# Al Methods for Building University Research Knowledge Base

# Scoring scientist using knowledge databases

Henryk Rybiński, Jakub Koperwas, <u>Łukasz Skonieczny</u>, Marek Kozłowski, Piotr Andruszkiewicz, Wacław Struk

Institute of Computer Science, Warsaw University Of Technology September 9, 2019

## Agenda



- From Institutional Repository...
- ...to University Knowledge Base.
- Al tools
  - semantic processing
    - information extraction
    - publication classifier
    - keywords extraction
    - sense indexing
  - Profiling researchers
    - person name disambiguation
    - finding and ranking an expert

### Institutional Repository (2010)

- Storing scientific information (meta data with full documents) – papers, theses, reports, patents, projects
- Exposing and sharing (with clean GUI, REST API and OAI)
- Advanced querying

 Key technology: Apache Jackrabbit
 (free implementation of Content Repository API for Java -JCR)

### Knowledge Base in 2013

#### $\Omega - \Psi^{R} = CRIS + IR + RPS$

### IR

- Run by libraries
- Long-term preservation
- **Providing access t**o full text
- Dublin Core
- Limited to OA
- Self-archiving

CRIS

- Run by adminstration
- Covers all the research activities
- Metadata (CERIF)
- Rather bibliogr. description
- Reporting
- Research assessment
- (e.g. IF vs costs)

### RPS

- Run by scientists
- Covers all the research activities
- Building profiles
- Vizualizes the activities and cooperation
- Social networks functions
- No evaluation

## University Knowledge Base

- Relations and connections between objects: authors, papers, books, projects, journals, conferences
- Auto-generated information pages for researchers and units (journals and conferences in the near future)
  - -with statistics regarding research activity
  - -with research area word-cloud
  - -with cooperation graph
- Social network elements
  - -keep your profile up-to-date
  - -share your scientific achievements
  - -find and contact researchers



|--|



Henryk Rybiński, PhD, DSc, Professor Professor

The Institute of Computer Science Faculty of Electronics and Information Technology

Email: H.Rybinski@ii.pw.edu.pl Phone: +48 22 234-7432, fax +48 22 234 6091 Room no: 204 Consultations: Monday 14.00-16.00

Researcher Report	
Publications	92
PhD theses	13
Participation in projects	30
Supervised BSc and MSc theses	2
Professional activity	9

#### h-index\*:12

MSc (1970), PhD (1974), DSc (1988), Tenured Professor (2001); Specialization: information systems, knowledge representation, data and text mining, databases, Professor, Director of the Institute (2008-), Head of Division of Information Systems (1994-2008), Co-ordinator of the Curriculum on Software Engineering and Information Systems (1994-2008), Co-



ordinator of the Subject Class "Databases and Information Systems" (1995-2001), voting member of ACM and SIGMOD (1989-), Affiliate Member of EEE (1990-1996); Member of several programme committees of international conferences and workshops, among others: IIS, ISMIS, IIPWM, AM, SWC, RSFDGRC, RSKT, TKE, PKDD, PAKDD, MCD; member of CREST Working Group; expert and consultant of many UN agencies and European Commission; Member of Informatics Committee at Polish Academy of Science (2011-); Editorial Board Member of the Journal of ntelligent Information Systems (2012-); Editorial Board Member of the International Journal of Social Network Mining (2012-); Chair of the Rector's Committee for the strategy of developing ICT infrastructure for WUT in 2013-2020.

\* presented value of the Hirsch index is approximative calculation obtained in the Repository based on the scientist's publications (including autocitations) in the Repository and Internet information analysis. The value is close to the value obtained with the Publish or Perish system. In general it is higher than the value given by the Scopus or Web of Science sites. In the case of undervalued number, first of all take care of completeness of the Repository.





Profile Publications	PhD Projects	BSc and MSc	Activities	Citations	Statistics	Cooperation			Edit
Publications	Henryk Rybińs Professor The Institute of Co Faculty of Electron Technology Email: H.Rybinski Phone: +48 22 23 Room no: 204 Consultations: Mo	ki, PhD, DSc omputer Science nics and Informa @ii.pw.edu.pl 4-7432, fax +48 onday 14.00-16.0	, <b>Professo</b> tion 22 234 6091	pr fi d	databa sm iormation retr synonyms proper nouns digital libu irequent items ecision rules idustry default logic Kr rougo	ise mana informa informa ary infor <b>dat</b> iowietic sets ass distributed con	gement tion ret mation a mi je disco ociation r	text analysis <b>ning</b> i learning algorit <b>systems</b> <b>ning</b> <b>byery</b> ules homonym ic resources	} hm S
Number of records: 92.									
🝸 🔯 🛅 Order	by: type/author	T					Export 0 as:	Ankieta 2013	T
Authored books									
Rybiński Henryk: Mo Google	odele baz danych, vo	l. 63, 1987, Centru	m Informacji N	aukowej, Te	chnicznej i Ek	onomicznej, 156 j	p.	(	0
Rybiński Henryk: Pr Google	oblem optymalizacji re	organizowania zł	oioru informacy	yjnego w sy:	stemie wyszu	ikiwania informac	iji, 1976, IINTE, 92	2 p. (	0
Edited books									
Bembenik Robert, S Computational Intelli Google	ikonieczny Łukasz, R igence, vol. 390, 2012	ybiński Henryk, Nie , Springer, ISBN 9	zgódka Marek 78-3-642-248(	( ( <i>eds</i> . ): Inte 08-5, 277 p.,	elligent Tools ( DOI:10.1007/	for Building a Scie 978-3-642-24809	entific Information 1-2	Platform, Studies in	n () (?)
Bembenik Robert, S Platform: Advanced 642-35647-6 Google	ikonieczny Łukasz, R I Architectures and S	ybiński Henryk, Kr olutions, Studies ir	yszkiewicz Ma i Computationa	arzena, Niez al Intelligence	gódka Marek e, vol. 467, 20	( <i>eds</i> . ): Intelligent 13, ISBN 978-3-6	Tools for Buildin 42-35646-9, 548	g a Scientific Inform p., DOI:10.1007/97	nation 8-3-
Bembenik Robert, S Platform: From Rese [978-3-319-04714-0 Google	ikonieczny Łukasz, R earch to Implementatic 0], 290 p., DOI:10.100	ybiński Henryk, Kr n, Studies in Com 7/978-3-319-0471	yszkiewicz Ma putational Intell 4-0	arzena, Niez ligence, vol.	gódka Marek 541, 2014, Sp	( <i>eds</i> . ): Intelligent ringer Internation	Tools for Buildin al Publishing, ISB	g a Scientific Inform № 978-3-319-04713	nation 3-3,
Kryszkiewicz Marz	ena, Rybiński Henryk	Skowron Andrze	i. Raś Zbioniev	w W. (eds.)	): Foundations	s of Intelligent Svs	stems. Lecture N	otes in Artificial	

Intelligence, vol. 6804, 2011, Springer, ISBN 978-3-642-21915-3, 746 p., DOI:10.1007/978-3-642-21916-0

Kryszkiewicz Marzena, Peters James F, Rybiński Henryk, Skowron Andrzej (eds.): Rough Sets and Intelligent Systems Paradigms, Lecture Notes in Artificial

0 🕕

Publications PhD Projects BSc and MSc Activities Citations Statistics Cooperation
---



Henryk Rybiński, PhD, DSc, Professor Professor The Institute of Computer Science Faculty of Electronics and Information Technology Email: H.Rybinski@ii.pw.edu.pl

Phone: +48 22 234-7432, fax +48 22 234 6091 Room no: 204 Consultations: Monday 14.00-16.00

hindex = 12, cited by total = 379



	ontologies	(artificial	intelligence)
itod	ontologies	(artificiar	intenigence

Edit

pub	title	cited	ontologies (artificial inte
	Finding Reducts in Composed Information Systems	70	19/04/2014
	Reducing information systems with uncertain attributes	31	19/04/2014
0	Computation of Reducts of Composed Information Systems	27	19/04/2014
0	On first order logic databases	27	03/02/2014
	Dataless transitions between concise representations of frequent patterns	24	19/04/2014
	Data Mining in Incomplete Information Systems from Rough Set Perspective	21	19/04/2014
	Discovering Synonyms Based on Frequent Termsets	17	19/04/2014
	Text onto miner - a semi automated ontology building system	14	19/04/2014
	Intelligent Tools for Building a Scientific Information Platform	14	22/04/2014
	Discovering Compound and Proper Nouns	12	19/04/2014
	Mining spatial association rules with no distance parameter	12	19/04/2014
	Knowledge sharing in default reasoning based multi-agent systems	12	19/04/2014
	Towards a unifying logic formalism for semantic data models	9	19/04/2014
	Discovering Word Meanings Based on Frequent Termsets	9	19/04/2014
0	Distributed default logic for multi-agent system	9	19/04/2014
0	Automatic Index Selection in RDBMS by Exploring Query Execution Plan Space	7	24/04/2014
0	Data Mining for Technical Operation of Telecommunications Companies: a Case Study	7	12/05/2014
	Incomplete database issues for representative association rules	6	19/04/2014
	Multilevel information system- Towards more flexible information retrieval systems	5	17/02/2014
	Word sense discovery for web information retrieval	5	23/03/2014
0	Learning Mechanism for distributed default logic based MAS - implementation considerations	4	17/02/2014
	Extending open source software solutions for CRM text mining	3	11/05/2014
0	A new approach to computing weighted attributes values in incomplete information systems	3	17/02/2014
	A Distributed Decision Rules Calculation Using Apriori Algorithm	3	18/02/2014
	Legitimate Approach to Association Rules under Incompleteness	3	03/02/2014
	Methods and Tools for Ontology Building, Learning and Integration – Application in the SYNAT Project	3	17/02/2014
	Regression - vet another clustering method	3	09/05/2014

Profile	Publications	PhD	Projects	BSc and MSc	Activities	Citations	Statistics	Cooperation			Edit
		Henry Profess The Ins Faculty Techno Email: Phone: Room Consul	<b>/k Rybiń:</b> sor stitute of Co of Electroi blogy H.Rybinski : +48 22 23 no: 204 ltations: Mo	ski, PhD, DSc omputer Science nics and Informa @ii.pw.edu.pl 4-7432, fax +48 onday 14.00-16.0	, <b>Profess</b> <i>tion</i> 22 234 609	or )1	inte tes tes tes tes tes tes tes tes tes t	Alligent stor <b>K M</b> <b>ta</b> <b>informa</b> <b>owledg</b> (artif- rithm <b>gh sets</b> aging (artif- (artif- (aging (artif- (aging (ag	a. idustry inin ation s e disc rmatio ssociati al library fre	reposi ng inter ng inter ng inter in	tory wiedge items mation retri- quent text pat. equent text pat. efault logic ecision rules synonyms sme homonyms



Cumulative O Annual

\* presented value of the Hirsch index is approximative calculation obtained in the Repository based on the scientist's publications (including autocitations) in the Repository and Internet information analysis. The value is close to the value obtained with the Publish or Perish system. In general it is higher than the value given by the Scopus or Web of Science sites. In the case of undervalued number, first of all take care of completeness of the Repository.

Back

#### Problem

Manual maintenance of such a semantic network is not easy, if possible at all.

### So AI methods are of crucial importance

**AI Tools** 

#### Semantic processing

Goals:

- enrich data by adding semantically meaningful descriptors
- 2. With enriched texts discover research areas (for authors, teams or institutions) and build vector representation of research
- 3. Given the vector representation visualize the research



AI Tools

#### Semantic processing

#### Methods

- 1. Tagging publications with OSJ categories
- 2. Keywords extraction
- 3. Sense indexing



# OSJ tagging

- OSJ ontology of scientific journals
- Created by Science-Metrix, Canada, 2011
- 15000 journals, hierarchical (3 levels) classification
- Ready to use for known journals but...
- a classifier is needed for other journals, conference papers, and book chapters

#### **Our solution: tree of bayesian classifiers, 85% accuracy**

### OSJ – level 2



### **Keyword Extraction from Document**

- TKE method was proposed based on RAKE (Rose et al. 2010) and KEA (Witten et al. 1999);
- Additionally Wikipedia is used for adding more general or synonymic keywords (senses, translations, summaries)
- 3. Extracted terms are enriched with context that eliminates disambiguity (SnS)

### **Keywords** extraction

The algorithm:

- identifies wikipedia entry corresponding to given candidate term
- 2. processes the entry in order to extract context: labels, senses, translation, summaries
- evaluates source document with respect to the context found in previous step
- 4. taging keywords with meaning dicovered by SenseSearcher algorithm

### **Keywords** extraction

IEEE.org | IEEE Xplore Digital Library | IEEE Standards | IEEE Spectrum | More Sites

IEEE Xplore®

Access provided by: **POLITECHNIKI WARSZAWSKIEJ** » Sign Out



Word meaning disambiguation has always been an important problem in many computer science tasks, such as information retrieval and extraction. One of the problems, faced in automatic word sense discovery, is the number of different senses a word can have. Often, senses are dominated by some other, more frequent ones. Discovering such dominated meanings can significantly improve quality of many text-related algorithms. In particular, Web search quality can be leveraged. In the paper, we present a novel approach for discovering word senses. The method is based on concise representations of frequent patterns. The method attempts to discover not only word senses that are dominating, but also senses that are dominated and under represented in the repository.

#### Published in:

Data Mining Workshops, 2008. ICDMW '08. IEEE International Conference on

Extracted terms: information, information retrieval, web search engine, computer science, word sense, World Wide Web

# **Profiling of researcher**

 $keywords(d) = OSJ(d) \cup EXTRACT(d)$ 

$$score(k) = \sum_{d \in D(p)} rel(k, d) \times (sif(d) + 1)$$

etisteiline information ceelin teriseinen teristeiline information ceelin teriseinen **COMPUTER Systems** retter interiseine rough sets in element information systems determining antik... association atos determining antik... association atos determining antik... association atos information systems determining antike interimentation synthesis <u>knowledge discovery</u> information interimentation antiketer atos atomatication interimentation atomatication interi

istributed systems

- rel(k, d) measures the relevance of keyword k with respect to document d; Usually it is the value of tf-idf; however, in the case of publications and technical reports the values of the OSJ keywords and the keywords provided by the authors are boosted;
- sif(d) is a scientific impact factor of the document d. For the journal papers this is a linear combination of the impact factor of the journal and the citations of the document. Arbitrary values are given to other objects, like conference papers, patents, supervised theses.

### **Profiling of an institution**

Given vectors of researchers affiliated at an institution we build a vector for the institution

Clearly, the same algorithm can be used for building an expertise vector for any subset of documents. So, given a query q we can obtain the set of documents D(q) and calculate the score vector the same way as D(p) for researcher p. This means that we can easily obtain an aggregated cloud of research interest for a faculty, department, as well as for a whole university. The profiles for a whole university and two different faculties can be seen in Fig. 4.

### Visualization of vectors Word clouds visualizing Unversity and 2 faculties: Chemistry, Physics



Fig. 4 Aggregated expertise clouds for a university (a), and two faculties (b) and (c)

#### Search for best experts in a domain

#### Step 1

all the knowledge base resources are searched with a specified search phrase q (formulated the same way as for searching publications, theses, etc.). The result of the query q is denoted by D(q);

#### Step 2

- a set of all persons P(q) related to items from D(q) is calculated as follows:

$$P(q) = \bigcup_{d \in D(q)} \{p : role(p, d) > 0\}$$

$$(3)$$

where the function role(p, d) provides a measure of relevance of role of person p in elaborating the document d; by role(p, d) > 0 we mean that p has some role in d, that is, p is an author of publication d, or is a supervisor of thesis d, or is a leader of project d, etc.; in particular, the function takes into account the roles which can be article author, book author, book editor, phd author, phd supervisor, master thesis author, master thesis supervisor, project member, project leader; in this way we can, e.g., value the role book author more than book editor;

#### Search for best experts in a domain

#### Step 3

- for each person  $p \in P(q)$  the person score measure, denoted by Pscore(p, q), is calculated:

$$Pscore(p,q) = \sum_{d \in D(q): role(p,d) > 0} score(p,d,q)$$
(4)

where score(p, d, q) is a function expressing the importance of d with respect to query q, and in relation to p; the function is calculated according to a selected ranking algorithm;

the set of persons P(q) is sorted in descending order by Pscore(p, q), and a list of top n persons is presented to the user.

#### Search for best experts in a domain

$$score(p, d, q) = rel(d, q) \times (sif(d) + 1) \times role(p, d)$$
(5)

#### where:

rel(d, q) is a measure of relevance of d with respect to query q; here we rely on the Lucene relevance score, which uses the cosine measure, with boosted values for the fields resulting from the semantic enrichment procedures (Section 5); sif(d) is a scientific impact factor of d; it is a linear combination of the impact factor of the journal (for the journal papers) and its citations.

### Literature

KOPERWAS, J.J., SKONIECZNY, Ł., KOZŁOWSKI, M., ANDRUSZKIEWICZ, P., RYBIŃSKI, H., AND STRUK, W. 2014. AI PLATFORM FOR BUILDING UNIVERSITY RESEARCH KNOWLEDGE BASE. LNAI, SPRINGER, 405–414.

KOPERWAS, J.J., SKONIECZNY, Ł., KOZŁOWSKI, M., ANDRUSZKIEWICZ, P., RYBIŃSKI, H., AND STRUK, W. 2017. Intelligent information processing for building university knowledge base. Journal of Intelligent Information Systems 48, 1, 141–163.

KOZŁOWSKI, M. AND RYBIŃSKI, H. 2017. Word Sense Induction with Closed Frequent Termsets. Computational Intelligence 33, 335–367.

RYBIŃSKI, H., SKONIECZNY, Ł., KOPERWAS, J.J., STRUK, W., STĘPNIAK, J., AND KUBRAK, W. 2017. Integrating IR with CRIS – a novel researcher-centric approach. *Program-*Electronic Library and Information Systems 51, 3, 298–321.

KOZŁOWSKI, M. AND RYBIŃSKI, H. 2018. Clustering of semantically enriched short texts. Journal of Intelligent Information Systems, 1–24.

### Thank you! http://omega-psir.ii.pw.edu.pl http://repo.pw.edu.pl http://wizzar.ii.pw.edu.pl/RepoPW

