

Automated Construction Of Domain Ontologies From Lecture Notes

Neelamadhav Gantayat

under the guidance of

Prof. Sridhar Iyer

Department of Computer Science and Engineering,
Indian Institute of Technology, Bombay
Powai, Mumbai - 400 076

June 28, 2011



Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation
- 7 Evaluation
- 8 Conclusion & Future Work
- 9 References

Motivation

- Courseware Repositories
 - MIT's OCW¹
 - NPTEL²
 - CDEEP³
- Searching in Repositories

¹<http://ocw.mit.edu/>

²<http://www.nptel.iitm.ac.in/>

³<http://www.cdeep.iitb.ac.in/>

Repositories

MITOPENCOURSEWARE MASSACHUSETTS INSTITUTE OF TECHNOLOGY

Home Courses Donate About OCW Help Contact Us

Home > Search [Email this page](#)

Search Results

MIT OpenCourseWare Highlights for High School

Threads [Advanced Search](#)

Results 1 - 10 of about 820 for **Threads**. Sort by **Date / Relevance**

Showing: All formats / PDF only / HTML only

[jwv1.6.087 Practical Programming in C, Lecture 12](#)
The main **thread** spawns multiple **threads**. The **thread** may communicate with one another. 8 Page 11. Not all **multi-threaded** code is safe
ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-087-practical-programming-in-c-january-sep-2010/lecture-notes/MIT6_087WP15 lec12.pdf - 2010-06-28

[jwv1.6.087 Practical Programming in C, Lecture 13](#)
of parallel processing with shared memory. Program organized to execute multiple

Figure: MIT's OCW search for "Threads"

MITOPENCOURSEWARE MASSACHUSETTS INSTITUTE OF TECHNOLOGY

Home Courses Donate About OCW Help Contact Us

Home > Search [Email this page](#)

Search Results

MIT OpenCourseWare Highlights for High School

Operating System Threads [Advanced Search](#)

Results 1 - 10 of about 215 for **Operating System Threads**. Sort by **Date / Relevance**

Showing: All formats / PDF only / HTML only

[jwv1.MiscKeywords](#)
influenced a number of commercial **operating systems**) has the fir parameter passing where ever possible: **systems** calls to get good overall **system** performance?
ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-828-operating-system-engineering-fall-2006/lecture-notes/15_miscKeyw.pdf - 2010-05-29

[jwv1.Threads_processes_and_context_switching](#)
The **thread** manager needs a method for deciding which **thread** to run if multiple **threads**

Figure: MIT's OCW search for "Operating system Threads"

NPTEL A JOINT VENTURE BY THE INDIAN INSTITUTES OF TECHNOLOGY & THE INDIAN INSTITUTE OF SCIENCE

Home | About NPTEL | FAQ | Contact us | Courses | AMBIC7 - Sakshat

NPTEL Courses

Basic Courses (Semesters I & II)

- Civil Engineering
- Computer Science & Engineering
- Electrical Engineering
- Electronics & Communication Engineering
- Metallurgical Engineering
- Chemical Engineering
- Biotechnology
- Mining Engineering
- Semester-wise suggested reading

Find Courses

Subject: Type:

Keyword:

Course Name:

Coordinator Name:

No matching courses found

Separate: (Please enable javascript in your browser and click on [Admin Login](#)) or view this site via: Mailto: npTEL@iitb.ac.in

Figure: NPTEL search for "Threads" under course OS

Searching Tools

A screenshot of a Google search for the term "Threads". The search bar contains the word "Threads" and the search button is visible. The results page shows three main entries:

- poni Poster - Slide 1**: File Format: PDF/Acrobat - Quick View. By P. Ajmani - Related articles. use of large number of parallel GPU threads. Contributions performance loss due to many threads accessing the same node ... www.cse.ibt.ac.in/~ibrahimajmani/pubs/0013008Poster.pdf
- poni Ruby Programming Language Threads and Processes**: File Format: PDF/Acrobat - Quick View. **Threads** Introduction. **Threads** in Ruby. Processes in Ruby. Ruby Programming Language ... returns: *task/active* depending upon **Thread** being in *Critical* region ... www.cse.ibt.ac.in/~cs703/ibrahimajmani/rubyThreads.pdf
- poni A Thread of One's Own**: File Format: PDF/Acrobat - Quick View. By S. Srinivasan - Cited by 8 - Related articles. They encapsulate data, code and a **thread** of control of their own and ... of lightweight **threads** of control, one per actor, we consider a **thread** lightweight ... www.cse.ibt.ac.in/~sathyNAC/05FinalNAC05-Srinivasan.pdf

There is also a link for [Servlet Tutorial](#).

Figure: Google tool search for "Threads"

A screenshot of a Google search for the term "Threads". The search bar contains the word "Threads" and the search button is visible. The results page shows several entries:

- Thread (computer science) - Wikipedia, the free encyclopedia**: In computer science, a **thread** of execution is the smallest unit of processing that can be scheduled by an operating system. It generally results from a fork ... [en.wikipedia.org/wiki/Thread_\(computer_science\)](http://en.wikipedia.org/wiki/Thread_(computer_science)) - Cached - Similar
- Threads - Wikipedia, the free encyclopedia**: **Threads** is a British television drama produced by the BBC in 1984. Written by Barry Hines and directed by Mick Jackson, it is a documentary-style account of ... en.wikipedia.org/wiki/Threads - Cached - Similar
- Threads (India) Limited - Sewing Threads and Braids of Nylon ...**: Manufacturer and exporter of Sewing **Threads** and Braids of Polyester, Cotton, Nylon, specially **Threads** and Polyester Cotton. www.threadsindia.com/ - Cached - Similar
- Videos for Threads**:
 - Threads - Nuclear War, 1984**: 1:09 min - 12 Nov 2006 video.google.com
 - Threads**: 8 min - 28 Oct 2006. Uploaded by Chandrabudhika www.youtube.com

Figure: simple google search for "Threads"

A screenshot of an advanced Google search for the term "operating system threads". The search bar contains the phrase "operating system threads" and the search button is visible. The results page shows several entries:

- Thread (computer science) - Wikipedia, the free encyclopedia**: **Operating system** schedule **threads** is one of two ways ... Kernel **threads** are preemptively multiplexed if the **operating system's** process scheduler is ... [en.wikipedia.org/wiki/Thread_\(computer_science\)](http://en.wikipedia.org/wiki/Thread_(computer_science)) - Cached - Similar
- Operating System:Threads-Gavin**: **Operating System:Threads-Gavin** - Presentation Transcript. THREADS SONALI CHALHAN SYBSC(IT UOI). INTRODUCTION: A thread is contained inside a process and ... www.slideshare.net/~operating-system/threads-presentation - United States - Cached - Similar
- Threads - Operating Systems Notes**: This is why **thread** needs its own stack. An **operating system** that has **thread** facility, the basic unit of CPU utilization is a **thread** ... www.personal.kent.edu/~mhaumera/OS/Systems/.../threads.htm - Cached - Similar
- Operating Systems Lecture Notes Lecture 2 Processes and Threads**: 25 Aug 1998 ... **Operating Systems** Lecture Notes Lecture 2 Processes and **Threads**. www.stallings.com/Extra/OS-Notes/L2.html - Cached - Similar
- poni OPERATING SYSTEMS Threads**

Figure: Advanced google search for "Operating system Threads"

Outline

- 1 Motivation
- 2 Problem Statement**
- 3 Ontology
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation
- 7 Evaluation
- 8 Conclusion & Future Work
- 9 References

Problem Statement

Given a set of lecture notes (pdf files) or PDF of a textbook from a course-ware repository,

- Provide user with the reading material, suggest some Basic and Advanced Topics.

KeyWord	<input type="text" value="Threads"/>	Search
Subject	<input type="text" value="Operating Systems"/>	
File Name	Pre-requisite Files	Follow-Up Files
Threads (Module 3)	* Process	* User Threads * Kernel Threads

Figure: Expected System

Problem Statement

- Dependency graph (Ontology) for a course.

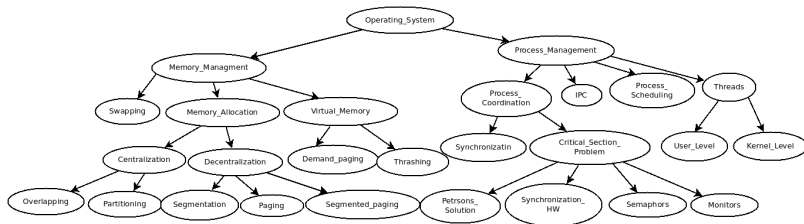


Figure: Dependency Graph (Ontology) for Operating System

Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology**
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation
- 7 Evaluation
- 8 Conclusion & Future Work
- 9 References

Ontology

It is borrowed from philosophy - the study of “The nature of being”.

It “consists of concepts, hierarchical organization of concepts”.

Domain Ontology

- Model which provides definitions and relationships of the concepts, and major theories, principles and activities in the domain.
- Domain ontologies provide shared and common understanding of a specific domain.

Applications of Ontology

- **Knowledge management:** Acquiring, maintaining, and accessing an organization's data.
 - "What is the birthplace of Gandhi?"
- **Web commerce:** On-line market places and auction houses.
- **E-learning:** Dependencies between the keywords of a topic in the repository.

Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology
- 4 Manual Ontology Generation**
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation
- 7 Evaluation
- 8 Conclusion & Future Work
- 9 References

Ontology for Operating System

- **Scope and Domain:** To find out the dependencies between the course ware repositories for operating system.
- **Reuse existing Ontology**
 - www.ksl.stanford.edu/software/ontolingua
 - www.daml.org/ontologies/
 - www.unspc.org, www.dmoz.org
 - www.roselternet.org
- **Important Keywords:**

Table: Keywords

Types of computing	Types of Systems
Memory Management	Process Management
Secondary Management	File Management
Memory Allocation	Virtual Memory
Disk Scheduling	Threads

Ontology Development for Operating System

- **Identify the classes:**



Figure: Classes

- **Define Properties:**

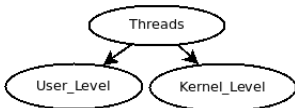


Figure: Types of Thread

Final Ontology using DOT⁴ Language

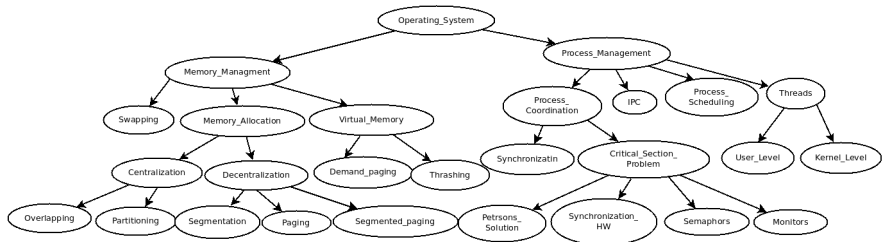


Figure: Ontology for Operating System

⁴http://en.wikipedia.org/wiki/DOT_language

Protégé⁵ Ontology

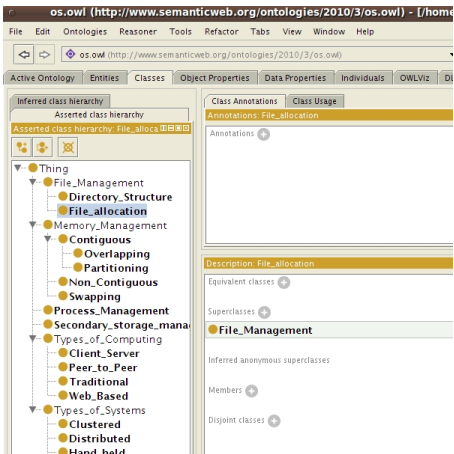


Figure: Protégé Ontology

⁵<http://protege.stanford.edu/>

Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation**
- 6 Solution & Implementation
- 7 Evaluation
- 8 Conclusion & Future Work
- 9 References

Difficulties with current systems[1]

- Requirement of an Expert.
- Manual processing of the data
- Markup languages and code fragments
- Assumption of More general Ontology.
- Availability of WordNet
- How terms are extracted from Text?

Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation**
- 7 Evaluation
- 8 Conclusion & Future Work
- 9 References

System - 1

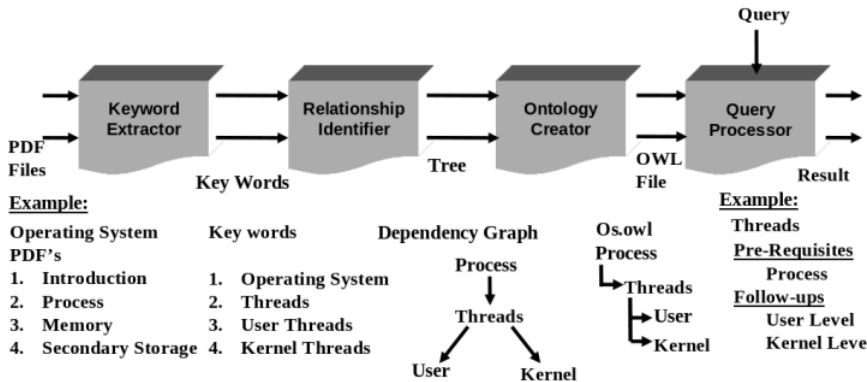
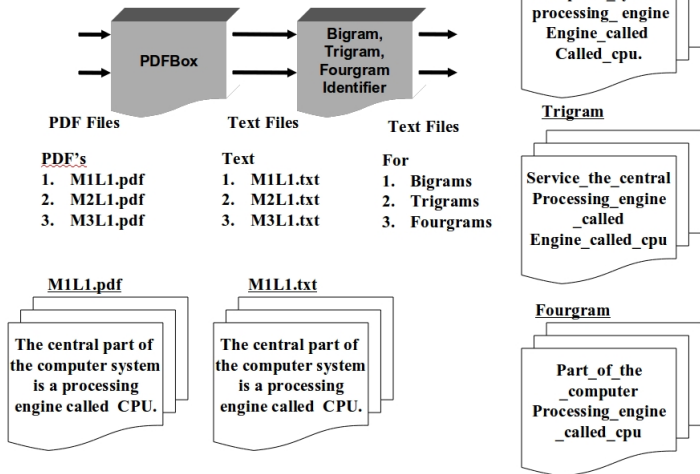


Figure: System overview of System - 1

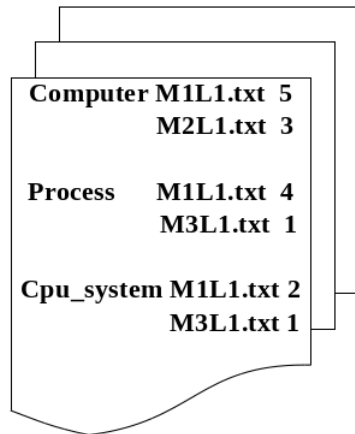
“ngram” Identification

Example:
Operating System

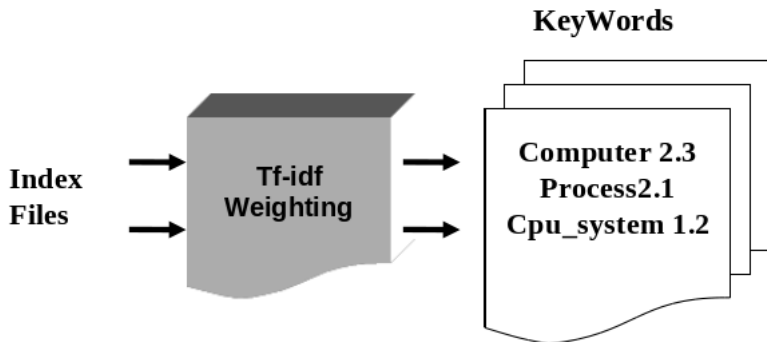


Indexing

Text Files
Bigram
Trigram
Fourgram



Keyword Extraction



Tf-Idf⁶

Given a document collection D , a word w , and an individual document $d \in D$

$$w_d = f_{w,d} * \log(|D|/f_{w,D})$$

$$Tf - Idfweight = \sum_{d=1}^{|D|} w_d$$

where

- $f_{w,d}$ Number of times w appears in the current document d ,
- $|D|$ Number of documents, and
- $f_{w,D}$ Number of documents in which w appears

⁶<http://en.wikipedia.org/wiki/Tf%E2%80%93idf>

Importance of Tf-Idf

- Highest when w occurs many times within a small number of documents;
- Lower when the term occurs fewer times in a document, or occurs in many documents;
- Lowest when the term occurs in virtually all documents.

Apriori Algorithm[5]

Apriori Algorithm is an algorithm for finding association rules.

- 1 Find Keywords (Tf-idf weights)
- 2 Find out the *frequent wordsets* with the given *support* and *confidence*, for all pairs of keywords.

Terminology

Association Rule $i \rightarrow j$ means “if a document contains i then it is likely to contain j ”.

Support The number of documents containing the words in w .

Confidence of this association rule is the probability of j given i .

Ontology for CN

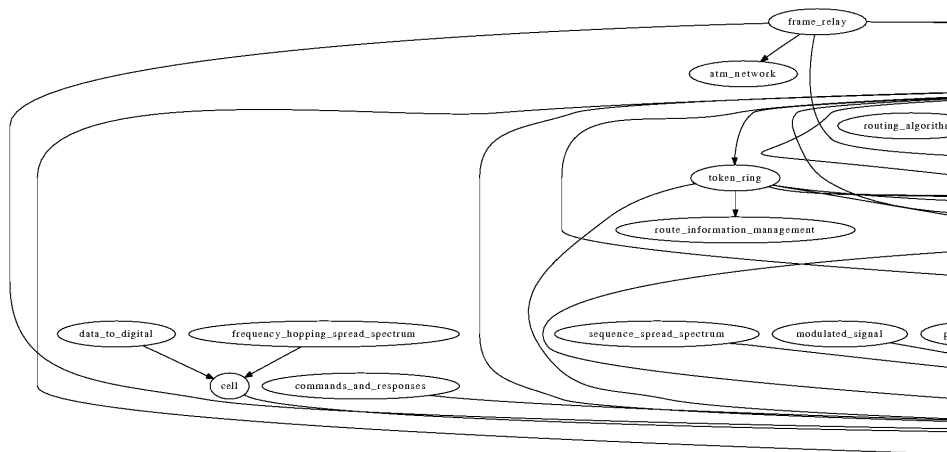


Figure: System output of Computer Networks

System -2

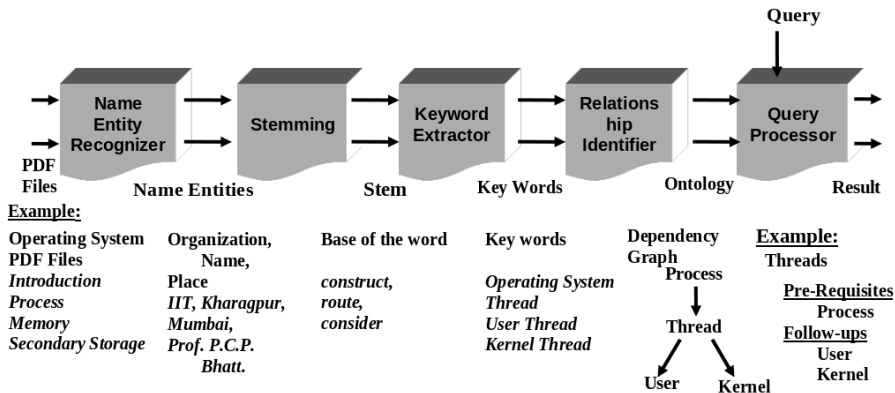


Figure: System overview of System - 2

Name Entity Recognition

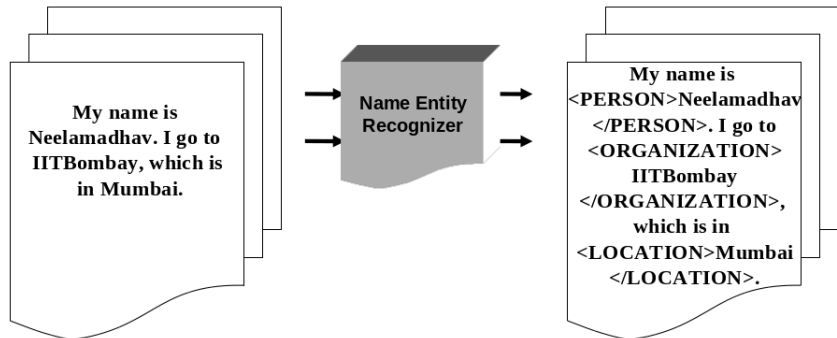


Figure: Name entity recognition

CRF-based NER [6]

We used Stanford NER mainly because of the following reasons

- The NER is trained on CoNLL, MUC and ACE English training data
- By default it recognizes the entities: Person, Location, Organization, Which we need for our experiment
- Finally the NER is trained on both British and American newswire, so robust across both domains

Stemming [7]

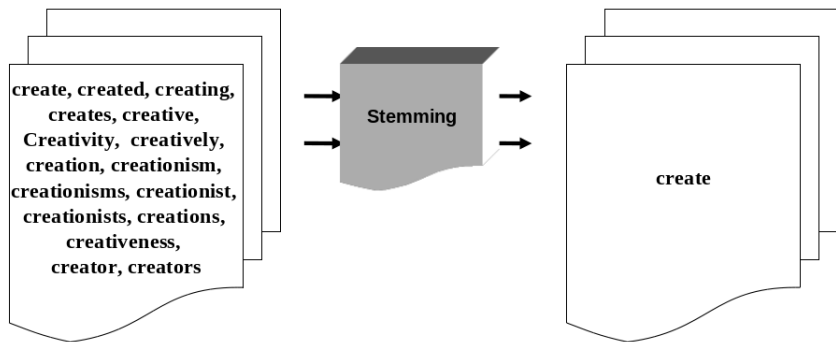


Figure: Stemming

Ontology Development

- We modified our Apriori algorithm in order to get good results.
- We considered a line as a dataset instead of the whole document.

Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation
- 7 Evaluation**
- 8 Conclusion & Future Work
- 9 References

Recall & Precision

We have compared results generated by our system with those of the expert generated results.

Recall (R) Ratio of the relevant results retrieved to the results suggested by the expert.

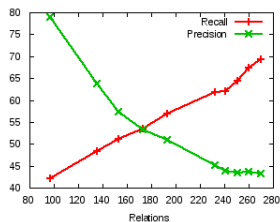
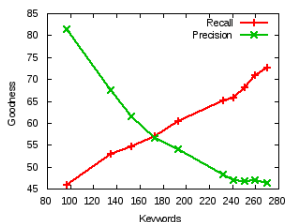
$$R_w = \frac{W_c}{W_e} * 100\% \qquad R_r = \frac{R_c}{R_e} * 100\%$$

Precision (P) Ratio of the relevant results retrieved to the total results identified by the system.

$$P_{pw} = \frac{W_c}{W_s} * 100\% \qquad P_r = \frac{R_c}{R_s} * 100\%$$

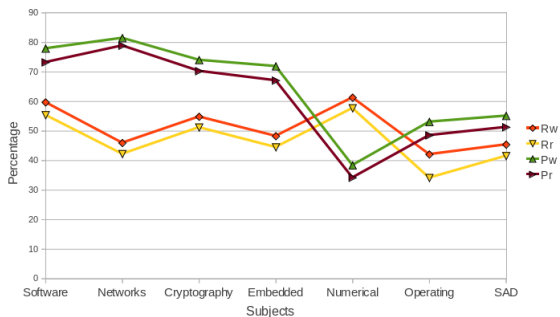
Results of Computer Networks

keywords	Recall(%)		Precision(%)	
	R_w	R_r	P_w	P_r
97	45.93	42.19	81.44	78.86
135	52.91	48.26	67.41	63.84
153	54.65	51.09	61.44	57.32
173	56.98	53.46	56.65	53.21
193	60.47	56.94	53.89	50.82
232	65.12	61.77	48.28	45.19
241	65.70	61.98	46.89	43.88
251	68.02	64.28	46.61	43.43
260	70.93	67.35	46.92	43.56
270	72.67	69.33	46.30	43.16



Recall & Precision for all the subjects

Subjects	keywords	Recall(%)		Precision(%)	
		R_w	R_r	P_w	P_r
Software	95	59.68	55.37	77.89	73.22
Networks	97	45.93	42.19	81.44	78.86
Cryptography	100	54.81	51.19	74	70.26
Embedded	96	48.25	44.47	71.88	67.07
Numerical	99	61.29	57.69	38.38	34.18
Operating	98	42.06	34.1	53.06	48.52
SAD	98	45.38	41.52	55.1	51.28



Confusion Matrix

		System Results	
		Positive	Negative
Expert Results	True	True Positive	False Negative
	False	False Positive	True Negative

Where

True Positive: Number of correct results that were correctly identified,

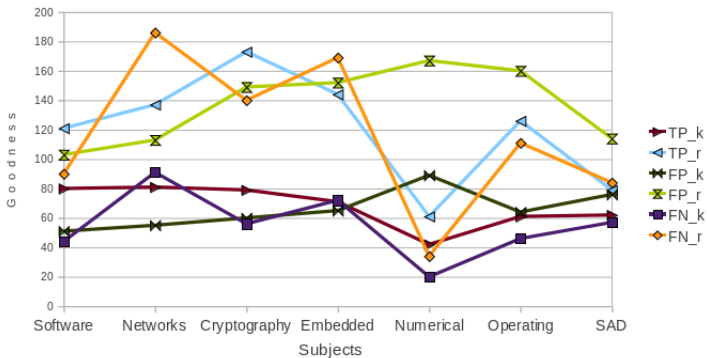
False Positive: Number of incorrect results that were incorrectly classified as positive,

True Negative: Number incorrect results that were identified as negative,

False Negative: Number of correct results that were incorrectly classified as negative.

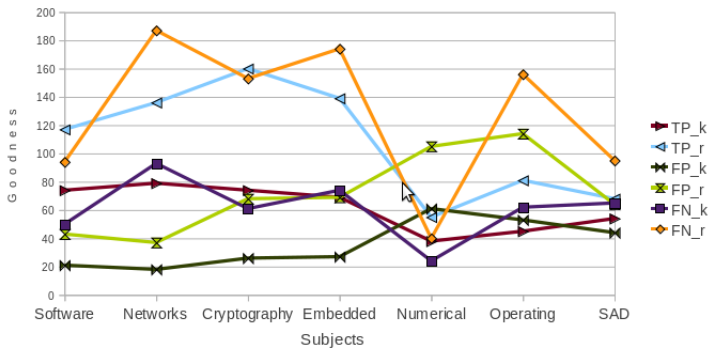
Confusion Matrix for System-1

Subject	System Generated		Expert identified		True Positive		False Positive		False Negative	
	Keywords	Relations	Keywords	Relations	Keywords	Relations	Keywords	Relations	Keywords	Relations
Software	131	224	124	211	80	121	51	103	44	90
Networks	136	250	172	323	81	137	55	113	91	186
Cryptography	139	322	135	313	79	173	60	149	56	140
Embedded	136	296	143	313	71	144	65	152	72	169
Numerical	131	228	62	95	42	61	89	167	20	34
Operating	125	286	107	237	61	126	64	160	46	111
SAD	138	193	119	163	62	79	76	114	57	84



Confusion Matrix for System-2

Subject	System Generated		Expert identified		True Positive		False Positive		False Negative	
	Keywords	Relations	Keywords	Relations	Keywords	Relations	Keywords	Relations	Keywords	Relations
Software	95	160	124	211	74	117	21	43	50	94
Networks	97	173	172	323	79	136	18	37	93	187
Cryptography	100	228	135	313	74	160	26	68	61	153
Embedded	96	208	143	313	69	139	27	69	74	174
Numerical	99	160	62	95	38	55	61	105	24	40
Operating	98	195	107	237	45	81	53	114	62	156
SAD	98	132	119	163	54	68	44	64	65	95



System-1 Vs. System-2

- *System-1* generated more number of False Positives than *System-2*
- Processing time of *System-2* is lesser than *System-1*
- Named Entity Recognition

Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation
- 7 Evaluation
- 8 Conclusion & Future Work**
- 9 References

Future Work

- We observed that Recall was maximum when the number of keywords was nearly equivalent to 95.
- The optimum value of Recall was obtained when the number of unigrams was 40, number of bigrams was 30, number of trigrams was 20, and number of fourgrams was 10.
- Efficiency of the system is mainly dependent on number of PDF files in case of multiple PDF files.
- And dependent on number of pages in case of a single pdf file.
- Relationship identification and keyword extraction algorithms can be modified for better results.

Outline

- 1 Motivation
- 2 Problem Statement
- 3 Ontology
- 4 Manual Ontology Generation
- 5 Semi Automatic Ontology Generation
- 6 Solution & Implementation
- 7 Evaluation
- 8 Conclusion & Future Work
- 9 References

References

- [1] Ivan Bedini and Benjamin Nguyen. *“Automatic ontology generation: State of the art.”* In PRISM Laboratory Technical Report. University of Versailles, 2007.
- [2] Natalya F. Noy and Deborah L. McGuinness, *“Ontology Development 101: A Guide to Creating Your First Ontology”*, Available from Internet:<http://protege.stanford.edu/publications/ontology_development/ontolgy101.html>
- [3] Oscar Corcho, Mariano Fernández-López, Asunción Gómez-Pérez, *“Methodologies, tools and languages for building ontologies. Where is their meeting point”*, Data & Knowledge Engineering 46(2003) 41-64.
- [4] Deborah L. McGuinness, Frank van Harmelen *“OWL Web Ontology Language Overview”*, <http://www.w3.org/TR/owl-features/>, W3C Recommendation 10 February 2004.
- [5] Rakesh Agrawal and Ramakrishnan Srikant. *“Fast algorithms for mining association rules in large databases.”* In Proceedings of the 20th International Conference on Very Large Data Bases, VLDB '94, pages 487-499, San Francisco, CA, USA, 1994. Morgan Kaufmann Publishers Inc.
- [6] Jenny Rose Finkel, Trond Grenager, and Christopher Manning. *“Incorporating non-local information into information extraction systems by gibbs sampling.”* In ACL, pages 363-370, 2005.
- [7] Martin Porter. *“An algorithm for suffix stripping.”* In Workshop on Multimedia Information Systems, 1980.