our system to explore the quality of the retrieved results and the performance of our system. Our quality evaluations showed that a combination of notational and structural features allow us to achieve significantly better results. Our performance results showed that the extraction of structural features through the Mathematica CAS imposes a heavy overhead of more than 10x on the indexing time, but as discussed previously, the extra work is worth investing in, given the improved set of results. The extra overhead at querying time is not insignificant but it is still acceptable for today's web usage standards. Our evaluation discarded the tree edit distance because the running time was 200 times more than the time needed by a simple angular distance used in standard vector space model implementations, making this approach impractical for a big system.
- Finally, during this work we identified areas of improvement which are worth looking into. They are presented in the next chapter.

Chapter 9

Future work

Several different work directions were identified during the development of this work.

9.1 Features extraction

The MathML specification include hundreds of different tags, which were not addressed in our work. They can affect drastically the semantics of a mathematical expression. A simple example is the fact that sometimes to stress that a variable represents a vector, it is standard to use a bold font. The same can be signified by an arrow on top of a variable. Dozens and even hundreds of these notations contain important information which is not taken into account by our system, since each one of them requires explicit treatment, therefore handling all of the cases is impractical. Also, there is the limitation of the extracting tools to handle complex constructs like matrices or embedded text.

9.2 Data Analysis

Different algorithms and techniques have been developed to overcome some of the shortcomings of the vector space model. Latent semantic indexing together with singular value decomposition, allow to automatically identify synonyms which can improve the recall of an IR system. During our work, we tested different projects such as Semantic Vectors[44] and Mallet[45] and Mahout. However, these tools are not yet ready for a big deployment and a considerable amount of time and effort was spent on trying to integrate the formats from the Lucene index to the input format of each tool, as well as debugging and fixing code of some of the tools to make them work. Several times we encountered errors due to the in-memory implementations of the algorithms. Further work on integrating the results of these techniques to the workflow of the system could provide better results.

In our work, we explored the clustering of the expressions and managed to find a reasonable number of clusters for k-means. However more insight from this should be developed and the information gained with the clustering step should be integrated in the workflow to improve both the results and the querying time.

9.3 End to end system

In our implementation, we explored and incorporated `LATEX`extraction to our system. For a full integration on CDS, PDF extraction becomes an important need since most of the records do not include the source files.

Bibliography

- [1] About cern, 2014. <http://home.web.cern.ch/about>.
- [2] The original proposal of the www, htmlized, 2014. <http://www.w3.org/History/1989/proposal.html>.
- [3] Cern document server, 2014. <http://cds.cern.ch/>.
- [4] Maja Gracco Jean Yves Le Meur Nicholas Robinson Tibor Simko Alberto Pepe, Thomas Baron and Martin Vesely. Cern document server software: the integrated digital library, 2005. <http://cds.cern.ch/record/1198695/files/CERN-THESIS-2009-057.pdf>.
- [5] Invenio, 2014. <http://invenio-software.org/>.
- [6] Modules overview, 2014. <http://invenio-demo.cern.ch/help/hacking/modules-overview>.
- [7] arxiv.org e-print archive, 2014. <http://arxiv.org>.
- [8] S. K. M. Wong and Vijay V. Raghavan. Vector space model of information retrieval: A reevaluation. In *Proceedings of the 7th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '84, pages 167–185, Swinton, UK, UK, 1984. British Computer Society.
- [9] Stephen E. Robertson and Karen Sparck Jones. Document retrieval systems. chapter Relevance Weighting of Search Terms, pages 143–160. Taylor Graham Publishing, London, UK, UK, 1988.
- [10] Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, April 2009.
- [11] Kaizhong Zhang and Dennis Shasha. Simple Fast Algorithms for the Editing Distance Between Trees and Related Problems. *Siam Journal on Computing*, 18:1245–1262, 1989.
- [12] Philip N. Klein. Computing the edit-distance between unrooted ordered trees. In *Proceedings of the 6th Annual European Symposium on Algorithms*, ESA '98, pages 91–102, London, UK, UK, 1998. Springer-Verlag.

-
- [13] Erik D. Demaine, Shay Mozes, Benjamin Rossman, and Oren Weimann. An Optimal Decomposition Algorithm for Tree Edit Distance. *ACM Transactions on Algorithms*, 6:146–157, 2007.
 - [14] Bruce Miller. Latexml a latex to xml converter. <http://dlmf.nist.gov/LaTeXML/>.
 - [15] Tex4ht: Latex and tex for hypertext. <https://www.tug.org/applications/tex4ht/mn.html>.
 - [16] The University of Edinburgh School of Physics and Astronomy. Snuggletex (1.2.2). <http://www2.ph.ed.ac.uk/snuggletex/documentation/overview-and-features.html>.
 - [17] A brief description of latexmathml, 2007. <http://math.etsu.edu/LaTeXMathML/>.
 - [18] Josef B. Baker, Alan P. Sexton, and Volker Sorge. Maxtract: Converting pdf to latex, mathml and text. In *Proceedings of the 11th International Conference on Intelligent Computer Mathematics*, CICM'12, pages 422–426, Berlin, Heidelberg, 2012. Springer-Verlag.
 - [19] Maxtract - scientific document analysis group, 2012. <http://www.cs.bham.ac.uk/research/groupings/reasoning/sdag/maxtract.php>.
 - [20] Masakazu Suzuki, Fumikazu Tamari, Ryoji Fukuda, Seiichi Uchida, and Toshihiro Kanahori. Infty: An integrated ocr system for mathematical documents. In *Proceedings of the 2003 ACM Symposium on Document Engineering*, DocEng '03, pages 95–104, New York, NY, USA, 2003. ACM.
 - [21] Inftyreader, 2013. <http://www.sciaccess.net/en/InftyReader/>.
 - [22] Heinrich Stamerjohanns; Deyan Ginev; Catalin David; Dimitar Misev; Vladimir Zamdzhiev; Michael Kohlhase. (mathml-aware article conversion from latex — a comparison study). 51:7–28, 2008.
 - [23] Mir 2012 workshop – mathematics information retrieval, july 8th, 2012, 2012. <http://cicm-conference.org/2012/cicm.php?event=mir>.
 - [24] National institute of informatics, 2014. <http://www.nii.ac.jp/en/>.
 - [25] Ntcir pilot task: Math task, 2013. <http://ntcir-math.nii.ac.jp/ntcir10-math/>.
 - [26] Petr Sojka and Martin Liška. Indexing and searching mathematics in digital libraries. In JamesH. Davenport, WilliamM. Farmer, Josef Urban, and Florian Rabe, editors, *Intelligent Computer Mathematics*, volume 6824 of *Lecture Notes in Computer Science*, pages 228–243. Springer Berlin Heidelberg, 2011.

-
- [27] Mrec — mathematical retrieval collection, 2012. <https://mir.fi.muni.cz/MREC/index.html>.
 - [28] Jozef Mišutka and Leo Galamboš. Extending full text search engine for mathematical content. *Towards Digital Mathematics Library*, pages 55–67, 2008.
 - [29] Jozef Mišutka and Leo Galamboš. System description: Egomath2 as a tool for mathematical searching on wikipedia.org. In *Proceedings of the 18th Calculemus and 10th International Conference on Intelligent Computer Mathematics*, MKM’11, pages 307–309, Berlin, Heidelberg, 2011. Springer-Verlag.
 - [30] NIST Digital Library of Mathematical Functions. <http://dlmf.nist.gov/>, Release 1.0.6 of 2013-05-06. Online companion to NIST Handbook of Mathematical Functions.
 - [31] F. W. J. Olver, D. W. Lozier, R. F. Boisvert, and C. W. Clark, editors. *NIST Handbook of Mathematical Functions*. Cambridge University Press, New York, NY, 2010.
 - [32] BruceR. Miller. Three years of dlmf: Web, math and search. In Jacques Carette, David Aspinall, Christoph Lange, Petr Sojka, and Wolfgang Windsteiger, editors, *Intelligent Computer Mathematics*, volume 7961 of *Lecture Notes in Computer Science*, pages 288–295. Springer Berlin Heidelberg, 2013.
 - [33] Michael Kohlhase, BogdanA. Matican, and Corneliu-Claudiu Prodescu. Mathwebsearch 0.5: Scaling an open formula search engine. In Johan Jeuring, JohnA. Campbell, Jacques Carette, Gabriel Reis, Petr Sojka, Makarius Wenzel, and Volker Sorge, editors, *Intelligent Computer Mathematics*, volume 7362 of *Lecture Notes in Computer Science*, pages 342–357. Springer Berlin Heidelberg, 2012.
 - [34] Peter Graf. Substitution tree indexing. In Jieh Hsiang, editor, *Rewriting Techniques and Applications*, volume 914 of *Lecture Notes in Computer Science*, pages 117–131. Springer Berlin Heidelberg, 1995.
 - [35] zbmath - the first resource for mathematics, 2014. <http://zbmath.org/about/>.
 - [36] Xuan Hu, Liangcai Gao, Xiaoyan Lin, Zhi Tang, Xiaofan Lin, and Josef B. Baker. Wikimirs: A mathematical information retrieval system for wikipedia. In *Proceedings of the 13th ACM/IEEE-CS Joint Conference on Digital Libraries*, JCDL ’13, pages 11–20, New York, NY, USA, 2013. ACM.
 - [37] Wolfram|alpha: Computational knowledge, 2014. <https://www.wolframalpha.com/>.
 - [38] Xah Lee. Unicode: Math symbols, 2010. <http://web.archive.org/web/20131223003502/>.

-
- [39] Online latex editor - create, integrate and download, 2014. <http://www.codecogs.com/latex/eqneditor.php>.
 - [40] Apache mahout: Scalable machine learning and data mining, 2014. <https://mahout.apache.org/>.
 - [41] Shahab Kamali and Frank Wm. Tompa. Retrieving documents with mathematical content. In *Proceedings of the 36th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '13, pages 353–362, New York, NY, USA, 2013. ACM.
 - [42] Tree edit distance using the zhang shasha algorithm, 2013. <https://github.com/timtadh/zhang-shasha>.
 - [43] Mateusz Pawlik and Nikolaus Augsten. Rted: A robust algorithm for the tree edit distance. *Proc. VLDB Endow.*, 5(4):334–345, December 2011.
 - [44] Package for creating and searching semantic vector indexes by wrapping apache lucene., 2014. <https://code.google.com/p/semanticvectors/>.
 - [45] Andrew Kachites. McCallum. Mallet: Machine learning for language toolkit, 2002. <http://mallet.cs.umass.edu/>.

Appendix A

Configuration files

A.1 Solr schema

The Solr schema file `schema.xml` indicates which are the fields that are going to be used in this instance, and what type of analyzer is going to be used for each one. One can also define field types in case different fields want to have the same behaviour but conceptually having different meanings.

The following snippet add the fields declarations:

```
<!-- Arthur Oviedo A0-->
<!-- Fields declaration using the field types declared below-->
<field name="math_notational_field"
  type="math_notational_type"/>
<field name="math_structural_field"
  type="math_structural_type"/>
<field name="math_normalized_notational_field"
  type="math_normalized_notational_type"/>
<field name="filename"
  type="string" indexed="true" stored="true"/>
<field name="number_occurrences"
  type="int" indexed="true" stored="true"/>
```

And the following one, defines the types and the configurations (Including the java class implementing the analyzer) for each one:

```
<fieldType name="math_normalized_notational_type"
class="solr.TextField"
positionIncrementGap="1"
indexed="true"
stored="true"
termVectors="true"
```

```

termPositions="true" >
<analyzer class="cern.ch.mathexplorer.lucene.analysis.analyzers.
SolrNormalizerNotationalAnalyzer"/>
</fieldType>

<!-- Field type for MathMl info
(Tokens that are extracted from the MathML representation)
over the normalized/simplified string using Full simplification-->
<fieldType name="math_full_normalized_notational_type"
class="solr.TextField"
positionIncrementGap="1"
indexed="true"
stored="true"
termVectors="true"
termPositions="true" >
<analyzer class="cern.ch.mathexplorer.lucene.analysis.analyzers.
SolrFullNormalizerNotationalAnalyzer"/>
</fieldType>

<!-- Field type for structural features extracted from the equation-->
<fieldType name="math_structural_type"
class="solr.TextField"
positionIncrementGap="1"
indexed="true"
stored="true"
termVectors="true"
termPositions="true" >
<analyzer class="cern.ch.mathexplorer.lucene.analysis.analyzers.
SolrStructuralAnalyzer"/>
</fieldType>

```

A.2 Solr configuration file

The `solrconfig.xml` file specifies additional configuration parameters of the instance, including what handlers to use for each different type of request. One specific type of handler is `queryParser`, and for using our specific query parser we include the following line of code:

```

<queryParser name="mathqueryparser"
class="cern.ch.mathexplorer.lucene.query.MathQueryParserPlugin" />

```

Appendix B

Test equations

| Id | Description | MathML code | Visual Form |
|-----|--|--|--|
| Eq1 | Time-independent Schrödinger equation | $\langle \text{mfenced close="" } " \text{ open="" } \langle \text{mrow} \rangle \langle \text{mo} \rangle \cdot \langle \text{mfrac} \rangle \langle \text{msup} \rangle \langle \text{mi} \rangle \langle \text{mn} \rangle 2 \langle \text{mn} \rangle \langle \text{msup} \rangle \langle \text{mrow} \rangle \langle \text{mn} \rangle 2 \langle \text{mn} \rangle \langle \text{msub} \rangle \langle \text{mi} \rangle m \langle \text{mi} \rangle e \langle \text{msub} \rangle \langle \text{mrow} \rangle \langle \text{mfrac} \rangle \langle \text{msup} \rangle \langle \text{mo} \rangle \nabla \langle \text{mo} \rangle \langle \text{mn} \rangle 2 \langle \text{mn} \rangle \langle \text{msup} \rangle \langle \text{mo} \rangle + \langle \text{mo} \rangle \langle \text{mi} \rangle V \langle \text{mi} \rangle \langle \text{mfenced close="" } " \text{ open="" } \langle \text{mi} \rangle r \langle \text{mi} \rangle \langle \text{mfenced} \rangle \langle \text{mrow} \rangle \langle \text{mfenced} \rangle \langle \text{mi} \rangle \langle \text{mfenced close="" } " \text{ open="" } \langle \text{mi} \rangle r \langle \text{mi} \rangle \langle \text{mfenced} \rangle \langle \text{mo} \rangle = \langle \text{mo} \rangle \langle \text{mi} \rangle E \langle \text{mi} \rangle \langle \text{mfenced close="" } " \text{ open="" } \langle \text{mi} \rangle r \langle \text{mi} \rangle \langle \text{mfenced} \rangle \langle \text{math} \rangle$ | $\left(-\frac{\hbar^2}{2m_e}\nabla^2 + V(\mathbf{r})\right)\psi(\mathbf{r}) = E\psi(\mathbf{r})$ |
| Eq2 | Atomic mass in terms of number of proton and number of neutrons | $\langle \text{math} \rangle \langle \text{mi} \rangle A \langle \text{mi} \rangle \langle \text{mo} \rangle = \langle \text{mo} \rangle \langle \text{mi} \rangle Z \langle \text{mi} \rangle \langle \text{mo} \rangle + \langle \text{mo} \rangle \langle \text{mi} \rangle N \langle \text{mi} \rangle \langle \text{math} \rangle$ | $A = Z + N$ |
| Eq3 | Radiation flux in term of the initial intensity, the absorption coefficient and the thickness of a substance | $\langle \text{math} \rangle \langle \text{mi} \rangle I \langle \text{mi} \rangle \langle \text{mo} \rangle = \langle \text{mo} \rangle \langle \text{msub} \rangle \langle \text{mi} \rangle I \langle \text{mi} \rangle \langle \text{mn} \rangle 0 \langle \text{mn} \rangle \langle \text{msub} \rangle \langle \text{msup} \rangle \langle \text{mi} \rangle e \langle \text{mi} \rangle \langle \text{mrow} \rangle \langle \text{mo} \rangle \cdot \langle \text{mo} \rangle \langle \text{mi} \rangle u \langle \text{mi} \rangle \langle \text{mi} \rangle x \langle \text{mi} \rangle \langle \text{mrow} \rangle \langle \text{msup} \rangle \langle \text{math} \rangle$ | $I = I_0 e^{-\mu x}$ |
| Eq4 | Higgs Boson decaying process into W^+ and W^- bosons. | $\langle \text{math} \rangle \langle \text{mi} \rangle H \langle \text{mi} \rangle \langle \text{mo} \rangle \rightarrow \langle \text{mo} \rangle \langle \text{msup} \rangle \langle \text{mi} \rangle W \langle \text{mi} \rangle \langle \text{mo} \rangle + \langle \text{mo} \rangle \langle \text{msup} \rangle \langle \text{mi} \rangle W \langle \text{mi} \rangle \langle \text{mo} \rangle \cdot \langle \text{mo} \rangle \langle \text{msup} \rangle \langle \text{math} \rangle$ | $H \rightarrow W^+ W^-$ |

| | | | |
|-----|--|---|--|
| Eq5 | Einstein Field equation | $G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2}Rg_{\mu\nu} = \frac{8\pi G}{c^4}T_{\mu\nu}$ | |
| Eq6 | Luminosity of a star in terms of its temperature and radius | $L = 4\pi R^2 \sigma T_{eff}^4$ | |
| Eq7 | Free-cyclotron motion of an electron in terms of its charge, mass and the magnetic field | $\omega_c = q\vec{B} /m$ | |
| Eq8 | Spin energy of a particle in presence of an external magnetic field | $E = -\vec{\mu} \cdot \vec{B}$ | |
| Eq9 | Quantum electrodynamics Lagrangian modelling the interaction of the -1/2 spin field with the electromagnetic field | $\mathcal{L}_{QED} = \bar{\psi}_i i ((\gamma^\mu D_\mu)_{ij} - m \delta_{ij}) \psi_j - \frac{1}{4} F_{\mu\nu} F^{\mu\nu}$ | |

| | | |
|------|----------------------------|---|
| Eq10 | Einstein-Hilbert action | $S = \frac{1}{2\kappa} \int R \sqrt{-g} d^4x$ |
|------|----------------------------|---|

Table B.1: Sample equations used during testing

Appendix C

Detailed system architecture diagrams

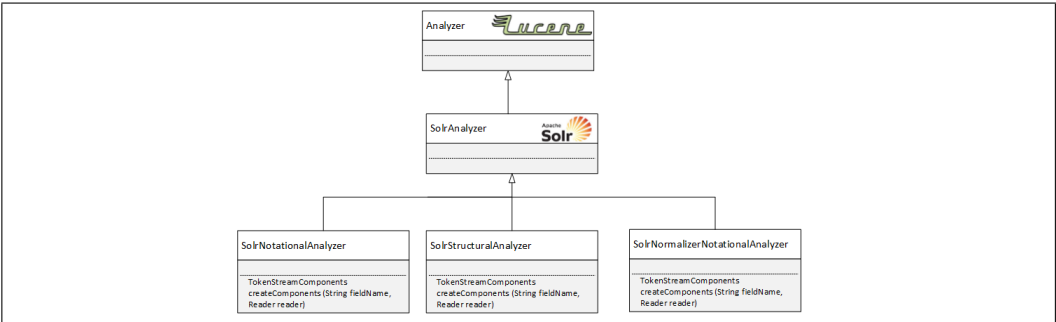


Figure C.1: Overview of `cern.ch.mathexplorer.lucene.analysis.analyzers` package

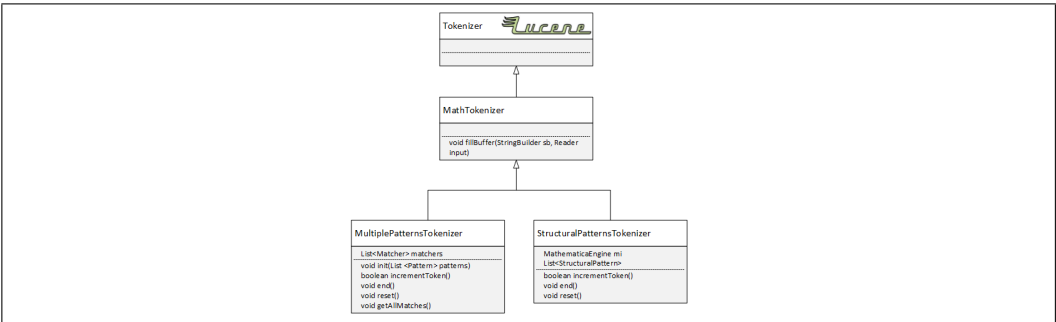


Figure C.2: Overview of `cern.ch.mathexplorer.lucene.analysis.tokenizers` package

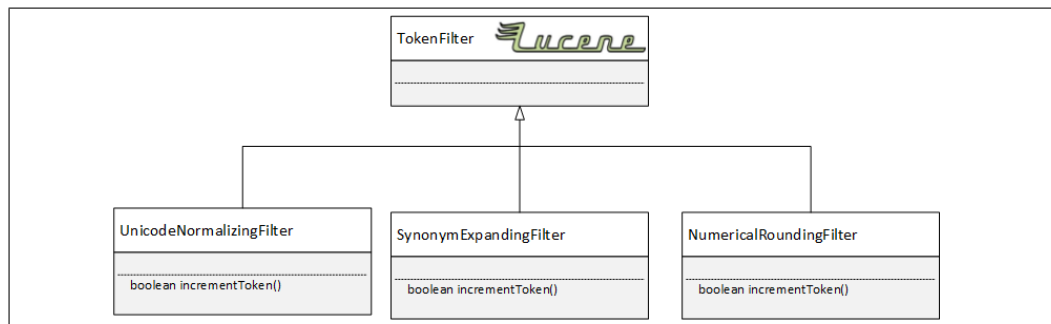


Figure C.3: Overview of `cern.ch.mathexplorer.lucene.analysis.filters` package

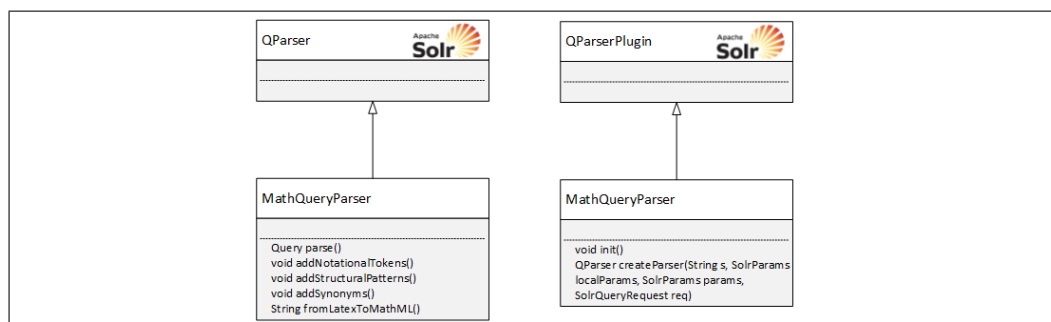


Figure C.4: Overview of `cern.ch.mathexplorer.lucene.query` package

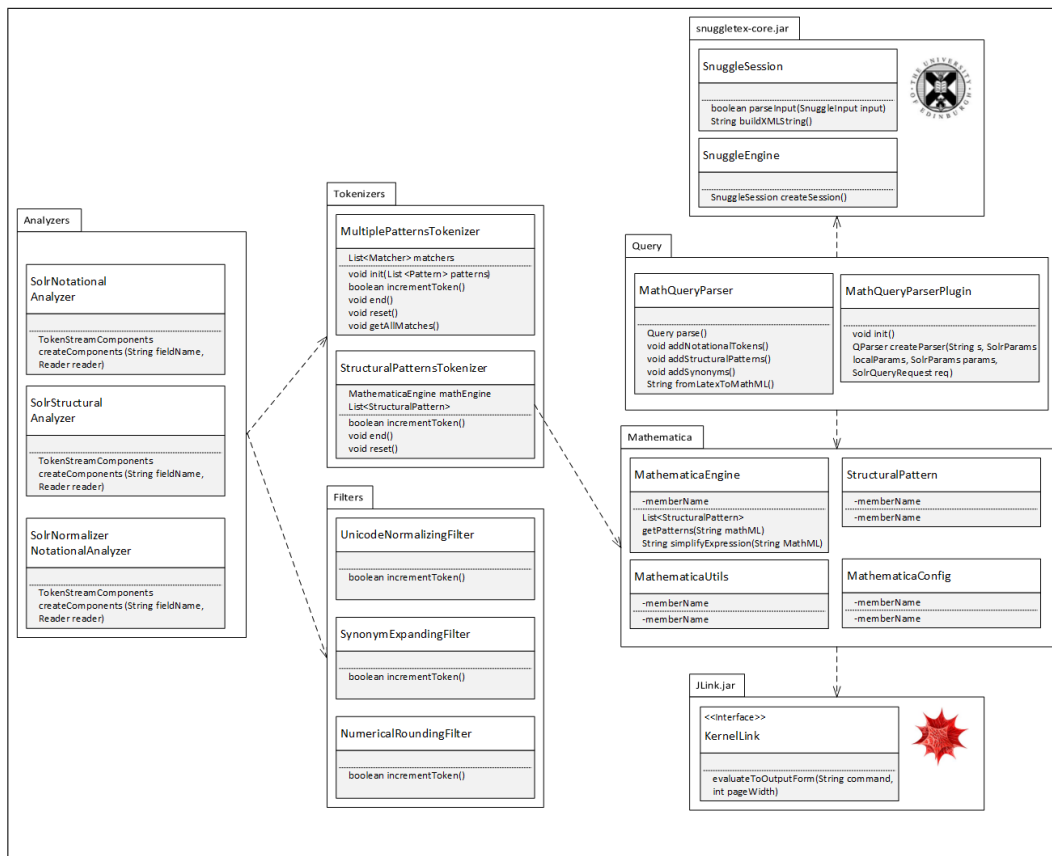


Figure C.5: Overview of the relations between the main packages