Overview
○○○○○○○○○○○○○○○○○○

Motivation
○○○○○○○○

DML-CZ
○○○○○○○○○

Math Indexer and Searcher
○○○○○○

Technologies
○○○○○○

Tools
○○○○○○○○○○○○○○○○○○

Summary
○○○○○

# Accessibility Issues in a Digital Mathematical Library

## Examples and Experience of DML-CZ and EuDML

### Petr Sojka

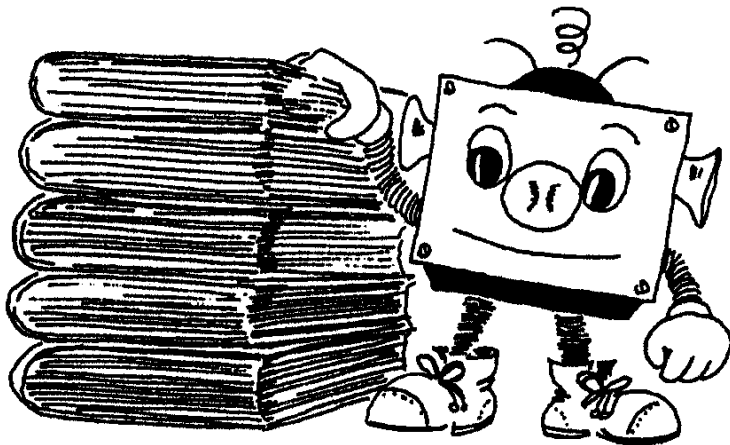Masaryk University, Faculty of Informatics, Brno, Czech Republic
<sojka@fi.muni.cz>

Universal Learning Design 2011, Brno, Czech Republic

February 10th, 2011, 9:50AM

## Outline and two take-home messages

1 Pictorial overview

2 Motivation, vision of PubMed Central for Mathematics

3 Complexity of digitization workflow of The Czech Digital Mathematics Library DML-CZ

4 Math Indexer and Searcher

5 Accessibility improving technologies and tools for DML-CZ and EuDML

6 Tools developed
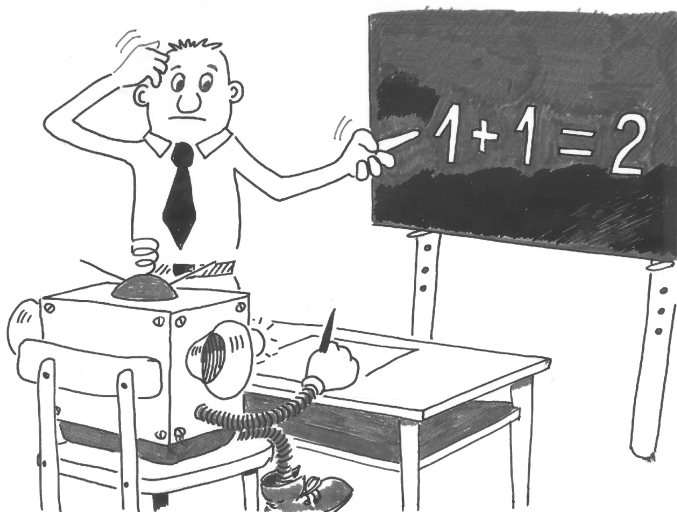
7 Summary, conclusions and future work

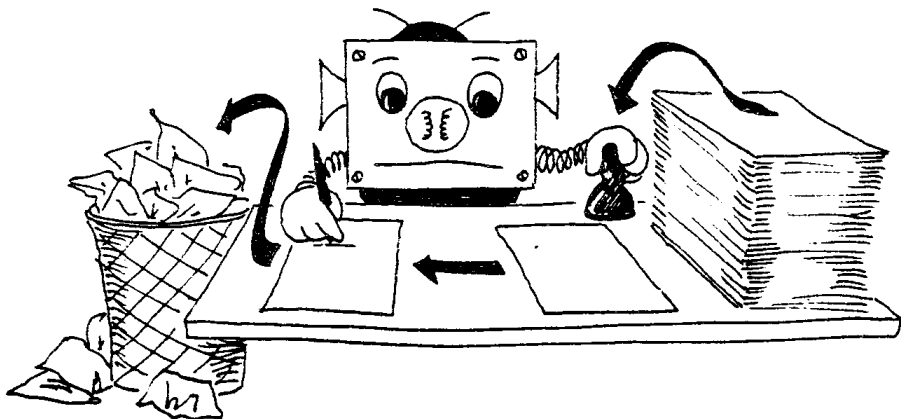# Accessibility issues in *digital* libraries

# Information overload in globalized *scientific* world

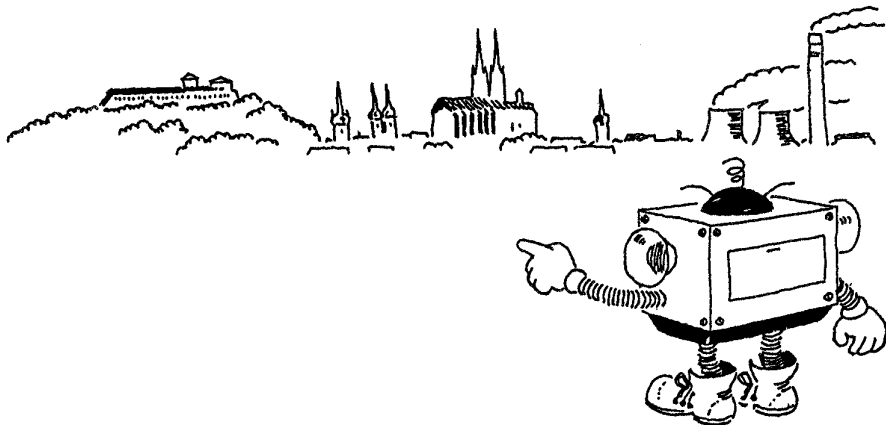# Universal learning design using digital *mathematics*

# From paper to digital *workflow*

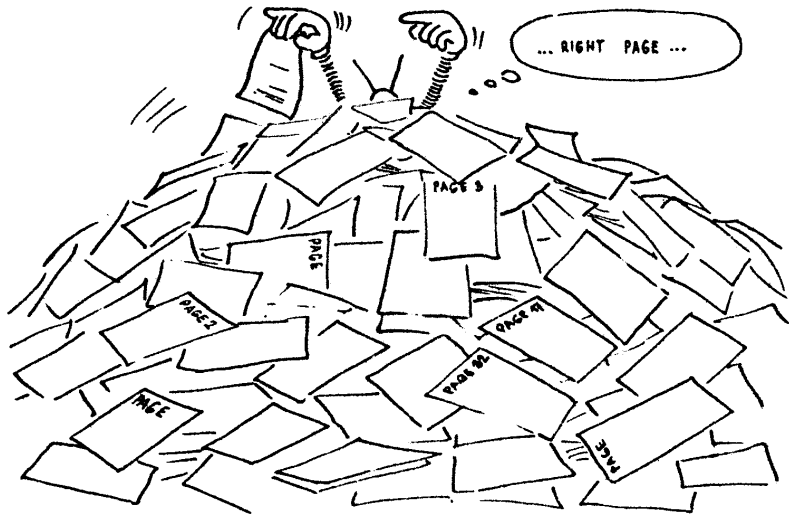# Retro-digitization, *accessible* digital library development

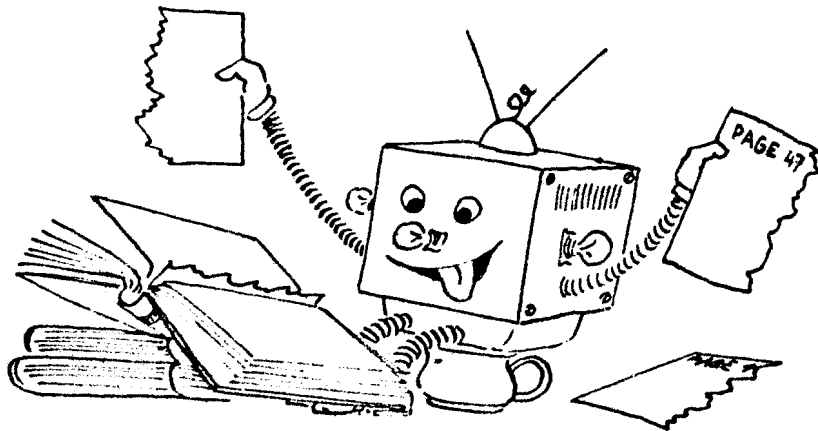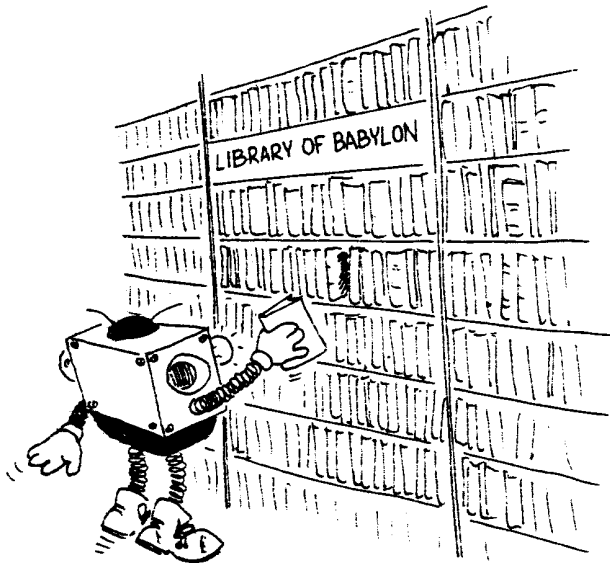# Experiences from project *DML-CZ* (Brno, CZ)

# DML-CZ: new *workflows* and math data indexing

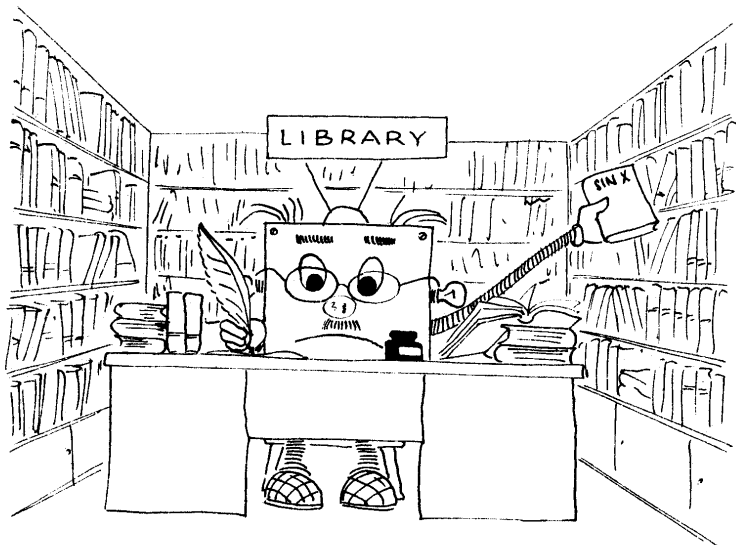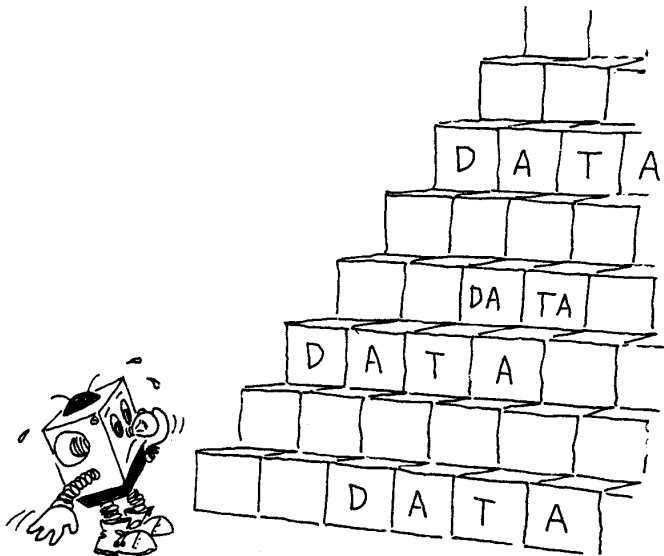## New approaches to *math document retrieval*

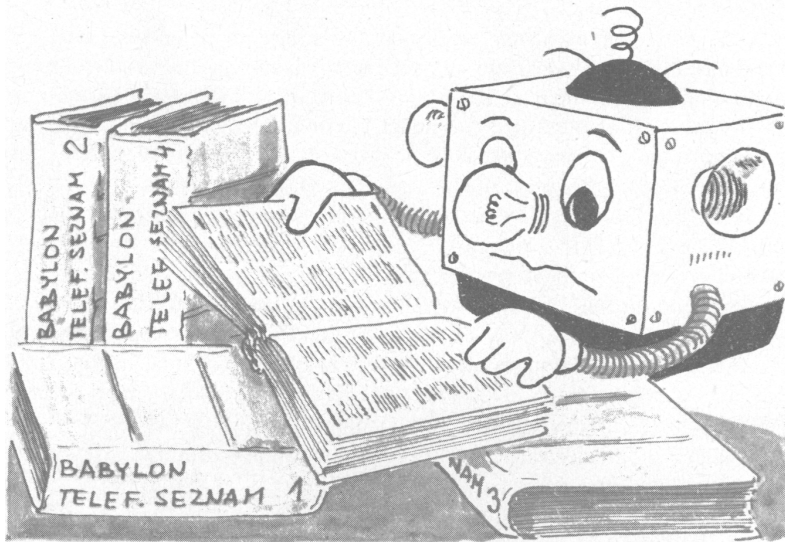# 'Bottom up' deployment towards EU or *worldwide scale*

# The European Digital Mathematics Library: *EuDML*

# EuDML: from local data collections to the virtual DL

## Tools for *automated math extraction* from PDF

# Yes, you can! You can have highly automated workflow to get accessible math, search, visibility, scalability,…

# End of talk overview

## Decade of the vision of WDML as PubMed 4 Math

In the beginning was vision of all mathematical knowledge, *peer reviewed, verified* (100,000,000 pages) and engineered into one-stop e-shop/DL.

Progress of IT, connectivity, cheap storage, new information retrieval technologies (Google).

AMS supported NSF preparation grant (in 2003) for WDML—Worldwide digital mathematics library, planned to be funded by de Moore foundation ($100,000,000 requested). Application was *not* successful.

## Vision of European Digital Mathematics Library

Even other attempts on the European level (FP5, FP6) were not successful. Finally three year project or *European Digital Mathematics Library, EuDML* (programme EU CIP-ICT-PSP, type Pilot B, EU contribution *(1.6 MEur, 50% of total budget only)* started from February 2010. The strategy of

**$\mathcal{E}u$DML**

*The* **EUROPEAN DIGITAL MATHEMATICS LIBRARY** is:

- to master the technology, develop tools and offer them;
- concept of *moving wall* to motivate and engage commercial publishers without Open Access bussiness model;
- to collect data (from existing local or publisher's) *digital libraries* into 'one-stop shop' and achieve critical mass in the domain → 'a must/me too' effect then as with PubMed.

## EuDML as a virtual library portal

EuDML will be a *virtual* library based on data from smaller data providers, DLs and publishers:

# European Digital Mathematics Library

## Bottom up—from building bricks of regional repositories

As DML content providers serve mostly publisher's or regional DML repositories as The Czech Digital Mathematics Library DML-CZ or NUMDAM, DML-PL, DML-PT, RusDML,…: aggregating content from local repositories to build the bigger (global?) DML.

Example of DML-CZ: up and running digital mathematics library `<http://dml.cz>` with 30,000+ papers (300,000+ pages). For more, see (who, what, browse, browse similar, how to search).

# From paper to digital processing, from local to the global DML

# DML-CZ workflow

## Challenges of Math handling: OCR, indexing, search…

# DML-CZ—data: scientific math published in CZ/SK

Proof. Let $\mathring{K}$ be a cube, $\mathring{K} \subset \mathring{O}$; put $K = \varphi^{-1}(\mathring{K})$. According to theorem 50 we have $K \in \mathfrak{A}$ and it follows from theorem 24 that

$$P(K, v) = \int_K f(x)\, dx. \tag{89}$$

The functional determinant $T$ of the mapping $v = \varphi^{-1}$ fulfils the relation $T(\varphi(x)) \cdot \det M(x) = 1$, so that

$$\int_K f(x)\, dx = \int_{\mathring{K}} f(\psi(y)) \cdot |T(y)|\, dy = \int_{\mathring{K}} \mathring{f}(y)\, dy. \tag{90}$$

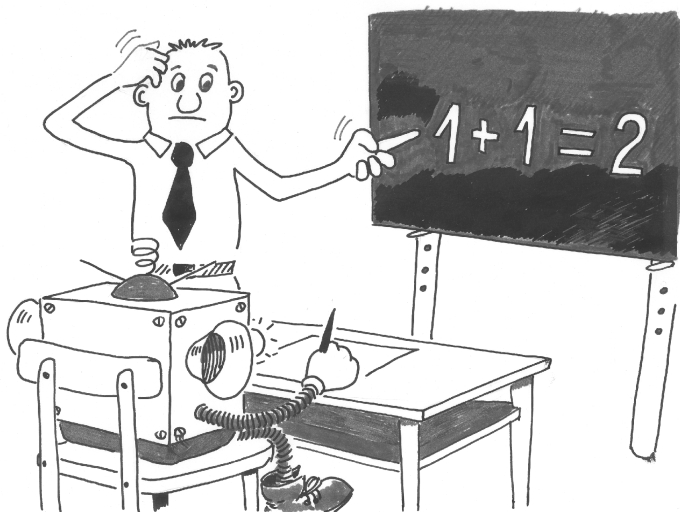From theorem 50 (and relation (86)) we see that $P(K, v) = P(\mathring{K}, \mathring{v})$; relations (89), (90) show therefore that $P(\mathring{K}, \mathring{v}) = \int_{\mathring{K}} \mathring{f}(y)\, dy$, which completes the proof.

Remark. The reader may compare this paper with [6].

REFERENCES

[1] *V. Jarník:* Diferenciální počet, Praha 1953.
[2] *V. Jarník:* Integrální počet II, Praha 1955.
[3] *J. Mařík:* Vrcholy jednotkové koule v prostoru funkcionál na daném polouspořádaném prostoru, Časopis pro pěst. mat., 79 (1954), 3—40.
[4] *Ян Маржик (Jan Mařík):* Представление функционала в виде интеграла, Чехословацкий мат. журнал, 5 (80), 1955, 467—487.
[5] *J. Mařík:* Plošný integrál, Časopis pro pěst. mat., 81 (1956), 79—82.
[6] *Ян Маржик (Jan Mařík):* Заметка к теории поверхностного интеграла, Чехословацкий мат. журнал, 6 (81), 1956, 387—400.
[7] *S. Saks:* Theory of the integral, New York.

Резюме

ПОВЕРХНОСТНЫЙ ИНТЕГРАЛ

ЯН МАРЖИК (Jan Mařík), Прага.

(Поступило в редакцию 10/X 1955 г.)

Пусть $m$ — натуральное число; пусть $E_m$ — $m$-мерное евклидово пространство. Для всякого ограниченного измеримого множества $A \subset E_m$ положим $\|A\| = \sup \int_{A} \sum_{i=1}^{m} \frac{\partial v_i(x)}{\partial x_i}\, dx$, где $v_1, \ldots, v_m$ — многочлены такие, что $\sum v_i^2(x) \leq 1$ для всех $x \in A$. Пусть $\mathfrak{A}$ — система всех ограниченных измеримых множеств $A$, для которых $\|A\| < \infty$. Теорема 18 тогда утверждает: *Пусть* $A \in \mathfrak{A}$; *пусть* $D$ — *граница множества* $A$. *Тогда на системе* $\mathfrak{B}$ *всех борелевских подмножеств множества* $D$ *существует мера* $\mu$ *и на*

557

ИОСИФ ВИССАРИОНОВИЧ СТАЛИН
1879—1953

# Document accessibility 4 DML processing challenges

Data heterogenity, plethora of formats, validation and conversions:

retro-digital period:  scanning, geometrical transformations (BookRestorer),
OCR (FineReader, InftyReader), *two-layer PDF*

retro-born-digital period:  not complete .tex or .dvi data, bad formats, bitmap
fonts of low resolution: PDF2Math (PDF2NLM?)

born-digital period:  typesetting by TeX with export of [meta]data into digital
library

## Document accessibility 4 DML processing challenges

Data heterogenity, plethora of formats, validation and conversions:

retro-digital period:  scanning, geometrical transformations (BookRestorer), OCR (FineReader, InftyReader), *two-layer PDF*

retro-born-digital period:  not complete .tex or .dvi data, bad formats, bitmap fonts of low resolution: PDF2Math (PDF2NLM?)

born-digital period:  typesetting by T$_E$X with export of [meta]data into digital library



**RETRO−DIGITISED**
Printed Document → Image Capture → Image → OCR → Printed Data

**RETRO−BORN DIGITAL**
PDF/PS → Extraction → Printed Data

Printed Data → Analysis → Meta−data and Document

Refinement → Meta−data and Document

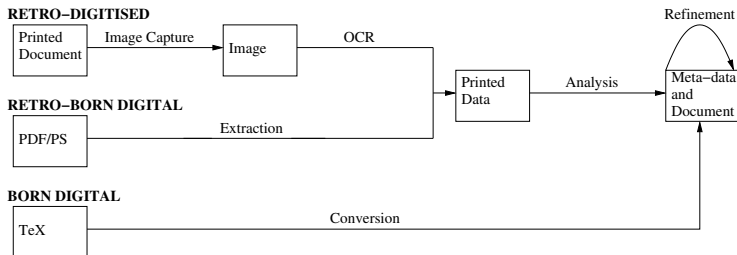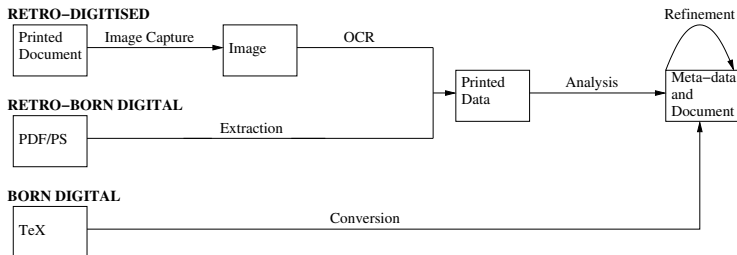**BORN DIGITAL**
TeX → Conversion → Meta−data and Document

## Document accessibility 4 DML processing challenges

Data heterogenity, plethora of formats, validation and conversions:

retro-digital period:  scanning, geometrical transformations (BookRestorer), OCR (FineReader, InftyReader), *two-layer PDF*

retro-born-digital period:  not complete .tex or .dvi data, bad formats, bitmap fonts of low resolution: PDF2Math (PDF2NLM?)

born-digital period:  typesetting by T$_E$X with export of [meta]data into digital library



**RETRO–DIGITISED**

| Printed Document | → Image Capture → | Image | → OCR → |

**RETRO–BORN DIGITAL**

| PDF/PS | → Extraction → |

| Printed Data | → Analysis → | Meta–data and Document |

Refinement

**BORN DIGITAL**

| TeX | → Conversion → |

# MathML or LaTeX? MathML and LaTeX!

Data heterogenity, plethora of formats, validation and conversions:

world of authors: LaTeX, TeX notation of mathematics

world of applications/data exchange: XML, *MathML*

Big volumes: $\rightarrow$ high *automation* to save costs

Exchange on the web—W3C standards: MathML, WAI-ARIA (Web Accessibility Initiative—Acessible Rich Internet Applications), WCAG (Web Content Accessibility Guidelines) 2.0 or Dogma W4.

By converting to MathML to allow discoverability and indexing (formulae fuzzy search).

# MathML or LaTeX? MathML and LaTeX!

Data heterogenity, plethora of formats, validation and conversions:

world of authors:  LaTeX, TeX notation of mathematics

world of applications/data exchange:  XML, *MathML*

Big volumes: $\rightarrow$ high *automation* to save costs

Exchange on the web—W3C standards: MathML, WAI-ARIA (Web Accessibility Initiative—Acessible Rich Internet Applications), WCAG (Web Content Accessibility Guidelines) 2.0 or Dogma W4.

By converting to MathML to allow discoverability and indexing (formulae fuzzy search).

# DML-CZ document engineering—data processing

## DML-CZ challenges and lessons learned

DML-CZ, the Czech Digital Mathematics Library, now serves almost *300,000 pages of 30,000 math papers*. Challenges were

- *migration of existing workflows (retro-digital, retro-digital and born-digital) into the repository*
- negotiations with Google Scholar towards better visibility
- math indexing and search
- alternative visualization
- space and processing demands
- ….

DML-CZ is according to The Ranking Web of World Repositories the best repository in CZ, 91. in EU and 203. in the world.

## Math Search and Indexing

- Usual way of seeking information via [Google] search

- Conventional searching approaches are not applicable for math

- Usage of existing mathematical search engines (MathDex, EgoMath, LaTeXSearch, LeActiveMath, MathWebSearch) problematic

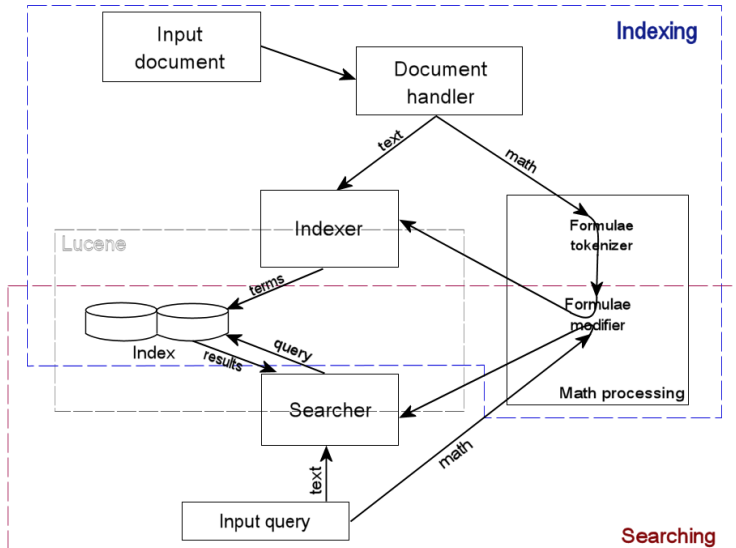- new Math Indexer and Searcher (MIaS) is being developed at MU

# Comparison

| | Input documents | Internal representation | Used converters | Approach | $\alpha$-eq. | Query language | Queries | Indexing core |
|---|---|---|---|---|---|---|---|---|
| MathDex | HTML, TEX/LaTEX, Word, PDF | Presentation MathML (text) | jtidy, blahtex, LaTeXML, Hermes, Word+Math-Type, pdf2tiff->Infty | syntactic | ✗ | ? | text, math, mixed | Apache Lucene |
| LeActiveMath | OMDoc, OpenMath | OpenMath (text) | - | syntactic | ✗ | OpenMath (palette editor) | text, math, mixed | Apache Lucene |
| LaTEXSearch | LaTEX | LaTEX(text) | - | syntactic | ✗ | LaTEX | titles, math, DOI | ? |
| MathWeb Search | Presentation MathML, Content MathML, OpenMath | Content MathML, OpenMath (substitution trees) | - | semantic | ✔ | QMath, LaTEX, Mathematica, Maxima, Maple, Yacas styles (palette editor) | text, math, mixed | Apache Lucene (for text only) |
| EgoMath | Presentation MathML, Content MathML, PDF | Presentation MathML (text) | Infty | mixed | ✗ | LaTEX | text, math, mixed | EgoThor |

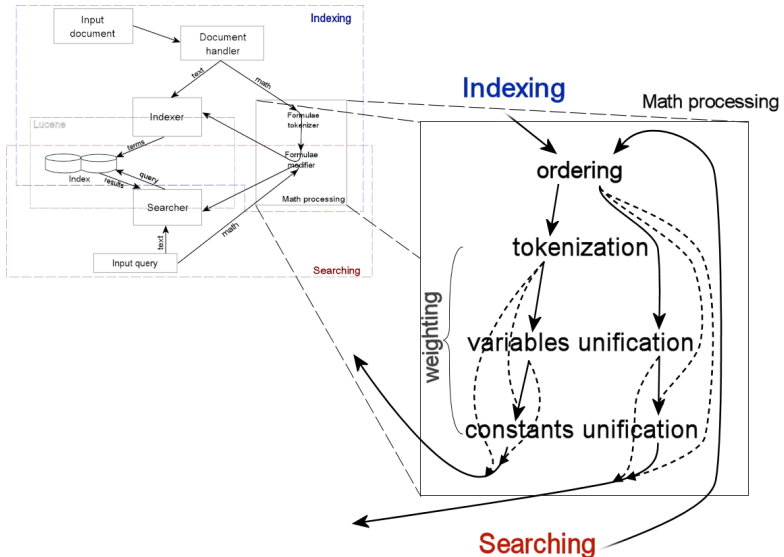## Math Indexer and Searcher — Features

- Inspired mostly by MathDex and EgoMath

- Based on full text core Apache Lucene

- Presentation MathML

- Allows similarity (not only exact match) between query and matched term
    - Commutativity
    - Unification of variables and constants
    - Subformulae matching

- Level of similarity calculation for expressions

- Mixed mathematical-textual queries

# Math Indexer and Searcher — Design

# Math Indexer and Searcher — Design II

## Formula Processing Example

input: $(a + b^{2+c}, 1)$

$\downarrow$ ("mi"+"mn" $\Rightarrow$ 2⊆c)

arranged: $(a + b^{c+2}, 1)$

tokenization: $(a, 0.5)$ $(+, 0.5)$ $(b^{c+2}, 0.5)$

$(b, 0.25)$ $(c + 2, 0.25)$

$(c, 0.125)$ $(+, 0.125)$ $(2, 0.125)$

variables unification: $(id_1 + id_2^{id_3+2}, 0.8)$ $(id_1^{id_2+2}, 0.4)$ $(id_1 + 2, 0.2)$

constants unification: $(a + b^{c+const}, 0.8)$ $(b^{c+const}, 0.4)$ $(c + const, 0.2)$

$(id_1 + id_2^{id_3+const}, 0.64)$ $(id_1^{id_2+const}, 0.32)$ $(id_1 + const, 0.16)$

## Implementation

- Java

- Lucene 3.0.1

- jTidy for text extraction

- Mathematical part implements Lucene's interface Tokenizer — able to integrate to any Lucene based system as YADDA, DSpace,…

## Evaluation

- Math corpus from arXMLiv in XHTML + MathML
    - 324,060 real math documents
    - Uncompressed corpus size 53 GB, ZIP compressed 6.7 GB
    - 112 million input formulae

- Indexed
    - Produced over 2 billion math index expressions
    - Index size 45 GB

- Simple demo web interface: WebMIaS
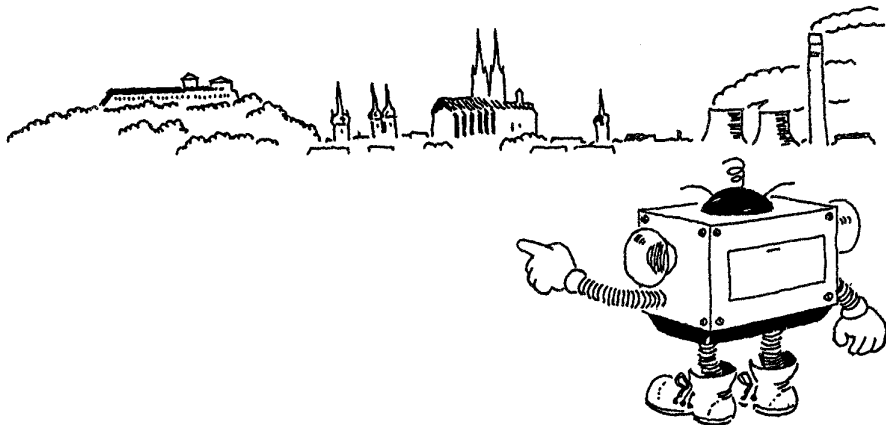
## Math Search Conclusion

MIaS is

- *text+math IR compatible* (fits mathematician's needs),
- *scalable* (index with 2 billions formulae tested),
- *Lucene/SOLR compatible* (transformable into SOLR-usable plugin) system,
- *easily deployable/integrable into EuDML core* (written in Java).

MIaS is *ready to be used in EuDML*!

## Accessibility improving technologies and tools

# 6+ years of local (Brno, CZ) document engineering

# Some of the verified and proven technologies (in DML-CZ)

- Scanned image processing and transformations (with BookRestorer) (BT Pulkrábek)

- Mathematical optical character recognition: OCR by combining FineReader (SDK 8.1) and Infty by prof. Suzuki (MT Panák, Mudrák, BT Vystrčil)

- Pre-MSC era papers' automated classification by MSC (with Radim Řehůřek)

- gensim framework: similarity article computations aka document clustering (Ph.D. research by Radim Řehůřek)

## Some of the verified and proven technologies (in DML-CZ)

- Scanned image processing and transformations (with BookRestorer) (BT Pulkrábek)

- Mathematical optical character recognition: OCR by combining FineReader (SDK 8.1) and Infty by prof. Suzuki (MT Panák, Mudrák, BT Vystrčil)

- Pre-MSC era papers' automated classification by MSC (with Radim Řehůřek)

- gensim framework: similarity article computations aka document clustering (Ph.D. research by Radim Řehůřek)

## Some of the verified and proven technologies (in DML-CZ)

- Scanned image processing and transformations (with BookRestorer) (BT Pulkrábek)

- Mathematical optical character recognition: OCR by combining FineReader (SDK 8.1) and Infty by prof. Suzuki (MT Panák, Mudrák, BT Vystrčil)

- Pre-MSC era papers' automated classification by MSC (with Radim Řehůřek)

- gensim framework: similarity article computations aka document clustering (Ph.D. research by Radim Řehůřek)

## Some of the verified and proven technologies (in DML-CZ)

- Scanned image processing and transformations (with BookRestorer) (BT Pulkrábek)

- Mathematical optical character recognition: OCR by combining FineReader (SDK 8.1) and Infty by prof. Suzuki (MT Panák, Mudrák, BT Vystrčil)

- Pre-MSC era papers' automated classification by MSC (with Radim Řehůřek)

- gensim framework: similarity article computations aka document clustering (Ph.D. research by Radim Řehůřek)

## Some of the verified and proven technologies (in DML-CZ) (cont.)

- Google Scholar partnership: interface to use our metadata instead of those parsed from landing pages' HTML

- Math retrieval: math formula indexing and search (MT Vítězslav Dostál, BT Martin Liška, BT Peter Mravec)

- Born-digital publishing system [for Archivum Mathematicum and for other 10 journals] and retro-born-digital paper conversions and enhancements (BT&MT Michal Růžička)

- Visualization and browsing interface (MT Zuzana Nevěřilová)

# Some of the verified and proven technologies (in DML-CZ) (cont.)

- Google Scholar partnership: interface to use our metadata instead of those parsed from landing pages' HTML

- Math retrieval: math formula indexing and search (MT Vítězslav Dostál, BT Martin Liška, BT Peter Mravec)

- Born-digital publishing system [for Archivum Mathematicum and for other 10 journals] and retro-born-digital paper conversions and enhancements (BT&MT Michal Růžička)

- Visualization and browsing interface (MT Zuzana Nevěřilová)

## Some of the verified and proven technologies (in DML-CZ) (cont.)

- Google Scholar partnership: interface to use our metadata instead of those parsed from landing pages' HTML
- Math retrieval: math formula indexing and search (MT Vítězslav Dostál, BT Martin Liška, BT Peter Mravec)
- Born-digital publishing system [for Archivum Mathematicum and for other 10 journals] and retro-born-digital paper conversions and enhancements (BT&MT Michal Růžička)
- Visualization and browsing interface (MT Zuzana Nevěřilová)

## Some of the verified and proven technologies (in DML-CZ) (cont.)

- Google Scholar partnership: interface to use our metadata instead of those parsed from landing pages' HTML

- Math retrieval: math formula indexing and search (MT Vítězslav Dostál, BT Martin Liška, BT Peter Mravec)

- Born-digital publishing system [for Archivum Mathematicum and for other 10 journals] and retro-born-digital paper conversions and enhancements (BT&MT Michal Růžička)

- Visualization and browsing interface (MT Zuzana Nevěřilová)

# Some of the verified and proven technologies (in DML-CZ) (cont.)

Metadata (in RDF) visualisation, browsing: Visual Browser tool (MT Zuzana Nevěřilová) for [Eu]DML GUI.

# Visual Browser to meet Universal Learning Design

One of ULD principles is to present information and content in different ways:

## Article processing



## Metadata extraction

## How Does It Work

- A lightweight set of LATEX macros in the form of a LATEX macro package.

  - Can be easily customized to meet needs of a particular journal document class / style file.
  - The LATEX macro package itself does not transform the LATEX source code to XML.
  - Literally exports selected parts of the LATEX document to an external file.
  - This file is subsequently processed by a journal-independent Tralics-based procedure.

## How Does It Work (cont.)

```
\documentclass[runningheads]{llncs}
\usepackage{dmlcommon}
\usepackage{dmlcz}

\begin{document}

\author{Petr Sojka}
\dmlaindex{Sojka}{Petr}
\dmltitle{Towards a Digital Mathematical Library}
...
\maketitle

\begin{dmlabstract}
The workshop's objectives were to formulate the strategy
and goals of a global mathematical digital library...
\end{dmlabstract}
...
```

## How Does It Work (cont.)

```
\documentclass{dmlczmeta}\begin{document}

\begin{xmlelement}{author}{Sojka, Petr
\XMLaddatt{order}{1}}\end{xmlelement}

\begin{xmlelement}{title}{Towards a Digital Mathematical
Library\XMLaddatt{lang}{eng}}\end{xmlelement}

\begin{xmlelement}{abstract}\XMLaddatt{lang}{eng}\bgroup
The workshop's objectives were to formulate the strategy
and goals of a global mathematical digital library...
\egroup\end{xmlelement}

\begin{xmlelement}{keyword}{OCR\XMLaddatt{lang}{eng}}
\end{xmlelement}

...
\end{document}
```

## How Does It Work (cont.)

- Tralics is a LaTeX to XML translator.
  - The most indispensable part of the system.
  - Its engine is able to process regular LaTeX code.
  - It is not necessary to
    - convert the LaTeX code to plain text directly,
    - nor deal with the LaTeX macro expansion or the complexity of its syntax.

- Tralics outputs a UTF-8 encoded XML file.

- This output is finally processed by the XLST processor furnishing DML-CZ metadata in its final form.

## How Does It Work (cont.)

```
<?xml version='1.0' encoding='UTF-8'?>
<!DOCTYPE std SYSTEM 'classes.dtd'>
<!-- Translated from latex by tralics 2.13.5,
     date: 2010/07/03-->
<std><p>
<author order='1'>Sojka, Petr</author>
<title lang='eng'>Towards a Digital Mathematical
Library</title>

<abstract lang='eng'>The workshop's objectives were to
formulate the strategy...</abstract>
<keyword lang='eng'>OCR</keyword>
<keyword lang='eng'>OpenMath</keyword>

<language>eng</language>
<abstractlanguage>eng</abstractlanguage>
...
</p></std>
```

## How Does It Work (cont.)

```
<?xml version="1.0" encoding="UTF-8"?>
<article>
  <title lang="eng">Towards a Digital
    Mathematical Library</title>

  <author order="1">Sojka, Petr</author>

  <language>eng</language>

  <keyword lang="eng">OCR</keyword>
  <keyword lang="eng">OpenMath</keyword>

  <summary lang="eng">The workshop's objectives
    were to formulate the strategy...</summary>
  ...
</article>
```

## Why It Is Useful

- It is easy to integrate this procedure to an existing journal processing workflow. It is thus acceptable to all the involved editors.

  - Current T$_E$X processing is used.
  - Platform independent.

    - The T$_E$X itself produces the source file.
    - XML generated using Tralics and XSLT.

  - No need for BibT$_E$X.

- It is safe.

  - At the same time as the final PDF document is created, the metadata is automatically generated based on the same source code.

- Since Tralics supports MathML we are able to translate mathematical expressions from the input L$^A$T$_E$X notation to this XML language.

## Maths, TEX, PDF

- PDF is widely adopted and very often used for electronic publications.
    - The DML-CZ project stores full texts of the articles as PDF files as do many other digital libraries.

- Thanks to pdfTEX, PDF is also the *de facto* standard output format of the modern TEX distributions.
- LATEX mathematical notation is well known and effective.
    - Used not only in LATEX documents, but also in a variety of other projects, such as Wikipedia.

- LATEX source code is usually a good choice for plain text representation of mathematical expressions.

## Standard PDF document



LATEX source code:

```
Text $\Pi(x) = \pi(x) +
\frac{1}{2}\pi(x^{1/2}) +
\frac{1}{3}\pi(x^{1/3}) + \cdots$
text.
```

# Standard PDF document



**PDF code:**

```
BT
/F16 9.9626 Tf 148.712 707.125 Td [(T)83(ext)]TJ/F17 9.9626 Tf 23.247 0 Td
[(\005\050)]TJ/F20 9.9626 Tf 11.346 0 Td [(x)]TJ/F17 9.9626 Tf 5.694 0 Td
[(\051)-278(=)]TJ/F20 9.9626 Tf 17.158 0 Td [(\031)]TJ/F17 9.9626 Tf 6.036 0 Td
[(\050)]TJ/F20 9.9626 Tf 3.875 0 Td [(x)]TJ/F17 9.9626 Tf 5.694 0 Td
[(\051)-222(+)]TJ/F18 6.9738 Tf 17.247 3.923 Td [(1)]TJ
ET
```

# Standard PDF document



Text obtained using Copy & Paste function of PDF reader:

```
Text  ( ) =  ( ) + 1
2 ( 1/2) + 1
3 ( 1/3) + · · · text.
```

# copymath-enabled PDF document



Text $\Pi(x) = \pi(x) + \frac{1}{2}\pi(x^{1/2}) + \frac{1}{3}\pi(x^{1/3}) + \cdots$ text.
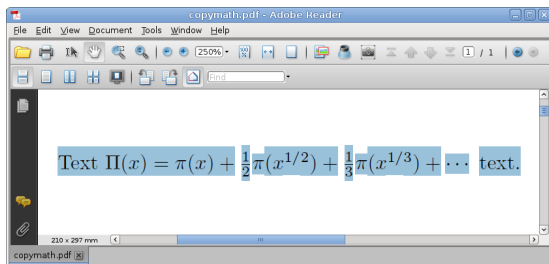
LaTeX source code:

```
Text $\Pi(x) = \pi(x) +
\frac{1}{2}\pi(x^{1/2}) +
\frac{1}{3}\pi(x^{1/3}) + \cdots$
text.
```
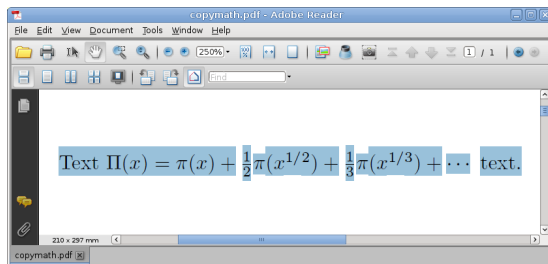
# copymath-enabled PDF document



**PDF code:**

```
BT
/F16 9.9626 Tf 148.712 707.125 Td [(T)83(ext)]TJ
ET
1 0 0 1 171.959 707.125 cm
/Span <<
/ActualText<245C506920287829203D205C706920287829202B205C66726163207B317D7B32
7D5C70692028785E7B312F327D29202B205C66726163207B317D7B337D5C70692028785E7B31
2F337D29202B205C63646F74732024> >> BDC
1 0 0 1 -171.959 -707.125 cm
BT
/F17 9.9626 Tf 171.959 707.125 Td [(\005\050)]TJ/F20 9.9626 Tf 11.346 0 Td
[(x)]TJ/F17 9.9626 Tf 5.694 0 Td [(\051)-278(=)]TJ/F20 9.9626 Tf 17.158 0 Td
[(\031)]TJ/F17 9.9626 Tf 6.036 0 Td [(\050)]TJ/F20 9.9626 Tf 3.875 0 Td
[(x)]TJ/F17 9.9626 Tf 5.694 0 Td [(\051)-222(+)]TJ/F18 6.9738 Tf 17.247 3.923
Td [(1)]TJ
ET
```
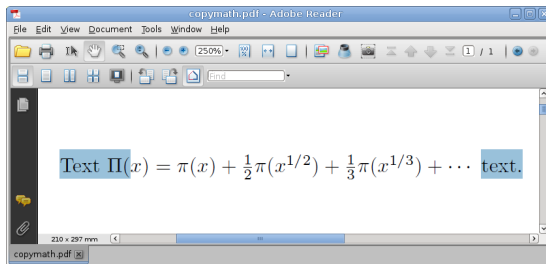
## copymath-enabled PDF document



Text obtained using Copy & Paste function of PDF reader:

```
Text $\Pi (x) = \pi (x) +
    \frac {1}{2}\pi (x^{1/2}) +
    \frac {1}{3}\pi (x^{1/3}) + \cdots $
text.
```

## Implementation

- The `ActualText` command of the PDF language is used to mark the region of the mathematical expression inside the PDF document.

- We want the package to be as user friendly as possible – users should not be forced to modify their mathematical expressions in any way, `\usepackage{copymath}` should cater for all their needs.

  - The implementation is not easy.
  - This requires nonstandard modifications of the LaTeX mathematical environments.

## Implementation (cont.)

- We need to add `\pdfliteral` at the beginning and end of every mathematical environment.

- The dollar sign ($) is activated and redefined.

- It is necessary to keep track of nested mathematical environments.

- Simple redefinition of $\mathcal{AMS}$-LATEX mathematical environments is not possible.

- It seems that not all PDF viewers respect contents of the `ActualText` command.

- Adobe Reader ignores the "_" sign inside `ActualText` provided another character is present.

- Possibility to be misused.

# Summary

- Verified complex DML-CZ digitization workflow and proven technologies and tools for math DL

- EuDML: Towards *accessible* worldwide digital mathematical library, based on DML-CZ know-how and tools

- DML workshop series, join us at DML 2011 c/o CICM Bertinoro, Italy, July 18th–23rd, 2011

Overview   Motivation   DML-CZ   Math Indexer and Searcher   Technologies   Tools   **Summary**

## Summary

- Verified complex DML-CZ digitization workflow and proven technologies and tools for math DL

- EuDML: Towards *accessible* wordwide digital mathematical library, based on DML-CZ know-how and tools

- DML workshop series, join us at DML 2011 c/o CICM Bertinoro, Italy, July 18th–23rd, 2011

Sojka: Accessibility Issues in a Digital Mathematical Library
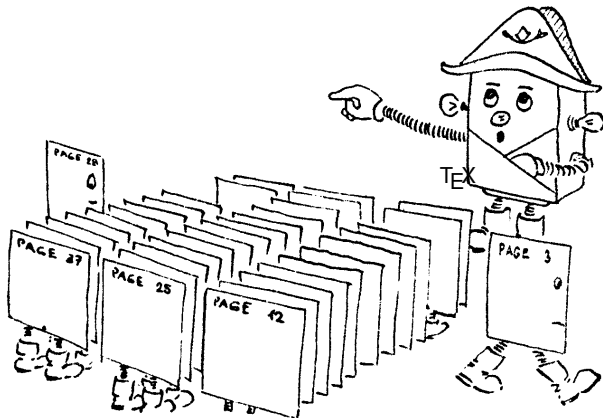Universal Learning Design, Brno, CZ, 10. 2. 2011

## Summary

- Verified complex DML-CZ digitization workflow and proven technologies and tools for math DL

- EuDML: Towards *accessible* wordwide digital mathematical library, based on DML-CZ know-how and tools

- DML workshop series, join us at DML 2011 c/o CICM Bertinoro, Italy, July 18th–23rd, 2011

# Yes, you can!

## Future work

- Robust Math OCR

- Robust PDF2Math conversion

- Design alternative, novel and accessible user interfaces for the digital library

- Other challenges: multilingual math retrieval, MathML indexing and search, math common sense, mathematical document classification, document similarity

## Future work

- Robust Math OCR
- Robust PDF2Math conversion
- Design alternative, novel and accessible user interfaces for the digital library
- Other challenges: multilingual math retrieval, MathML indexing and search, math common sense, mathematical document classification, document similarity

## Future work

- Robust Math OCR
- Robust PDF2Math conversion
- Design alternative, novel and accessible user interfaces for the digital library
- Other challenges: multilingual math retrieval, MathML indexing and search, math common sense, mathematical document classification, document similarity
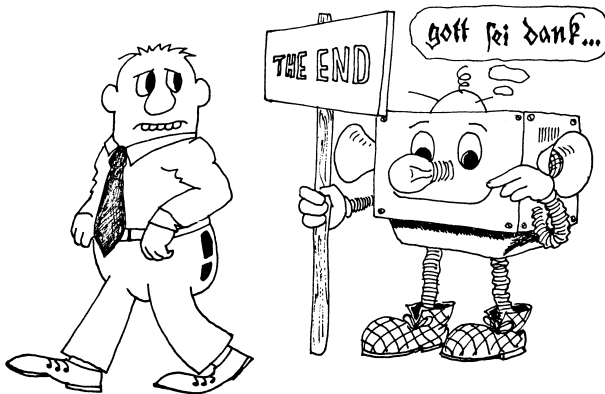
# Future work

- Robust Math OCR
- Robust PDF2Math conversion
- Design alternative, novel and accessible user interfaces for the digital library
- Other challenges: multilingual math retrieval, MathML indexing and search, math common sense, mathematical document classification, document similarity

## Take-home messages and acknowledgements

- Even though some problems are not resolved and complexity of dealing with math is high,

- there seems to be light at the end of the tunnel!

Many thanks to all members of DML-CZ and EuDML team, especially to those whose tools and work had been reported: Michal Růžička, Martin Líška, Martin Šárfy, Radim Řehůřek, Volker Sorge, Thierry Bouche et al.

## End of the talk



Questions? Comments? Cooperation offers?

Czech Digital Mathematics Library [online].

[online, cit. 2011-02-10].
Available from WWW: <http://dml.cz/>.

EuDML: The European Digital Mathematics Library [online].

This page was last modified on 20 January 2010, at 08:09. [online, cit. 2011-02-10].
Available from WWW: <http://www.eudml.eu/>.

Bouche, T.:

A pdfLATEX-based automated journal production system.
In Proceedings of EuroTEX 2006, TUGboat **27**(1) (2006) 45–50.

Centre de diffusion de revues académiques mathématiques [Center for diffusion of mathematic journals] [online].

[online, cit. 2011-02-10].
Available from WWW: <http://www.cedram.org/>.

Růžička, M.:

Automated Processing of TEX-Typeset Articles for a Digital Library.
In: Sojka Petr (editor): *DML 2008 – Towards Digital Mathematics Library*, Birmingham, UK, July 27th, 2008, 167–176.

Archivum Mathematicum [online].

Masaryk University, Brno, Czech Republic.
Last modified December 18, 2009 [online, cit. 2011-02-10].
Available from WWW: <http://www.emis.de/journals/AM/>.

Grimm, J.:

Tralics, a LATEX to XML Translator.
In Proceedings of EuroTEX, TUGboat **24**(3) (2003) 377–388.

Tralics: a LaTeX to XML translator [online].

Last modified $Date: 2009/11/24 17:17:03 $ [online, cit. 2011-02-10].
Available from WWW: <http://www-sop.inria.fr/apics/tralics/>.

Infty Project: Research Project on Mathematical Information Processing [online].

[online, cit. 2011-02-10].
Available from WWW: <http://www.inftyproject.org/en/>.

Suzuki, M.; Kanahori, T.; Ohtake, N.; Yamaguchi, K.:
An Integrated OCR Software for mathematical Documents and Its Output with Accessibility.
*Computers Helping people with Special Needs, 9th International Conference ICCHP 2004*, Paris, July 2004, Lecture Notes in Computer Sciences 3119, Springer (2004) 648–655.

EuDML at MU team.
*EuDML at MU project info* [online, cit. 2011-02-10].
<http://nlp.fi.muni.cz/projekty/eudml/> or
<http://www.muni.cz/research/projects/10067>.