# Acquisition of Data for Plasma Simulation by Automated Extraction of Terminology from Article Abstracts

**L. Pichl, M. Suzuki[1], M. Murata[2], A. Sasaki[3], D. Kato[4] and I. Murakami[4]**

*International Christian University, Osawa 3-10-2, Mitaka, Tokyo 181-8585, Japan*
*1) School of Systems Science, Arkansas Tech University, Russellville, Arkansas 72801, USA*
*2) National Institute of Information and Communications Technology, Kyoto 619-0289, Japan*
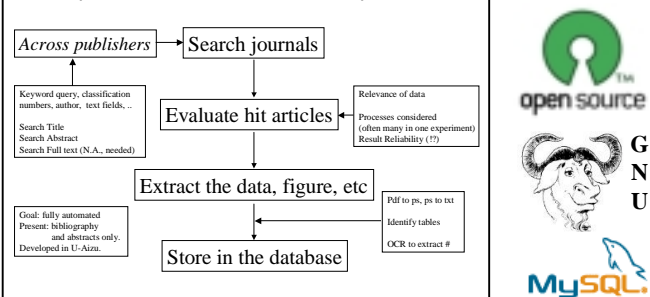*3) Japan Atomic Energy Agency, 8-1 Umemidai, Kizu-cho, Kyoto 619-0215, Japan*
*4) National Institute for Fusion Science, Oroshi-cho, Toki, Gifu 509-5292, Japan*

## Motivation

❏ Maintaining data center DBs costs time and money
❏ Data-searching and data-input are very low qualified, manual, stereotype activities → should be automated
❏ Commercial DB system solutions: costly & rigid, if a design change is need later.

The process to select and input data

*Across publishers* → Search journals

Keyword query, classification numbers, author, text fields, ..
Search Title
Search Abstract
Search Full text (N.A., needed)

Evaluate hit articles

Relevance of data
Processes considered (often many in one experiment)
Result Reliability (!?)

Extract the data, figure, etc

Goal: fully automated
Present: bibliography and abstracts only.
Developed in U-Aizu.

Pdf to ps, ps to txt
Identify tables
OCR to extract #

Store in the database

❏ Present work automates abstract download & processing

## Database System Design

❏ FCII linux OS, MySQL DB management system
HTML and logic layer with PHP scripting.

➢ Pre-designed queries are sent to the multiple publisher databases by using unix wget in command line mode.
➢ HTML output is analyzed for relevant data input fields using string matching search in PHP for each provider.
➢ Access rate is interrupted by intermissions following N(T,S) (or time histogram).
➢ DB input format is HTML; figures for special characters are also input into the DB.

String matching & Data input

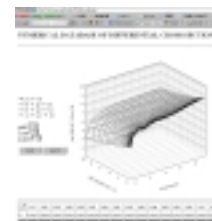## Abstract search and text classification (Dr. Murata, NIICT)

## Information extraction

✓ Relevance judgment is done based on the extraction of abstracts
✓ Special terminology of A+M physics is automatically recognized (species, states, elements, processes)
✓ Satisfactory recall and precision are obtained based on
✓ simultaneous process/species matching
✓ Further improvement is possible by the analysis of table and figure captions

Evaluation:

| Efficiency of rule-based method | |
|---|---|
| Total No. of papers | 348 |
| No. of text-formatted pdf files | 167 |
| No. of relevant articles | 64 |
| No. of irrelevant articles | 103 |
| No. of keyword | 92 (61 proc., 31 species) |
| Precision 97%, Recall 100% (TL, not 64) | |

## Work in development

▪ Numerical database of differential cross sections
  - prototype on the left, also crdb.nifs.ac.jp

▪ Online database builder

## Conclusion

Created an automated-input bibliography database for fusion plasma at NIFS.
Free-software open-source abstract and fulltext DB system development is finished.
A differential cross-section database with GUI features is tested on the same footing.
Interfaced with online article search engine, which includes a module for recognition of A+M terminology (Dr. Murata) in a joint project with Dr. Sasaki.