


## Efficient Information Retrieval Using Measures of Semantic Similarity

**Krishna Sapkota**  
Center for Research in Social Defense Technology, Language Technology Group, nec  
**Student Affiliate**


**Laxman Thapa**  
Nepal Engineering College (nec)

**Shailesh Pandey**  
Center for Research in Social Defense Technology, Language Technology Group, nec  
**Research Engineer**



## Introduction


- Information in WWW are scattered and diverse in nature
- Users frequently fails to describe the information to retrieve
- Traditional search techniques are constrained by keyword based matching techniques
  - Hence low precision and recall is obtained



## Limitations of Traditional Search


- Miss to retrieve synonymy terms
- Users must be intelligent
- Do not retrieve conceptual terms

Hence there is the need for semantic feature in search



## Foundations


- Semantic similarity and relatedness
  - Degree of closeness of meaning of words
- Semantic Similarity
  - For example car and bicycle are more similar than car and human
- Semantic Relatedness
  - Car and bicycle are similar but hot and cold are more related
- Similarity applies to Noun but relatedness covers all category



## Approaches for Computing Semantic Similarity

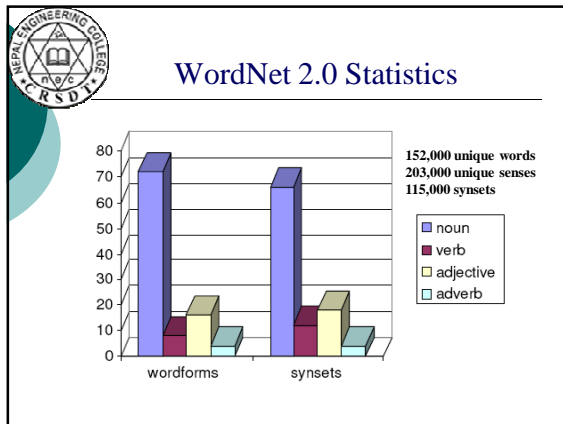
- Dictionary based
- Thesaurus based
- Semantic Network
  - Using path or node
  - Using Information Content

Word Pairs	Humans	Sim	Sim <sub>in</sub>	Sim <sub>jac</sub>
fruit-furnace	0.05	1.85	0.14	0.05
monk-slave	0.57	2.53	0.21	0.05
coast-hill	1.26	6.19	0.53	0.09
magician-oracle	1.82	13.5	0.96	1.00
brother-lad	2.41	2.53	0.23	0.06
food-fruit	2.69	1.50	0.22	0.09
furnace-stove	3.11	1.85	0.13	0.04
boy-lad	3.82	8.29	0.72	0.18
automobile-car	3.92	8.62	1.00	maximum

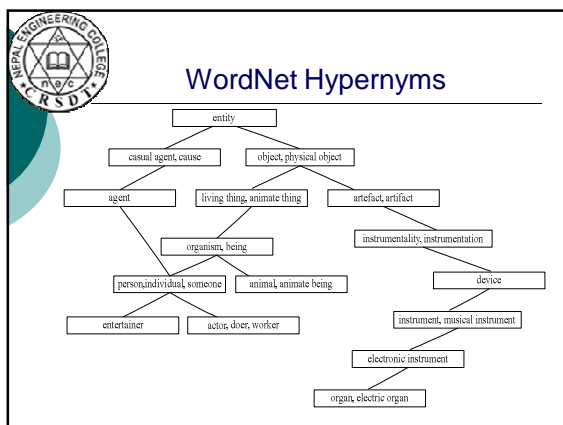


## WordNet Semantic Network

- Developed at the cognitive science laboratory of Princeton University
- English online lexical database
- Based on psycholinguistic principle
- Represents World Knowledge
- Four syntactic categories noun, verb, adverb, adjective are organized in to different relationships

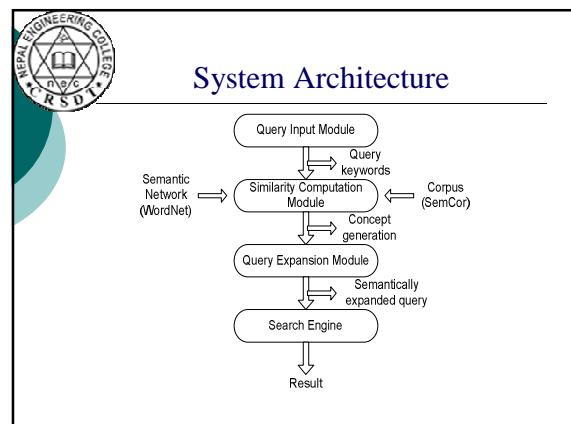


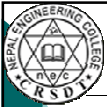
- ### WordNet Relations
- Semantic relations
    - Hypernymy/hyponymy or IS-A
    - Meronymy/holonymy or HAS-A
    - Entailment
  - Lexical relation
    - Synonymy
    - Antonymy



- ### Extending the Existing Search Engine
- Low cost & less time rather than to built from scratch
  - Our system uses Google API to perform search
  - Google API can be accessed using SOAP from different programming language
  - The users original queries are expanded and fed to Google Search Engine
    - Hence our system lies in between users and Google's interface

- ### Our Algorithm
- Original query as "Bank Interest"
    - Bank has 10 concept and interest has 7 concept in WordNet 2.0
  - To get optimized query
    - calculate similarity score of each pair of concepts as  $10 \times 7 = 70$  similarity score is obtained
    - Pick up the concept with highest value of similarity score
    - Replace the original query with synonyms and hypernyms of concept that has highest similarity score






## Search Results

Topic Retrieved	WWW Address
The National Neopian Bank	<a href="http://www.neopets.com/bank.html">http://www.neopets.com/bank.html</a>
Bankrate.com	<a href="http://www.bankrate.com">http://www.bankrate.com</a>
Certificate of deposit of interest rates: Compare the best rate	<a href="http://www.bankrate.com/bm/rate/deposits">http://www.bankrate.com/bm/rate/deposits</a>
Bank Interest Calculator	<a href="http://www.digita.com/itsali/calculators/bankinterestcalculator/">http://www.digita.com/itsali/calculators/bankinterestcalculator/</a>
Personal Banking System, savings, bank interest rate, tax	<a href="http://www.thisismoney.co.uk/saving">http://www.thisismoney.co.uk/saving</a>
Indian Bank-Interest rates	<a href="http://www.indian-bank.com/interest.htm">http://www.indian-bank.com/interest.htm</a>
National Australian Bank	<a href="http://www.national.com.au/business-solution/02253300.htm">http://www.national.com.au/business-solution/02253300.htm</a>
Bank Interest	<a href="http://www.ato.gov.au/content/48327.htm">http://www.ato.gov.au/content/48327.htm</a>
Infochoice Banking	<a href="http://www.infochoice.com.au/banking/default.asp">http://www.infochoice.com.au/banking/default.asp</a>
National Bank Interest Rate Graph	<a href="http://www.nbnz.co.nz/economics/interest/">http://www.nbnz.co.nz/economics/interest/</a>

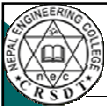
*First 10 results from Google search for 'bank interest'*



## Search Results Contd.


Topic Retrieved	WWW Address
Scholarly articles (for all expanded term)	<a href="http://www.google.com">http://www.google.com</a>
Bloomberg.com: Financial Glossary	<a href="http://www.bloomberg.com/invest/glossary/bg/glost.htm">http://www.bloomberg.com/invest/glossary/bg/glost.htm</a>
FCAC-Glossary	<a href="http://www.fcac-acfc.gc.ca/eng/glossary.asp">www.fcac-acfc.gc.ca/eng/glossary.asp</a>
Guide to organize a new state bank in Florida	<a href="http://www.flfc.com/banking/howtoorg.htm">http://www.flfc.com/banking/howtoorg.htm</a>
PSI- Performance Solution International	<a href="http://www.goto-psi.com/glossary.htm">http://www.goto-psi.com/glossary.htm</a>
[pdf] How should financial institution and market should be structured	<a href="http://www.iadh.org/res/publications/pubfiles/">http://www.iadh.org/res/publications/pubfiles/</a>
Women's wallstreet.com - glossary	<a href="http://www.wimenswallstreet.com/tools-resources/glossary/f.htm">http://www.wimenswallstreet.com/tools-resources/glossary/f.htm</a>
Operational risk poses challenges to financial institution	<a href="http://knowledge.wharton.upenn.edu/article.cfm?articleid=582">http://knowledge.wharton.upenn.edu/article.cfm?articleid=582</a>
FDIC:FDIC Banking Review	<a href="http://www.fdic.gov/bank/analytical/banking/2004nov/article1/index">www.fdic.gov/bank/analytical/banking/2004nov/article1/index</a>
BKD LLP: Financial Services	<a href="http://www.bkd.com/industry/financial-services">http://www.bkd.com/industry/financial-services</a>

*First 10 results from our extended Google for bank interest*




## Programming Languages and Tools

- Used PHP as it provides easy access to internet service
- Used MySQL format of WordNet 2.0
- All 186 files of semantic concordance were converted to MySQL format
- Used nuSOAP a implementation of SOAP architecture in PHP by nuSphere corporation
- Used Google API service




## Further Enhancements

- Word Sense Disambiguation (WSD)
  - Process of assigning correct sense to the word
    - The chair for principal was not occupied. (position)
    - This chair is made up of good wood. (furniture)
- Measures for more syntactic category
- Customizable interface to search



## Conclusion

- Users queries were replaced by our analyzed terms
- User's must not be intelligent
- We provide Concept based search
- For better precision and recall incorporate WSD



## THANK YOU

## QUESTIONS?

