

# Similarity Search for Mathematics: Masaryk University team at the NTCIR-10 Math Task

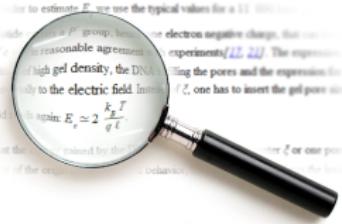
Martin Líška, Petr Sojka, Michal Růžička

Masaryk University, Faculty of Informatics, Brno, Czech Republic

June 21st, 2013

In order to estimate  $E_c$  we use the typical values for a  $\text{DNA}$  molecule. Inside a  $\text{P}_\text{v}$  group, because of electron negative charge, the concentration of  $\text{DNA}$  is about  $1 \text{ mg}/\text{ml}$ . It's in reasonable agreement with experiments/[22, 23]. The expression for  $E_c$  is valid only for high gel density, the DNA molecule is completely adsorbed onto the pores and the expression for  $E_c$  is valid only for the electric field. Instead of  $E_c$ , one has to insert the gel pore size and again:  $E_c \simeq 2 \frac{k_\text{B} T}{q \ell}$ .

that the pores are caused by the gel pores. After  $\xi$  or one pore size of the original  $\text{DNA}$  molecule, the density of the pores is



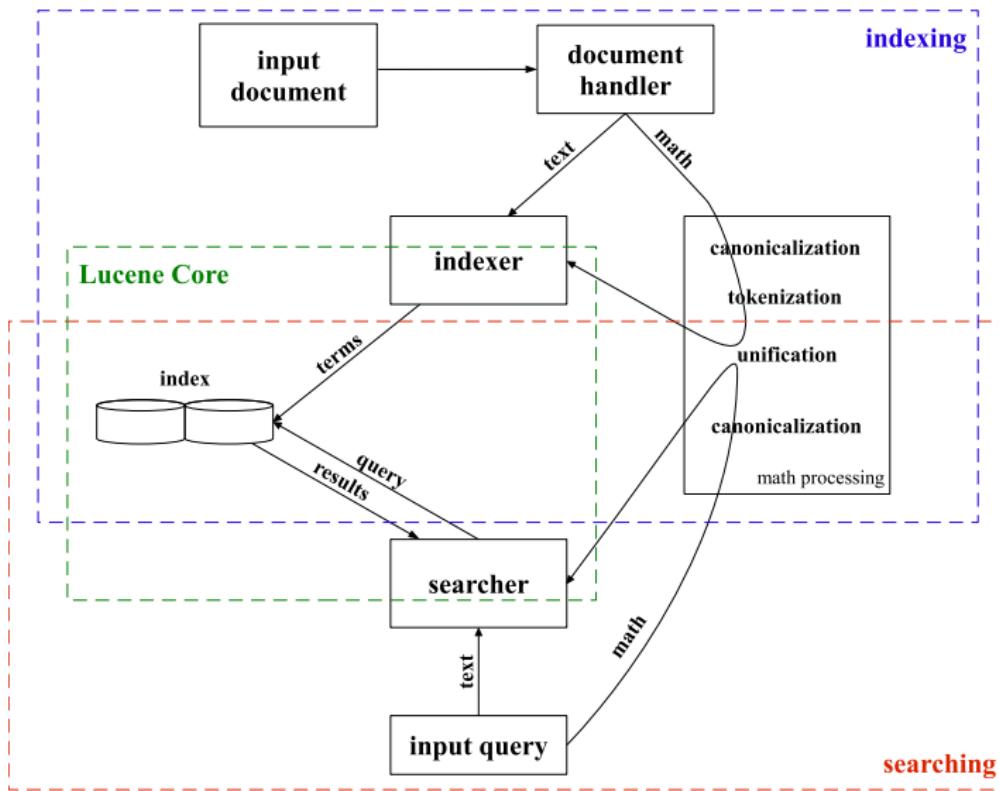
# Motivation

- Digital Mathematics Libraries – EuDML, DML-CZ, ...
- Knowledge in mathematical expressions
- Example
- MlaS evaluation
  - How well does it work?

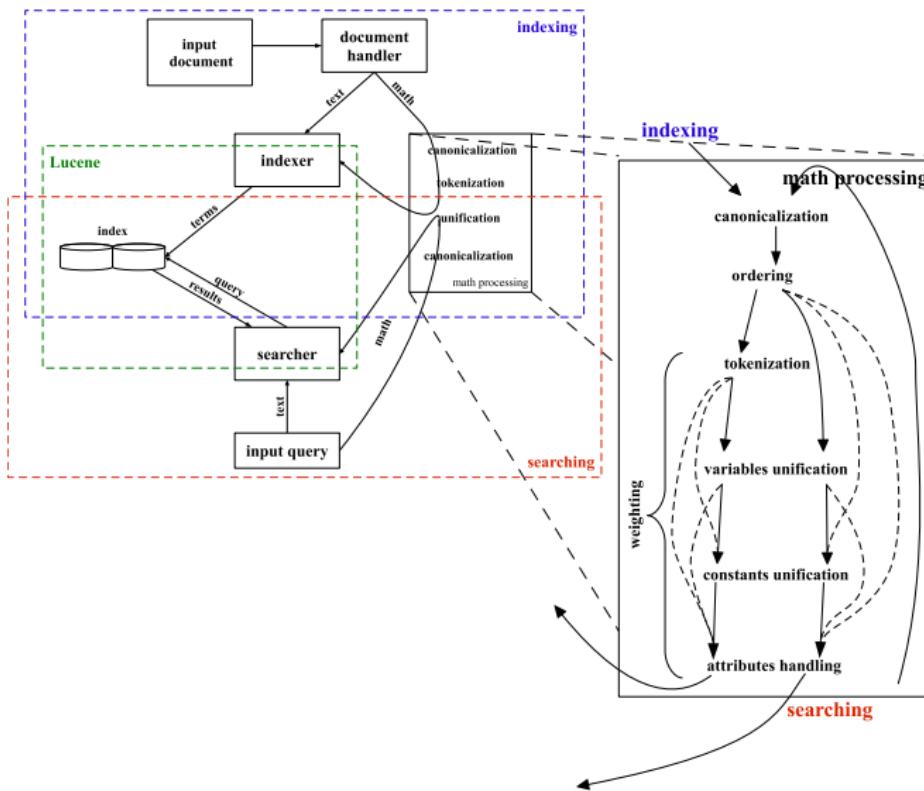
# Math Indexer and Searcher (MiAS)

- Similarities of math expression trees (*both* structural (PMath) and semantic (CMath))
- Full-text based, document oriented, math-aware search engine
- Mixed text-math indexing and searching
- Query language: MathML,  $\text{\LaTeX}$
- Results ordering based on similarity level
- Snippets with hit highlighting
- On-the-fly query formula rendering

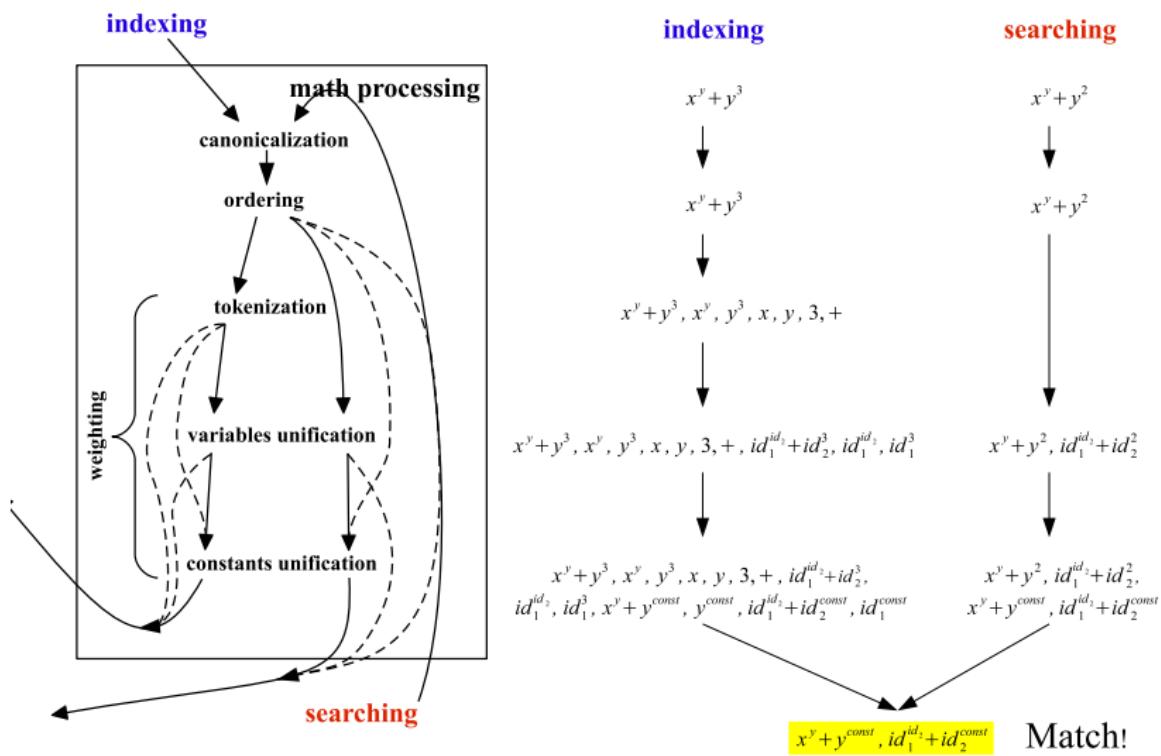
# Design



# Design II



# Example



# MlaS – Math Task

- Formula search and Full-text search
  - 4 runs submitted – differ in query language
    - PMath – Run #1
    - CMath – Run #2
    - PCMath – Run #3
    - T<sub>E</sub>X – Run #4
- Open Information Retrieval
  - 1 run submitted – T<sub>E</sub>X+text mixed queries

# Querying

- Automatic querying scripts, generated:
  - trec\_eval format
  - XML
  - Clickable HTML for investigation and fine tuning

## Querying

Results for 'IToX' run with query ID 'INTCIR10\_55\_31' in Index '31' (XML response)

**Results for 'RMath' run with query ID 'INTCIR10-ET-1' in index '0' (XML response)**

[View Details](#) | [Edit](#) | [Delete](#) | [Print](#) | [Email](#) | [Share](#) | [Report](#)

This XML file does not appear to have any style information associated with it. The document tree is shown below.

Similarity Search for Mathematics: Masaryk University team at the NTCIR-10 Math Task

# Results

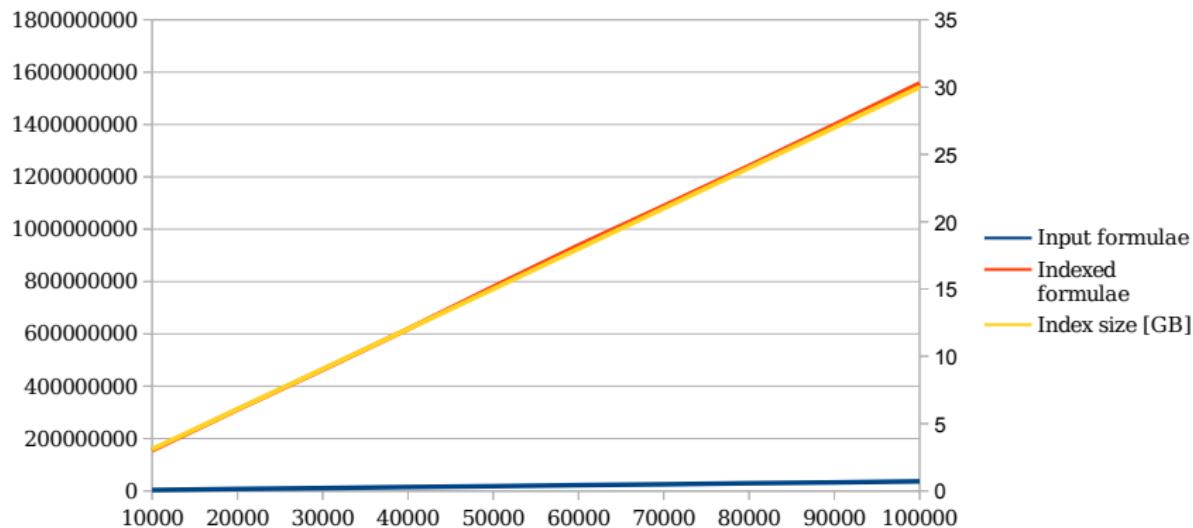
Table 1: Result metrics for submitted runs in Formula Search with Relevance Level  $\geq 3$  (Relevant)

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

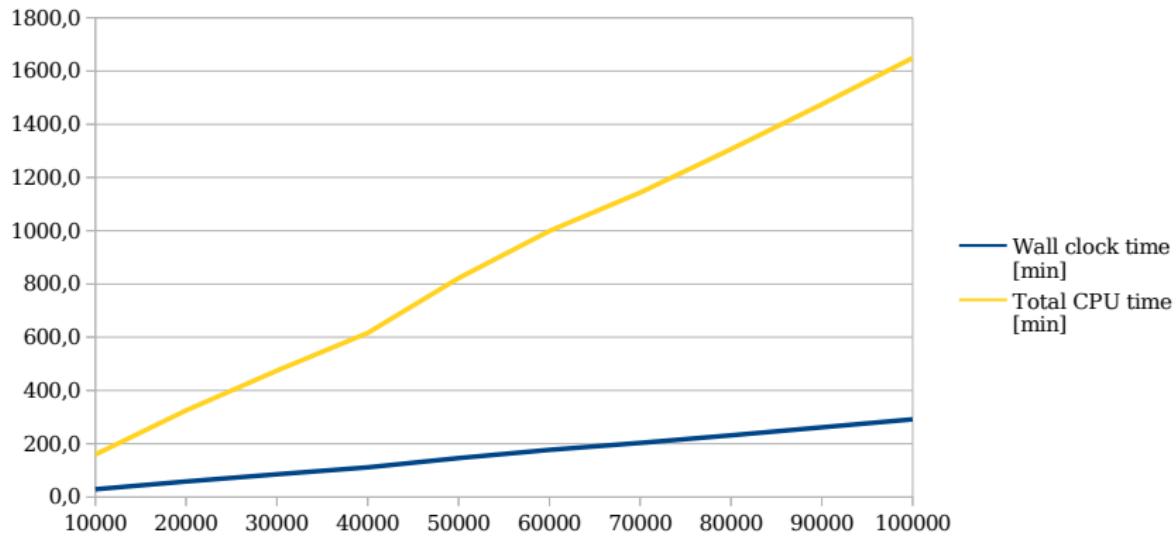
Table 2: Result metrics for submitted runs in Formula Search with Relevance Level  $\geq 1$  (Partially Relevant)

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

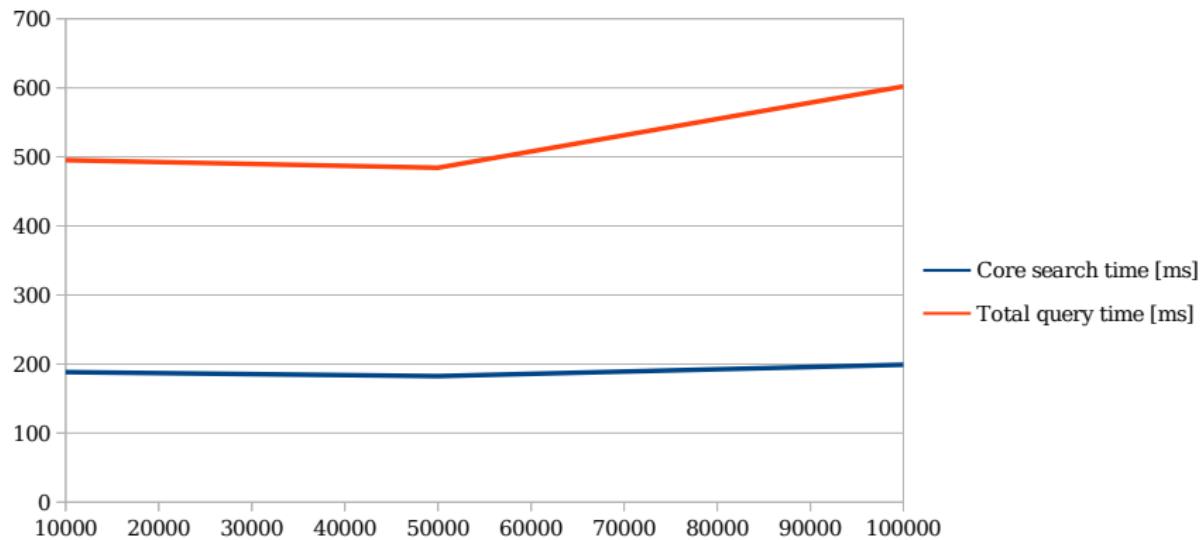
# Efficiency



# Efficiency



# Efficiency



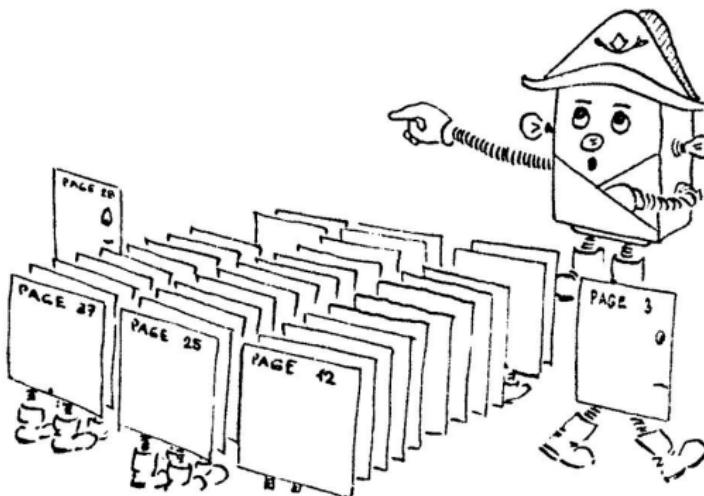
# Conclusions

- Evaluation is driving the development
- Our participation to NTCIR shown useful
- $\text{\TeX}$  queries (both PMath and CMath used) most precise
- Deeper investigation of the results

## Future work

- Own evaluation framework (evaluation driven development)
- Canonicalization
- Experiment with weight factors and formulae relevance weighting
- Content (aka semantics) focused search
- Analysis from running system (EuDML) and its logs

# Questions?



-  Liška, Martin and Petr Sojka and Michal Růžička. Similarity Search for Mathematics: Masaryk University team at the NTCIR-10 Math T task. In Proceedings of the 10th NTCIR Workshop Meeting on Evaluation of Information Access Technologies: Math Pilot Task
-  D. Formánek, M. Liška, M. Růžička, and P. Sojka. Normalization of digital mathematics library content. In J. Davenport, J. Jeuring, C. Lange, and P. 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.
-  Sojka, Petr and Martin Liška. 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, 4 pp. ISBN 978-1-4503-0863-2. <<http://dx.doi.org/10.1145/2034691.2034703>>
-  Sojka, P., Liška, M.: Indexing and Searching Mathematics in Digital Libraries – Architecture, Design and Scalability Issues. In: Davenport, J.H., Farmer, W., Urban, J., Rabe, F., (eds.) Proceedings of CICM Conference 2011 (Calculemus/MKM). Lecture Notes in Artificial Intelligence, LNAI, vol. 6824, pp. 228–243. Springer-Verlag, Berlin, Germany (Jul 2011), <[http://dx.doi.org/10.1007/978-3-642-22673-1\\_16](http://dx.doi.org/10.1007/978-3-642-22673-1_16)>
-  Sojka, P. (ed.): Towards a Digital Mathematics Library. Masaryk University, Paris, France (Jul 2010), <<http://www.fi.muni.cz/sojka/dml-2010-program.html>>
-  Sylwestrzak, W., Borbinha, J., Bouche, T., Nowiński, A., Sojka, P.: EuDML—Towards the European Digital Mathematics Library. In: Sojka [5], pp. 11–24, <<http://dml.cz/dmlcz/702569>>
-  Martin Liška, Petr Sojka, Michal Růžička, and Petr Mravec.  
**Web Interface and Collection for Mathematical Retrieval.**  
In Petr Sojka and Thierry Bouche, editors, *Proceedings of DML 2011*, pages 77–84, Bertinoro, Italy, July 2011. Masaryk University. <<http://www.fi.muni.cz/sojka/dml-2011-program.html>>.