Overview
○○○○○○○○○○○○○○○
Motivation
○○○○
DML-CZ
○○○○○○○○○○○○○○○
EuDML
○○○○○○○○○○○○○○
Summary
○○○○

# Specifics of Open Access Publishing and Retrodigitization in Mathematics: An Experience from DML-CZ and EuDML Projects
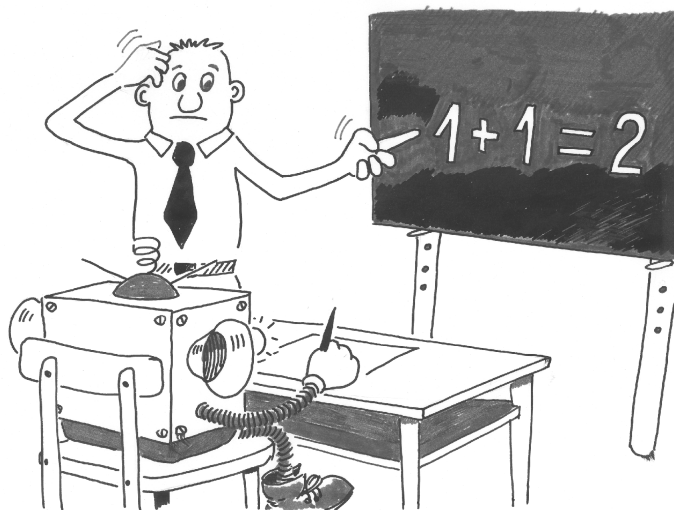
Petr Sojka

<sojka@fi.muni.cz> (Faculty of Informatics, Masaryk University, Brno)
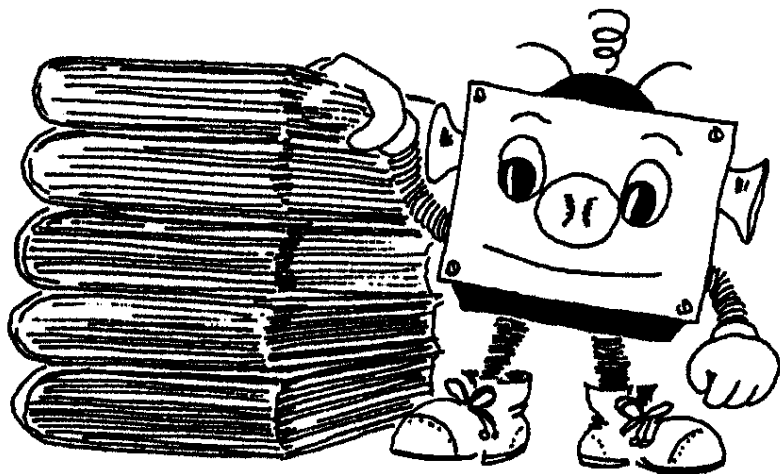
COASP 2010, Prague, CZ, August 24th, 2010, 9:40 a.m.

# Specifics of mathematics: Mathematics is the Queen of Science and Arithmetics is the Queen of Mathematics (C.F. Gauss)
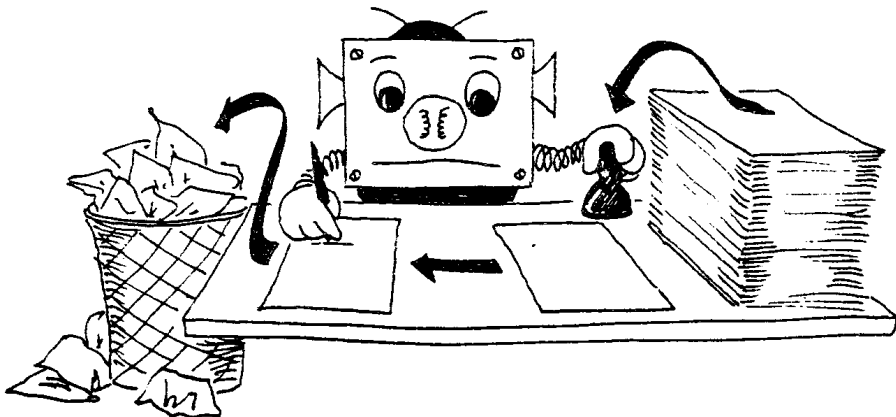
# Digital processing is full of challenges

# Amount of math ever published is about 100.000.000 pages

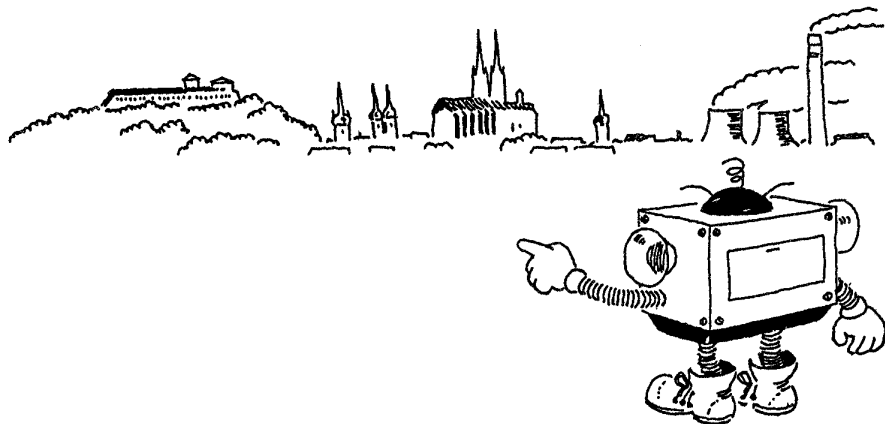# Specific publishing workflow and document engineering (tough reviewing, formulae handling)

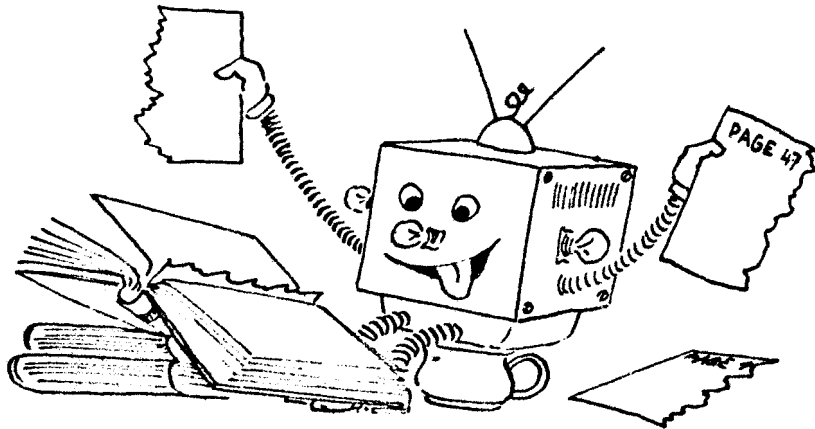# Specific problems of digitization (OCR, digital library development)

# Challenges of math formulae search and indexing

# Challenges tackled during DML-CZ in Brno, CZ (pdfTeX's born town)
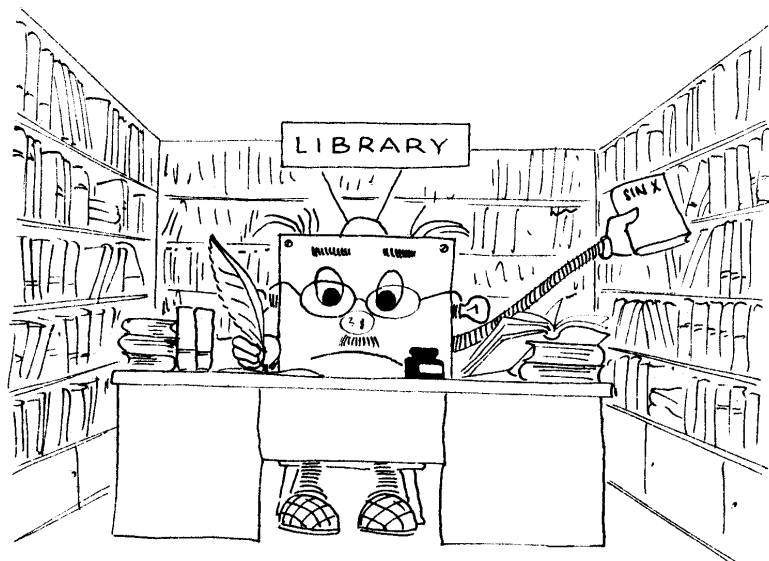
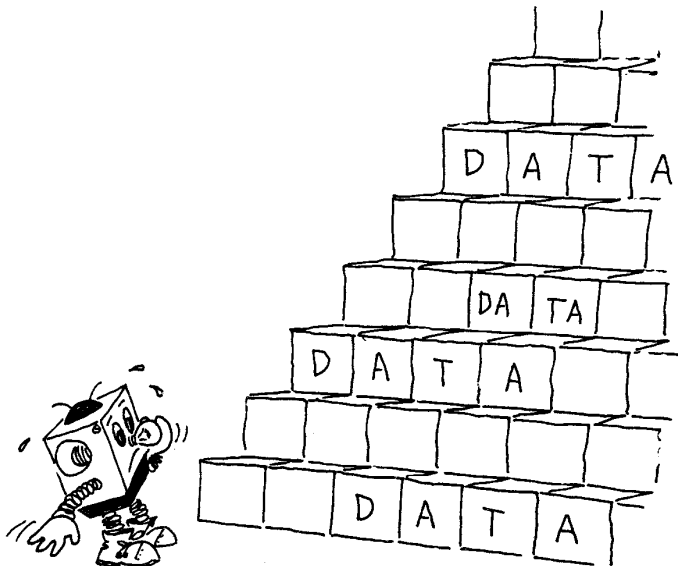# DML-CZ document engineering—best practices and tools

# Towards European Virtual Library of Mathematics: bottom up DML processing towards EU (or worldwide) scale
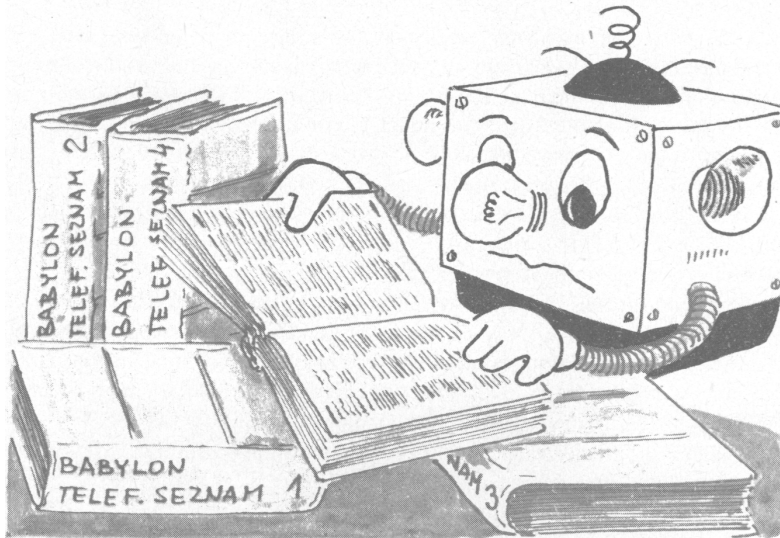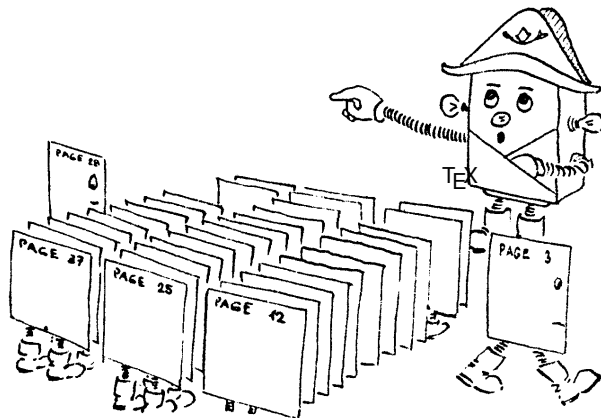
# European Digital Mathematics Library

## EuDML—from local data collections to the virtual digital library

# EuDML document engineering—scalable tools development

# Yes, you can!

# End of talk overview, quiz!

Overview
0000000000000000

Motivation
●000

DML-CZ
0000000000000000

EuDML
000000000000000

Summary
0000

## Open-ended story of a Vision

At the beginning there was a vision of all mathematical knowledge, *peer reviewed and verified* (100,000,000 pages) in 'one stop shop' (in the digital form), articulated at the beginning of millenium: reference mathematics library, grounding truths for not only mathematitions, but also other sciences.

Progress of IT, cheap disk space, 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 (some nine digit sum in $'s requested). Application was not successful–top down approach failed.

Overview
○○○○○○○○○○○○○○○○

Motivation
○●○○

DML-CZ
○○○○○○○○○○○○○○○○

EuDML
○○○○○○○○○○○○○○○

Summary
○○○○

# Vision of (European) Digital Mathematics Library

Several attempts to fund development of DML on European level (FP5, FP6) also was not successful.

Now, it starts to be realized: three year EU project EuDML (programme EU CIP-ICT-PSP, type Pilot B, EU contribution 1.6 MEur) from February 2010



(MU and MU AV).

The strategy:

- to master the technology, develop tools and offer them;

- concept of *moving wall* to motivate and engage commercial publishers.

- to collect (bottom up) [virtual] *digital library*, 'one-stop shop' and achieve critical mass in the domain → 'me too' effect then.

# From paper to digital processing, from local to the whole

Overview
○○○○○○○○○○○○○○○

Motivation
○○○●

DML-CZ
○○○○○○○○○○○○○○○

EuDML
○○○○○○○○○○○○○

Summary
○○○○

# Bottom up—from building bricks of regional repositories

As building bricks current (regional or publisher's) Digital Mathematics Library (DML) repositories as DML-CZ or NUMDAM, DML-PL, DML-PT, RusDML,…(from local repositories bottom-up to build the final thing—to realize the vision).

Example of DML-CZ: up and running digital mathematic library with nearly 30,000 papers. For more, see (who, what, browse, browse similar, how to search).

Live project—all comments to DML-CZ welcome!

Overview
0000000000000

Motivation
0000

DML-CZ
●000000000000

EuDML
000000000000

Summary
0000

## DML-CZ: main facts

- Czech Academy of Sciences grant (program Information Society) 2005–2009, *full* (retro)digitization of 50,000 pages of mathematical literature per year, 8M CZK in total.

- Research part: **1)** gradual enhancement of the digital material by 'knowledge enhancing' filters on markup-rich XML data. **2)** New methods for (semantic) text processing tested on the available data. **3)** The design of the work-flow aiming at mathematical knowledge stored in digital library.

- IPR part: sharing/delivery (economic models for knowledge sharing due to interests of content owners/publishers).

- Production part: dig. center Jenštejn, overestimated costs.

Overview
○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○●○○○○○○○○○○○○○○

EuDML
○○○○○○○○○○○○○○

Summary
○○○○

## DML-CZ data: what is there

- 12 mathematics journals, 7 OA, 5 with moving wall.

- 8 publishers of math in CZ and SK, close cooperation.

- 27,000+ articles with rich metadata (full bibliographic record, similarities, MSC classifications, translated titles).

- heavily multilingual (7+ languages).

Links on who, what, browse, browse similar, how to search).

# DML-CZ – data: scientific math published in Czech and Slovak

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

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

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

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

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 теме поверхностного интеграла, Чехословацкий мат. журнал, 6 (81), 1956, 387—400.
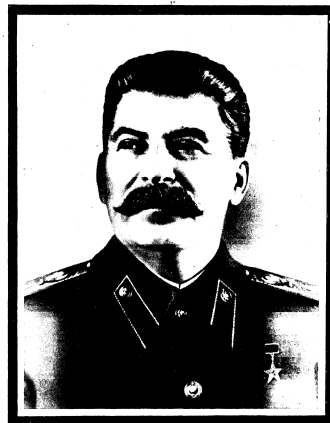[7] S. Saks: Theory of the integral, New York.

### Резюме

#### ПОВЕРХНОСТНЫЙ ИНТЕГРАЛ

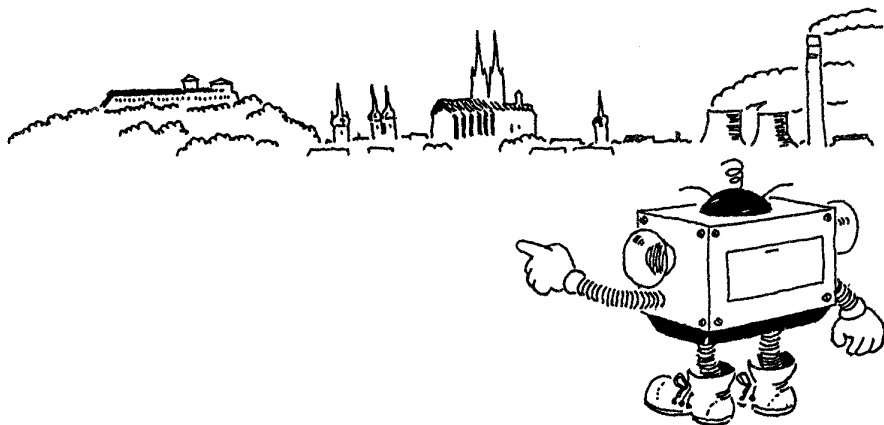ЯН МАРЖИК (Jan Mařík), Прага.
(Поступило в редакцию 10/X 1955 г.)

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 $\sum_{i=1}^m v_i^2(x) \le 1$ для всех $x \in A$. 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 $p$ и на

557

ИОСИФ ВИССАРИОНОВИЧ СТАЛИН
1879—1953

Overview
○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○○○●○○○○○○○○○○○

EuDML
○○○○○○○○○○○○○

Summary
○○○○

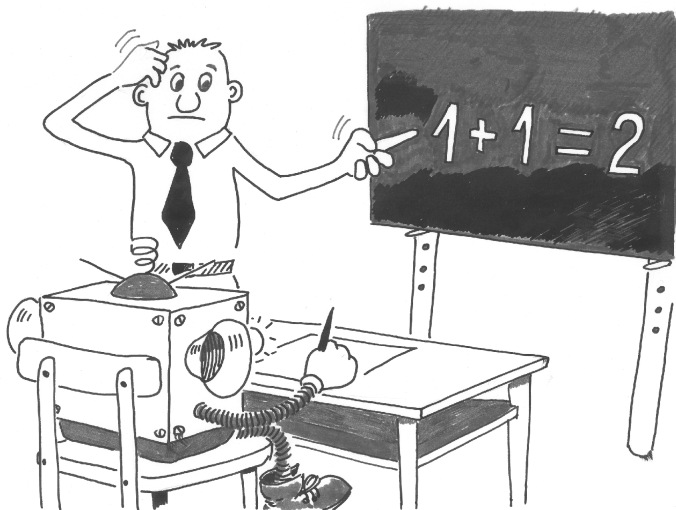# Bottom up processing—local (Brno, CZ) document engineering

## The approach used in DML-CZ

A succesfully built repository (e.g. set of *workflows*) needs a *coordinated* effort of *librarians*, *IT specialists* and representatives of users—*content specialists*: (D+M+L)=success 'equation'.
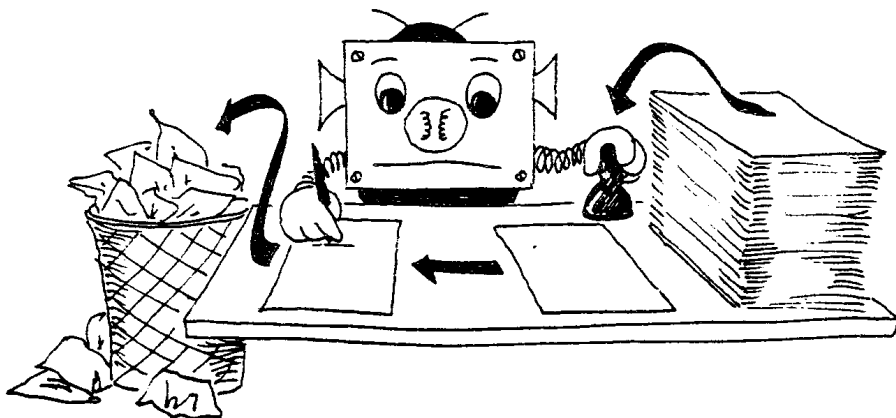
*Design, technical* and *political decisions* behind building the *Czech Digital Mathematics Library DML-CZ* (<http://dml.cz>) in the context of other thematical community projects (PubMed Central, ADS, INSPIRE, SCOAP3 and EuDML) have been solved. *No wheel reinvention.*

Our framework integrates workflow for the articles scanned from a paper (*math OCR*), for documents from retro-born digital period (data available in some type of electronic form) and for born-digital ones. As much automated (using robots :-) as possible.
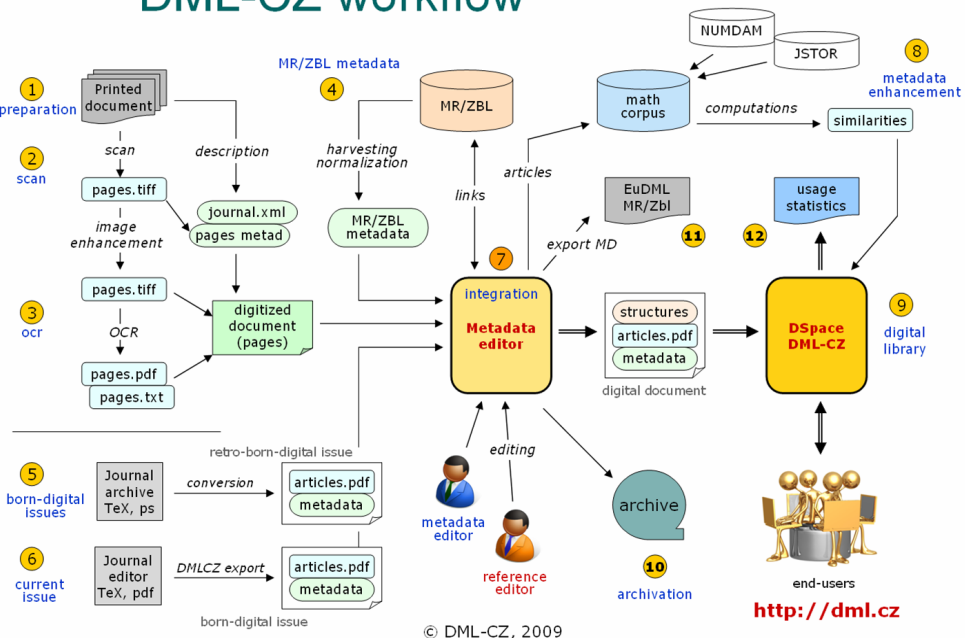
# Math handling poses challenges—math OCR, math indexing,...

# Document engineering—from paper to digital *workflow*

# DML-CZ workflow

# DML-CZ document engineering—data processing



... RIGHT PAGE ...

DML-CZ now serves about *275,000 pages of math papers*.

Problems of *migration of existing workflows (born-digital, retro-digital) into the repository*. Negotiations with Google Scholar towards better visibility, indexing and search, and problems of copyright and sustainability issues, visualization, space and processing demands,….

# Document engineering—digitization, digital library development

## New tools and best practices for [meta]data processing

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

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

world of authors: LaTeX, TeX notation of mathematics

world of applications/data exchange: XML, MathML

Overview
○○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○○○○○○○○○○○○○●○

EuDML
○○○○○○○○○○○○○○

Summary
○○○○

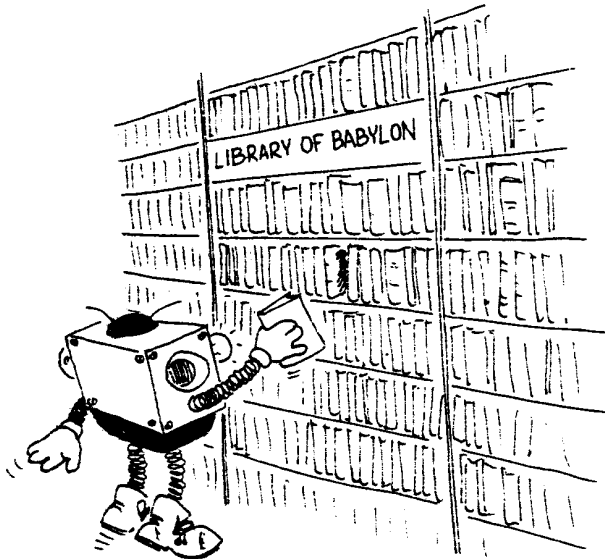## Verified and proven technologies (in DML-CZ)

- scanned image processing and transformations, mathematical optical character recognition: OCR

- digital signature of PDF: pdfsign

- web-based long distance metadata editing: web application metadata editor

- optimization and recompression of PDF: downsizing of more than 50% without quality loss

- article similarity computations, demo.

- retroborn paper automated classification by MSC.

Overview
0000000000000000

Motivation
0000

DML-CZ
0000000000000●

EuDML
000000000000000

Summary
0000

## Verified and proven technologies (cont.)

- born-digital publishing system [for Archivum Mathematicum and other 4 journals] and conversions.

- retro-born-digital paper conversions and enhancements.

- data vizualization, browsing: adaptation of Visual Browser

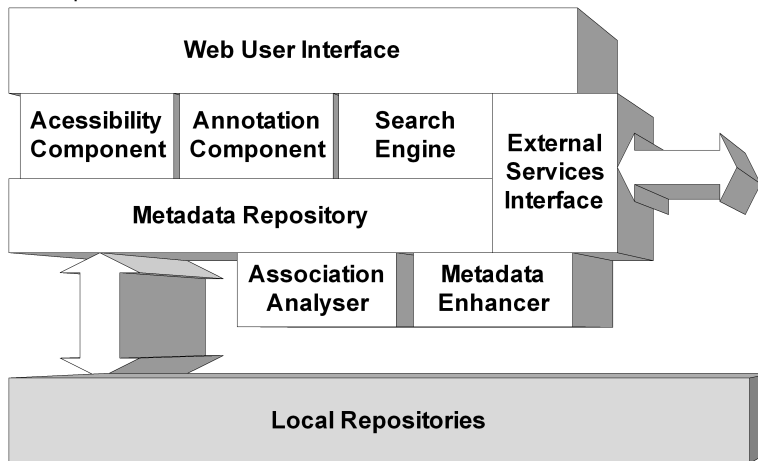- math retrieval: math formula indexing and search

- citation linking: CiteCrawl

many open areas/challenges: multilingual retrieval?. MathML indexing,…

# Bottom up processing towards EU or worldwide scale

# 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

Overview
○○○○○○○○○○○○○○

Motivation
○○○○

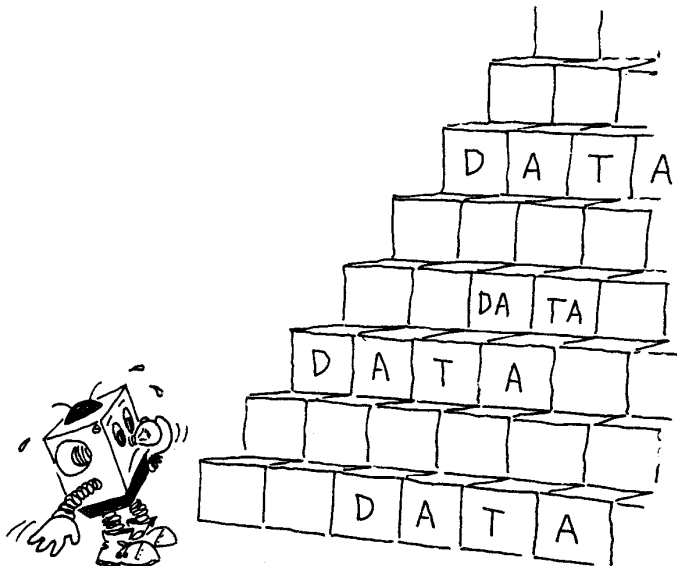DML-CZ
○○○○○○○○○○○○○○

EuDML
○○○●○○○○○○○○○

Summary
○○○○

## EuDML – data: legacy scientific math

- By 2013, EuDML should integrate *12 repositories*, have content from *200 integrated collections* (journals, book series, conference proceedings,…), more than *160,000 digital items* (papers, book chapters), *500,000 links between database objects*.

- It should be 'live' DL, having more than *1,000 users* contributing annotations, and more than *10,000 annotation* by 2013.

- Concept of *moving wall*: legacy data even from commercial publishers.

But how to actually implement it?

Experience from project partners from current digital library development.

Overview
○○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○○○○○○○○○○○○○○○○

EuDML
○○○○●○○○○○○○○○

Summary
○○○○

# EuDML—from data collection to the virtual digital library

Overview
○○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○○○○○○○○○○○○○○○○

EuDML
○○○○○○●○○○○○○○

Summary
○○○○

# EuDML service based architecture

# Alternative GUIs—visualization research

Overview
oooooooooooooooo

Motivation
ooooo

DML-CZ
oooooooooooooooo

EuDML
ooooooooo●ooooooo

Summary
ooooo

# Visual Browser development (DML-CZ)

# EuDML document engineering—scalable tools development

Overview
○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○○○○○○○○○○○○○○○○

EuDML
○○○○○○○○○○○●○○○

Summary
○○○○

## EuDML service based architecture II



based on existing YADDA (used in Driver, Driver II) and REPOX (used in EuropeanaLocal, Telplus) projects – both are verified and mature platforms (implemented in Java)

math specifics needed to develop (TeX to MathML converter, math OCR, math in metadata support,…

MU offers: Metadata editor and other tools and expertize, mainly to be used in *WP7 Metadata Enhancements*

# PDF Re-compression

New tools developed (Radim Hatlapatka) to re-compress [bitonal] PDF files:

| | Original PDF | After using PDF re-compressor | After using pdfsizeopt.py | After both |
|---|---|---|---|---|
| Size of whole PDF | 100% | 74.61% | 50.02% | 40.23% |
| Size of image and other objects | 69.46% | 37.14% | 45.14% | 35.36% |

May be used for any PDF 1.4 (since Acrobat 5 released in 2001) file—JBIG2 compression.

Overview
○○○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○○○○○○○○○○○○○○○

EuDML
○○○○○○○○○○○○○●○

Summary
○○○○

## Metadata Editor http://editor.dml.cz

Web-based client-server tool, developed (MU) for [meta]data import, editing, validation and checking. For testing, try <http://editor.dml.cz:9129>, admin/admin

Overview
0000000000000000
Motivation
0000
DML-CZ
0000000000000000
EuDML
000000000000●0
Summary
0000

# Metadata Editor localization (open source development)

# Yes, you can!

## Summary

Publishing and retrodigitization of mathematics poses many challenges, some yet to solve.

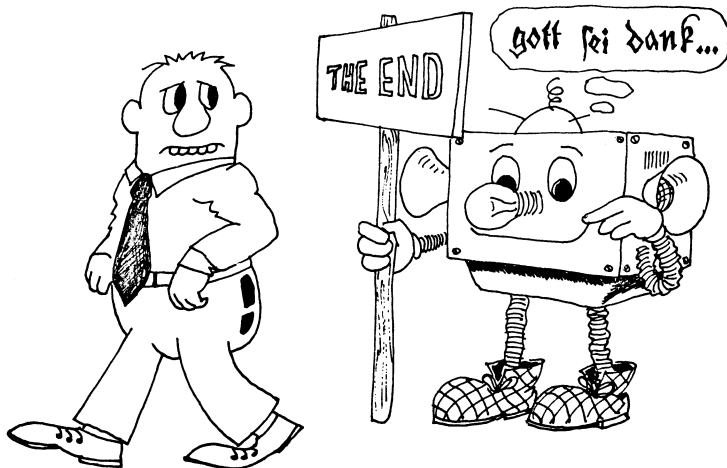DML-CZ: up and running, most open access.

EuDML: work in progress, some DML-CZ experience and tools developed to be (re)used.

Invitation: EuDML meeting with content providers in Prague, October 15th: math publishers wanted to cooperate.

Further reading: DML workshop series proceedings (Openly accessed in DML-CZ).

Comments, cooperation offers welcome! We are Open and AND-minded (no black&white thinking :-).

Overview
○○○○○○○○○○○○○○○○

Motivation
○○○○

DML-CZ
○○○○○○○○○○○○○○○○

EuDML
○○○○○○○○○○○○○○

Summary
○○●○

# End of the talk



Questions?

Overview
○○○○○○○○○○○○○○○
Motivation
○○○○
DML-CZ
○○○○○○○○○○○○○○○
EuDML
○○○○○○○○○○○○○○○
Summary
○○○●

# References, links

📄 DML-CZ team.
*Materials about DML-CZ, project publications* [online, cit. 2010-08-24].
<http://project.dml.cz/documents.html>.

📄 EuDML team.
*EuDML project info* [online, cit. 2010-08-24].
<http://ec.europa.eu/information_society/apps/projects/factsheet/index.cfm?project_ref=250503>.

📄 EuDML team.
*EuDML webpage* [online, cit. 2010-08-24].
<http://eudml.eu/>.

📄 EuDML at MU team.
*EuDML at MU project info* [online, cit. 2010-08-24].
<http://nlp.fi.muni.cz/projekty/eudml/> or <http://www.muni.cz/research/projects/10067>.