

Maths Information Retrieval for Digital Libraries

Michal Růžička

Masaryk University, Faculty of Informatics, Brno, Czech Republic
mruzicka@mail.muni.cz

<https://mir.fi.muni.cz/>



Illustrations by Jiří Franek.



Outline

1 Motivation

2 MathML Canonicalizer

3 My Goals

4 Evaluation

5 Future Works

6 Summary

Motivation



Motivation

- *Czech Digital Mathematics Library (DML-CZ)*, <http://www.dml.cz/>
- *Centre de diffusion de revues académiques mathématiques (CEDRAM)*,
<http://www.numdam.org/>
- *Numérisation de documents anciens mathématiques (NUMDAM)*, <http://www.cedram.org/>
- *Göttingen Göttinger Digitalisierungszentrum (GDZ)*, <http://gdz.sub.uni-goettingen.de/>
- *Electronic Research Archive for Mathematics (ERAM)*, <http://www.emis.de/projects/JFM/>
- *The Electronic Library of Mathematics (ELibM)*,
<http://siba-sinmemis.unile.it/ELibM.html>
- *Journal STORage (JSTOR)*, <http://www.jstor.org/>
- *Project Euclid*, <http://projecteuclid.org/>
- *Russian Digital Mathematics Library (RusDML)*, <http://www.rusdml.de/>
- *Polish Digital Mathematical Library (DML-PL)*, <http://pldml.icm.edu.pl/>
- *Biblioteca Digital Española de Matemáticas (DML-E)*, <http://dmle.cindoc.csic.es/>
- *Japanese Digital Mathematics Library (DML-JP)*,
<http://sparc1.math.sci.hokudai.ac.jp/dmljp/>
- *Riviste Elettroniche Italiane di Matematica (REIM)*, <http://siba2.unile.it/sinm/reim/>
- *Biblioteca Digitale Italiana di Matematica (bdim)*, <http://www.bdim.eu/>

Motivation (cont.)

Q: 'What functionality and incentives would made a working mathematician to login and use a modern DML as EuDML?'

A: '**Math formulae search.**'

Prof. James Davenport, CEIC member, MKM 2011 PC chair, on panel at DML 2011 workshop in Bertinoro as a reply.



Motivation (cont.)

- DML without maths-aware search support is an oxymoron.



- Simple search based on text keywords is not appropriate or sufficient for mathematical contents.

Outline

1 Motivation

2 MathML Canonicalizer

3 My Goals

4 Evaluation

5 Future Works

6 Summary

Sources of MathML in Digital Libraries

- ‘Hand made’ ‘Hand made’ $x^2 + y^2$ MathML
- Tralics
- LATEXML
- InftyReader
- MaxTract
- MATLAB
- Wolfram Alpha
- ...

```
<math xmlns='http://www.w3.org/1998/Math/MathML'>
  <msup>
    <mi>x</mi><mn>2</mn>
  </msup>
  <mo>+</mo>
  <msup>
    <mi>y</mi><mn>2</mn>
  </msup>
</math>
```

Sources of MathML in Digital Libraries

- 'Hand made'
- Tralics
- LATEXML
- InftyReader
- MaxTract
- MATLAB
- Wolfram Alpha
- ...

Matlab $x^2 + y^2$ MathML

```
generate::MathML(x^2 + y^2,
                  Content = FALSE, Annotation = FALSE)
<math xmlns='http://www.w3.org/1998/Math/MathML'>
  <mrow xref='No7'>
    <msup xref='No3'>
      <mi xref='No1'>x</mi>
      <mn xref='No2'>2</mn>
    </msup>
    <mo>+</mo>
    <msup xref='No6'>
      <mi xref='No4'>y</mi>
      <mn xref='No5'>2</mn>
    </msup>
  </mrow>
</math>
```

Sources of MathML in Digital Libraries

- 'Hand made'
- Tralics
- L^AT_EXML
- InftyReader
- MaxTract
- MATLAB
- Wolfram Alpha
- ...

L^AT_EXML $x^2 + y^2$ MathML

```
<math xmlns="http://www.w3.org/1998/Math/MathML"
      alttext="x^{2}+y^{2}" display="block">
  <semantics>
    <mrow>
      <msup><mi>x</mi></msup>2
      <mo>+</mo>
      <msup><mi>y</mi></msup>2
    </mrow>
    <annotation encoding="application/x-tex">
      x^{2}+y^{2}
    </annotation>
  </semantics>
</math>
```

Sources of MathML in Digital Libraries

- 'Hand made'
- Tralics
- LATEXML
- InftyReader
- MaxTract
- MATLAB
- Wolfram Alpha
- ...

InftyReader $x^2 + y^2$ MathML

```
<math xmlns="http://www.w3.org/1998/Math/MathML">
  <msup>
    <mi mathvariant="italic">x</mi>
    <mrow>
      <mn mathvariant="normal">2</mn>
    </mrow>
  </msup>
  <mo mathvariant="normal">+</mo>
  <msup>
    <mi mathvariant="italic">y</mi>
    <mrow>
      <mn mathvariant="normal">2</mn>
    </mrow>
  </msup>
</math>
```

MathML Canonicalizer

- Our own MathML canonicalization tool.
- The main design imperatives:
 - Modularity,
 - simplicity,
 - extensibility, and
 - flexibility.
- The speed of the canonicalization application is also a critical parameter.
 - In our MREC corpora there is 168,000,000 formulae to canonicalize.

MathML Canonicalizer Use Cases

<mp phantom> Omission

```
<mfrac>
  <mrow>
    <mi> x </mi>
    <mo> + </mo>
    <mi> y </mi>
    <mo> + </mo>
    <mi> z </mi>
  </mrow>
  <mrow>
    <mi> x </mi>
    <mp phantom>
      <mo> + </mo>
      <mi> y </mi>
    </mp phantom>
    <mo> + </mo>
    <mi> z </mi>
  </mrow>
</mfrac>
```

```
<mfrac>
  <mrow>
    <mi> x </mi>
    <mo> + </mo>
    <mi> y </mi>
    <mo> + </mo>
    <mi> z </mi>
  </mrow>
  <mrow>
    <mi> x </mi>
    <mo> + </mo>
    <mi> z </mi>
  </mrow>
</mfrac>
```

MathML Canonicalizer Use Cases

Unnecessary Attributes

```
<mfrac linethickness="2"  
       bevelled="true">  
  <mi> a </mi>  
  <mi> b </mi>  
</mfrac>  
  
<mfrac>  
  <mi> a </mi>  
  <mi> b </mi>  
</mfrac>
```

MathML Canonicalizer Use Cases

<mrow> Minimizing

```
<msqrt>                                <msqrt>
  <mrow>
    <mo> - </mo>                      <mo> - </mo>
    <mn> 1 </mn>                      <mn> 1 </mn>
  </mrow>
</msqrt>                                </msqrt>
```

MathML Canonicalizer Use Cases

Unifying Fences

```
<mfenced open="[">
  <mi> x </mi>
  <mi> y </mi>
</mfenced>
<mrow>
  <mo> [ </mo>
<mrow>
  <mi> x </mi>
  <mo> , </mo>
  <mi> y </mi>
<mrow>
  <mo> ) </mo>
</mrow>
```

MathML Canonicalizer Use Cases

Sub-/Superscripts Handling

```
<msubsup>
  <mi> x </mi>
  <mn> 1 </mn>
  <mn> 2 </mn>
</msubsup>
```

```
<msubsup>
  <msup>
    <mi> x </mi>
    <mn> 1 </mn>
  </msup>
  <mn> 2 </mn>
</msubsup>
```

MathML Canonicalizer Use Cases

Applying Functions

```
<mi> f </mi>
<mo> &#x2061; </mo>
<mrow>
  <mo> ( </mo>
  <mi> x </mi>
  <mo> ) </mo>
</mrow>
```

```
<mi> f </mi>
<mrow>
  <mo> ( </mo>
  <mi> x </mi>
  <mo> ) </mo>
</mrow>
```

MathML Canonicalizer Use Cases

Applying Functions

```
<mi> sin </mi>
```

```
<mo> &#x2061; </mo>
```

```
<mi> x </mi>
```

```
<mi>sin</mi>
```

```
<mrow>
```

```
<mo>(</mo>
```

```
<mi>x</mi>
```

```
<mo>)</mo>
```

```
</mrow>
```

MathML Canonicalizer Web Evaluation Application

- JUnit testing seems not to be enough.
- Collaboration-enabled evaluation system needed:
 - Visualization of test data collection.
 - History of canonicalization results on the same data with different version of the Canonicalizer.
 - Annotations on the results.
 - Statistics.
 - Coverage of all the mark-up of the MathML standard.

MathML Canonicalizer Web Evaluation Application

MathMLCanEval

Jmeno

Heslo

Přihlásit

Vzorec

Detaily vzorce	
ID	7519
Uživatel	admin
Zdrojový dokument	textr-10
Konverzní program	LaTeXML
Poznámky	

MathMLCanEval

Jmeno

Heslo

Přihlásit

Kanonikálizovány vzorec

Nejít podobné

Podrobnosti	
ID	6568
Původní vzorec	7519
Doba behu	2

Anotace	

MathML kód	
Kanonikální kód	Původní kód

```
<math display="block">\frac{e^++e^-}{2}
```

```
smath smmlm="http://alpha.planetmath.org/orita/mse"
smmlm:dc="http://purl.oclc.org/dc/terms/"
smllm:rcs="http://alpha.planetmath.org/orita/mse"
smllm:version="20140927T092114Z"
smllm:mime="text/html"
smllm:id="1dp3881904"
smllm:content="<math display='block'>{\frac{{e}^{+}+{e}^{-}}{2}}</math>
<math display='block'>{\frac{{e}^{+}+{e}^{-}}{2}}</math>
```

PDF view generated by pdfkit

Outline

1 Motivation

2 MathML Canonicalizer

3 My Goals

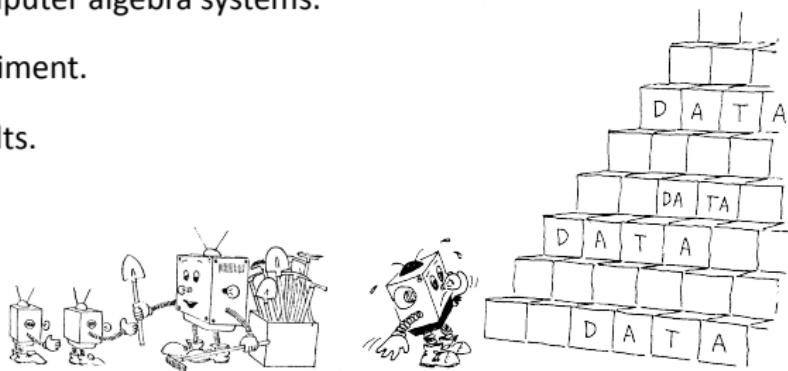
4 Evaluation

5 Future Works

6 Summary

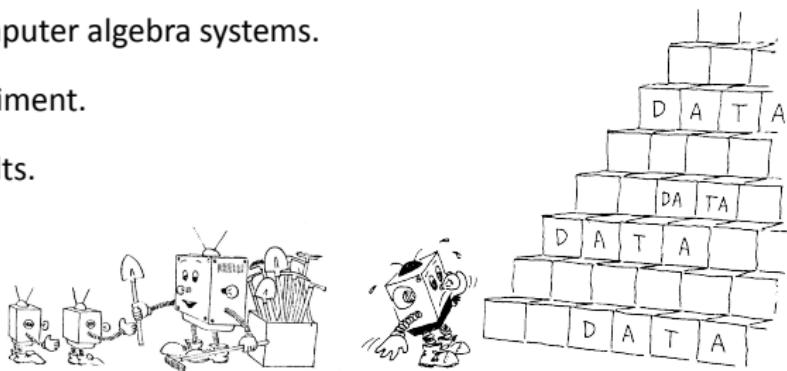
Objectives

- Improvements of the relevance of the results of our math-aware search engine.
 - MathML Normalization.
 - Classification of identifiers.
 - Context driven search.
 - Involvement of computer algebra systems.
 - Image search experiment.
 - Ranking of the results.



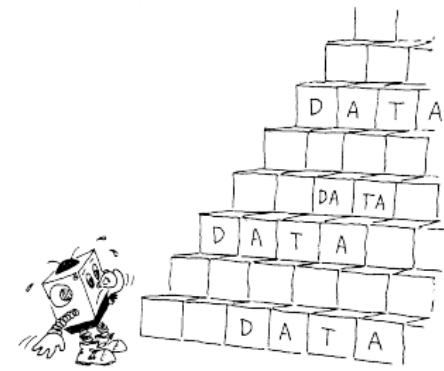
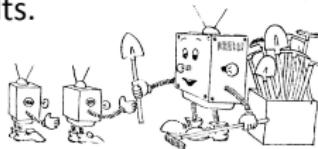
Objectives

- Improvements of the relevance of the results of our math-aware search engine.
 - MathML Normalization.
 - Canonicalization of both Presentation and Content MathML.
 - Classification of identifiers.
 - Context driven search.
 - Involvement of computer algebra systems.
 - Image search experiment.
 - Ranking of the results.



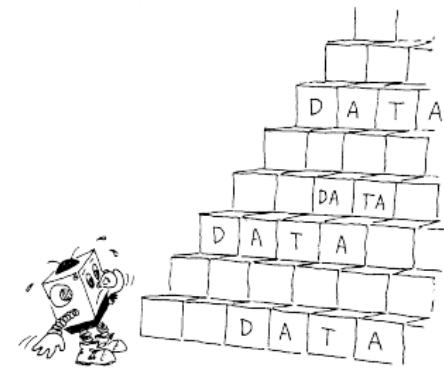
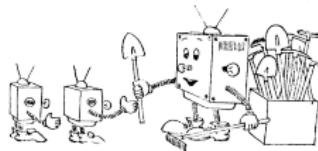
Objectives

- Improvements of the relevance of the results of our math-aware search engine.
 - MathML Normalization.
 - Classification of identifiers.
 - To mark particular identifiers as variable name, function name, and so on.
 - Start from the metadata available for the documents.
 - Context driven search.
 - Involvement of computer algebra systems.
 - Image search experiment.
 - Ranking of the results.



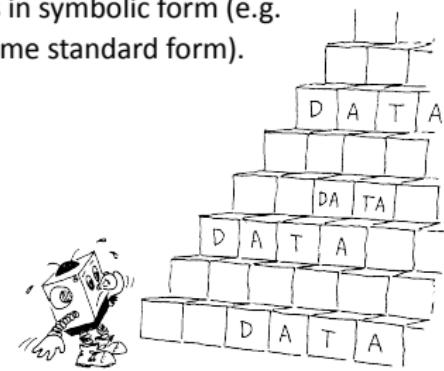
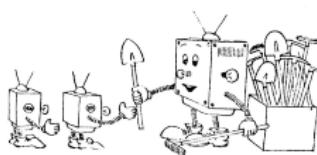
Objectives

- Improvements of the relevance of the results of our math-aware search engine.
 - MathML Normalization.
 - Classification of identifiers.
 - Context driven search.
 - Exploitation of the metadata available for the documents.
 - Exploitation of the users' inputs.
 - Involvement of computer algebra systems.
 - Image search experiment.
 - Ranking of the results.



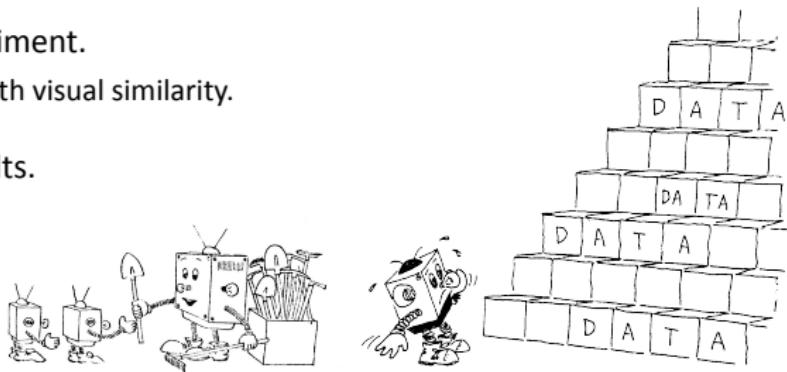
Objectives

- Improvements of the relevance of the results of our math-aware search engine.
 - MathML Normalization.
 - Classification of identifiers.
 - Context driven search.
 - Involvement of computer algebra systems.
 - Manipulation of mathematical expressions in symbolic form (e.g. simplification to a smaller expression or some standard form).
 - Image search experiment.
 - Ranking of the results.



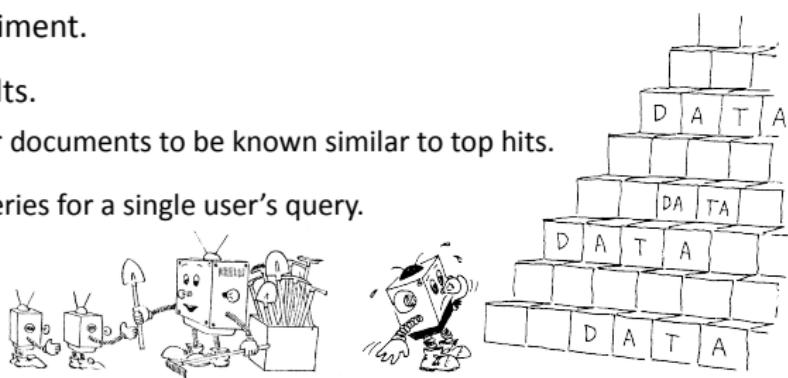
Objectives

- Improvements of the relevance of the results of our math-aware search engine.
 - MathML Normalization.
 - Classification of identifiers.
 - Context driven search.
 - Involvement of computer algebra systems.
 - Image search experiment.
 - Experiments with visual similarity.
 - Ranking of the results.



Objectives

- Improvements of the relevance of the results of our math-aware search engine.
 - MathML Normalization.
 - Classification of identifiers.
 - Context driven search.
 - Involvement of computer algebra systems.
 - Image search experiment.
 - Ranking of the results.
 - Score bonus for documents to be known similar to top hits.
 - Internal subqueries for a single user's query.



Outline

1 Motivation

2 MathML Canonicalizer

3 My Goals

4 Evaluation

5 Future Works

6 Summary

NTCIR-10 Math Task

- The first (pilot) year of the math task event last year (i.e. 2013).
- Formula search and Full-text search.
 - 4 runs submitted – differ in query language.
 - PMath – Run #1.
 - CMath – Run #2.
 - PCMath – Run #3.
 - T_EX – Run #4.
- Open Information Retrieval.
 - 1 run submitted – T_EX + text mixed queries.

NTCIR-10 Math Task Results

Table 1: Result metrics for submitted runs in Formula Search with Relevance Level ≥ 3 (Relevant)

Metric	Run 1	Run 2	Run 4
P-10 avg	0.105	0.191	0.219
P-5 avg	0.133	0.229	0.276
MAP avg	0.060	0.112	0.127
Precision	0.109 (64/589)	0.185 (92/496)	0.123 (96/778)

Table 2: Result metrics for submitted runs in Formula Search with Relevance Level ≥ 1 (Partially Relevant)

Metric	Run 1	Run 2	Run 4
P-10 avg	0.143	0.214	0.267
P-5 avg	0.181	0.267	0.343
MAP avg	0.066	0.081	0.100
Precision	0.148 (87/589)	0.232 (115/496)	0.161 (125/778)

NTCIR-11 Math Task

- A greater number of participants.
 - Increase from 6 to 8.
- Only one type of queries.
 - 50 queries, each
 - 1–4 formulae,
 - 1–4 keyphrases.
- Our results submitted for judgement in June 2014.

NTCIR-11 Math Task Approach

- Query expansion & strip-merging of subresults.
 - Query expansion.

query 1 (the original query):	f_1	f_2	k_1	k_2	k_3
query 2:	f_1	f_2	k_1	k_2	
query 3:	f_1	f_2	k_1		
query 4:	f_1	f_2			
query 5:			k_1	k_2	k_3
query 6:			k_1	k_2	k_3

NTCIR-11 Math Task Approach

- Strip-merging of subresults.
 - Example on three subqueries (the original one and two derived subqueries).

Results of the original query:

1: r_1 original
2: r_2 original
3: r_3 original
4: r_4 original
5: r_5 original
6: r_6 original
7: r_7 original
8: r_8 original
9: r_9 original
10: r_{10} original
11: r_{11} original

Results of the subquery 1:

1: r_1 subquery 1
2: r_2 subquery 1
3: r_3 subquery 1
4: r_4 subquery 1
5: r_5 subquery 1

Results of the subquery 2:

1: r_1 subquery 2
2: r_2 subquery 2
3: r_3 subquery 2
4: r_4 subquery 2
5: r_5 subquery 2

The final result list:

1: r_1 original
2: r_2 original
3: r_3 original
4: r_1 subquery 1
5: r_2 subquery 1
6: r_1 subquery 2
7: r_4 original
8: r_5 original
9: r_6 original
10: r_3 subquery 1
11: r_4 subquery 1
12: r_2 subquery 2
13: r_7 original
14: r_8 original
15: r_9 original
16: r_5 subquery 1
No more results from subquery 1.
17: r_3 subquery 2
18: r_{10} original
19: r_{11} original
No more results from the original query.
20: r_4 subquery 2
21: r_5 subquery 2
No more results from subquery 2.
22: r^1 random
23: r^2 random
...
1000: r^{979} random

NTCIR-11 Math Task Results

Table: Formulae count statistics

Documents	Formulae	
	Original	Indexed
8,301,545	59,647,566	3,021,865,236

Table: Index statistics

Indexing times [min]		Index size [GiB]
Wall Clock	CPU	
1940.0	3413.55	68

NTCIR-11 Math Task Results (cont.)

Table: Results of submitted runs with Relevance Level ≥ 3 (Relevant)

	PMath	CMath	PCMath	TeX
MAP avg	0.3073	0.3630	0.3594	0.3357
P-10 avg	0.3040	0.3520	0.3480	0.3380
P-5 avg	0.5120	0.5680	0.5560	0.5400

Table: Results of submitted runs with Relevance Level ≥ 1 (Partially Relevant)

	PMath	CMath	PCMath	TeX
MAP avg	0.2557	0.2807	0.2799	0.2747
P-10 avg	0.5020	0.5440	0.5520	0.5400
P-5 avg	0.8440	0.8720	0.8640	0.8480

NTCIR-11 Math Task: Our Investigation of Our Results

Improper query conversion

Index: \operatorname{Im}P^{+}_{\Gamma}=C_{\mu}^{+}(\Gamma)

Query: ImP^{+}_{\gamma}=C^{+}_{\mu}(\gamma)

Index

Query

...
<mrow>
 <mo>Im</mo>
 <mo></mo>
 <msup>
 ...
 </msup>

...
<mrow>
 <mi>I</mi>
 <mi>m</mi>
 <msup>
 ...
 </msup>

NTCIR-11 Math Task: Our Investigation of Our Results

Substructure difference tolerance should be improved

```
...
<mrow>
  QUERY-FORMULA-SUBPART-1
  [[ INDEX
    <mrow>
      <mo>∫ </mo>
    <mrow>
  || QUERY
    <mi>o</mi>
  ]]
  QUERY-FORMULA-SUBPART-2
  INDEX-FORMULA-SUPPLEMENT
    </mrow>
  </mrow>
  </mrow>
</mrow>
...

```

NTCIR-11 Math Task: Our Investigation of Our Results

\qvar{} handling

Original task query:

```
\qvar{S}=-\qvar{T}_{\qvar{p}}\int\qvar{d}^{\qvar{p}+1}\qvar{x}\sqrt{\qvar{g}}
```

Index: $S = -T_{p} \int d^{p+1} x \sqrt{-g}$

Query: $S = -T_{p} \int d^{p+1} x \sqrt{g}$

\sqrt{g} does not match $\sqrt{-g}$.

NTCIR-11 Math Task: Our Investigation of Our Results

Unification would be helpful

Query:

$$\text{\qvar\{x\}} \text{\frac\{\qvar\{y\}\}{\qvar\{z\}}} - \text{\qvar\{u\}} \text{\frac\{\qvar\{v\}\}{\qvar\{w\}}}$$

$$x \frac{y}{z} - u \frac{v}{w}$$

Matches:

$$\begin{aligned} \{q_s, q_r\} &= \int dx \int dy \{A(x, \mu)^p, A(y, \nu)^q\} \Big|_{\mu^{s+1} \nu^{r+1}} \\ &= pq \mu \nu \int dx A(\mu)^p (A(\nu)^q)' \left[\left\{ \frac{s}{p} \nu - \frac{r}{q} \mu \right\} \frac{1}{\mu - \nu} + \frac{1}{h} \frac{rs}{pq} \right] \Big|_{\mu^{s+1} \nu^{r+1}} \end{aligned}$$

Does not match:

$$\zeta \sim c_1 \frac{\delta \rho_\sigma}{\rho_\sigma} - c_2 \frac{\delta H_{\text{osc}}}{H_{\text{osc}}}$$

NTCIR-11 Math Task: Our Investigation of Our Results

- Combination of both formulae and text keywords in one query is important.
- Multiple subqueries derived from the original query with result lists merging turned out to be very useful.
 - One-by-one removal of the keywords and formulae.
 - “Strip-merging” of the results of the subqueries.
 - Further investigation of the best strategies needed.
 - Subqueries with subformulae?
 - Different strategy for merging of results?

Outline

1 Motivation

2 MathML Canonicalizer

3 My Goals

4 Evaluation

5 Future Works

6 Summary

Future Works

- Long-term goals:
 - All the objectives mentioned earlier. ☺
 - Whole arXiv indexing for MlaS search.
 - Evaluation.
- Short-term goals:
 - Improved strategies for MlaS internal subqueries derived from a single user's query.
 - Integration of MlaS with DML-CZ DSpace.
 - Exploitation of Gensim-Math-computed document similarities to improve of ranking results in MlaS.
 - Evaluation, Evaluation, Evaluation.

Future Works – Indexing-Time Simple Unification

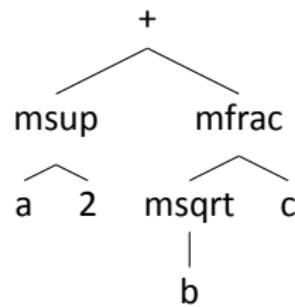
- The core idea is to implement a special identifier, let denote it \underline{U} , working as a single universal unifying element through the whole index.
- The symbol \underline{U} would be used to derive a set of unified versions of formulae from the original formula.
The derived versions are generated ‘layer-by-layer’ according to the MathML tree structure by substituting subtrees for \underline{U} . For example:

Future Works – Indexing-Time Simple Unification (cont.)

- Original formula:

$$a^2 + \frac{\sqrt{b}}{c}$$

$$a^2 + \frac{\sqrt{b}}{c}$$

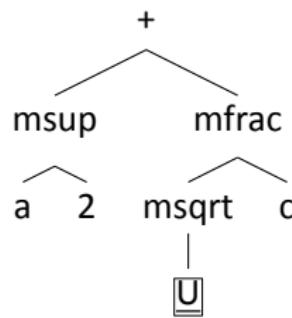


Future Works – Indexing-Time Simple Unification (cont.)

- Sequence of the unified formulae derivation – step 1:

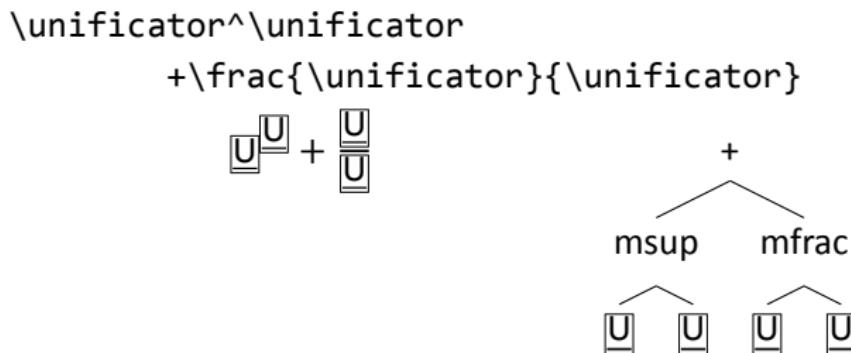
$a^2 + \frac{\sqrt{U}}{c}$

$$a^2 + \frac{\sqrt{U}}{c}$$



Future Works – Indexing-Time Simple Unification (cont.)

- Sequence of the unified formulae derivation – step 2:



Future Works – Indexing-Time Simple Unification (cont.)

- Sequence of the unified formulae derivation – step 3:

\unifier + \unifier

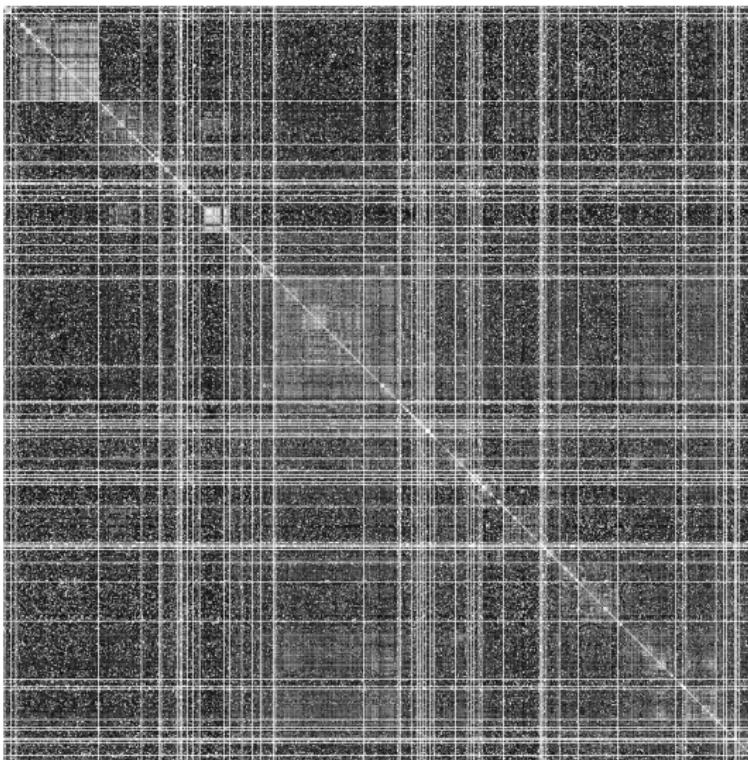


This expansion of the original formula would be applied to every formula during the indexing of the database of documents. The same method would also be used to expand every formula in the user's query to extend the original query. Hits on unified formulae would be rated with gradually decreased score according to the unification level of the formula.

Future Works – Gensim Math Document Similarities

- Gensim by Radim Řehůřek is “the most robust, efficient and hassle-free piece of software to realize unsupervised semantic modelling from plain text”: <http://radimrehurek.com/gensim/>
- We are experimenting with correlation of document similarities based on Mathematics Subject Classification (MSC) vs. document similarities based on formulae and other prominent parts (title, authors, abstract...) of the documents.
- Visualized similarity matrices.

Future Works – Gensim Math Document Similarities



Method: TfIdf-LSI; Weighted MTerms: true; MTerm Weight Conversion: 1

Future Works – Gensim Math Document Similarities

- MSC sorted documents in columns/rows, white lines separate MSC codes with different two characters (top category).
- All the documents compared each other.
 - Grayscale level indicates similarity of the document on the row to the document in the column.

white Very similar.

black Very different.

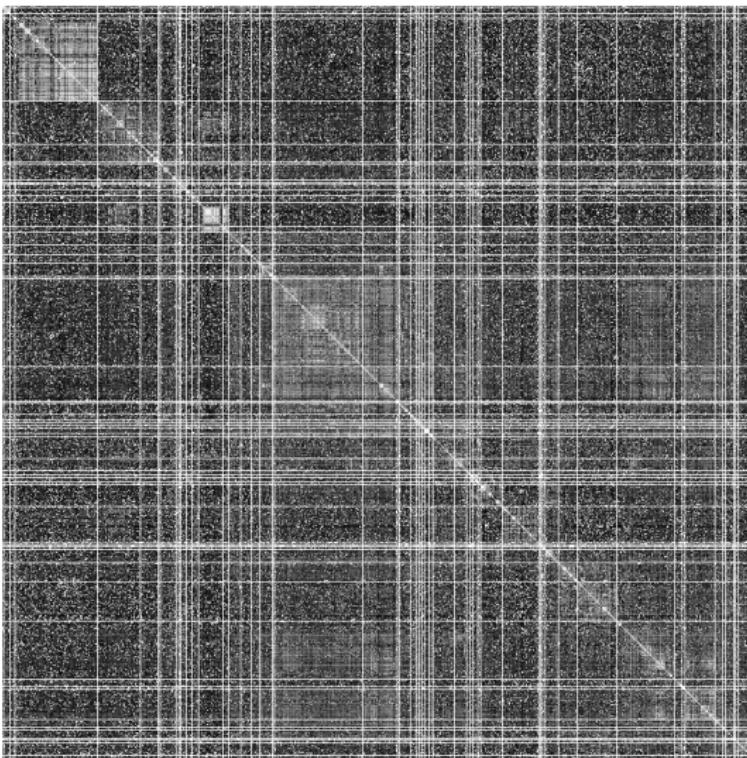
- Different weighting (or even omitting) of tokens from various metadata fields.
- Vector space model transformation method:

TfIdF-LSI Inverse Document Frequency wrapped by Latent Semantic Indexing

LDA Latent Dirichlet Allocation

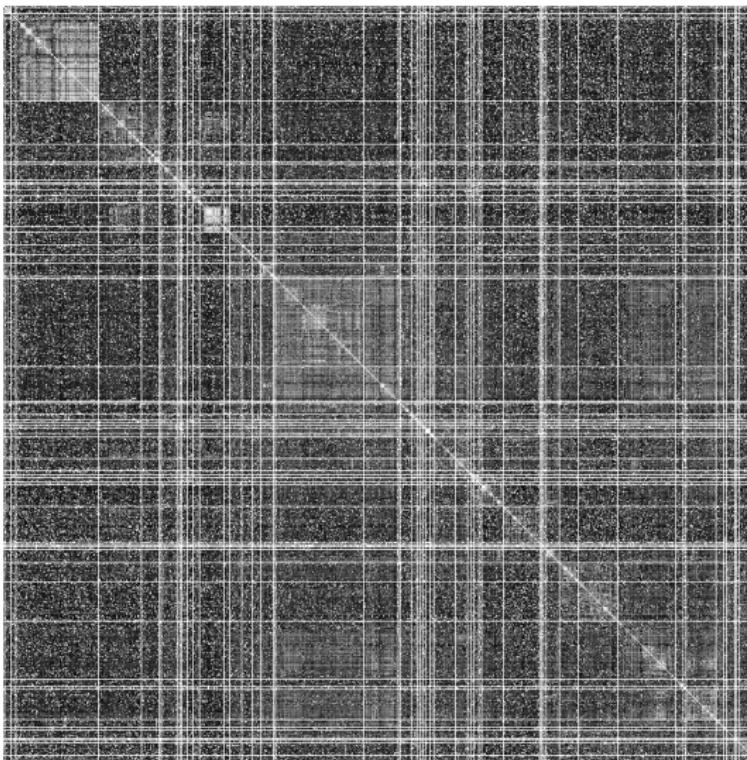
... ...

Future Works – Gensim Math Document Similarities



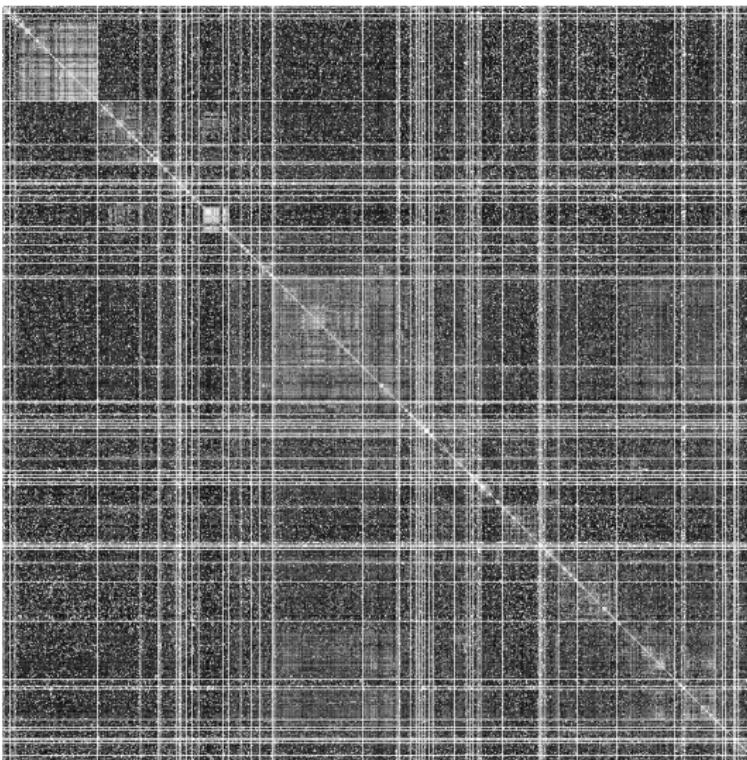
Method: TfIdf-LSI; Weighted MTerms: true; MTerm Weight Conversion: 1

Future Works – Gensim Math Document Similarities



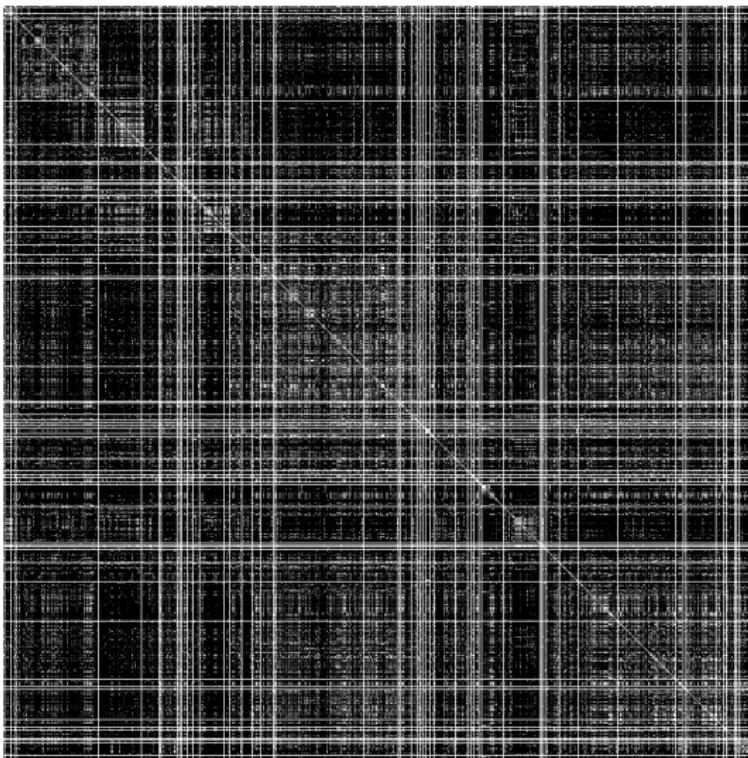
Method: TfIdf-LSI; Weighted MTerms: true; MTerm Weight Conversion: trunc(3.9 * mtermWeight)

Future Works – Gensim Math Document Similarities



Method: TfIdf-LSI; Weighted MTerms: false

Future Works – Gensim Math Document Similarities



Method: LDA; Weighted MTerms: false

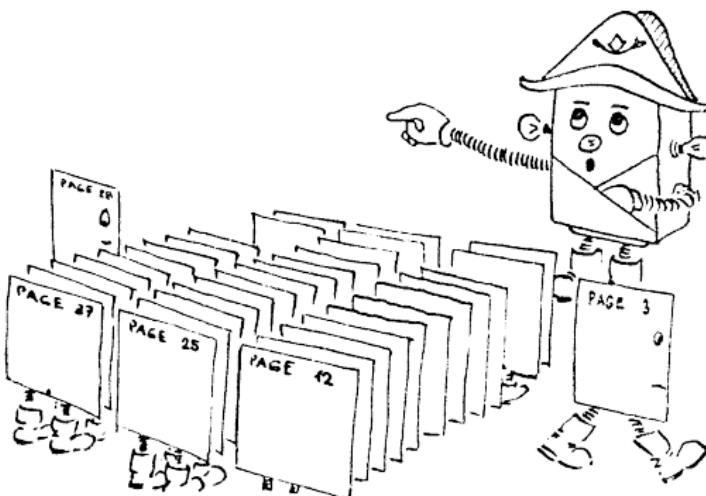
Future Works – Gensim Math Document Similarities

- Gensim by Radim Řehůřek is “the most robust, efficient and hassle-free piece of software to realize unsupervised semantic modelling from plain text”: <http://radimrehurek.com/gensim/>
- We are experimenting with correlation of document similarities based on Mathematics Subject Classification (MSC) vs. document similarities based on formulae and other prominent parts (title, authors, abstract...) of the documents.
- Visualized similarity matrices.
 - How to compute similarity of these matrices rigorously?
 - Variance of the normalized similarity inside each of the MSC regions + weighted average over the matrix?

Summary

- We have our own math-aware search engine.
- We have a lot of possible improvements in our minds.
- We are interested in evaluation a lot.

Questions?





Illustrations by Jiří Franek.



SOJKA, Petr and Martin LÍŠKA. The Art of Mathematics Retrieval. In Matthew R. B. Hardy, Frank Wm. Tompa. Proceedings of the 2011 ACM Symposium on Document Engineering. Mountain View, CA, USA: ACM, 2011. p. 57–60. ISBN 978-1-4503-0863-2. doi:10.1145/2034691.2034703.



LÍŠKA, Martin, Petr SOJKA and Michal RŮŽIČKA. Similarity Search for Mathematics: Masaryk University team at the NTCIR-10 Math Task. In Noriko Kando, Kazuaki Kishida. Proceedings of the 10th NTCIR Conference on Evaluation of Information Access Technologies. Tokyo: National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo 101-8430 Japan, 2013. s. 686–691, 6 s. ISBN 978-4-86049-062-1.



LÍŠKA, Martin, Petr SOJKA, Michal RŮŽIČKA and Peter MRAVEC. Web Interface and Collection for Mathematical Retrieval : WebMLaS and MREC. In Petr Sojka, Thierry Bouche. DML 2011: Towards a Digital Mathematics Library. Brno: Masaryk University, 2011. p. 77–84. ISBN 978-80-210-5542-1.



FORMÁNEK, David, Martin LÍŠKA, Michal RŮŽIČKA and Petr SOJKA. Normalization of Digital Mathematics Library Content. CEUR Workshop Proceedings, Aachen, 2012, vol. 921, October, p. 91–103. ISSN 1613-0073.



LÍŠKA, Martin, Petr SOJKA, Michal RŮŽIČKA and Peter MRAVEC. Web Interface and Collection for Mathematical Retrieval : WebMLaS and MREC. In Petr Sojka, Thierry Bouche. DML 2011: Towards a Digital Mathematics Library. Brno: Masaryk University, 2011. p. 77–84. ISBN 978-80-210-5542-1.



ŘEHŮŘEK, Radim and Petr SOJKA. Software Framework for Topic Modelling with Large Corpora. In Proceedings of LREC 2010 workshop New Challenges for NLP Frameworks. Valletta, Malta: University of Malta, 2010. p. 46–50. ISBN 2-9517408-6-7.