

# The DBLP Computer Science Bibliography: Evolution, Research Issues, Perspectives

SPIRE 2002  
Lisbon, Portugal

# Thank you for the invitation

- Many of you use DBLP
- This talk gives some background information about the service
- You are invited to use the DBLP data to test and evaluate for your algorithms ...

# About me ...

- born March 1959
- 1986 diploma in informatics from Aachen University of Technology
- 1993 Ph.D. from University of Trier
- since 1993 lecturer at U Trier
  - Programming for 1st/2nd year students
  - DB implementation, Digital Libraries

# Outline

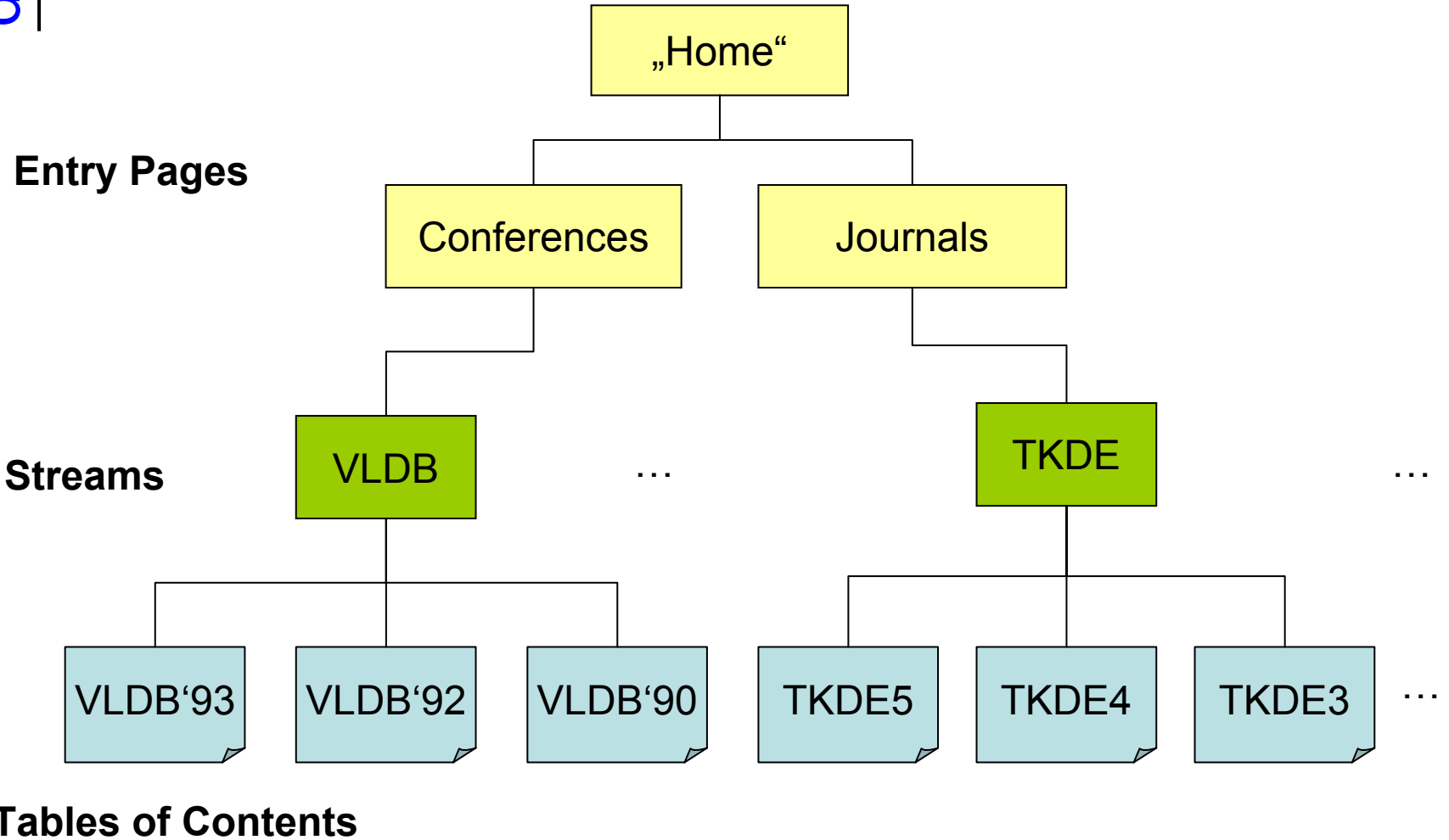
1. History
2. Technical Background
3. Perspectives & Research Issues

# 1. History

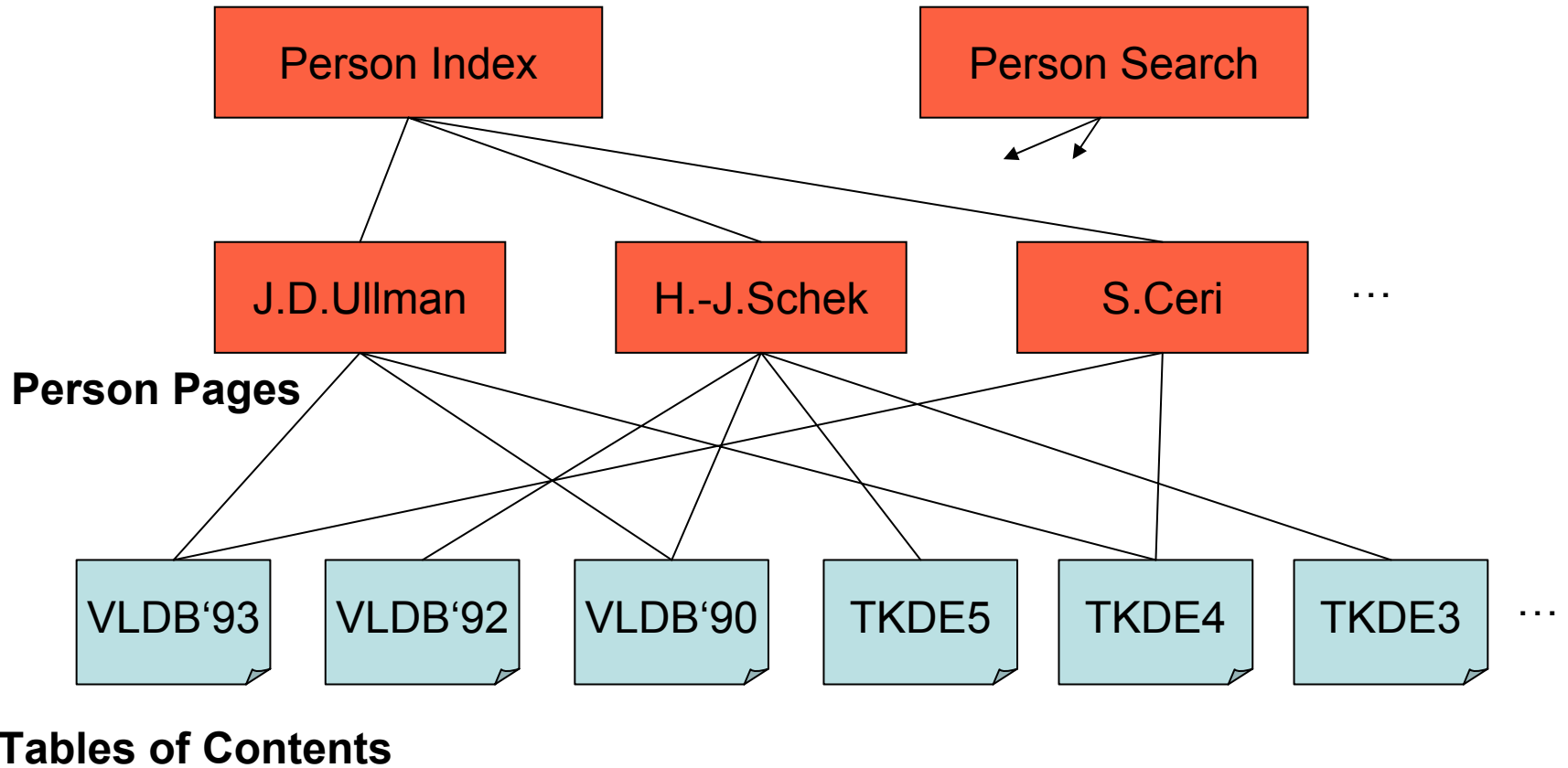
- The Beginning
- Person Pages
- Service
- Early Recognition
- ACM SIGMOD Anthology
- Sponsor
- Growth of DBLP
- Software Labs

# The Beginning

- End of 1993 - simple test of Web technology: Xmosaic, NCSA HTTP server
- Tables of contents:
  - Journals and proceedings
  - DataBase systems / Logic Programming



# Person-Publication Network





# Person-Publication Network

- Adding a hyperlink from each person name to a page which enumerates this person's publications
- Result of the complex social network behind research
- DB, IR: simply a view, a „canned“ query

# DBLP is a service, not a research project

- limited resources
- trade-off:
  - development of new features & software
  - vs. entering & maintaining contents

# No DBMS used until today

- prototype implementations (diploma theses): DBLP with SHORE, DB2
  - + better tools to maintain consistency
  - software too large for maintenance
- trade-off: disk-space vs. CPU speed

# Early Recognition

- 1997:
  - ACM SIGMOD Service Award
  - VLDB Endowment Special Recognition Award
- Helped to make DBLP a more „official“ project & to get a small initial fund

# ACM SIGMOD Anthology

- SIGMOD had made some profits with it's conferences
- Idea by Rick Snodgrass:
  - use the money to scan in „historical“ DB publications
  - Combine these full texts and an improved version of DBLP to a digital library

# Anthology: Contents

## Journals, Newsletters:

- TODS
- TKDE
- VLDB Journal
- Distr. & Parallel DB
- Data Engineering
- SIGMOD Record
- SIGKDD Expl.
- SIGIR Forum
- Data Base

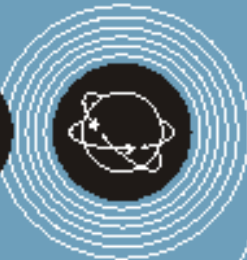
Proceedings:

- ACM DL
- ACM GIS
- ADBIS
- CIKM
- COOPIS
- DASFAA
- DBPL
- DOLAP
- EDBT
- ER
- Hypertext
- ICDT
- KRDB
- MFDDBS
- NPIV
- PDIS
- PODS
- POS
- SIGIR
- SIGMOD
- SIGFIDET
- SSD
- SSDBM
- VLDB
- XP
- + several Workshops

## Books:

- Abiteboul/Hull/Vianu: Foundations of DBs
- Bernstein/Hadzilacos/Goodman: Concurrency Control and Recovery in Database Systems
- Maier: Theory of Relational Databases
- Gray: The Benchmark Handbook
- Stonebraker: The INGRES Papers
- Wiederhold: Database Design (2nd Ed.)
- Snodgrass: The TSQL2 Temporal Query L.





The ACM SIGMOD



# Anthology

21 CDROMs / 2 DVDs

>150000 pages full text

	size	# files	# PDF files
DVD 1	7.5 G	10339	10025
DVD 2	6.31 G	201902	4284

# Anthology: Citation Links

References:  
[1] ...  
[2] B: xxx. ←  
[3] ...  
...

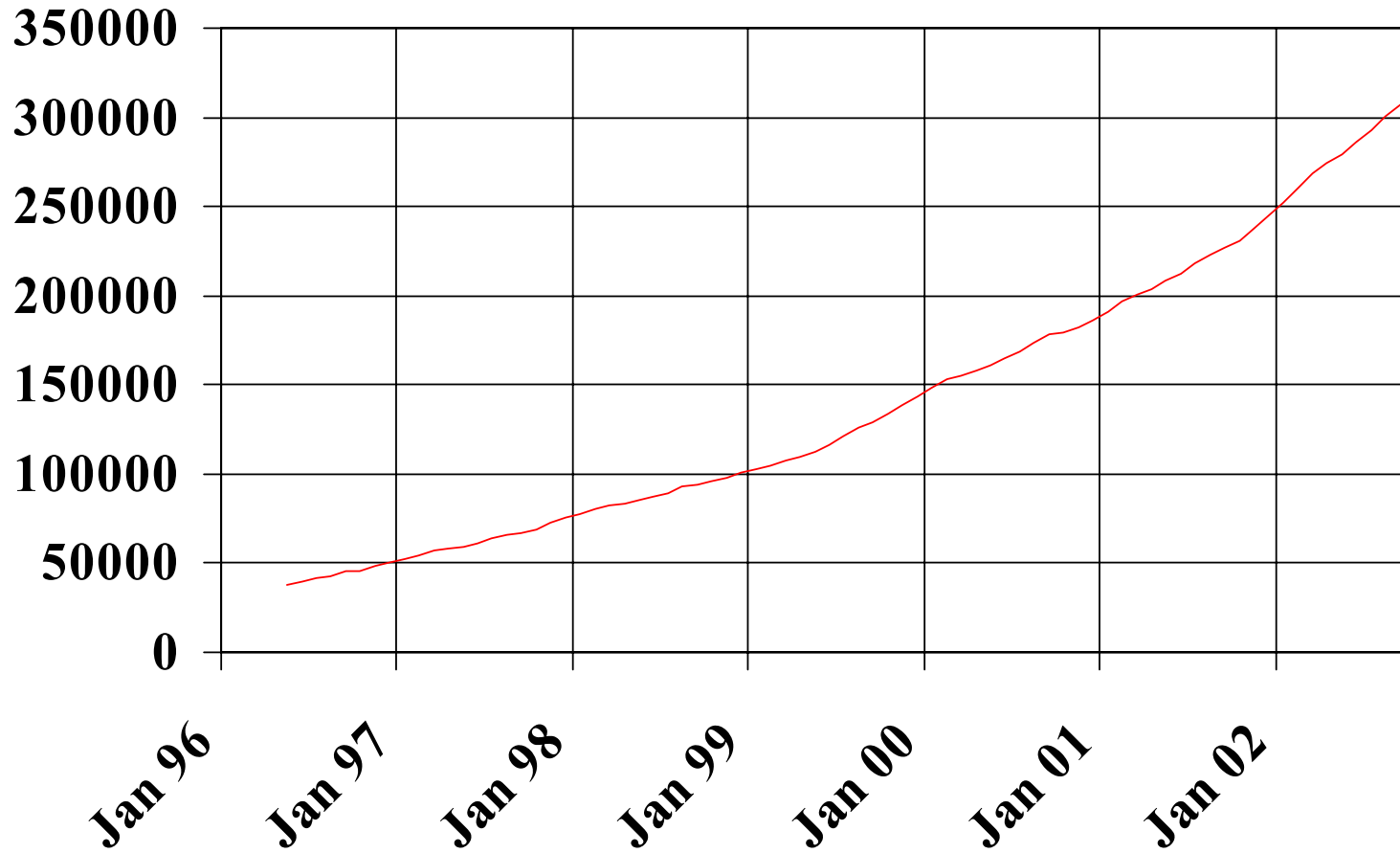
References:  
[1] ...  
...  
Referenced by:  
• A: yyy.  
• ...

# Sponsor Found / Expansion

- help by MS Research & Jim Gray made it possible to expand DBLP to cover most areas of computer science
- students were hired to enter data



# Growth of DBLP



# Teaching Java ...

- At U Trier we teach Java as the first programming language
- Assignments: DBLP XML data are used
  - graph algorithms
  - user interfaces
  - simple search engines
  - ...

## 2. Technical Background

- Initial Design
- Person Pages
- Mirrors
- Simple Search
- XML-Records
- BHT-Files
- MG: advanced search
- Anthology search

# Initial Design

- Entry pages
- Collection of HTML tables of contents
- TOCs were parsed to generate „Person Pages“ (customized xmosaic parser)
- Person Index

## . SPIRE 1998: Santa Cruz, Bolivia

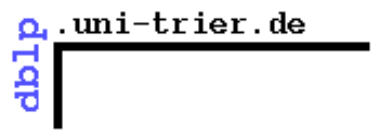
ing Processing and Information Retrieval: A South American Symposium, September 9-11, 1998, Santa Cruz de la Sierra Bolivia. IEEE Computer Society, 1998

- [Mauricio Ayala-Rincón](#), [Paulo D. Conejo](#): A Linear Time Lower Bound on Updating Algorithms for Suffix Trees. 1-6
- [Ricardo A. Baeza-Yates](#), [Jesús Vegas](#), [Gonzalo Navarro](#), [Pablo de la Fuente](#): A Model and a Visual Query Language for Structured Text. 7-13
- [Ricardo A. Baeza-Yates](#), [Gonzalo Navarro](#): Fast Approximate String Matching in a Dictionary. 14-22
- [Ricardo A. Baeza-Yates](#): Searching the Web: Challenges and Partial Solutions (Invited Paper). 23-31
- [Lasse Bergroth](#), [Harri Hakonen](#), [Timo Raita](#): New Approximation Algorithms for Longest Common Subsequences. 32-40
- [Aurelio López-López](#), [Sung H. Myaeng](#): Evidence Accumulation with Competition in Information Retrieval. 41-49
- [Ruy Luiz Milidú](#), [Artur Alves Pessoa](#), [Eduardo Sany Laber](#): In-Place Length-Restricted Prefix Coding. 50-59
- [Andrey A. Mironov](#), [Pavel A. Pevzner](#): SST versus EST in Gene Recognition (Invited Paper). 60-64
- [Matthew Montebello](#): Information Overload - An IR Problem? 65-74
- [Mario A. Nascimento](#), [Adriano C. R. da Cunha](#): An Experiment Stemming Non-Traditional Text. 75-80
- [Arlindo L. Oliveira](#), [João P. Marques Silva](#): Efficient Search Techniques for the Inference of Minimum Size Finite Automata. 81-89
- [Edleno Silva de Moura](#), [Gonzalo Navarro](#), [Nivio Ziviani](#), [Ricardo A. Baeza-Yates](#): Direct Pattern Matching on Compressed Text. 90-95
- [Maria Emilia M. T. Walter](#), [Zanoni Dias](#), [Joao Meidanis](#): Reversal and Transposition Distance of Linear Chromosomes. 96-102
- [Leandro Krug Wives](#), [Stanley Loh](#): Hyperdictionary: A Knowledge Discovery Tool to Help Information Retrieval. 103-110

BLP: [[Home](#) | Search: [Author](#), [Title](#) | [Conferences](#) | [Journals](#)]

opyright © Wed Jan 16 12:22:43 2002 by [Michael Ley](#) ([ley@uni-trier.de](mailto:ley@uni-trier.de))





# Justin Zobel

Online Version of this page: [ACM SIGMOD](#) - [VLDB Endowment](#) - [Uni Trier](#)

[Home Page](#)

## 2001

72	<a href="#">W. Bruce Croft</a> , <a href="#">David J. Harper</a> , <a href="#">Donald H. Kraft</a> , Justin Zobel: SIGIR 2001: Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, September 9-13, 2001, New Orleans, Louisiana, USA. <a href="#">ACM 2001</a>
71	<a href="#">Dennis Bahle</a> , <a href="#">Hugh Williams</a> , Justin Zobel: Compaction Techniques for Nextword Indexes. <a href="#">SPIRE 2001</a> : 33-45
70	<a href="#">Marcin Kaszkiel</a> , Justin Zobel: Effective ranking with arbitrary passages. <a href="#">JASIS 52(4)</a> : 344-364 (2001)
69	<a href="#">Hugh E. Williams</a> , Justin Zobel, <a href="#">Steffen Heinz</a> : Self-adjusting trees in practice for large text collections. <a href="#">Software - Practice and Experience 31(10)</a> : 925-939 (2001)

## 2000

68	<a href="#">EE</a> <a href="#">Daryl J. D'Souza</a> , <a href="#">James A. Thom</a> , Justin Zobel: A Comparison of Techniques for Selecting Text Collections. <a href="#">Australasian Database Conference 2000</a> : 28-32
----	--

## Author-Index: Root Page

- [-- Abber](#)
- [Abbey -- Abd](#)
- [Abe -- Aboa](#)
- [Aboe -- Abreu, P](#)
- [Abreu, S -- Ackerman, A](#)
- [Ackerman, J -- Adams, A. A](#)
- [Adams, A. G -- Addison, M.](#)
- [Addison, Ma -- Adriae](#)
- [Adrian -- Agas](#)
- [Agat -- Agrawal, P.](#)
- [Agrawal, Pr -- Aha](#)
- [Ahe -- Ahn, R](#)
- [Ahn, S -- Aiken, M](#)
- [Aiken, P -- Akav](#)
- [Akaz -- Akm](#)
- [Akn -- Al-Qaq](#)
- [Al-Qas -- Albanese](#)
- [Albanesi -- Alb](#)
- [Alc -- Alexander, K.](#)
- [Alexander, Ke -- Ali, M.](#)
- [Ali, Ma -- Allen, C.](#)
- [Allen, Ch -- Allman](#)
- [Allmar -- Alonso, M.](#)
- [Alonso, M. -- Altham](#)
- [Althan -- Alves, Marce](#)
- [Alves, Marco -- Ambroise, C](#)

...

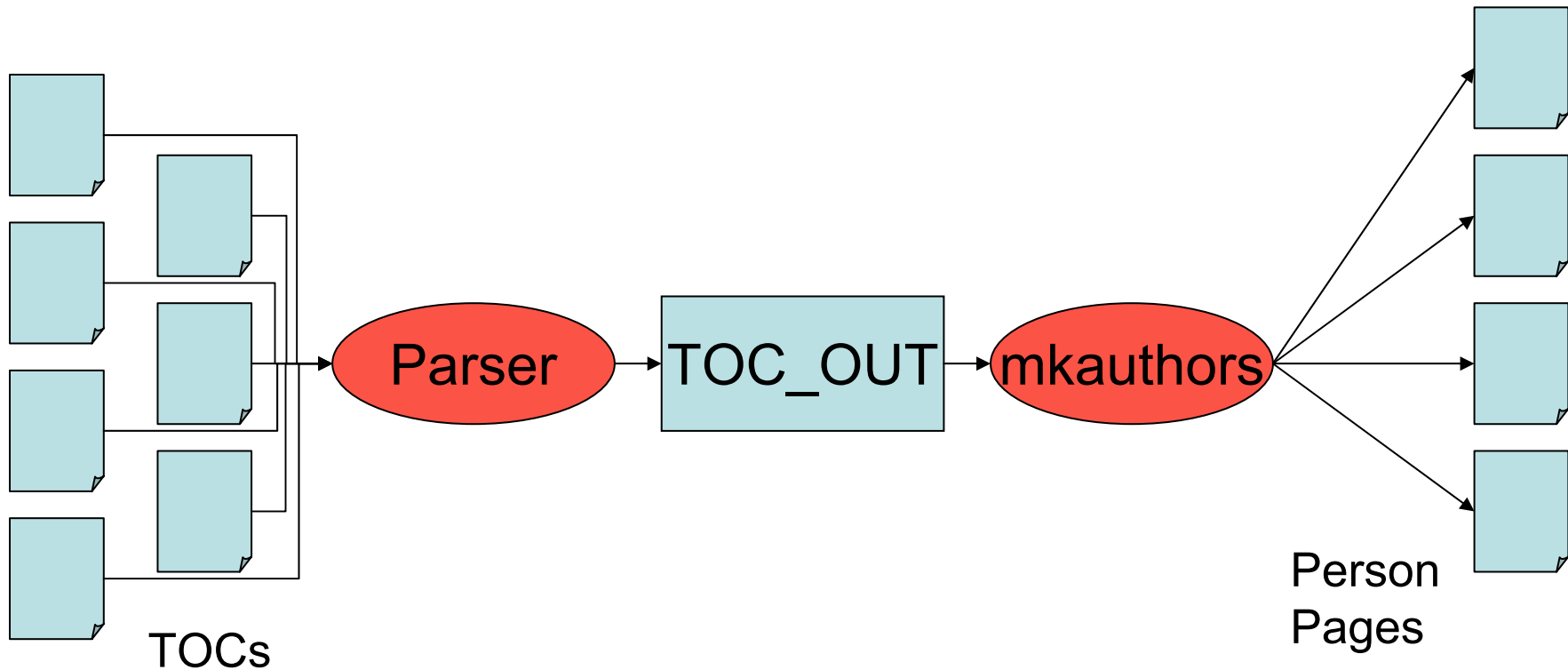
## Author-Index

[previous page](#) - [index root](#)

- [Lewis, S.](#)
- [Lewis, S. M.](#)
- [Lewis, Scott](#)
- [Lewis, Scott M.](#)
- [Lewis, Steven A.](#)
- [Lewis, Steven D.](#)
- [Lewis, Suzanna](#)
- [Lewis, T. A.](#)
- [Lewis, T. G.](#)
- [Lewis, T. H.](#)
- [Lewis, T. S.](#)
- [Lewis, Ted G.\\*](#)
- [Lewis, Vasily](#)
- [Lewis, Wendy](#)
- [Lewis II, Philip M.](#)
- [Lewit, A. F.](#)
- [Lewkowicz, Myriam](#)
- [Lewontin, Steve](#)
- [Lewski, Frank H.](#)
- [Lewyckyj, N.](#)
- [Lexa, Chuck](#)
- [Lexis, Clayton](#)
- [Lext, Jonas](#)
- [Ley, Dominik\\*](#)
- [Ley, Michael\\*](#)
- [Ley, Susan H.](#)

...

# TOC Parser / Generation of Person Pages



# Mirrors

- University of Trier had a slow internet connections
- DBLP should have a high availability
- Technique:
  - transfer all TOCs + entry pages + TOC\_OUT in a tar.gz file
  - run mkauthors on the mirror

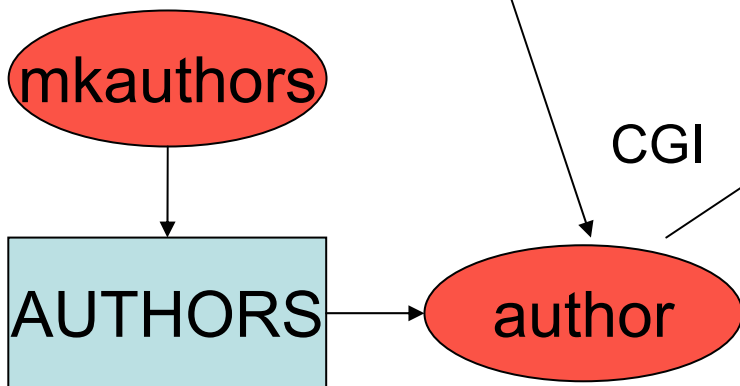
# Search Author

Name:    -> [Help](#)

## Index

DBLP: [[Home](#) | [Search: Author, Title](#) | [Conferences](#) | [Journals](#)]

*Michael Ley (ley@uni-trier.de) Fri Aug 23 15:10:24 2002*



# Search Results for 'Gray'

- [Graydon Barz](#)
- [Gray Clossman](#)
- [V. Grayson Cuglock-Knopp](#)
- [Graydon Davison](#)
- [C. Gray Girling](#)
- [A. Gray](#)
- [A. J. Gray](#)
- [Alex Gray](#)
- [Alexander Gray](#)
- [Alexander G. Gray](#)
- [Andrew R. Gray](#)
- [Brett Gray](#)
- [C. Thomas Gray](#)
- [Cary G. Gray](#)
- [D. Gray](#)
- [D. I. Gray](#)
- [David Gray](#)
- [David N. Gray](#)
- [E. M. Gray](#)
- [F. Gail Gray](#)
- [George Gray](#)
- [Harry J. Gray](#)
- [Heather Gray](#)
- [Helen Frances Gray](#)

...

# Search Title

Keyword:



DBLP: [[Home](#) | Search: [Author](#), [Title](#) | [Conferences](#) | [Journals](#)]  
*Michael Ley (ley@uni-trier.de) Fri Aug 23 15:10:24 2002*

TOC\_OUT

title

CGI

11 September 2002

# Search Results for 'data cube'

- [Yi-Leh Wu](#), [Divyakant Agrawal](#), [Amr El Abbadi](#): Using Wavelet Decomposition and Approximate Range-Sum Queries over Data Cubes. [CIKM 2001](#): 467-476
- [Jeffrey Scott Vitter](#), [Min Wang](#), [Balakrishna R. Iyer](#): Data Cube Approximation Using Wavelets. [CIKM 1998](#): 96-104
- [David Wai-Lok Cheung](#), [Bo Zhou](#), [Ben Kao](#), [Hongjun Lu](#), [Tak Wang](#): Requirement-Based Data Cube Schema Design. [CIKM 1999](#): 162-171
- [Hua-Gang Li](#), [Tok Wang Ling](#), [Sin Yeung Lee](#), [Zheng Xuan Loh](#): Range-Max/Min Queries over Data Cubes. [CODAS 2001](#): 79-86
- [Y. Chung](#), [M. Kim](#), [W. Park](#), [M. Kim](#): Fractionalized View Materialization. [VLDB 2001](#): 156-157
- [Mirek Riedewald](#), [Divyakant Agrawal](#), [Amr El Abbadi](#), [Renato Pajares](#): Data Cube Compression in Dynamic Environments. [DaWaK 2000](#): 24-33
- [Pedro Furtado](#), [Henrique Madeira](#): Data Cube Compression with Query Support. [VLDB 2001](#): 167
- [Hua-Gang Li](#), [Tok Wang Ling](#), [Sin Yeung Lee](#): Range-Max/Min Queries over Data Cubes. [VLDB 2000](#): 467-476
- [ZhongWei Luo](#), [Tok Wang Ling](#), [Chuan -eng Ang](#), [Sin Yeung Lee](#): Range-Max/Min Queries in OLAP Sparse Data Cubes. [DEXA 2001](#): 678-687
- [Mauricio Minuto Espil](#), [Alejandro A. Vaisman](#), [Leonardo Terribile](#): Data Cube Compression: a rule-based perspective. [DMDW 2002](#): 72-81
- [Chris Jermaine](#), [Renée J. Miller](#): Approximate Query Answering in High-Dimensional Data. [SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery 2002](#): 10-19
- [Yuping Yang](#), [Mukesh Singhal](#): Accessing Data Cubes along Compressed Dimensions. [VLDB 2002](#): 10-19

# Bibliographic Records

citation linking, reviews, annotated bibliographies, ...

- assign an unique ID to each publication
- make it accessible by this ID
- store the information in classical bibliographic records

# Bibliographic Records

- You may download them from <http://dblp.uni-trier.de/xml/>  
(uncompressed ~120MByte)
- Simple DTD
- Idea: BibTeX++ in XML syntax
- Examples ...



```
<article key="journals/wi/OberweisS91">
<author>Andreas Oberweis</author>
<author>Wolffried Stucky</author>
<title>Die Behandlung von Ausnahmen in
Software-Systemen: Eine
Literatur&uuml;bersicht.</title>
<pages>492-502</pages>
<year>1991</year>
<volume>33</volume>
<journal>Wirtschaftsinformatik</journal>
<number>6</number>
<url>db/journals/wi/wi33.html#OberweisS91
</url>
</article>
```

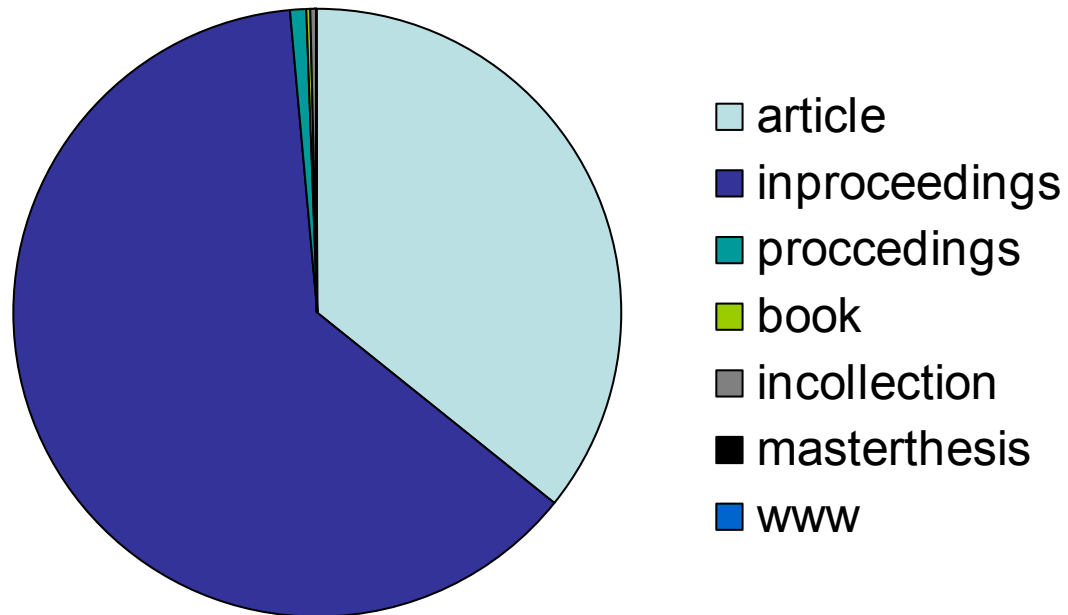
```
<inproceedings key="conf/gi/OberweisPS01">  
<author>Andreas Oberweis</author>  
<author>Oliver Paulzen</author>  
<author>Hagen J. Sexauer</author>  
<title>Ein wissensbasiertes Vorgehensmodell  
zur Gestaltung von CRM-Systemen.</title>  
<pages>429-436</pages>  
<year>2001</year>  
<booktitle>GI Jahrestagung (1)</booktitle>  
<url>db/conf/gi/gi2001-1.html#OberweisPS01  
</url>  
</inproceedings>
```

```
<inproceedings key="conf/er/JaeschkeOS93">  
<author>Peter Jaeschke</author>  
<author>Andreas Oberweis</author>  
<author>Wolffried Stucky</author>  
<title>Extending ER Model Clustering by  
      Relationship Clustering.</title>  
<pages>451-462</pages>  
<year>1993</year>  
<booktitle>ER</booktitle>  
<url>db/conf/er/er93.html#JaeschkeOS93</url>  
<crossref>conf/er/93</crossref>  
<cdrom>er93/ER93-P447.pdf</cdrom>  
<ee>db/conf/er/JaeschkeOS93.html</ee>
```

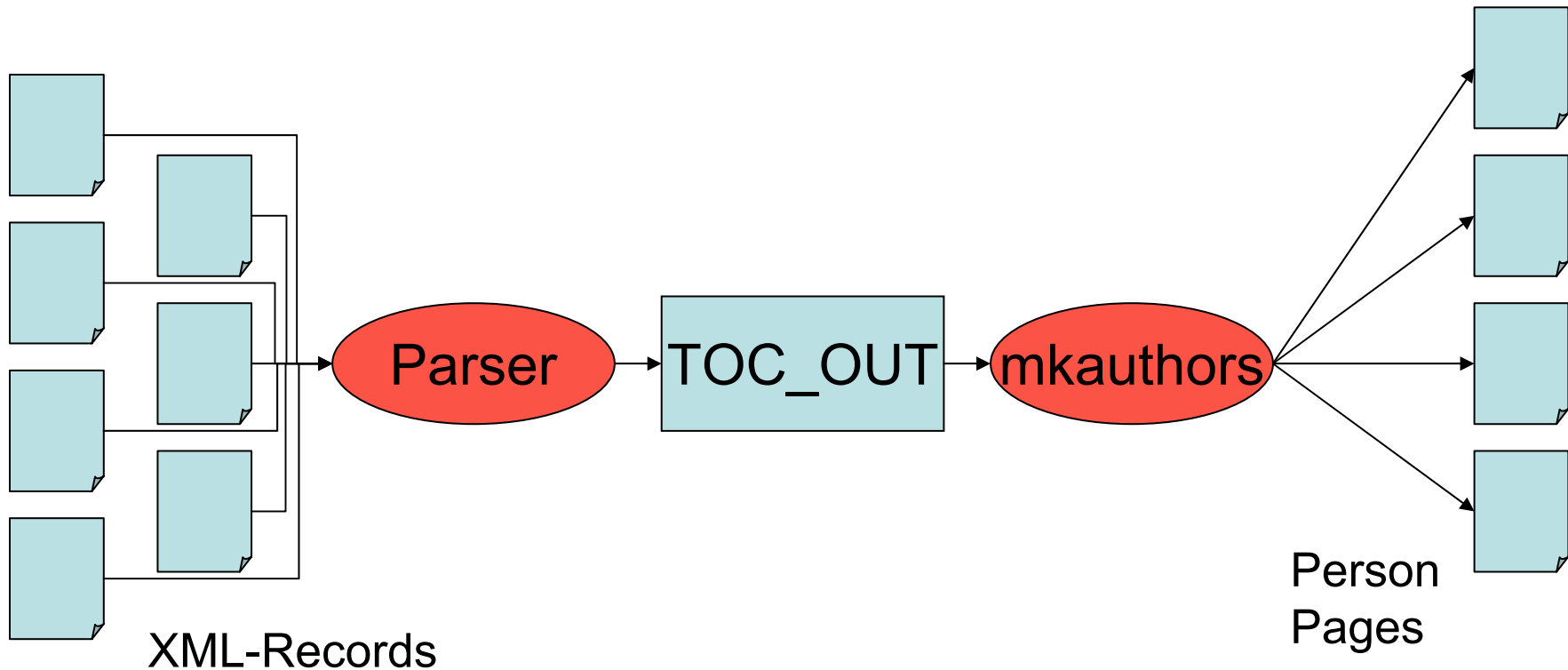
```
<cite label="CaJA89">conf/er/CarlsonJA89</cite>
<cite label="Chen76">journals/tods/Chen76</cite>
<cite label="FeMi86">journals/cj/FeldmanM86</cite>
<cite label="Mart89">...</cite>
<cite label="Mist91">...</cite>
<cite label="RaSt92">conf/er/RauhS92</cite>
<cite label="ScSt83">...</cite>
<cite label="ScSW79">conf/er/ScheuermannSW79</cite>
<cite label="TeYF86">journals/csur/TeoreyYF86</cite>
<cite label="TWBK89">journals/cacm/TeoreyWBK89</cite>
</inproceedings>
```

```
<proceedings key="conf/er/93">
<editor>Ramez Elmasri</editor>
<editor>Vram Kouramajian</editor>
<editor>Bernhard Thalheim</editor>
<title>Entity-Relationship Approach - ER'93,
  12th International Conference on the
  Entity-Relationship Approach, Arlington,
  Texas, USA, December 15-17, 1993,
  Proceedings</title>
<booktitle>ER</booktitle>
<series href="db/journals/lncs.html">Lecture
  Notes in Computer Science</series>
<volume>823</volume>
<publisher>Springer</publisher>
<year>1994</year>
<isbn>3-540-58217-7</isbn>
<url>db/conf/er/er93.html</url>
</proceedings>
```

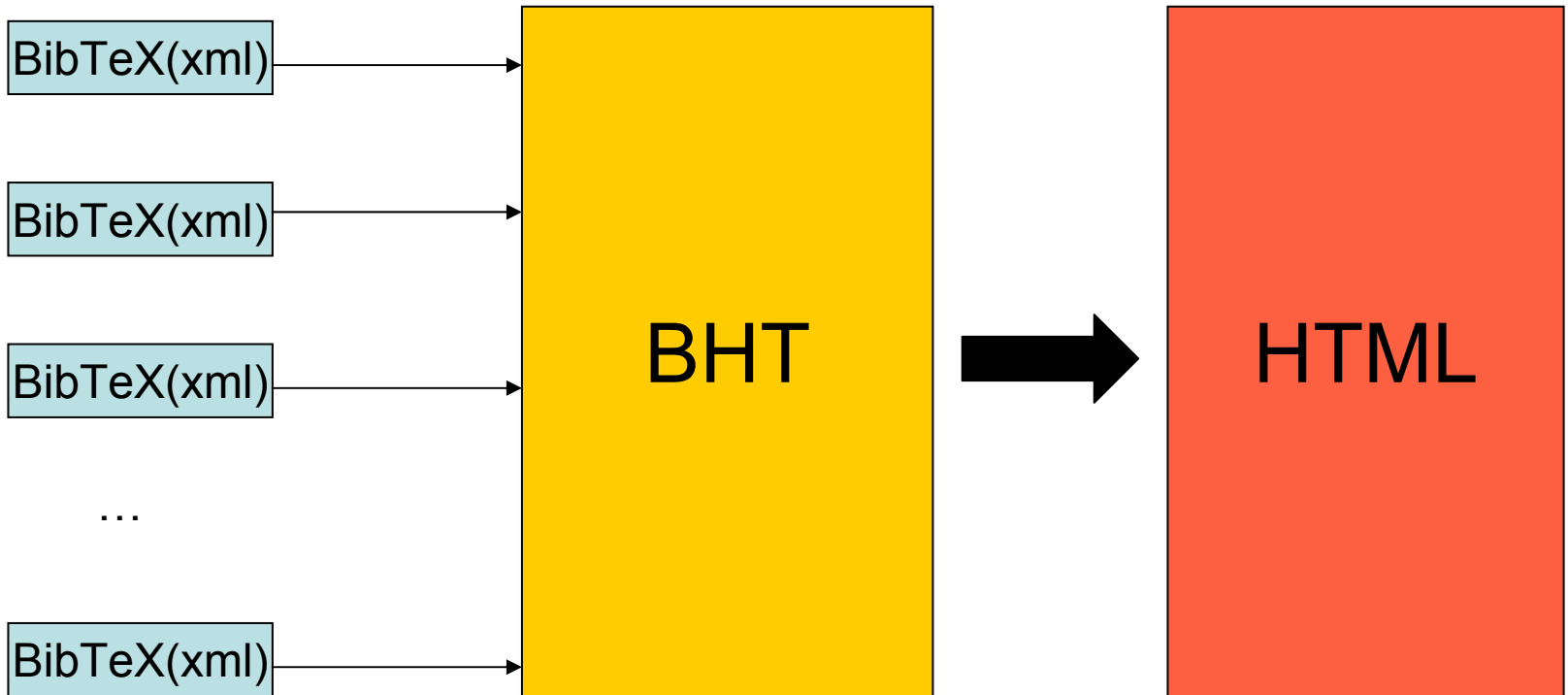
# Publication Types



# XML Parser / Generation of Person Pages



# Tables of Contents, ...





# „Bibliography HyperText“ (BHT)

- include mechanism  
(`<cite key='...' style='... '>`)
- several additional tags:  
`<logo>`, `<footer>`, `<ref href="...">`, ...
- HTML

```
<html><head><title>25. SIGIR 2002: Tampere, Finland</title></head>
<body bgcolor="#ffffff" text="#000000" link="#000000"><logo style="sigir">
<h1>25. <ref href="db/conf/sigir/index.html">SIGIR</ref> 2002: Tampere,
Finland</h1><hr>
<cite key="conf/sigir/2002"><cite key="conf/sigir/2002" style="bibtex">
<center>

</center>
<ul>
<li><cite key="conf/sigir/Rijsbergen02" style=ee>
</ul>
<h2>Web Information Retrieval</h2>
<ul>
<li><cite key="conf/sigir/AnhM02" style=ee>
<li><cite key="conf/sigir/ParkPGK02" style=ee>
<li><cite key="conf/sigir/SiC02" style=ee>
<li><cite key="conf/sigir/KraaijWH02" style=ee>
</ul>
<h2>Information Retrieval Theory</h2>    ...
```

# Search

Authors	<input type="text" value="gray"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Title	<input type="text" value="cube"/>	Year <input type="text"/>	Page <input type="text"/>	
Conference	<input type="text"/>	ID <input type="text"/>		
Journal	<input type="text"/>	Volume <input type="text"/>	Number <input type="text"/>	
<input type="button" value="Search"/>	<input type="button" value="Reset"/>	Maximum of <input type="text" value="100"/> matches		

This search engine is powered by the [MG Software](#).

ACM SIGMOD Anthology - DBLP: [[Home](#) | [Search: Author](#), [Title](#) | [Conferences](#) | [Journals](#)]

## Search Result

Query: author = "gray", title = "cube"

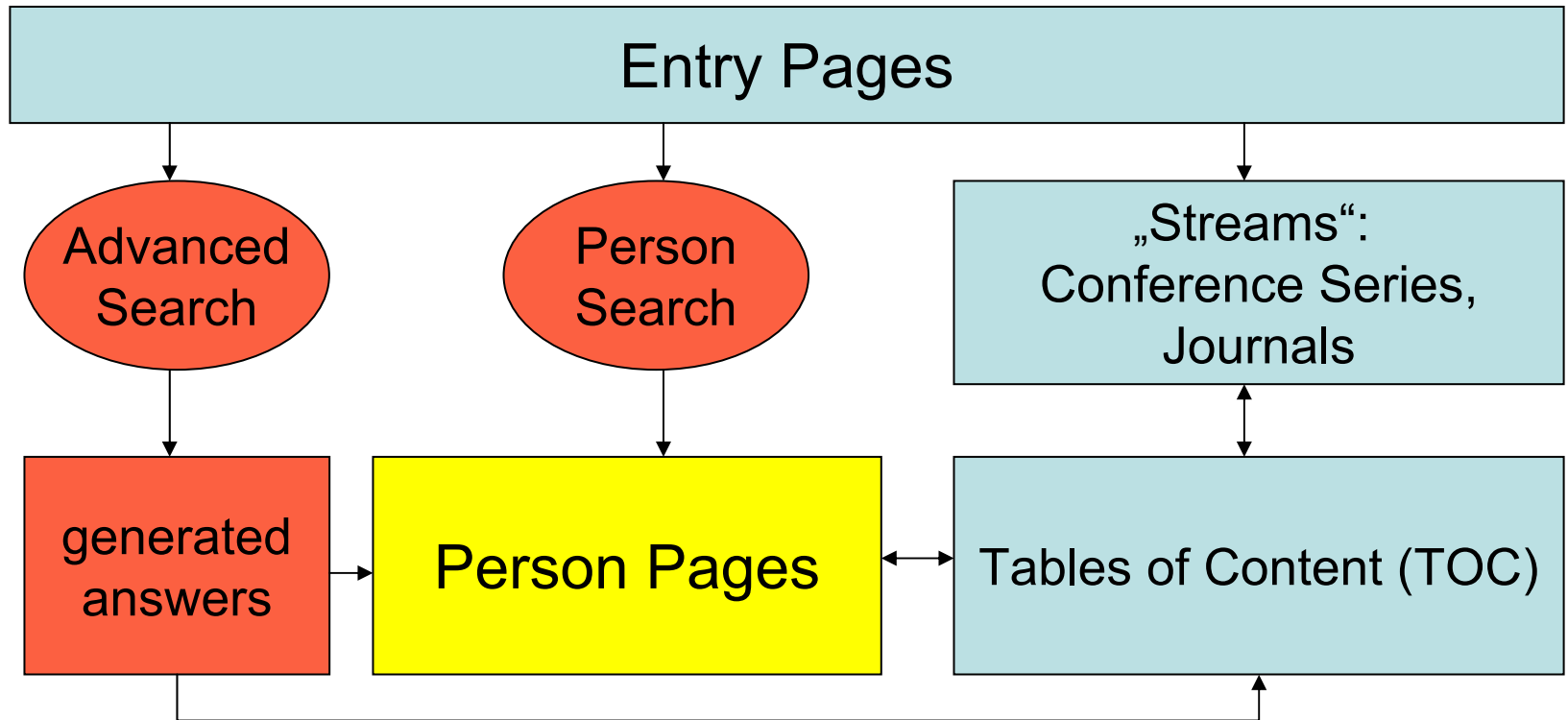
- |    |   |
|----|---|
| EE | <p><u>Jim Gray</u>, <u>Adam Bosworth</u>, <u>Andrew Layman</u>, <u>Hamid Pirahesh</u>: Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tab, and Sub-Total. <u>ICDE 1996</u>: 152-159 [DBLP:conf/icde/GrayBLP96]</p>   |
|    | <p><u>J. R. Armstrong</u>, <u>F. Gail Gray</u>: Fault Diagnosos in a Boolean <math>n</math> Cube Array of Microprocessors. <u>IEEE Transactions on Computers</u> 30(8): 587-590 (1981) [DBLP:journals/tc/ArmstrongG81]</p>  |
|    | <p><u>Jim Gray</u>, <u>Surajit Chaudhuri</u>, <u>Adam Bosworth</u>, <u>Andrew Layman</u>, <u>Don Reichart</u>, <u>Murali Venkatrao</u>, <u>Frank Pellow</u>, <u>Hamid Pirahesh</u>: Data Cube: A Relational Aggregation Operator Generalizing Group-by, Cross-Tab, and Sub Totals. <u>Data Mining and Knowledge Discovery</u> 1(1): 29-53 (1997) [DBLP:journals/datamine/GrayCBLRVPP97]</p> |

ACM SIGMOD Anthology - DBLP: [[Home](#) | [Search](#) | [Conferences](#) | [Journals](#)]

# „Advanced Search“: MG

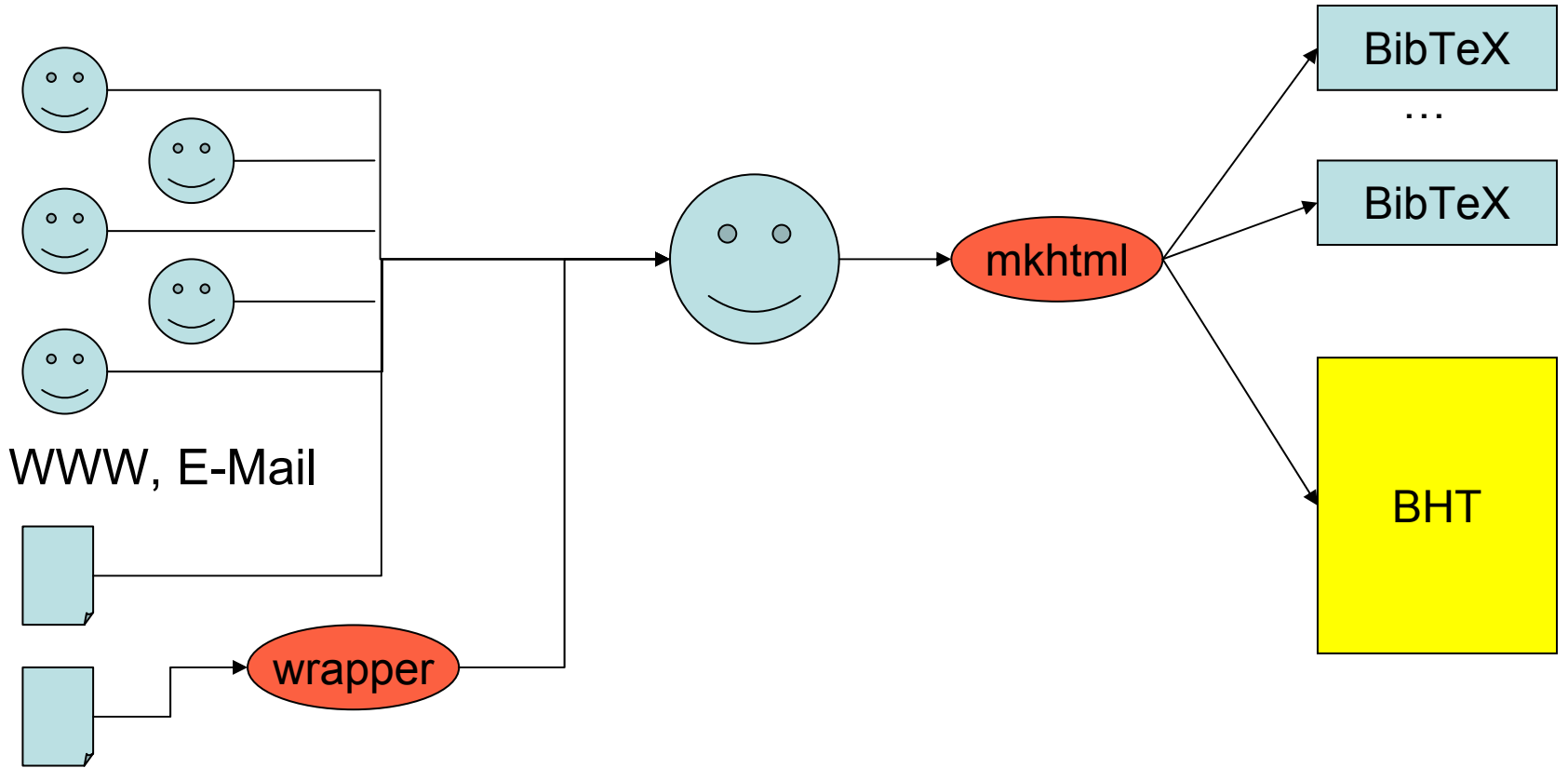
- „Managing Gigabytes“ Software by Witten, Moffat, Bell
- DBLP XML Records → MG Documents
- Filter: matching terms in the required field

# DBLP Architecture



# Entering Data

Students



# Anthology: Offline Search

- Simple search engine for all standard platforms → Java Applet (Java 1.0)
- Only author search
- Applet reads simple tree data structure from files:



root:  
separators

...  
curran,john H  
cutting,d  
czarnecki,h  
dabbaghc  
dagmar br  
daigr  
dale ro  
dambr  
dan cart  
dan sm  
dang,z  
daniel d. fu  
daniel j. bur  
daniel mend  
daniel sa  
daniela mer  
dannenber,c  
darad  
...

leaves: names  
(permutations  
of name parts)

...  
Mann,Robert I.  
Mann,Sally Fahrenholz-  
Mann,Samuel  
Mann,Stefan  
Mann,Stephen  
...  
Mann-Ho Lee  
Mann-May Yau  
Manna,I.  
Manna,M. La  
Manna,Serena La  
Manna,Zohar  
Mannai,Dhamir N.  
Mannarino,Gabriela Susana  
Mannava,Phanindra K.  
Manne,Fredrik  
Manne,Srilatha  
Manneback,Pierre  
...

Search Author - Microsoft Internet Explorer

Datei Bearbeiten Ansicht Favoriten Extras ?

Zurück Zurück Suchen Favoriten Medien

Adresse D:\db\indices\A-tree\index.html Wechseln zu Links >>

Anthology ACM SIGMOD dblp.uni-trier.de

## Search Author

### Applet Search

Name:  Search Reset

open Page (select with double click):

Marcelo Borges Ribeiro  
Berthier A. Ribeiro-Neto  
João Bosco Ribeiro do Val

Applet AuthorSearch started Arbeitsplatz

# Anthology: Full Text Search

- Acrobat „catalog“: index construction
- „capture“: OCR for scanned texts
- manual entry of title fields was necessary
  - insufficient software support by Adobe
- Acrobat Reader „with search“
  - not available for Linux



The ACM SIGMOD



# Anthology

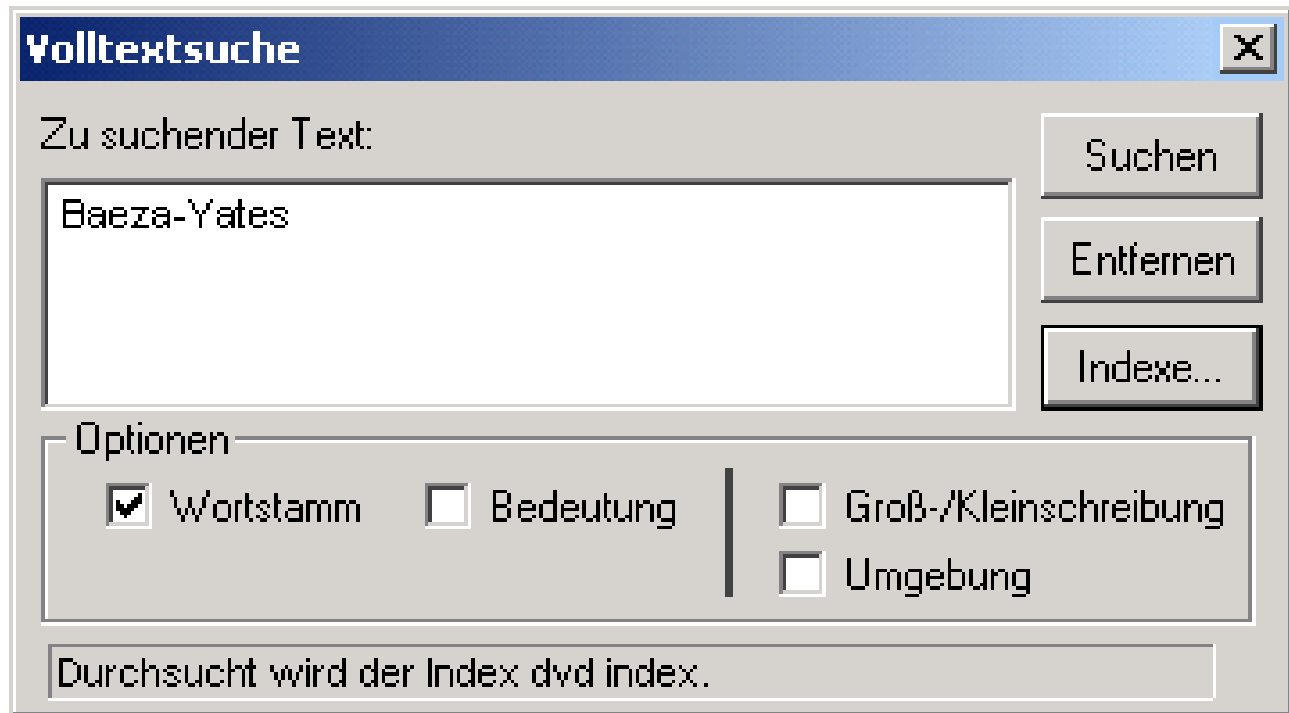
DVD 2, 2001

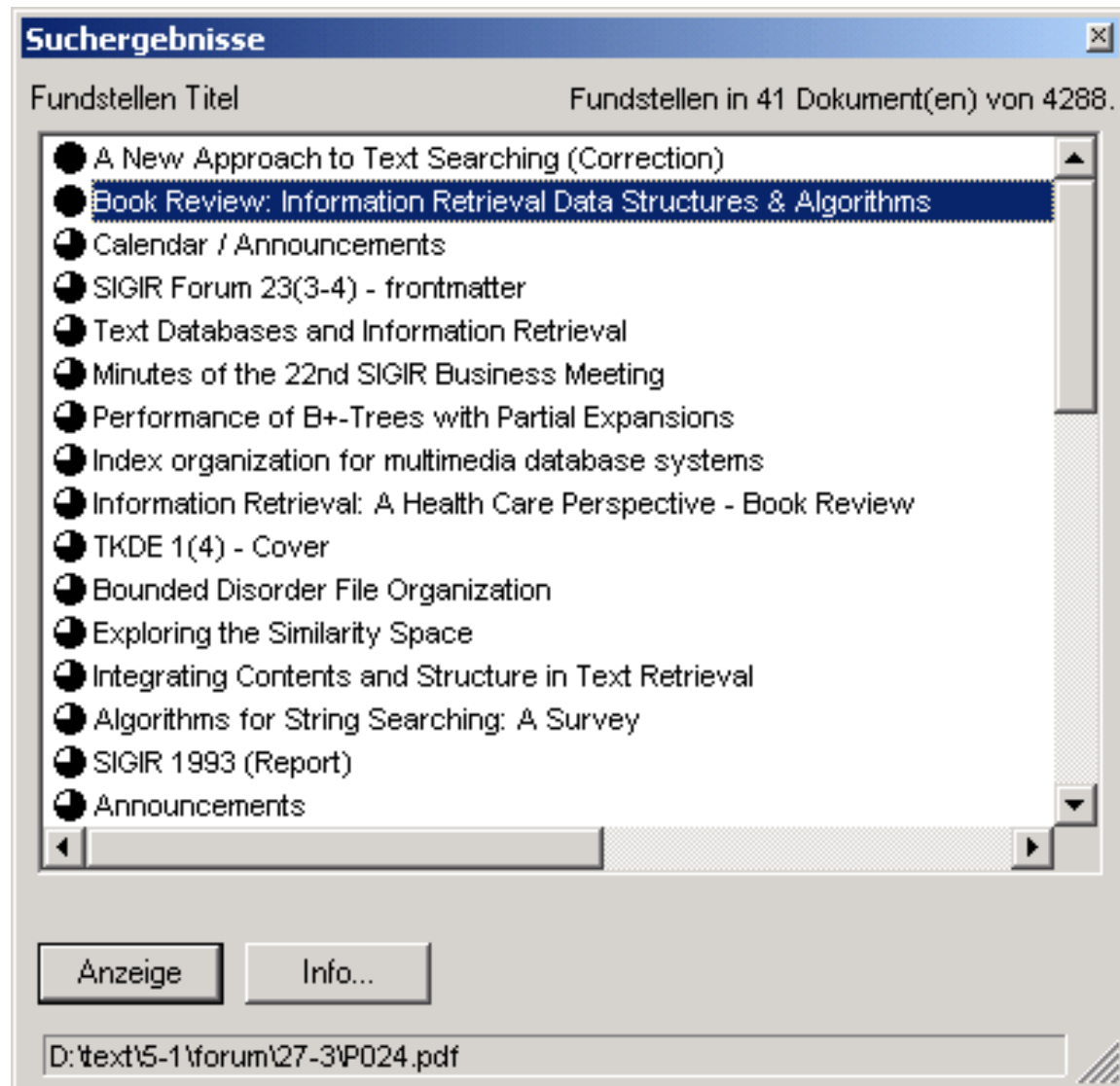
Table of Contents

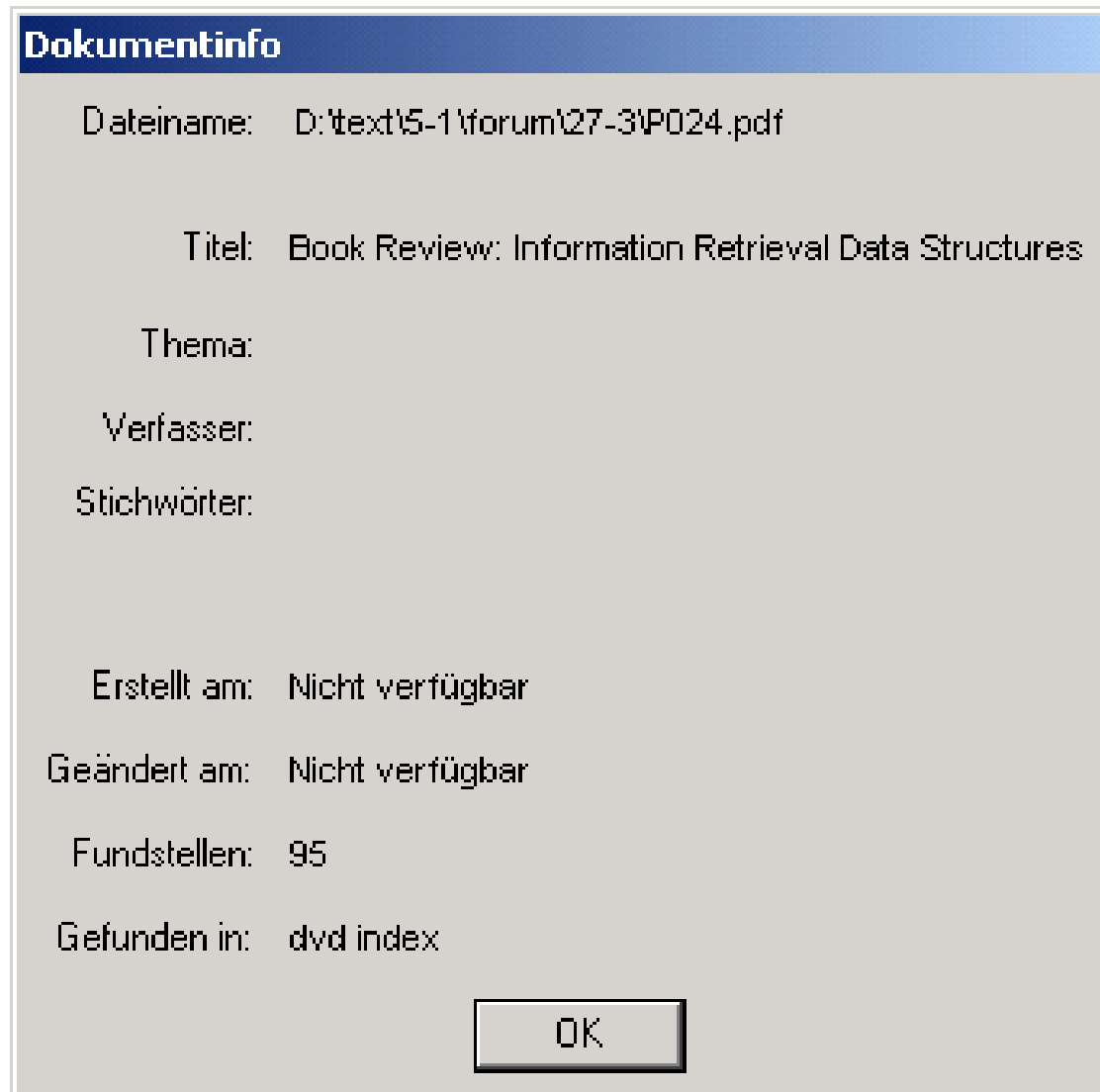
Full Text Index

Starts Web Browser

requires Acrobat Reader  
WITH SEARCH







William B. Frakes (Software Engineering Guild), Ricardo **Baeza-Yates** (University of Chile), editors. **Information Retrieval Data Structures & Algorithms**. Prentice Hall, Englewood Cliffs, New Jersey, 1992, 504 pp., \$50.00 (Retail Price), ISBN 0-13-463837-9.

This is an extremely welcome addition to the Information Retrieval (IR) literature. Because of its technical approach it is much different from most of the available books on IR. The book consists of five sections containing eighteen chapters. The chapters are written by different authors.

The *introduction* section is made-up of two chapters and each chapter is written by one of the editors. The chapters acquaint the reader with various IR related concepts and provide pointers for the other chapters of the book.

The first chapter (by D. Harman, E. Fox, R. A. **Baeza-Yates**, W. Lee) of the *file structures* section provides a survey of techniques that can be used in the construction of inverted files then gives the details of two techniques that can be used for large document databases. The next chapter (by C. Faloutsos) provides a classification of signature-based methods and a short survey of the related studies. An alternative to inverted and signature files, the relatively new data structures, PAT trees and PAT arrays, is the topic of Chapter 5 (by G. H. Gonnet, R. A. **Baeza Yates**, T. Snider). This



## 3. Perspectives & Research Issues

- Funding & collaborations
- Person Names
- DBLP Browser
- Visualization
- Classification / Clustering

# Informatics Portal

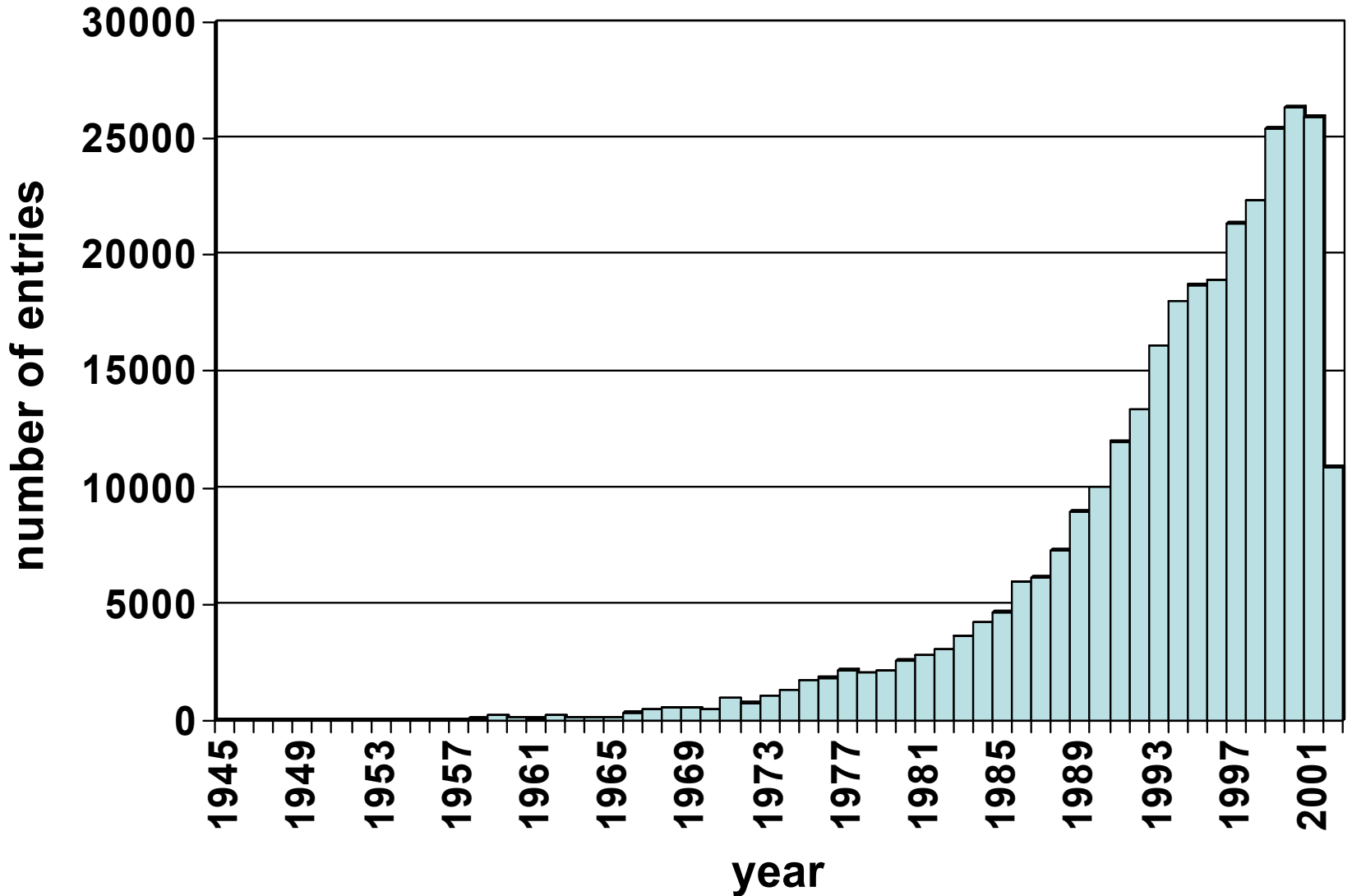
- DBLP will be a central part of a larger Informatics Portal
- Other parts:
  - CompuScience (FIZ Karlsruhe)
  - Collection of CS Bibliographies (Achilles/Vollmer, Karlsruhe)
  - LeaBib (Mayr, Munich)

## Informatics Portal (2)

- Funded by the German Federal Ministry of Education and Research, 3 years
- In Cooperation with the „Gesellschaft für Informatik“ (German Computer Society)
- 2 open positions at Trier:  
service+research

# Additional Contents ...

- complete coverage of LNCS
- more ACM & IEEE publications
- cooperation with IFIP, Usenix, ...
- cooperation with the libraries at Dagstuhl & Max-Planck-Institut für Informatik (Saarbrücken)
- several German series (LNI, IA, IFB...)



# Challenge: How to manage the exponential growth of (computer) science publications ?

- maintenance problem
- „lost in publication space“

# Person Names

- widely accepted method to identify persons
- names may not be unique
- a person may change her/his name (marriage, emigration to other cultural environment)
- variations of person names ...

- abbreviations: Jeffrey D. Ullman, J. D. Ullman, J. Ullman, Jeff Ullman, ...
- nicknames: Michael / Mike, William / Bill, Joseph / Joe
- permutations: Liu Bin / Bin Liu
- different transcriptions: Andrei / Andrej / Andrey
- accents: Stephane / Stéphane
- umlauts: Muller / Müller / Mueller



- ligatures: Weiß / Weiss, Åström / Aastrom
- case: Al-A'Ali / Al-A'ali
- hyphens: Hans-Peter / Hans Peter
- composition: MaoLin / Mao Lin, Kenichi / Ken-ichi / Ken'ichi
- postfixes: Karel Culik II, Jr. / Sr.
- typos

# Person Names: Problems

- there should be a 1:1 mapping between persons and person pages
- how to search persons best ?
- how to normalize different spellings ?

name normalization costs more than 60% of our time ...

# Person Name Normalization

- for each new entry we try to locate the authors/editors in the existing collection
- if spellings differ, but we are confident that they are variations of the same person's name → make them equal
- write out most parts of the name
- person's preferred spelling ?

- for persons with many publications the name spelling usually converges to a stable & correct state
- for persons with a few known publications it is more likely that there are duplicate person pages, incorrect or incomplete spellings

# Heuristics in the decision process

- coauthor relationship gives strong indications for the identity of persons
- streams (journals/conference series): condensation points for communities, weaker indication
- same keywords in titles
- time frame ?
- a lot of background knowledge ...

# Wanted: Tool to make DBLP maintenance more efficient

- specialized weight functions for the name normalization problem
  - query: new entry = list of names, title, publication stream, year, ...
- „DBLP browser“: framework for experiments

# DBLP Browser: Roles

- maintenance tool
- bibliographic tool for users of DBLP
  - composition of reference lists
  - export in popular formats like BibTeX
- platform for experiments in visualization (SemIPort project ...)

# DBLP Browser

- main memory IR system
  - compressed representation of the bibliographic records
- convenient graphical user interface
- Java / Swing
- first prototype implemented



## Welcome to DBLPbr

DBLPbr is an experimental browser for the DBLP computer science bibliography ([dblp.uni-trier.de](http://dblp.uni-trier.de)). It loaded 308223 bibliographic records into it's main memory database and provides fast access to the "author pages". The current version does not contain the complete table of contents information with session titles from the Web version of DBLP.

- [main memory database statistics](#)
- [top authors](#)
- [publications per year](#)
- [journals](#)
- [conferences](#)

---

© Copyright 2002, Michael Ley ([ley@uni-trier.de](mailto:ley@uni-trier.de))

Dblpbr is free software; you can redistribute it and/or modify it under the terms of the [GNU General Public License](#) as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version. Dblpbr is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with dblpbr; if not, write to the [Free Software Foundation, Inc.](#), 59 Temple Place, Suite 330, Boston, MA 02111-1307 USA

---

back forward home clear Name: i witten

a=A a=

Ian H. Witten (132)

by year

- 2001
- 2000
- 1999
- 1998
- 1997
- 1996
- 1995
- 1994
- 1993
- 1992
- 1991
- 1990
- 1989
- 1988
- 1987
- 1985
- 1984
- 1983
- 1982
- 1981
- 1980
- 1977
- 1976
- 1975

## Ian H. Witten: 2001

101	EE	Ian H. Witten, <a href="#">David Bainbridge</a> , <a href="#">Stefan J. Boddie</a> : Greenstone: Open-source DL software. <i>CACM 44</i> (5): 47 (2001) * [ <a href="#">journals/cacm/WittenEB01</a> ]
100	EE	Ian H. Witten, <a href="#">Michel Loots</a> , <a href="#">Maria F. Trujillo</a> , <a href="#">David Bainbridge</a> : The promise of digital libraries in developing countries. <i>CACM 44</i> (5): 82-85 (2001) * [ <a href="#">journals/cacm/WittenLTB01</a> ]
99	EE	<a href="#">Mark D. Apperley</a> , <a href="#">Sally Jo Cunningham</a> , <a href="#">Te Taka Keegan</a> , Ian H. Witten: NIUPEPA: a historical newspaper collection. <i>CACM 44</i> (5): 86-87 (2001) * [ <a href="#">journals/cacm/ApperleyCKW01</a> ]
98	EE	Ian H. Witten, <a href="#">David Bainbridge</a> , <a href="#">Stefan J. Boddie</a> : Greenstone: Open-Source Digital Library Software. <i>D-Lib Magazine 7</i> (10): (2001) * [ <a href="#">journals/dlib/WittenEB01</a> ]
97	EE	<a href="#">Stuart Yeates</a> , Ian H. Witten, <a href="#">David Bainbridge</a> : Tag Insertion Complexity. <i>Data Compression Conference 2001</i> : 243-252 [ <a href="#">conf/dcc/YeatesWB01</a> ]
96	EE	<a href="#">Gordon W. Paynter</a> , Ian H. Witten: A Combined Phrase and Thesaurus Browser for Large Document Collections. <i>ECDL 2001</i> : 25-36 [ <a href="#">conf/ercimd/PaynterW01</a> ] <sup>Ⓢ</sup>
95	EE	<a href="#">David Bainbridge</a> , <a href="#">George Buchanan</a> , <a href="#">John McPherson</a> , <a href="#">Steve Jones</a> , <a href="#">Abdelaziz Mahoui</a> , Ian H. Witten: Greenstone: A Platform for Distributed Digital Library Applications. <i>ECDL 2001</i> : 137-148 [ <a href="#">conf/ercimd/BainbridgeBMJMW01</a> ] <sup>Ⓢ</sup>
94		<a href="#">Te Taka Keegan</a> , <a href="#">Mark D. Apperley</a> , <a href="#">Sally Jo Cunningham</a> , Ian H. Witten: The Niupepa Collection: Opening the Blinds on a Window to the Past. <i>ICHIM (1) 2001</i> : 347-356 [ <a href="#">conf/ichim/KeeganACW01</a> ]
93	EE	<a href="#">Malcolm Ware</a> , <a href="#">Eibe Frank</a> , <a href="#">Geoffrey Holmes</a> , <a href="#">Mark Hall</a> , Ian H. Witten: Interactive machine learning: letting users build classifiers. <i>International Journal of Human Computer Studies 55</i> (2): 291-303 (2001) *

Name: 



Ian H. Witten (132)

 by year

 by coauthor

- Robert M. Akscyn
- Stuart Anderson
- Mark D. Apperley
- R. H. Atkin
- Joscha Bach
- David Bainbridge
- Timothy C. Bell
- Graham M. Birtwistle
- Stefan J. Boddie
- Mike Bonham
- Bob Bramwell
- Zane Bray
- George Buchanan
- Chang Chui
- John G. Cleary
- Darrell Conklin
- Chris Corbett
- George Coulouris
- Sally Jo Cunningham
- John J. Darragh
- Eibe Frank
- Brian R. Gaines

### Ian H. Witten: joint work with Sally Jo Cunningham

99	EE	<a href="#">Mark D. Apperley</a> , <a href="#">Sally Jo Cunningham</a> , <a href="#">Te Taka Keegan</a> , Ian H. Witten: NIUPEPA: a historical newspaper collection. <i>CACM</i> 44 (5): 86-87 (2001) * [ <a href="#">journals/cacm/ApperleyCKW01</a> ]
94		<a href="#">Te Taka Keegan</a> , <a href="#">Mark D. Apperley</a> , <a href="#">Sally Jo Cunningham</a> , Ian H. Witten: The Niupepa Collection: Opening the Blinds on a Window to the Past. <i>ICHIM</i> (1) 2001: 347-356 [ <a href="#">conf/ichim/KeeganACW01</a> ]
89	EE	<a href="#">Gordon W. Paynter</a> , Ian H. Witten, <a href="#">Sally Jo Cunningham</a> , <a href="#">George Buchanan</a> : Scalable browsing for large collections: a case study. <i>ACM DL</i> 2000: 215-223 [ <a href="#">conf/dl/PaynterWCB00</a> ]
75		Ian H. Witten, <a href="#">Rodger J. McNab</a> , <a href="#">Steve Jones</a> , <a href="#">Mark D. Apperley</a> , <a href="#">David Bainbridge</a> , <a href="#">Sally Jo Cunningham</a> : Managing Complexity in a Distributed Digital Library. <i>IEEE Computer</i> 32 (2): 74-79 (1999) * [ <a href="#">journals/computer/WittenMJABC99</a> ]
67		Ian H. Witten, <a href="#">Craig G. Nevill-Manning</a> , <a href="#">Rodger J. McNab</a> , <a href="#">Sally Jo Cunningham</a> : A Public Library Based on Full-text Retrieval. <i>CACM</i> 41 (4): 71-75 (1998) * [ <a href="#">journals/cacm/WittenNMC98</a> ]
53	EE	<a href="#">Rodger J. McNab</a> , <a href="#">Lloyd A. Smith</a> , Ian H. Witten, <a href="#">Clare L. Henderson</a> , <a href="#">Sally Jo Cunningham</a> : Towards the Digital Music Library: Tune Retrieval from Acoustic Input. <i>Digital Libraries</i> 1996: 11-18 [ <a href="#">conf/dl/McNabSWHC96</a> ]
51		Ian H. Witten, <a href="#">Sally Jo Cunningham</a> , <a href="#">Mahendra Vallabh</a> , <a href="#">Timothy C. Bell</a> : A New Zealand Digital Library for Computer Science Research. <i>DL</i> 1995: 0- [ <a href="#">conf/dl/WittenCVB95</a> ]

back forward home clear Name: i witten

a=A a=

Ian H. Witten (132)

- by year
- by coauthor
- by journal
  - ACM Computing Surveys
  - AI & Society
  - CACM
  - Computational Linguistics
  - D-Lib Magazine
  - IEEE Computer
  - IEEE Transactions on
  - Information Processing
  - Information and Control
  - Int. J. on Digital Libraries
  - Interacting with Computers
  - International Journal of
  - International Journal of
  - JACM
  - JAIR
  - JASIS
  - Knowledge and Information Systems
  - Machine Learning
  - Multimedia Tools and Applications
  - Operating Systems Reviews
  - SIGART Bulletin

### Ian H. Witten: publications in IEEE Computer

75	<p>Ian H. Witten, <a href="#">Rodger J. McNab</a>, <a href="#">Steve Jones</a>, <a href="#">Mark D. Apperley</a>, <a href="#">David Bainbridge</a>, <a href="#">Sally Jo Cunningham</a>: Managing Complexity in a Distributed Digital Library. <i>IEEE Computer</i> 32 (2): 74-79 (1999) * [<a href="#">journals/computer/WittenMJABC99</a>]</p>
59	<p><a href="#">David Maulsby</a>, Ian H. Witten: Teaching Agents to Learn: From User Study to Implementation. <i>IEEE Computer</i> 30 (11): 36-44 (1997) * [<a href="#">journals/computer/MaulsbtW97</a>]</p>
44	<p><a href="#">Harold W. Thimbleby</a>, <a href="#">Stuart Inglis</a>, Ian H. Witten: Displaying 3D Images: Algorithms for Single-Image Random-Dot Stereograms. <i>IEEE Computer</i> 27 (10): 38-48 (1994) * [<a href="#">journals/computer/ThimblebyTW94</a>]</p>
27	<p><a href="#">John J. Darragh</a>, Ian H. Witten, <a href="#">Mark L. James</a>: The Reactive Keyboard: A Predictive Typing Aid. <i>IEEE Computer</i> 23 (11): 41-49 (1990) * [<a href="#">journals/computer/DarraghWJ90</a>]</p>

back forward home clear Name: i witten

a=A a-

Ian H. Witten (132)

- by year
- by coauthor
- by journal
- by conf
- ACM DL
- ADL
- AISB (ECAI)
- CPM
- DL
- Data Compression Conf
- Digital Libraries
- ECDL
- EWSL
- ICHIM (1)
- IFIP Congress
- IJCAI
- IJCAI (2)
- JCDL
- Learning for Natural Lang
- ML
- New Results and New Tr
- PAKDD
- PRICAI Workshop on Tex
- SIGIR

## Ian H. Witten: publications in Data Compression Conference

97	EE	Stuart Yeates, Ian H. Witten, David Bainbridge: Tag Insertion Complexity. <i>Data Compression Conference 2001</i> : 243-252 [ <a href="#">conf/dcc/YeatesWB01</a> ]
86	EE	Eibe Frank, Chang Chui, Ian H. Witten: Text Categorization Using Compression Models. <i>Data Compression Conference 2000</i> : 555 [ <a href="#">conf/dcc/FrankCW00</a> ]
85	EE	Stuart Yeates, David Bainbridge, Ian H. Witten: Using Compression to Identify Acronyms in Text. <i>Data Compression Conference 2000</i> : 582 [ <a href="#">conf/dcc/YeatesBW00</a> ]
78	EE	Ian H. Witten, Zane Bray, Malika Mahoui, W. J. Teahan: Text Mining: A New Frontier for Lossless Compression. <i>Data Compression Conference 1999</i> : 198-207 [ <a href="#">conf/dcc/WittenBMT99</a> ]
77	EE	Craig G. Nevill-Manning, Ian H. Witten: Protein is Incompressible. <i>Data Compression Conference 1999</i> : 257-266 [ <a href="#">conf/dcc/Nevill-ManningW99</a> ]
76	EE	Joscha Bach, Ian H. Witten: Lexical Attraction for Text Compression. <i>Data Compression Conference 1999</i> : 516 [ <a href="#">conf/dcc/BachW99</a> ]
66	EE	Craig G. Nevill-Manning, Ian H. Witten: Phrase Hierarchy Inference and Compression in Bounded Space. <i>Data Compression Conference 1998</i> : 179-188 [ <a href="#">conf/dcc/Nevill-ManningW98</a> ]
60		Craig G. Nevill-Manning, Ian H. Witten: Linear-time, Incremental Hierarchy Inference for Compression. <i>Data Compression Conference 1997</i> : 3-11 [ <a href="#">conf/dcc/Nevill-ManningW97</a> ]
55		Craig G. Nevill-Manning, Ian H. Witten, Dan R. Olsen: Compressing Semi-Structured Text using Hierarchical Phrase Identification. <i>Data Compression Conference 1996</i> : 63-72 [ <a href="#">conf/dcc/Nevill-ManningW096</a> ]
54		Stuart Inglis, Ian H. Witten: Bi-level Document Image Compression using Layout Information.

[back](#) [forward](#) [home](#) [clear](#) Name: [a=A](#) [a=](#)

## journals

 by name by # of papers 5293..937 885..536 528..300 299..201 200..141 140..95 95..62 61..41 41..19 18..6 6..1






















1. [CACM](#) (5293)
2. [TCS](#) (4737)
3. [IEEE Transactions on Computers](#) (4477)
4. [Information Processing Letters](#) (4396)
5. [TSE](#) (2400)
6. [IEEE Computer](#) (2288)
7. [JACM](#) (2265)
8. [The Computer Journal](#) (2098)
9. [SIAM J. Comput.](#) (2010)
10. [Software - Practice and Experience](#) (1988)
11. [Artificial Intelligence](#) (1673)
12. [Information and Control](#) (1626)
13. [JCSS](#) (1595)
14. [IEEE Software](#) (1384)
15. [IEEE Transactions on Pattern Analysis and Machine Intelligence](#) (1337)
16. [JASIS](#) (1271)
17. [IEEE Transactions on Parallel and Distributed Systems](#) (1198)
18. [JSC](#) (1145)
19. [IEEE Journal on Selected Areas in Communications](#) (1133)
20. [Information and Computation](#) (1105)
21. [Acta Informatica](#) (1100)
22. [IBM Systems Journal](#) (1040)
23. [International Journal of Man-Machine Studies](#) (1027)
24. [ACM Computing Surveys](#) (1019)
25. [J. Algorithms](#) (1005)
26. [TKDE](#) (999)
27. [Information Sciences](#) (974)

[back](#)[forward](#)[home](#)[clear](#)

Name:

[a=A](#)[a-](#)

top authors

 [267..156](#) [152..126](#) [126..115](#) [115..105](#) [104..98](#) [98..93](#) [93..89](#) [89..86](#) [86..83](#) [83..80](#) [80..77](#) [77..75](#) [75..72](#) [72..71](#) [71..69](#) [69..68](#) [68..67](#) [67..65](#) [65..63](#) [63..62](#) [62..61](#) [61..59](#) [59..59](#) [59..58](#) [59..57](#)

1. [Grzegorz Rozenberg](#) (267)
2. [Philip S. Yu](#) (244)
3. [Hector Garcia-Molina](#) (238)
4. [Jeffrey D. Ullman](#) (227)
5. [Christos H. Papadimitriou](#) (221)
6. [Moshe Y. Vardi](#) (220)
7. [Kang G. Shin](#) (215)
8. [Micha Sharir](#) (206)
9. [Elisa Bertino](#) (202)
10. [Oded Goldreich](#) (192)
11. [Robert Endre Tarjan](#) (191)
12. [Kurt Mehlhorn](#) (189)
13. [Michael Stonebraker](#) (187)
14. [Sushil Jajodia](#) (175)
15. [Moti Yung](#) (175)
16. [Joseph Y. Halpern](#) (173)
17. [Jeffrey Scott Vitter](#) (172)
18. [Abraham Silberschatz](#) (171)
19. [Oscar H. Ibarra](#) (170)
20. [Avi Wigderson](#) (167)
21. [Nancy A. Lynch](#) (166)
22. [Donald F. Towsley](#) (165)
23. [Hermann A. Maurer](#) (165)
24. [Marek Karpinski](#) (164)
25. [Leonidas J. Guibas](#) (163)
26. [Amir Pnueli](#) (163)
27. [Ugo Montanari](#) (161)

```
<article key="journals/ai/SchreyeBV89">  
<author>Danny De Schreye</author>  
<author>Maurice Bruynooghe</author>  
<author>Kristof Verschaetse</author>  
<title>On the Existence of Nonterminating Queries  
for a Restricted Class of PROLOG-Clauses.</title>  
<pages>237-248</pages>  
<year>1989</year>  
<volume>41</volume>  
<journal>Artificial Intelligence</journal>  
<number>2</number>  
<url>db/journals/ai/ai41.html#SchreyeBV89</url>  
</article>  
...
```



## Representation of <title>-Fields:

- construct canonical Huffman codes on word level [MG-Book]
- degree of tree-nodes: 213 (0x2a-0xff)
- lexicon: 3-in-4 front coding
- Publications are sorted by journal/booktitle, year

# Lexicon of title words

0='\*', 1='+', 2=',', 3='-',  
4='.', 5='/', 6='0', 7='1', ...

...  
7Comprehending32sibility5-ons0sion  
2Compress21ibility3+n3onless  
6Compressions3,ve2.lets2or  
5Compressors/-ise2+d2s  
4Comprising/5omisability3,ed4s  
6Compromising0+r..sing.ter  
2Compters0,ur/-ing/on  
2Comptuer.4uatational1-ion4al  
5Compuations20tional//iting/log  
4Compulsory/,nd0.ting/s  
0Comput03abilities2-les3y  
7Computacional1-ion4,al1ting  
7Computationali6.lism5,el4ons  
...

## Lexicon: title words

```
<article key="journals/ai/SchreyeBV89">  
<author>Danny De Schreye</author>  
<author>Maurice Bruynooghe</author>  
<author>Kristof Verschaetse</author>  
<title>C1C2C3 ...</title>  
<pages>237-248</pages>  
...  
</article>  
...
```

# Representation of <author>- and <editor>-Fields:

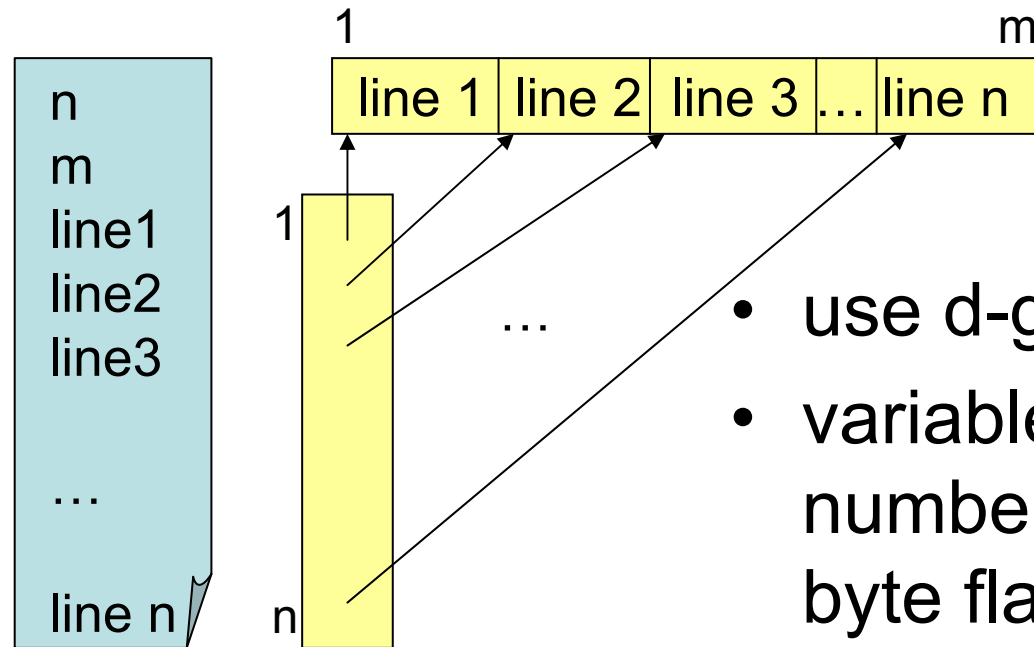
- construct list of all persons
- sort them by number of publications
- add inverted lists (publication numbers) to all persons

```
Peter Smith 17 36 1003 7790 8800
```

```
Johanna Mayer 36 2077 9002
```

```
...
```

# Representation of inverted indexes



- use d-gaps [MG-book]
- variable-byte coding (base 111 numbers,  $\approx 7$ bit value + last byte flag) [Scholer et al., SIGIR 2002]

...

Alberto H. F. Laender)s»™"bê~□5°žÛ\üç"@Í&ñα\$¥"ð8i?Øμ\-"ÜÖà·\$œα  
"iï\$-¯"c—UÚLÌ"?™"<°#Há"nÆWÕh·â'#:Ò"=£'

Melvin A. Breuer)#^È"j¥%ù"û"β"—"-÷"W¯"0à"†ò\$È"Â¶£"ž"n·\Û\*Û&ª'&·í\$Ÿ  
"Ò\$,,Äö&»#-™"²#½"¼¹ø"ÁyÅ&pFÌ

Jean-Louis Lassez)0Ñ%₀×qØ¹\$4Ò=ÆgÅ";'áM¹5Ä±#R÷Ò"Í#RÅ#&è#ø"ký¯  
-Â\*~"ô"·q™)j)Î.-]ò8βfŸ"Ì(\$\$;- \$õGí"æ#™"aÒ'

Weiyi Meng)P¾¼1à:¼¶#^ž½Ê"5ÉâEî'8,®#Ú&ø%Xž#×j"ÇÔfô~ú'\$O"}ì²QÜ  
<Ã<©#uÚ"<û#¯-Ÿ\$`'GÉ2¯JÑ°Vå

Walter L. Ruzzo)"SšŽ""åO«V©r,ž\$4-'>"m—\$Æ%jà"<Æ&Ö~òHö#α|\_#ë"Ÿ®  
%'©Õ' μfÝ%°À?ø"É'\$§ÅñçQÚš"~□ú

James E. Rumbaugh)o' # % ñ # Í \$ D ® & D ã Q õ % \ i ! Ú ! ! © ! > ž š □ š š > œ ç ç š ~ œ š >™™  
□>#\$žö\$" ÈÅ&G«"1ø

Abdul Sattar)8é°,Ó5ò'£'x^"®À »—"Sα'>ÒÜa>&&,0Ö=Ú0ÊI—"mØ#F¹pë#yÃ  
=š? Jª÷•™\$,—%>ì©Ç²

...

Lexicon: title words

Person table

```
<article key="journals/ai/SchreyeBV89">  
<author>a1</author> <author>a2</author>  
<author>a3</author> <title>c1c2c3 ...</title>  
<pages>237-248</pages>  
<journal>Artificial Intelligence</journal>  
...  
</article>  
...
```

# Representation of <journal>- and <booktitle>-Fields:

- construct lists of all journals and booktitles
- add position of first publication with this journal/booktitle and count

```
Artificial Intelligence 5000 603  
TODS 7890 400
```

...



331

10509

CACM)#^ãPÜ

TCS)7IçKÛ

IEEE Transactions on Computers),[ž|µ

Information Processing Letters)/KàHÓ

TSE)89Ô6Û

IEEE Computer)+□ã5Ô

JACM)0|~5½

The Computer Journal)8UÒ3ô

SIAM J. Comput.)5@-3œ

Software - Practice and Experience)6zÛ2õ

Artificial Intelligence)"&lt;Ê0~

Information and Control)/Ž÷/Ø

JCSS)12Î/¹

IEEE Transactions on Pattern Analysis and Machine Intelligence),□Ò.α

...

Lexicon: title words

Person table

Journal table

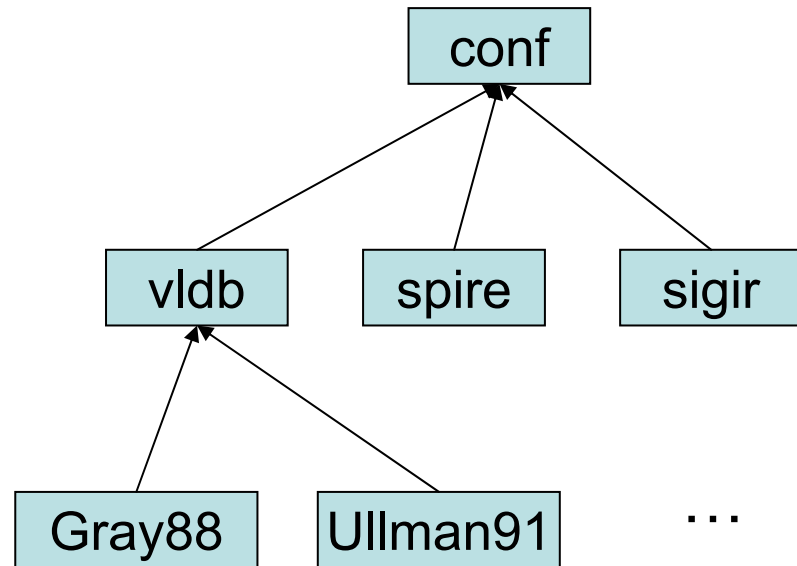
Booktitle table

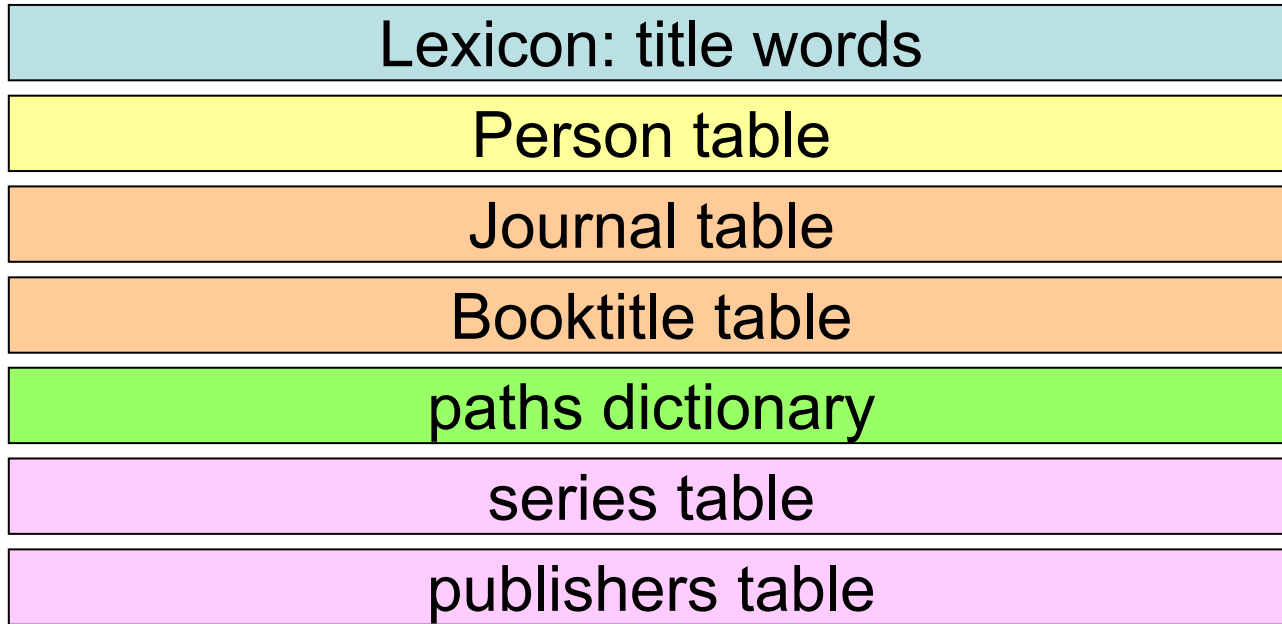
```
<article key="journals/ai/SchreyeBV89">  
<author>a1</author> <author>a2</author>  
<author>a3</author> <title>c1c2c3 ...</title>  
<pages>237-248</pages> <journal>j</journal>  
<number>2</number>  
<url>db/journals/ai/ai41.html#SchreyeBV89</url>  
</article>
```

...

# Representations of Paths: key, <ee>, <url>

construct paths  
dictionary:  
bottom-up tree of  
path elements





...  
Av'SchreyeBV89)T<sub>1/2</sub>fuf: 'föf\*ÚñfL•fófifNÿf3;föfHpl3Km)g#ÿœuv Y □ pš' b "#û""'+  
...

$A_{p_1}$ SchreyeBV89)T<sub>1/2...</sub>)g<sub>n<sub>1</sub>n<sub>2</sub></sub>u<sub>p<sub>2</sub></sub>Yypjvnb<sub>3</sub>a<sub>1</sub>a<sub>2</sub>a<sub>3</sub>

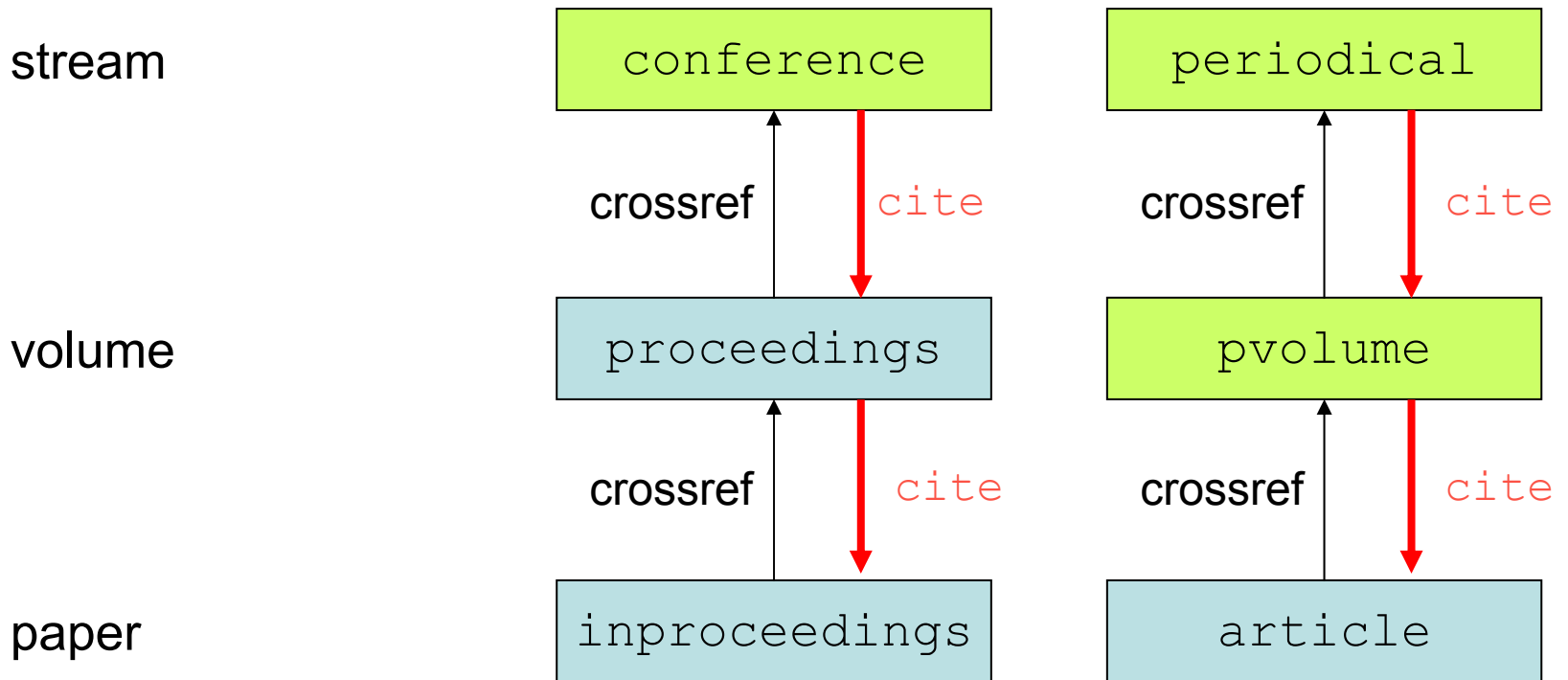
# DBLP Browser: Compression

XML-file (Sept 3, 2002)	123328 KB
zip(XML-file)	22100 KB
Compressed representation	26664 KB
zip(compressed representation)	14732 KB

# DBLP (Browser): next step

- search
- visualization
- BHT files, but standard cases should not require BHT files: extension of the BibTeX data model ...

# Extension of the BibTeX data model



```
<conference key=„conf/vldb/@“>
```

```
<title>Conference on Very Large Databases</title>
```

```
<booktitle>VLDB</booktitle>
```

```
<url>db/conf/vldb/index.html</url>
```

```
<toc>
```

```
...
```

```
<cite>conf/vldb/2002</cite>
```

```
...
```

```
<cite>conf/vldb/2001</cite>
```

```
...
```

```
</toc>
```

```
</conference>
```



# User Interfaces for Bibliographic Collections

- conventional wisdom: provide a good IR system
- DBLP: browsing oriented interface + very simple search tools

Why is DBLP nevertheless used?

„...there was little comprehensive searching by faculty researchers because they worked within specialized forums where they could more efficiently find the materials they needed. This precluded comprehensive searching to identify an exhaustive set of materials in a wider variety of forums. [...] They also used browsing to examine electronic announcements of conferences, and tables of contents of journals ...“

[Lisa Covi, PhD thesis 1996]

# Problems of the many search-oriented Web-Portals

- How to query ... ?
- Description / characterization of the collection unknown

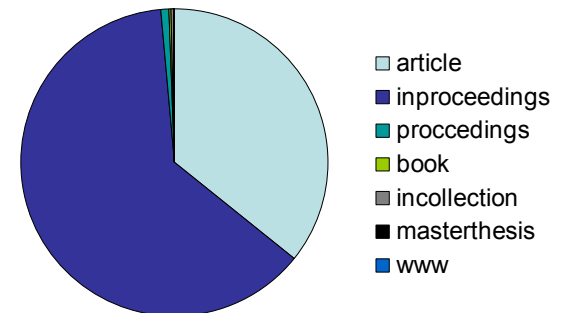
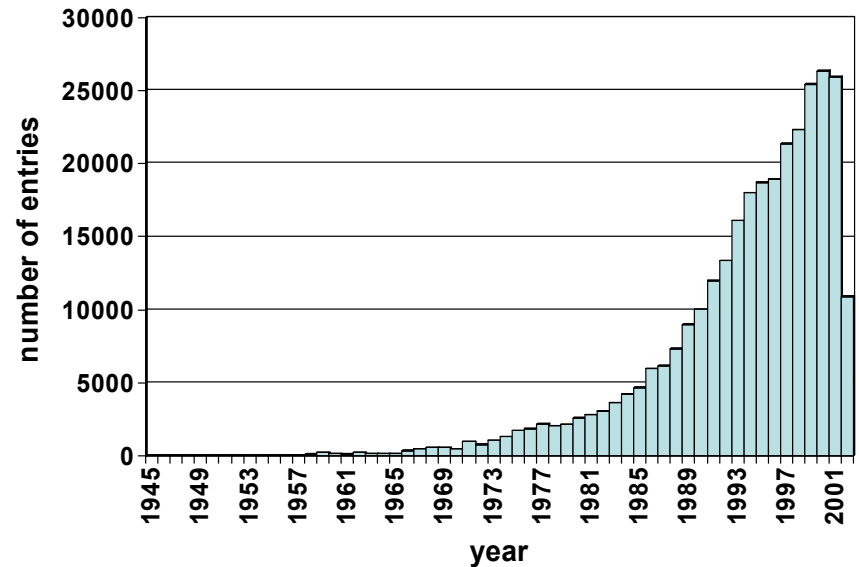
# views – visualizations – diagrams

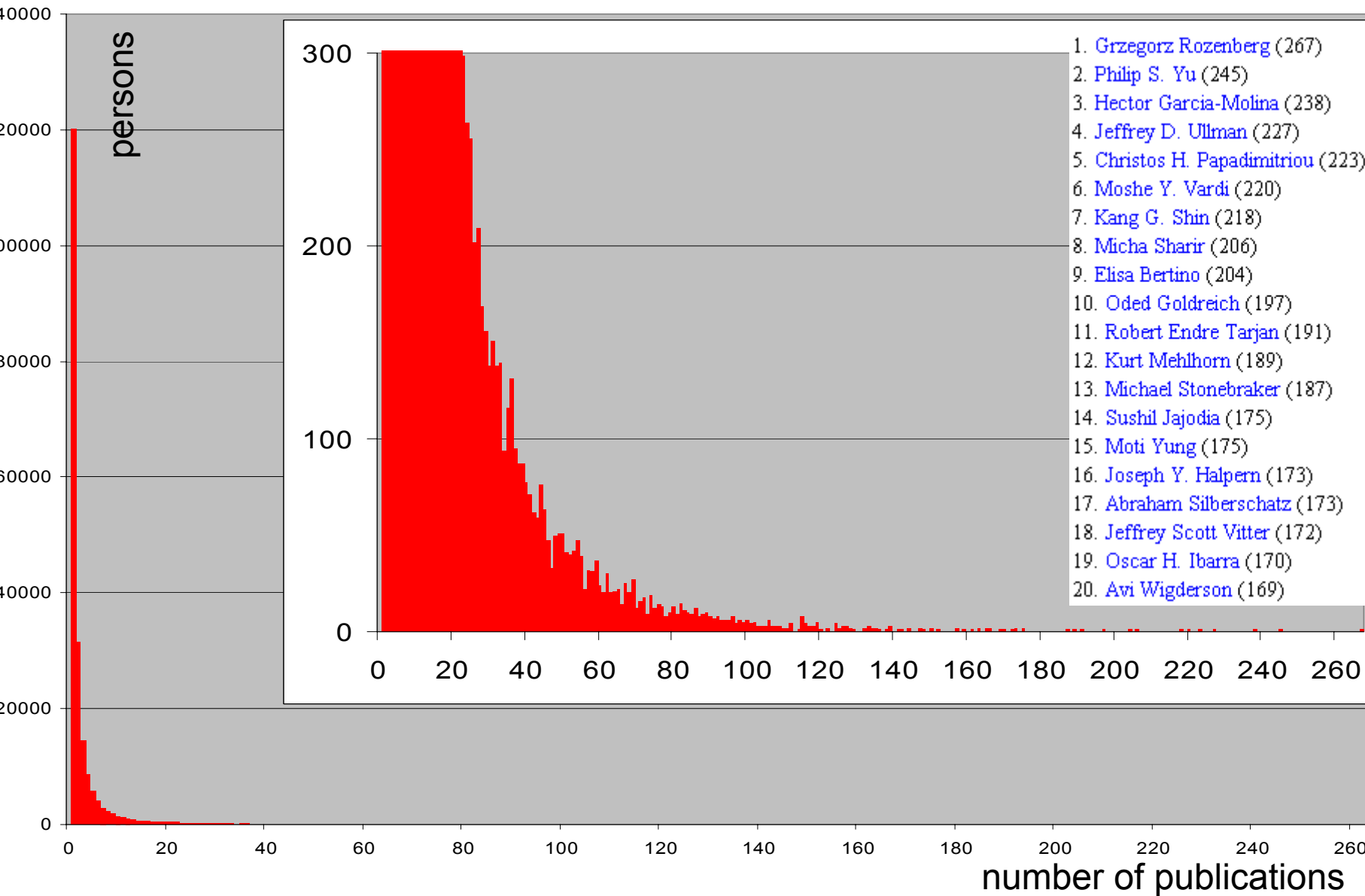
- global characterizations of the collection (coarse navigation)
- stream level diagrams
- person page visualizations

# global characterizations ...

- 1. CACM (5361)
- 2. TCS (4737)
- 3. IEEE Transactions on Computers (4477)
- 4. Information Processing Letters (4396)
- 5. TSE (2407)
- 6. IEEE Computer (2288)
- 7. JACM (2265)
- 8. The Computer Journal (2098)
- 9. SIAM J. Comput. (2010)
- 10. Software - Practice and Experience (1988)
- 11. Artificial Intelligence (1673)
- 12. Information and Control (1626)
- 13. JCSS (1595)
- 14. IEEE Transactions on Pattern Analysis and Machine Intelligence (1584)
- 15. IEEE Software (1384)
- 16. JASIS (1271)
- 17. IEEE Transactions on Parallel and Distributed Systems (1265)
- 18. JSC (1145)
- 19. IEEE Journal on Selected Areas in Communications (1140)
- 20. Information and Computation (1105)
- 21. Acta Informatica (1100)
- 22. IBM Systems Journal (1040)

- 1. IJCAI (2863)
- 2. DAC (2478)
- 3. INFOCOM (2345)
- 4. ISCAS (2042)
- 5. FOCS (2037)
- 6. STOC (1811)
- 7. AAI (1747)
- 8. VLDB (1728)
- 9. ICALP (1664)
- 10. Winter Simulation Conference (1483)
- 11. SIGMOD Conference (1460)
- 12. ICSE (1454)
- 13. MFCS (1434)
- 14. HICSS (1424)
- 15. ICDE (1390)
- 16. IFIP Congress (1388)
- 17. GI Jahrestagung (1318)
- 18. ICDCS (1187)
- 19. Int. CMG Conference (1173)
- 20. DAGM-Symposium (1131)
- 21. DEXA (1122)
- 22. SIGIR (1111)





# Classification

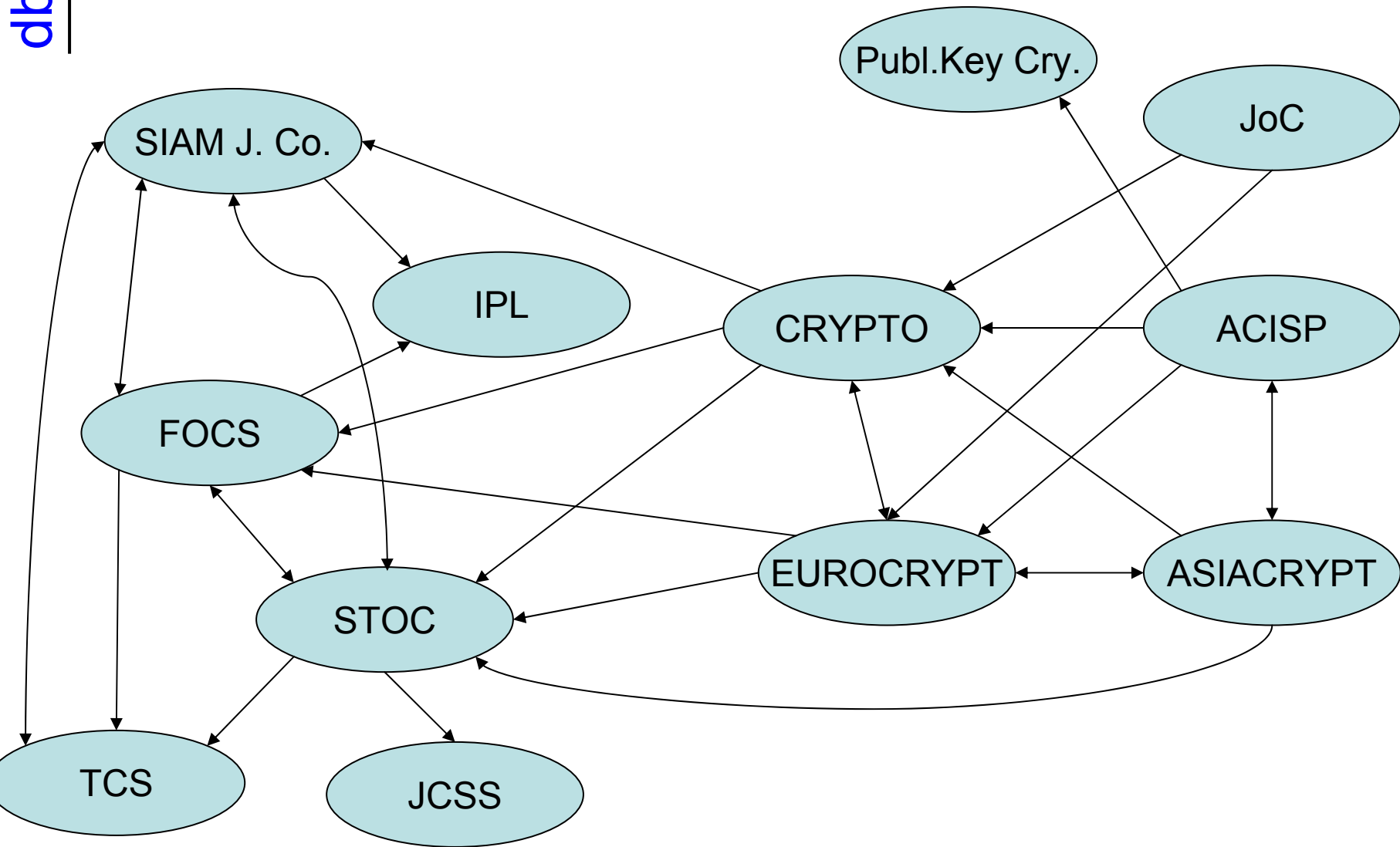
- ACM CR System
  - classifications only available for small subset of computer science literature
  - often criticized as too inflexible / too coarse
- intellectual classification (on the paper level) is too expensive

# Clustering on stream level

related conferences or journals:

- in which other streams did the authors of this stream publish most frequently ?
- todo:
  - what is „most frequently“?
  - what is the neighborhood of a stream ?
  - ...





Thank you for your attention



Porta Nigra“ – the most famous landmark of Trier

[www.trier.de](http://www.trier.de)