INSTITUTE OF
COMPUTER SCIENCE
Masaryk University

# Mathematics Index and Search in DSpace

**Vlastimil Krejčíř,** `krejcir@ics.muni.cz`

OAI 10 DSpace User Group Meeting, Geneva, 20th June 2017

# Motivation

In 2005, we started to build (upon DSpace + XMLUI) the **Czech Digital Mathematics Library (DML-CZ)**. Later on, the question *"What about a mathematics formulae search?"* appeared.

- Simple search based on text keywords is not sufficient for mathematical content because
- representation of thoughts in math = formulae, so
- math formulae search brings a great benefit to mathematicians (and related science disciplines).

Lets look what we struggle with…

# Simple Text Search

- Plain text indexing and searching is 'easy'
- Known tools and techniques
- Searching for DSpace $\Rightarrow$ the query is *'dspace'*
  - or any substring of this query…
- Matching is done character by character

…but for math formulae the situation is much more complicated!

# Mathematics Formulae Index and Search

Questions/thoughts:

- How to represent a formula?
- How to index a formula?
- How to write a query?
- How to match a query and weight results?

Consider that:

- symbols and graphics heavily used
- big complexity of formulae
  - subformulae usually have sense
- variables, constants, …
- syntax ambiguity
  - from the 'search and index' point of view

# Example: Ambiguity

These make the things really hard:

$$0.5 = \frac{1}{2} = 2^{-1}$$

$$\sqrt{8} = 2\sqrt{2}$$

It is easy to find a lot of similar examples…

# Example 2: Pythagorean Theorem

**Pythagorean Theorem**

$$a^2 + b^2 = c^2$$

# Example 2: Pythagorean Theorem

**Pythagorean Theorem**

$$a^2 + b^2 = c^2$$

is equivalent to

$$b^2 + a^2 = c^2$$

# Example 2: Pythagorean Theorem

**Pythagorean Theorem**

$$a^2 + b^2 = c^2$$

is equivalent to

$$b^2 + a^2 = c^2$$

is equivalent to

$$x^2 + y^2 = z^2$$

# Example 2: Pythagorean Theorem

**Pythagorean Theorem**

$$a^2 + b^2 = c^2$$

is equivalent to

$$b^2 + a^2 = c^2$$

is equivalent to

$$x^2 + y^2 = z^2$$

and is special case of **Fermat's Last Theorem**

$$a^n + b^n = c^n$$

*Besides, see book: Simon Singh: Fermat's Last Theorem*

# MathML

MathML - XML formulae representation

```
<math>
  <mfrac>
    <mn>1</mn>
    <msup>
      <mi mathvariant="bold">x</mi>
      <mn>2</mn>
    </msup>
  </mfrac>
</math>
```

$$\frac{1}{x^2}$$

# How to get MathML

To get MathML out of existing articles – very hard task...

...our real pain :-(.

There are some tools:

- InftyReader (OCR tool)
- LaTexML
- MATLAB
- 'hand made'
- ...

# MIaS: Mathematics Index and Search tool

**MIaS** is the Java tool that provides the necessary job:
MathML $\Rightarrow$ Lucene:

- canonicalization
- ordering
- tokenization
- variables and constant unification

The result is *M-term*, MIaS processed and plain text coded formula:

$$F(N(1)J(I[V{=}B](1)N(2)))$$

# MIaS: schema

# Math in DSpace (DML-CZ)

Assume we have already prepared MathML in time of ingest.

- extra metadata registry for math
    - *dmlcz.math*: MathML formulae storage
- configure SOLR to process *dmlcz.math*
    - we use *search* core
    - index *dmlcz.math* using **MIaS** analyzer
        - increases the size of index (approx. 100 times)
        - in DML-CZ: 264 MB $\rightarrow$ 28 GB (40 thousand items)
    - search *dmlcz.math* using **MIaS** analyzer again + **MIaS Payload Similarity** module
        - MIaS Payload Similarity takes care of results ranking

# Math in DSpace (cont.)

- integrate 'user friendly' formulae search in DSpace UI
  - in our case XMLUI
- separate form for math search
  - MathML or LaTeX notation
  - on the fly rendered and displayed using **MathJax**
    - JavaScript library
  - LaTeX converted to MathML query via LaTeXML
    - written in Perl

# Math formula search example in DML-CZ

# Math formula search example in DML-CZ

# Math formula search example in WebMIaS



Data (MathML formulea) taken from ArXiv.org

# Math formula search example in WebMIaS



Data (MathML formulea) taken from ArXiv.org

# Special thanks to (in alphabetical order):

**Martin Líška** (Maths Information Retrieval team)

**Michal Růžička** (Maths Information Retrieval team)

**Petr Sojka** (Maths Information Retrieval team)

**Dominik Szalai** (DSpace integration)

# References

MIR team homepage:
`https://mir.fi.muni.cz/`

DML-CZ DSpace at GitHub (branch *dspace5-dmlcz*):
`https://github.com/empt-ak/DSpace`

DSpace and MIaS integration tech report:
`https://empt-ak.gitbooks.io/dmlcz/content/dml.html`

# Questions?