

MASARYK UNIVERSITY
FACULTY OF INFORMATICS



Normative Textual Representation of Mathematical Formulae

MASTER'S THESIS

Maroš Kucbel

Brno, 2013

Declaration

Hereby I declare, that this paper is my original authorial work, which I have worked out by my own. All sources, references and literature used or excerpted during elaboration of this work are properly cited and listed in complete reference to the due source.

Maroš Kucbel

Advisor: doc. RNDr. Petr Sojka, Ph.D.

Acknowledgement

@todo Thanks!

Abstract

@todo Abstract

Keywords

MathML, StAX @todo more

Contents

1	Introduction	3
2	Mathematical Markup Language - MathML	4
2.1	<i>Structure of MathML</i>	4
2.1.1	Presentation MathML	5
2.1.2	Content MathML	6
2.2	<i>Creating MathML Markup</i>	8
2.2.1	L ^A T _E X _{ML}	9
2.2.2	Tralics	9
2.2.3	MaxTract	10
2.2.4	INFTY Reader	10
2.2.5	T _E X4ht	10
2.3	<i>Canonicalization of MathML</i>	11
3	Status Quo of Current Tools	12
3.1	DAISY	12
3.2	MathPlayer	13
3.3	AsTeR	13
3.4	MathTalk	14
3.5	Lambda	14
3.6	Web Browser Support	14
4	Analysis and Solution Proposal	16
4.1	<i>Load Document</i>	18
4.2	<i>Create Tree</i>	19
4.3	<i>Convert Tree</i>	21
4.3.1	Presentation MathML	21
4.3.2	Content MathML	23
4.3.3	Operators and Symbols	25
5	Implementation	26
5.1	<i>DOM-like Representation of MathML</i>	27
5.2	<i>Input Processing</i>	27
5.2.1	XML Parser	27
5.3	<i>Conversion</i>	30
6	Results	33

6.1	<i>Mathematical Retrieval Collection MREC</i>	33
6.2	<i>Changing StAX Implementation</i>	33
7	Conclusion	34

Chapter 1

Introduction

Chapter 2

Mathematical Markup Language - MathML

World Wide Web pages and their main publishing language, HTML¹, provide many ways to present desired information to the user. However, presenting mathematics is not so easy and straightforward. There aren't any special tags in the HTML specification for including mathematics. In most cases mathematical equations are presented in the form of an image. On one hand this approach renders the equation the same way in every web browser, on the other hand there is no way to copy the equation for further use, not to mention editing the equation.

MathML² was specifically created to circumvent this obstacle. It was therefore designed to be used in web pages alongside HTML. It follows that MathML is an application of XML with special set of tags used to capture the content, structure and even presentation details of mathematical equations. We will take closer look at how this is achieved in the following sections. In April 1998 MathML became the recommendation of the W3C working group³ for including mathematics into web pages.

2.1 Structure of MathML

MathML as an application of XML consists of a tree of nodes. Each node is either empty, has textual value or has a list of descendant nodes. To provide additional information, a node can have an arbitrary amount of attributes, key – value pairs. Each MathML tree has to have a root node named `math` that belongs to the MathML namespace. Also it is important not to forget to include XML declaration and MathML doctype declaration.

MathML comes two distinct sets of tags. For the visual form of equations there are Presentation MathML tags, Content MathML tags focus on the

1. HyperText Markup Language

2. <http://www.w3.org/Math/>

3. The World Wide Web Consortium (W3C) is the main international standards organization for the World Wide Web

semantics and meaning of formulas. Presentation tags are being used primarily by web browsers to display mathematical expressions to the users, Content tags are important for further processing of expressions by specialized mathematical programs.

2.1.1 Presentation MathML

The main focus of the Presentation MathML is displaying the equation. For this purpose there are around thirty elements, all starting with the prefix “m”. Elements are divided into two classes, first of which is called tokens. It consist of elements that represent individual content and do not contain other nested elements. These include:

- `<mi>x</mi>` - identifiers,
- `<mn>2</mn>` - numbers,
- `<mo>+</mo>` - operators, fences, separators,
- `<mtext>free text</mtext>` - text.

The content of the token can of course be expressed in more than one character (`<mo>sin</mo>`) or with an XML and HTML character entities, for example `>` and `>` have the same meaning as `>`, greater than. It is completely up to the user which notation he will choose.

Even though nesting other elements in tokens is not allowed, there are exceptions. For example HTML5 allows almost any HTML inline tags inside the `mtext` element. So `<mtext>free text</mtext>` would be rendered with the bold word free. The other class of elements is layout schemata. This collection of elements is further divided into following groups:

- General Layout
 - `<mrow>` - general horizontal grouping,
 - `<mfrac>` - fractions and binomial numbers,
 - `<msqrt>`, `<mroot>` - radicals
- Script and Limit - superscripts, subscripts,
- Tabular Math - tables and matrices,
- Elementary Math - notation for lower grades mathematics.

We can thought of the layout schemata as a form of expression constructors, that specify the way in which sub-expression are constructed into larger expressions, starting with tokens and ending with the math element. Therefore elements belonging to layout schemata class do not contain any characters, only other nested elements (layout schemata or tokens).

```
<math>
  <mrow>
    <mi>e</mi>
    <mo>=</mo>
    <mi>m</mi>
    <mo>*</mo>
    <msup>
      <mi>c</mi>
      <mn>2</mn>
    </msup>
  </mrow>
</math>
```

Figure 2.1: Expression $e = m \cdot c^2$ in Presentation MathML

Presentation MathML also provides over fifty attributes for fine tuning of expressions, like setting colors, dimensions, alignments and many others.

2.1.2 Content MathML

Mathematical notation, rigorous as it is, is still not standardized around the world and there are many different ways of writing mathematical expressions based on cultural customs. Even simple multiplication operation can be written as $x * y$, xy , $x \text{ times } y$, $x \times y$. In many situations there's no need to actually render the expression, only the underlying meaning is important. For this reason Content MathML provides a framework and a markup language that capture the semantics of mathematical expressions.

The structure of Content MathML tree is based on a very simple principle – applying an operator to sub-expressions. For example the quotient x/y can be thought of as applying the division operator to two arguments x and y . It follows that the cornerstone of Content MathML is the function application element `<apply>`. The first child of the `<apply>` element signifies the operator, which can be an `<apply>` element again, and the rest of the `<apply>` element direct descendants are arguments to which the oper-

ator is applied. Token elements `<ci>` and `<cn>` are available to represent numbers and variables respectively.

Content MathML provides two ways of defining the operator. The first one is a set of more than 100 elements, each corresponding to some mathematical operator. For example for the addition operator there is `<plus>`, for logical xor `<xor>` and so on. Then there is the second way, the element `<csymbol>`. This element contains a textual representation of the operator, that is bound to a definition in a content dictionary referenced by either the attribute `cd` or `definitionURL`. External content dictionaries are important for communication between agents and there exist several public repositories of mathematical definitions. Most notable is the OpenMath Society repository.

```
<math>
  <apply>
    <eq/>
    <ci>e</ci>
    <apply>
      <times/>
      <ci>m</ci>
      <apply>
        <power/>
        <ci>c</ci>
        <cn>2</cn>
      </apply>
    </apply>
  </apply>
</math>
```

Figure 2.2: Expression $e = m \cdot c^2$ in Content MathML

In MathML 3, a subset of Content MathML elements is defined: Strict Content MathML. This uses only minimal, but sufficient, amount of elements, most importantly allows only the usage of `<csymbol>` element for defining operators. Strict Content MathML is designed to be compatible with OpenMath standard.

```
<math>
  <apply>
    <csymbol cd="dict">equals</csymbol>
    <ci type="real">e</ci>
    <apply>
      <csymbol cd="dict">times</csymbol>
      <ci type="real">m</ci>
      <apply>
        <csymbol cd="dict">power</csymbol>
        <ci type="real">c</ci>
        <cn type="integer">2</cn>
      </apply>
    </apply>
  </apply>
</math>
```

Figure 2.3: Expression $e = m \cdot c^2$ in Strict Content MathML

2.2 Creating MathML Markup

Creating MathML documents is very simple. All that is needed is a word processor available on every operating system. However, this approach is prone to errors, be it the wrong letter case or unclosed tags. Since MathML is an XML language, the usage of some sophisticated word processor that supports tags highlighting, code completion and XML validation will greatly improve the efficiency of creating MathML documents. But the whole process is still very time consuming and impractical. As is the case with most XML documents, even writing simple mathematical expressions takes up a lot of space and time. Fortunately there are many dedicated MathML editors, that provide some degree of abstraction from the actual MathML markup.

Multiple software products designated for working with mathematical expressions provide an option to output the results of calculations in MathML format. These include popular web service Wolfram Alpha, Mathematica, Maple or Matlab. MathML can also be used as a data exchange format or input format for aforementioned programs.

Another way of creating MathML documents is a conversion from different formats. Among scientist, mathematicians particularly, \TeX is the format of choice. It then comes as no surprise that there are many sources of

mathematical texts written in \TeX . However, not every publication comes with the source files. Often only print-ready PDF files are released [3]. Also many academic writings, especially older ones, are available only in printed form. These have to be scanned using an Optical Character Recognition (OCR) software.

2.2.1 \LaTeX XML

The lack of a suitable tool for converting \LaTeX to XML prompted the participants of the Digital Library of Mathematical Functions project to develop their own solution. Main goals of \LaTeX XML [11] design contain the aspiration to faithfully emulate \TeX behavior, the ease of extensibility and the preservation of both semantic and presentation aspects of original documents. To this end \LaTeX XML provides two main commands, `latexml` and `latexmlpost`. `latexml` converts the initial \TeX document to XML based on a set of \LaTeX XML-bindings files, that define the mapping of \TeX macros to XML. `latexmlpost` then processes resulting XML output by converting mathematics and graphics, cross-referencing and applying an XSLT stylesheet.

Both commands come with a multitude of parameters that allow users to customize the whole process, such as loading user-defined bindings or choosing the output format. Based on the requested output format, mathematics is converted to graphics (png images for HTML) or Presentation MathML in case of XHTML or HTML5.

\LaTeX XML is freely available online as an installation package for Linux systems as well as Windows and MacOS. An extensive and detailed manual is available online or as a PDF document.

2.2.2 Tralics

Tralics [6][7] is a freeware software designed to translate \LaTeX sources into XML documents, that can be further converted into either PDF or HTML. The generated XML is conforming to the local ad-hoc DTD (a simplification of the TEI) with mathematical formulas conforming to the Presentation MathML 2.0 recommendations.

Similarly to \LaTeX XML, Tralics provides many ways for customization of resulting documents, among them the possibility to change element and attribute names in the XML file. Besides Presentation MathML, mathematical formulae can be translated to \LaTeX -like elements as well.

Tralics is readily available online in the form of source files or binaries

for Linux, Windows and MacOS operating systems. An extensive documentation regarding customization and usage is also available online.

2.2.3 MaxTract

MaxTract [4] is a tool that through spacial analysis of symbols and fonts in PDF document reverse-engineers its source files in the form of \LaTeX or XHTML + MathML documents. The conversion process requires valid PDF files to work correctly as it needs the information about symbols, font encoding and width of objects contained in the PDF file. PDF documents created via \LaTeX fulfill these requirements, therefore MaxTract is able to process most of the scientific and mathematical material.

2.2.4 INFTY Reader

INFTY [18] is an Optical Character Recognition (OCR) software especially created for mathematical documents. INFTY reads scanned pages (images) and yields their character recognition results in multiple formats, including \LaTeX and MathML. It does so in four steps: layout analysis, character recognition, structure analysis of mathematical expressions, and manual error correction (optional).

The most important phase, character recognition, is responsible for distinguishing mathematical expressions and running a character recognition engine originally developed for mathematical symbols together with non-specialized commercial OCR engine.

INFTY Reader is available free of charge for limited amount of time, then it is necessary to purchase a license key.

2.2.5 \TeX 4ht

\TeX 4ht [8] is a system for producing structured output in a markup language from sources written in the \TeX -based family of languages. It is highly extensible and configurable, most common configurations include the conversion from \LaTeX to HTML with MathML, braille or DocBook targets.

The conversion is invoked by running the `htlatex` command. To the user the process seems similar to producing standard DVI or PDF outputs. \TeX 4ht uses hooks within the \LaTeX constructs and associates configurations to them. By modifying default configuration files, the user can change the resulting HTML document to his liking. \TeX 4ht comes with the support of

Cascading style sheets⁴, which provides further possibilities of customizing the output files.

2.3 Canonicalization of MathML

With a wide variety of options of creating MathML documents it is necessary to modify and unify input documents in such a way that will decrease the number of possible ambiguities (especially in the Presentation MathML), in the best case completely eliminating all ambiguities.

The team standing behind the Universal Maths Conversion library (UMCL) [1] proposed a notion of a unified structure - Canonical MathML [2]. Canonical MathML is valid Presentation MathML and recommends a set of rules that should be followed to make Presentation MathML documents unambiguous. A correct use of `msup` (`msub`) for superscripts (subscripts); parentheses defined in the `mo` element (`mfenced` element is not used); or unified way of writing summations, integrals and products to name a few.

UMCL incorporates this notion and provides a module for converting Presentation MathML markup into Canonical MathML. The canonicalization module employs the use of XSL stylesheets and transformations ?? to create the desired Canonical MathML. However XSL Transformations tend to be quite slow and the UMCL canonicalization is rather error-prone and sometimes even changes the semantics of mathematical formulae, as was shown in M. Jarmar's thesis [9, chapter 5].

These shortcomings of UMCL prompted a team at Masaryk University to design and implement their own canonicalization solution [5] for use in the (Web)MIaS project [14]. Although the implementation of this solution is still under development, the core functionality is already working and provides better results than UMCL (from the point of performance as well as the structure of converted documents).

4. <http://www.w3schools.com/css/>

Chapter 3

Status Quo of Current Tools

Making mathematical content available to visually impaired users or those with dyslexia has been a focus of many projects and researches over the last two decades. As a result several products (commercial or open source) have been developed that help the users with working with mathematical formulae; by providing text-to-speech solutions, easy to use and comfortable editors, support for browsing and searching documents that contain mathematical equations.

In this chapter we will present some solutions that are unique in this field and had impact on another similar products, either by the theoretical research, that fuels such solutions or by general aspects that make them stand out.

3.1 DAISY

DAISY (Digital Accessible Information System) is a technical standard for digital audio books, periodicals and computerized text. It is specifically intended for people that have problems reading a printed text including blindness, dyslexia or impaired vision. DAISY is based on combination of XML and MP3 and has various functions that traditional audio books do not provide. These include line by line navigation, searching in the text, adding bookmarks or adjusting the speed of the speech. It also provides support for embedded content such as images, graphics and MathML.

DAISY books (books conforming to the DAISY standard) can be listened to on designated DAISY readers, computers with installed DAISY software, mobile phones or MP3 players. There are many implementation of software players, commercial, free or even open source. They range from standalone applications to addons to internet browsers (especially Mozilla Firefox). Some software players are capable of displaying HTML pages with embedded MathML markup and subsequently read the mathematical expressions (this functionality was tested on a commercial product Dolphin

EasyReader that supports Czech language as well).

3.2 MathPlayer

MathPlayer is an application developed by Design Science¹ that enables users of Microsoft Internet Explorer web browser to display mathematical equations written in MathML markup. It uses Microsoft's internal HTML engine (MSHTML) on which Internet Explorer is based. Also for any application that makes use of MSHTML to display formatted content MathPlayer is able to display MathML content. This may include email clients, alternative browsers or RSS readers.

Besides displaying MathML in the browser MathPlayer comes equipped with a wide range of MathML-related functions. It enables equations to be copied to the clipboard as MathML markup and pasted to any MathML-aware software, be it simple text editor or more sophisticated application. Drag-and-drop functionality is also supported. Another important feature of the MathPlayer is the ability to speak expressions on the web page.

MathPlayer is shipped as a free to use software with English localization only, but is still a close source product with a license that limits its use to only computers owned by the user and also prohibits translation and reverse engineering of the application.

3.3 AsTeR

AsTeR (Audio System for Technical Readings) [12] is an exceptional piece of software developed by T.V. Raman as a part of his dissertation at Cornell University. AsTeR is capable of producing audio renderings of technical documents, even those containing higher mathematical expressions, written in \TeX or \LaTeX . However, processing different markup languages does not pose a problem for AsTeR. All that is required is a recognizer for given markup language.

The logical structure of the document is transformed to internal representation, which is then rendered in audio using a collection of rendering rules. This enables AsTeR to provide different views of the document; a user can listen to the whole document or select a portion of the document for listening. An important feature of AsTeR makes use of voice intonation to read long and complicated formulae. Such formulae are divided into log-

1. <http://www.dessci.com/en/products/mathplayer/>

ical parts (like a single summation or a numerator of a fraction) and read with different intonation than the surrounding part of the expression.

3.4 MathTalk

MathTalk² [16] is a commercial text-to-software that enables visually impaired people to read algebra notation in a quick manner. MathTalk was designed and created as a part of doctoral dissertation of R. D. Stevens [17] and employs the use of prosody in the synthetic voice to decrease the mental workload of the listener. MathTalk enables the users to browse the document and change the speed of reading which gives them comfortable control of the reading process.

3.5 Lambda

Lambda [13] (Linear Access to Mathematic for Braille Device and Audio-synthesis) aims at solving the problem of mathematics text management by blind users. It consists of two sections: the Lambda code and the editor.

The Lambda code directly derives from MathML, it is automatically convertible into equivalent MathML and through it into the most popular mathematical formats. The Lambda code was designed to have explicit meaning with no ambiguities; to provide full Braille output of mathematical equations; to preserve peculiarities of national Braille; to have a compact linear representation (to minimize the movement while reading Braille).

The Lambda editor allows to write and manipulate mathematical expressions. It is essentially a text editor designed to read and write Lambda code. Mathematical elements are grouped into blocks, that can be collapsed or expanded for easier reading. Blocks can also be deleted or copied to a different place in the document. The output of the Lambda editor can be written (displayed) in Braille, spoken by speech synthesis or both at the same time. The speech synthesis exploits the block structure which enables it to change the speech speed and insert brakes to better communicate the meaning of the mathematical formula.

3.6 Web Browser Support

Strictly speaking web browsers by themselves are not considered an accessibility software, but they play a key role in sharing information and

2. <http://www.mathtalk.com/index.htm>

making mathematical content available to many people across the globe.

MathML was designed as a standard for including mathematical formulae in web pages. It stands to reason that at the most well-known and widespread web browsers should have good support of MathML, at least the Presentation MathML. Unfortunately the situation is not that good, developers of these browsers are mainly focused on different areas and MathML stand very low on the ladder of priorities.

Mozilla Firefox Firefox has the best native support for MathML rendering out of the most popular web browsers and is capable of displaying most elements of the Presentation MathML. Firefox can be enhanced by an extension called FireMath that adds a rich MathML editor capable of generating complicated mathematical expressions quite easily and right into MathML markup.

Google Chrome and Safari Both Google Chrome as well as Safari are based on the WebKit layout engine that has a development version of MathML. The support for MathML is available in latest versions of Safari (since version 5.1). Unfortunately Chrome does not have native support for MathML, but uses the combination HTML and CSS to render mathematical expression entered in MathML.

Internet Explorer Internet Explorer does not have any native support for MathML, but MathML rendering capability and some additional functionality can be added by the MathPlayer extension, that is described in more detail in section 3.2. However MathPlayer's support only extends to Internet Explorer 8, support in version 9 is rather spotty and in version 10 it does not work at all.

Opera Opera was one of the first browsers to include native support for MathML. The rendering of MathML is not as good as in Firefox, Opera has issues with positioning of elements in more complicated constructs.

Browser for Hand-held Devices At the time of writing this thesis none of the browsers for hand-held devices has a satisfactory native support for MathML. However they are capable of rendering expressions using CSS.

Chapter 4

Analysis and Solution Proposal

The goal of this thesis is the development of a conversion tool, that will be able to process XML document or documents and transform each occurrence of MathML markup into the appropriate textual representation. For example $5 \times \alpha = x + 3$ should be transformed into "five times alpha equals x plus three". In case of a visual (not semantic) change of the operator the tool should produce the same result. In our example $5 \times \alpha = x + 3$, $5 \cdot \alpha = x + 3$ and $5\alpha = x + 3$ have the same meaning, so both of them should be transformed into identical strings. For better results the input can be first canonicalized (as described in section 2.3 on page 11). From these requirements we can devise a basic top-level workflow diagram 4.1 on the next page of the application.

Before the input documents are loaded into in-memory structure for representing XML documents, a preprocessing phase will take place. This includes determining the file structure on the hard-drive and retrieving all files from the input document and all its descendants. This structure will be preserved and output documents will have the same structure, just in the user-specified directory. Also the input files can be packed using the .zip archive file format (as is the case with documents in the MREC 6.1) and such files have to be unpacked before further processing. Then, if requested, we will canonicalize input files. Since the canonicalization comes from external module (library), we have to be careful of possible errors or the time efficiency of this process.

One important factor about the input documents, that was not taken into account, is their quantity. The application must be able to process large corpora of mathematical documents in reasonable time and with efficient memory usage. Fortunately modern computation systems provide a way of running the application in multiple threads. The application just has to ensure that all threads work with the same settings. This can be accomplished by providing a single point of entry to the settings instance. A globally visible class (object) that implements the Singleton design pattern is the best solution for accessing settings across the whole application. Also

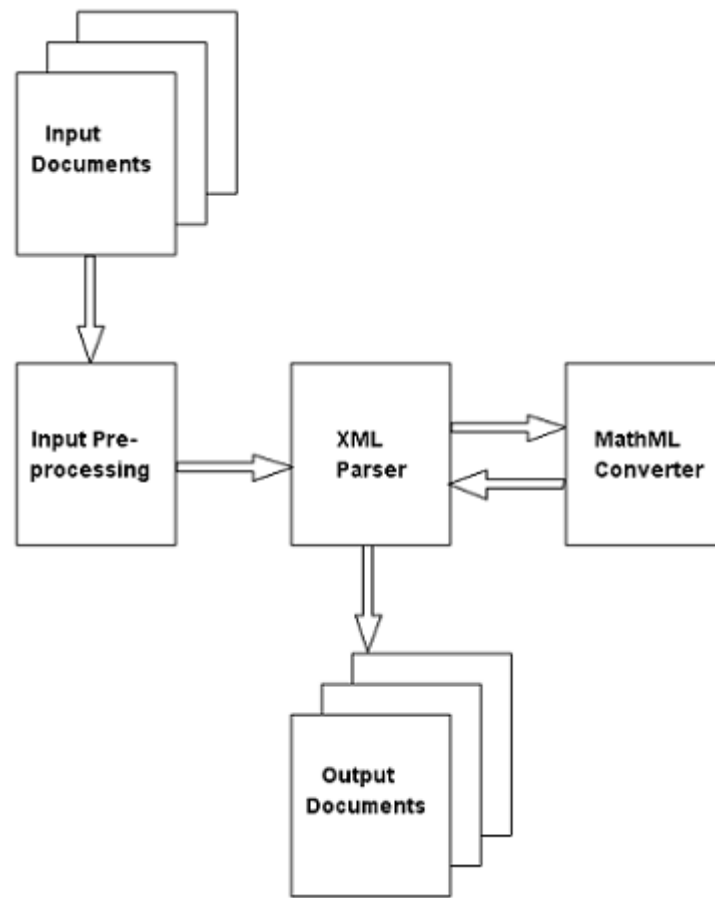


Figure 4.1: Basic diagram of application workflow

retrieving external resources for each single thread is very time consuming, therefore an in-memory cache of resources common for all threads will be implemented.

We can now create a more complex and detailed diagram of internal workings and dataflows inside the application. See diagram 4.2 on the following page.

Three activities (processes) still remain to be defined:

1. Load Document
2. Create Tree
3. Convert Tree

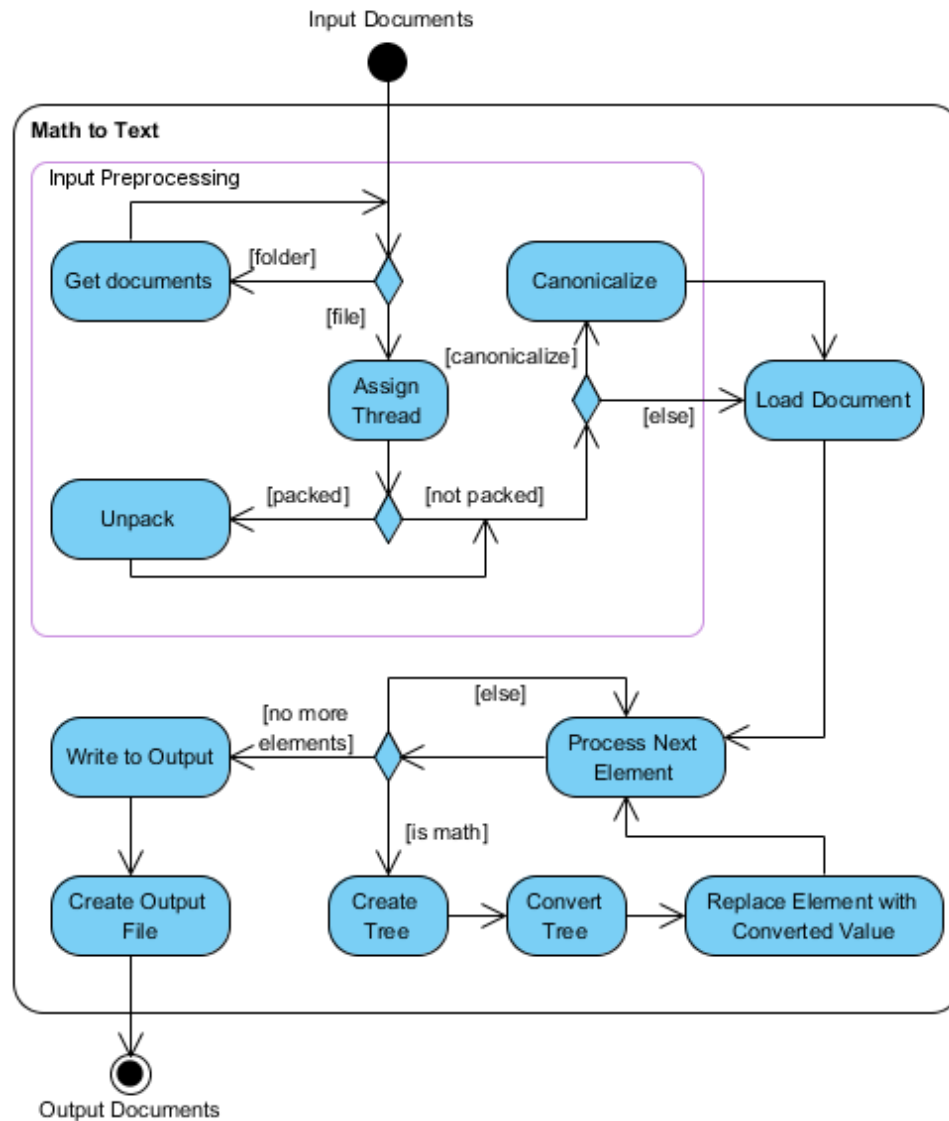


Figure 4.2: Activity diagram of the application

4.1 Load Document

When it comes to processing XML documents, there are four common approaches:

1. Document Object Model¹ (DOM) retains the tree structure of the XML document and is composed of nodes. Node represent elements in the XML structure and form a hierarchical structure.
2. Extensible Stylesheet Language² (XSL) refers to a family of languages used to transform and render XML documents.
3. Streaming XML is an event-based, sequential and unidirectional approach of processing XML documents.
4. XML data binding is a process of mapping XML documents to objects in computer memory (deserialization).

	Processing speed	Memory requirements
DOM	slow	high
XSL	slow	high
Streaming XML	fast	low
Data binding	average	high

Table 4.1: XML processing options comparison

The memory requirements shown in table 4.1 stem from the maximum amount of data that each method uses to process documents. For Streaming XML method, that reads only one event at a time and retains only minimal processing information, these requirements are low. The rest of compared methods need to load the whole document into computer memory and therefore their memory requirements are high. The processing speeds comparison comes from experience of the author with working with these methods.

Based on the comparison in table 4.1 we choose the fastest and most memory-efficient method for processing our input documents - the Streaming XML method. Algorithm 1 on the next page shows a way to process input documents, using an in-memory tree structure defined in section 4.2.

4.2 Create Tree

The tree in this case means an in-memory representation of MathML with tree structure. DOM seems like a good option for this application, however

1. <http://www.w3.org/DOM/DOMTR>

2. <http://www.w3.org/Style/XSL/>

Algorithm 1 Process input algorithm

```

1: procedure PROCESSINPUT
2:   clear tree ▷ a DOM-like MathML tree
3:   while next event exists do
4:     event ← next event
5:     if event is a start of math element then
6:       tree ← create a new tree
7:     else if event is an end of math element then
8:       convert tree and write the result to the output
9:       clear tree
10:    else if event is inside math element then
11:      insert a new node, value or attribute in the tree
12:    else
13:      write event to the output
14:    end if
15:  end while
16: end procedure

```

a fully-fledged DOM representation of MathML markup is not required. The information this model provides takes up resources (time and memory) and a big part of it would be discarded. However we need a structure that will provide a comfortable traversal; moving from parent to children and vice versa.

We have therefore designed a simplified DOM tree, with just the information we actually need. Since we only use it to build a tree representation of MathML markup, every node in this tree has a special property, that denotes its type - name of the element and part of MathML it belongs to (Presentation or Content).

Besides the type property our simplified model contains a list of children, pointer to the parent, a text value, a list of XML attributes (key-value pairs) and a attribute that signifies whether this node has already been processed (we want to process each node only once).

Our very simple tree has one more advantage, besides simplicity - it allows the application to provide an option to change the method of loading XML documents.

4.3 Convert Tree

We are starting the conversion process with our tree representation of the MathML markup. At the beginning of the conversion we need to determine whether the tree consist of only the Presentation markup, Content markup or both, since each requires a slightly different approach to the conversion. In most cases the Presentation MathML resides directly in the `math` tree (is a direct descendant of `math` element) or in the element `semantics`, while the Content MathML is often found enclosed in the `annotation-xml` element with an attribute `encoding` set to the value `MathML-Content`. Or we can simply traverse the tree till we find an element from either the Presentation of Content MathML markup and continue based on our finding.

The conversion process starts at the top level element, `math`, and then recursively continues to traverse the tree, converting the Presentation and Content elements based on their own specifics. Each converted node is marked as processed to prevent it to be processed twice, since in some cases the conversion requires to look ahead and jump out of the recursion pattern to process a sibling node (or any other node).

The basic outline of the conversion process can be seen in figure 4.3 on the next page and will be described in detail in the next sections.

4.3.1 Presentation MathML

Every element of the Presentation MathML requires individual approach to processing. But there are still groups of elements with similar characteristics.

The first group consists of token elements. The conversion is straightforward and depends only on the value of the element:

- `mn` - convert number to string or take the original value,
- `mi` - take the original value if it is in Latin script, otherwise convert the value if possible (e.g. Greek alphabet),
- `mo` - determine the operation defined by the `mo` element value ($\Sigma \rightarrow$ sum, $\int \rightarrow$ integral, ...).

It is important to remember that one mathematical operation can be expressed with various symbols and the conversion has to unify them into the same string representation.

The second group contains element with specific purpose. In other words these elements clearly state (with their names) what is their intended

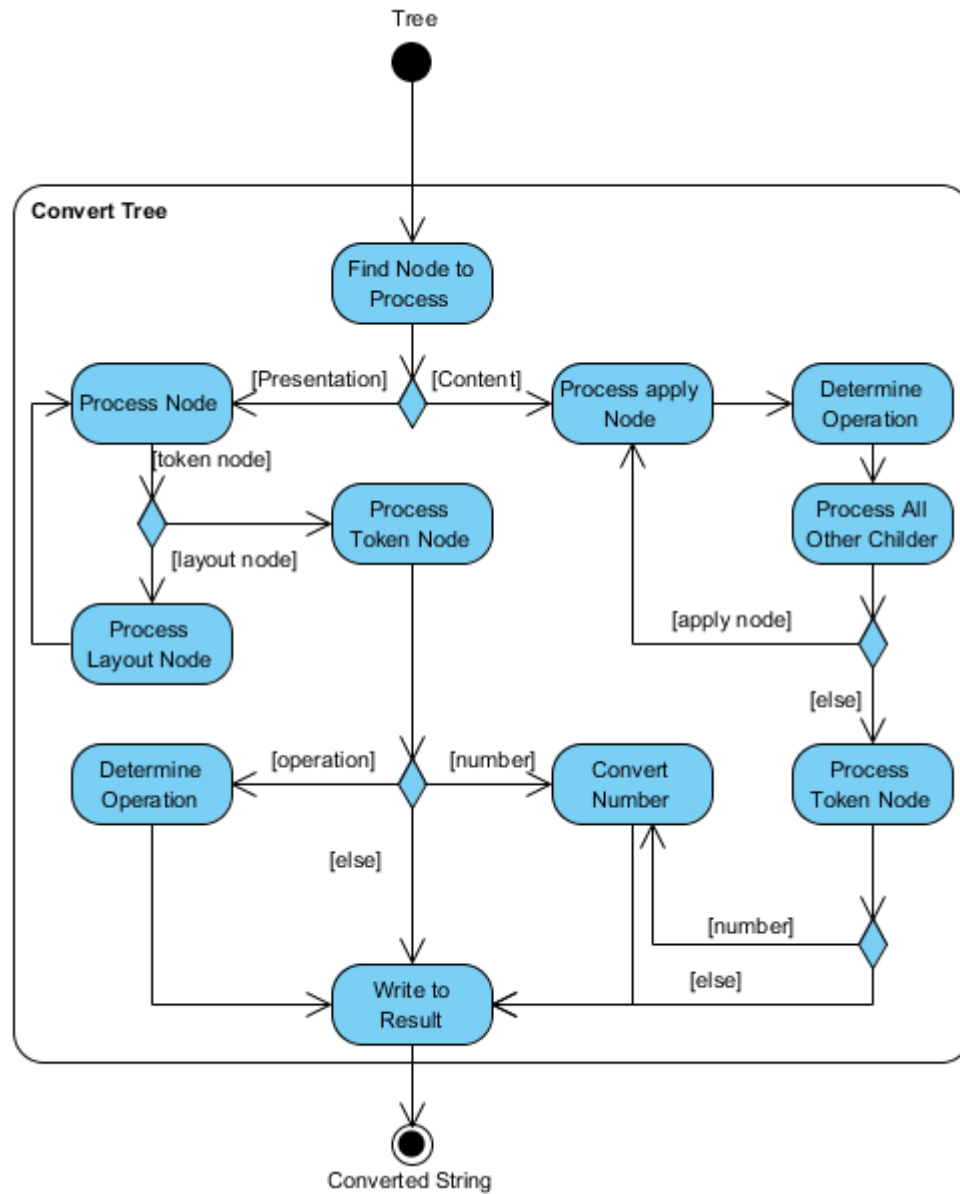


Figure 4.3: Tree conversion process

function. The conversion uses this fact and is therefore very simple. For example:

- `mfrac` - fractions or binomial numbers,

- `msqrt` - square root,
- `mfenced` - subexpression is enclosed in parenthesis.

The third group is composed of elements that specify only the layout of the subexpression. To determine the exact function the author wanted to convey we need additional information provided by child elements and in some cases it may also depend on sibling or parent elements.

```
<math>
  <munder>
    <mo>lim</mo>
    <mrow>
      <mi>x</mi>
      <mo>&rarr;</mo>
      <mn>0</mn>
    </mrow>
  </munder>
  <mrow>
    <msqrt>
      <mi>x</mi>
    </msqrt>
  </mrow>
</math>
```

Figure 4.4: Expression $\lim_{x \rightarrow 0} \sqrt{x}$ in Presentation MathML

As we can see in figure 4.4 the `munder` element has to be interpreted in the context of its child elements (especially the first child) and its first sibling element. Similar rules apply to elements `mover`, `munderover`, `msub`, `msup` and `msubsup`, that describe various mathematical constructs such as summations, integrals, products, logarithms and many more.

The last group is formed by elements that have purely presentation function: spacing, padding, and so on. These can be ignored altogether since they do not provide any information about the meaning of presented expressions.

4.3.2 Content MathML

The cornerstone of the Content MathML is without a doubt the element `apply`, function application. It describes the application of its first child

element on the rest of child elements, all of which can be `apply` elements themselves. The `apply` element will therefore serve as a hub, assigning the processing of expression or subexpression based on the operator (the first child element).

To be able to convert the expression correctly we need a different approach than for the Presentation MathML, where the elements are ordered for display purposes and can be converted basically from top to bottom (with a few exceptions). In the Content MathML a function can be applied to multiple elements, but it will be declared only once as seen in figure 4.5 (in the Presentation MathML the plus symbol would be declared twice, between x and y , y and z).

```
<math>
  <apply>
    <plus/>
    <ci>x</ci>
    <ci>y</ci>
    <ci>z</ci>
  </apply>
</math>
```

Figure 4.5: Expression $x + y + z$ in Content MathML

Therefore we can not determine word order of the expression from the position of elements in the document (as is the case in the Presentation MathML), but we have to specify the desired word order for each operation separately. Fortunately operations can be sorted into logical groups based on the word order we use when presenting them.

- Infix form - the operator is used between pairs of inputs, starting with the first and the second input, then the second and the third and so on. These include operators for division ($x/y/z$), multiplication ($x \cdot y \cdot z$), comparison ($x = y = z, \geq, <$) and many more.
- Prefix form - the operator is used at the beginning, preceding inputs. In this case the operator is used only once at the beginning, followed by converted inputs. Examples include absolute value ($|x|$), negation (\neg) or floor ($\lfloor x \rfloor$).
- Prefix form with multiple inputs - a special case of the prefix form, where there are multiple inputs. In this instance the operator is also used

only once, but inputs have to be divided by commas, with the last two inputs divided by the word "and". Function $\min(x, y, z)$ should be converted to string: minimum of x , y and z . Another examples might be sets, lists, functions greatest common divisor, maximum or lowest common multiple.

- Quotient and remainder - $\text{rem}(x, y)$ should be converted to: reminder of x divided by y . Similarly the quotient operator.
- Others - operators that do not belong to any of abovementioned categories or belong to more than one (like plus or minus, which can be used in both the prefix and infix form) have to be treated in a separate way.

4.3.3 Operators and Symbols

As we mentioned before many mathematical operations can be expressed using more than one operator or symbol. Our task lies in identifying operators and symbols belonging to each operation, creating a storage structure for this data and developing a method or methods for searching in this structure, i.e. finding an operation for a given operator.

This way all operators that denote the same operation are grouped together and we can easily unify the way each operation is converted.

Chapter 5

Implementation

This chapter describes various implementation aspects that occurred during the creation of the conversion application. The application is written in Java programming language, mainly because of personal preference of the author. Also the existence and availability of many frameworks and tools for Java language (for working with XML documents among others) are a big advantage of using this language.

The application uses Apache Maven for build automation, distribution management and dependency management.

The structure of the application can be seen on the component diagram 5.1 below.

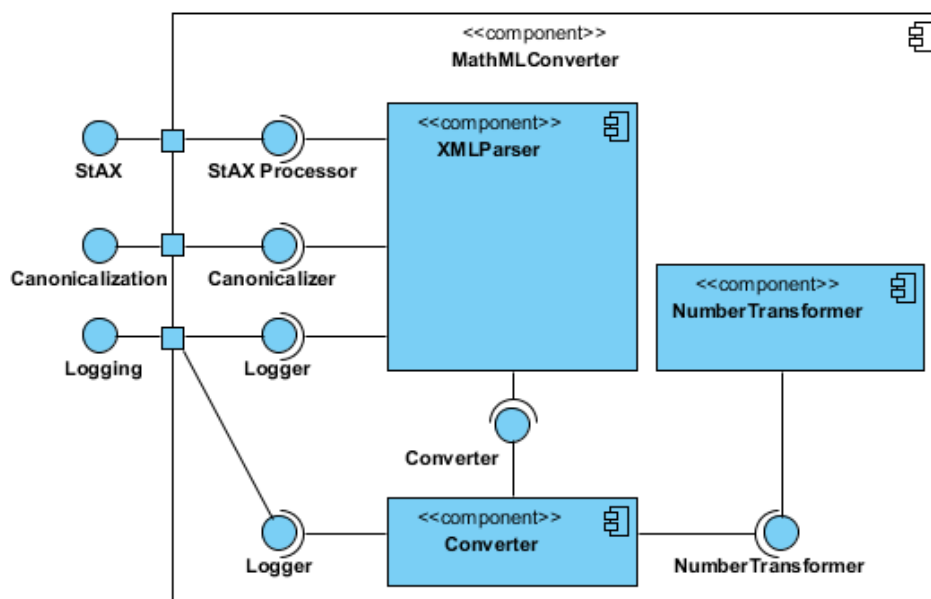


Figure 5.1: Component diagram of the application

5.1 DOM-like Representation of MathML

The implementation of the tree designed in section 4.2 on page 19 can be seen in figure 5.2 on the next page.

5.2 Input Processing

As we outlined in section 4.1 on page 18, the streaming approach with custom DOM-like internal representation of MathML tree is used.

In Java two major APIs¹ for streaming processing of XML documents are popular among developers:

- SAX (Simple API for XML) - SAX implements the push principle, reporting events as it encounters them. From the point of view of this application a better approach would be
- StAX (Streaming API for XML) - pull principle; the application requests events from the StAX processor. (StAX is a specification defined by the JSR 173².)

This application uses default StAX implementation that is shipped with Java Standard Edition 6 runtime, but there are also other implementations available, such as Woodstox³ or Aalto⁴, and there is a possibility to request one of these two implementations to use instead of the default implementation.

StAX offers two ways of traversing documents: Cursor API and Iterator API. We will use the first, Cursor API, because it should be faster and more memory-efficient⁵. Two main interfaces are available in the Cursor API: `XMLStreamReader` for accessing all possible information retrievable from XML documents and `XMLStreamWriter` which in turn provides methods for outputting this information.

5.2.1 XML Parser

The implementation of the `XmlParser` interface, `XmlParserStAX`, is responsible for processing input XML documents using the StAX Cursor API

1. Application programming interface

2. <http://www.jcp.org/en/jsr/detail?id=173>

3. <http://woodstox.codehaus.org>

4. <http://wiki.fasterxml.com/AaltoHome>

5. http://docs.oracle.com/cd/E17802_01/webservices/webservices/docs/1.6/tutorial/doc/SJSXP3.html


```
public final class MathMLNode {
    /**
     * Type of this node.
     */
    private MathMLElement type;
    /**
     * List of all children of this node. Empty if
     * there are none. In this case value must be
     * set.
     */
    private List<MathMLNode> children = new
        ArrayList<MathMLNode>();
    /**
     * Text value of this node. {@code null} if
     * there are some child nodes.
     */
    private String value;
    /**
     * Parent node.
     */
    private MathMLNode parent;
    /**
     * Was this node already processed? Useful when
     * you have to "look ahead" and process
     * element sooner.
     */
    private boolean processed = false;
    /**
     * Set of attributes.
     */
    private Set<XmlAttribute> attributes = new
        HashSet<XmlAttribute>();

    ... getters
    ... setters
}
```

Figure 5.2: Simplified DOM for internal representation of MathML markup

```
public interface XMLStreamReader {
    public int next();
    public boolean hasNext();
    public String getText();
    public String getLocalName();
    public String getNamespaceURI();
    ...
}

public interface XMLStreamWriter {
    public void writeStartElement(String localName)
        ;
    public void writeEndElement();
    public void writeCharacters(String text);
    ...
}
```

Figure 5.3: Example of methods in XMLStreamReader and XMLStreamWriter interfaces of the Cursor API

and at the same time creating output documents in the same run through the file. The internal processing follows steps outlined in algorithm 1 on page 20. The parser also provides a method for processing string input and producing string output. In this case the parser processes the document defined by this string the same way as any other input document, however the output contains only converted strings, without any of original documents elements outside the MathML markup.

```
public class XmlParserStAX implements XmlParser {
    public String parse(final String inputString,
        final Locale language);
    public File parse(final File file, final Locale
        language);
    public List<File> parse(final List<File> files,
        final Locale language);
}
```

Figure 5.4: Public methods of XmlParserStAX class

The parser is capable of accepting a folder as an input (instead of a document) and processing all files in this folder, preserving the file structure on the path entered by the user.

A thread pool of a user-specified size is initialized at the beginning and for every file from the input a `java.util.concurrent.Callable` instance is created. Using the `java.util.concurrent.ExecutorService` with our thread pool these callables are concurrently invoked (whenever there is an free thread in the thread pool next callable is invoked).

5.3 Conversion

When the parser creates a complete tree, it immediately initializes the conversion of that tree by calling the `convert(tree, language)` method of the `MathMLConverter` class. The converter then decides which element of the input tree will be used as the root element of the conversion.

In case that the tree contains Content MathML markup, this markup has precedence over the Presentation markup.

The converter consists of several classes that represent elements of MathML markup; a separate class for each element of the Presentation MathML markup and three classes (`Apply`, `Cn`, `Ci`) for the Content MathML markup. Every class has only one static method, `process(node, settings)`, that is called whenever the traversal of the input tree encounters corresponding element. The run through the tree is done via the class `Node` and its static method `process(node, settings)`. This method works as a switch for every type of node and delegates the processing to a specific class. It also marks each encountered node as processed (see figure 5.5 for illustration). The resulting string is gradually built from partial results of individual nodes.

Before we get to the explanation of the conversion inside these classes, we have to define how mathematical operators and operations are defined in the application. The `Operation` class is an enumeration of all operations known to this application. Each operation is assigned a unique key, that is also used to retrieve appropriate textual representation of this operation from the localization files. Next there is a type of the operation as defined in section 4.3.2. The last property is an array of possible operators or symbols that might be used to represent this operation in documents. These can be standard HTML character entities entered by their name (`−`) or unicode code point (decimal or hexadecimal) or a name of the correspond-

```

public class Node {
    public static String process(final MathMLNode
        node, final ConverterSettings settings) {
        final StringBuilder builder = new
            StringBuilder();
        switch (node.getType()) {
            ...
            case MN: {
                builder.append(Mn.process(node,
                    settings));
                break;
            }
            case MFRAC: {
                builder.append(Mfrac.process(node,
                    settings));
                break;
            }
            ...
        }
        node.setProcessed();
    }
}

```

Figure 5.5: Excerpt from the Node class

ing element in the Content MathML markup as seen in figure 5.6.

```

public enum Operation {
    MINUS_PLUS("minus_plus", OperationType.INFIX,
        "&#8723;", "&#x2213;", "±");
}

```

Figure 5.6: An example of an operation

The conversion of Presentation MathML elements is very simple and derives from the names of respective nodes. We offer special treatment for layout elements (such as `munder`, `msup`, ...) and distinguish a few notable operations, that occur often in mathematical texts, like limits, integrals or summations.

On the other hand the conversion of Content MathML elements is more complicated. Everything important is happening inside the class `Apply`. The first child of the `apply` element denotes the operation. The operation can be defined by the appropriate Content MathML element, in which case we use the element name to retrieve the operation from the enumeration; by the `csymbol` element, then the value of the element is used; or by another `apply` element, for which we invoke the `Node.process()` method.

When we have determined the exact operation, we proceed by processing it based on its type, i.e. based on its word order when spoken. Operations with a type `OperationType.SPECIAL` have to be processed individually since they require a distinctive approach (e.g. logarithm) or there are multiple ways of expressing them (e.g. summation, integral).

Lastly the conversion of numbers to strings ($11 \rightarrow \text{eleven}$) is optionally executed and contents of identifiers (elements `mi` and `ci`) are verbatim copied to the output.

Chapter 6

Results

6.1 Mathematical Retrieval Collection MREC

MREC [10] is a large corpus of mathematical texts (numbering close to 500 thousand documents). MREC¹ is based on documents downloaded from arXMLiv² [15], which in turn is created by transforming T_EX documents from arXiv³ - a huge library of freely available texts from multiple scientific fields, including Physics, Mathematics and Computer science. Using L^AT_EXML 2.2.1 these texts are converted into XML format with mathematics being represented in MathML. Not all documents from the arXiv are part of MREC - only documents for which the conversion yielded results in categories successful and complete with errors are included.

6.2 Changing StAX Implementation

-
1. <https://mir.fi.muni.cz/MREC/index.html>
 2. <http://kwarc.info/projects/arXMLiv/>
 3. <http://arxiv.org/>

Chapter 7

Conclusion

1 page

Bibliography

- [1] Dominique Archambault, Donal Fitzpatrick, Gopal Gupta, Arthur I Karshmer, Klaus Miesenberger, and Enrico Pontelli. Towards a universal maths conversion library. In *Computers Helping People with Special Needs*, pages 664–669. Springer, 2004.
- [2] Dominique Archambault and Victor Moço. Canonical MathML to simplify conversion of MathML to Braille mathematical notations. In *Computers Helping People with Special Needs*, pages 1191–1198. Springer, 2006.
- [3] Josef B Baker, Alan P Sexton, and Volker Sorge. Towards reverse engineering of PDF documents. *Towards a Digital Mathematics Library. Bertinoro, Italy, July 20-21st, 2011*, pages 65–75, 2011. <http://dml.cz/dmlcz/702603>.
- [4] Josef B Baker, Alan P Sexton, and Volker Sorge. MaxTract: Converting PDF to \LaTeX , MathML and Text. In *Intelligent Computer Mathematics*, pages 422–426. Springer, 2012.
- [5] David Formánek, Martin Líška, Michal Růžička, and Petr Sojka. Normalization of Digital Mathematics Library Content. In James Davenport, Johan Jeuring, Christoph Lange, and Paul Libbrecht, editors, *24th OpenMath Workshop, 7th Workshop on Mathematical User Interfaces (MathUI), and Intelligent Computer Mathematics Work in Progress*, number 921 in CEUR Workshop Proceedings, pages 91–103, Aachen, 2012.
- [6] José Grimm. Tralics, a \LaTeX to XML Translator. *TUGboat*, 24(3):377–388, 2003.
- [7] José Grimm. Producing MathML with Tralics. In Petr Sojka, editor, *Proceedings of DML 2010*, pages 105–117, Paris, France, July 2010. Masaryk University. <http://dml.cz/dmlcz/702579>.
- [8] Eitan M Gurari. TEX4ht: HTML Production. *TUG-boat*, 25(1):39–47, 2004.

-
- [9] Martin Jarmar. *Conversion of Mathematical Documents into Braille*. PhD thesis, Master's thesis, Faculty of Informatics (Jan 2012), https://is.muni.cz/th/172981/fi_m/?lang=en, 2012.
- [10] Martin Líška, Petr Sojka, Michal Růžička, and Petr Mravec. Web Interface and Collection for Mathematical Retrieval: WebMIaS and MREC. *Towards a Digital Mathematics Library. Bertinoro, Italy, July 20-21st, 2011*, pages 77–84, 2011.
- [11] Bruce Miller. LaTeXML: A L^AT_EX to xml converter. *Web Manual at <http://dlmf.nist.gov/LaTeXML/>, seen April 2013*, 2013.
- [12] TV Raman. *Audio System for Technical Readings*. PhD thesis, Cornell University, 1994.
- [13] Waltraud Schweikhardt, Cristian Bernareggi, Nadine Jessel, Benoit Encelle, and Margarethe Gut. LAMBDA: A European system to access mathematics with Braille and audio synthesis. In *Computers Helping People with Special Needs*, pages 1223–1230. Springer, 2006.
- [14] Petr Sojka and Martin Líška. Indexing and Searching Mathematics in Digital Libraries. In *Intelligent Computer Mathematics*, pages 228–243. Springer, 2011.
- [15] Heinrich Stamerjohanns, Michael Kohlhase, Deyan Ginev, Catalin David, and Bruce Miller. Transforming Large Collections of Scientific Publications to XML. *Mathematics in Computer Science*, 3:299–307, 2010. <http://dx.doi.org/10.1007/s11786-010-0024-7>.
- [16] Robert Stevens and Alistair Edwards. Mathtalk: The design of an interface for reading algebra using speech. In *Computers for Handicapped Persons*, pages 313–320. Springer, 1994.
- [17] Robert David Stevens. *Principles for the design of auditory interfaces to present complex information to blind people*. PhD thesis, University of York, 1996.
- [18] 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*, pages 95–104. ACM, 2003.