# Classifications and lightweight ontologies

**Knowdive Group**
**Fausto Giunchiglia, Ilya Zaihrayeu**

**First presented at Jilin University, Changchun, 19 Nov 2007**
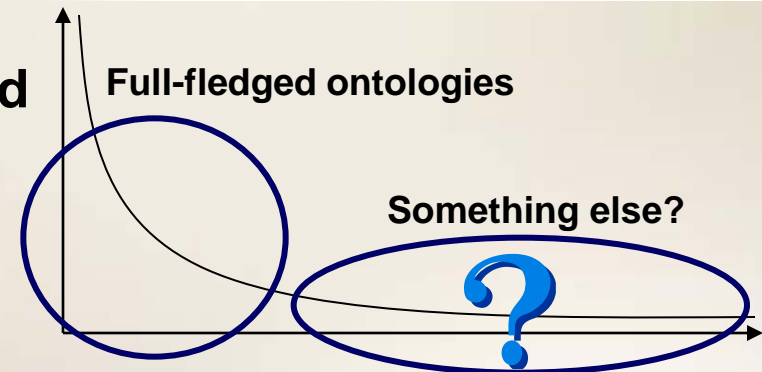
# Index

# Motivation

- The semantic web (SW) chicken-and-egg problem: users will not mark up their data unless they perceive an added value from doing so, and tools to demonstrate this value will not be developed unless a "critical mass" of annotated data is achieved (Hendler 2001)



- Natural language processing (NLP) has advanced to the point where it can break the impasse and open up the possibilities of the Semantic Web (Barney Pell, invited talk at ISWC 2007)



- People (computer scientists) tend to use simple ontologies, simple class hierarchies, simple rules. What is on the long tail? (Chris Welty, invited talk at ISWC 2007)
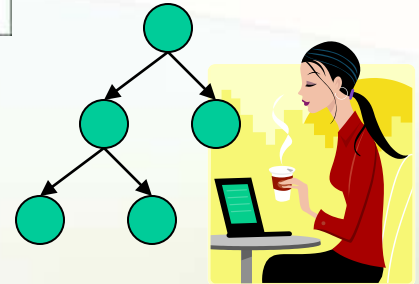


Full-fledged ontologies

Something else?

# Motivation

•**Taxonomies, thesauri, business catalogues, faceted classifications, web directories, user classifications are widely used as means to encode knowledge and organize data on the web and in personal document collections**
**Is not it the long tail?**

•**However, since their labels are written in natural language (NL), they are very hard to be reasoned about by automated software agents and represent annotations of little use for SW apps.**
**Is not it where NLP can help?**

Arts and entertainment

Production of man-made fibers

Cars, Boats, Vehicles & Parts

Vegetables & Vegetarian

Pictures from vacation in Busan

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- **Disambiguating Labels**
  - **Named Entity Locating**
  - **Part-of-Speech Tagging**
  - **Word Sense Disambiguation**
- **Disambiguating Edges**
- **Applications**
- **Conclusions**

# Classifications vs. Ontologies

- **Classifications**
  - **Rooted trees where nodes are assigned natural language labels**
  - **Easy to understand (for humans), pervasively used**
    - Industry and standards: DMOZ, DDC, Amazon, BBC
    - Personal: favorites, email folders, file system. A handful of classifications at each PC!
- **Ontologies**
  - **Can be complex graph-like structures described in a formal language**
  - **Can only be operated by the (very narrow) community of ontology engineers**
    - How many ontologies do you have on your PC?

# Classifications vs. Ontologies, cont'd

- ## Classifications
  - **Labels are ambiguous (because of natural language)**
  - **No semantics for edges (not necessarily ontological relations such as is-a, part-of)**
  - **Automated reasoning about them is very hard!**
- ## Ontologies
  - **No ambiguity due to the use of formal language**
  - **When represented as a graph-like structure, all structure elements have a well-defined semantics**
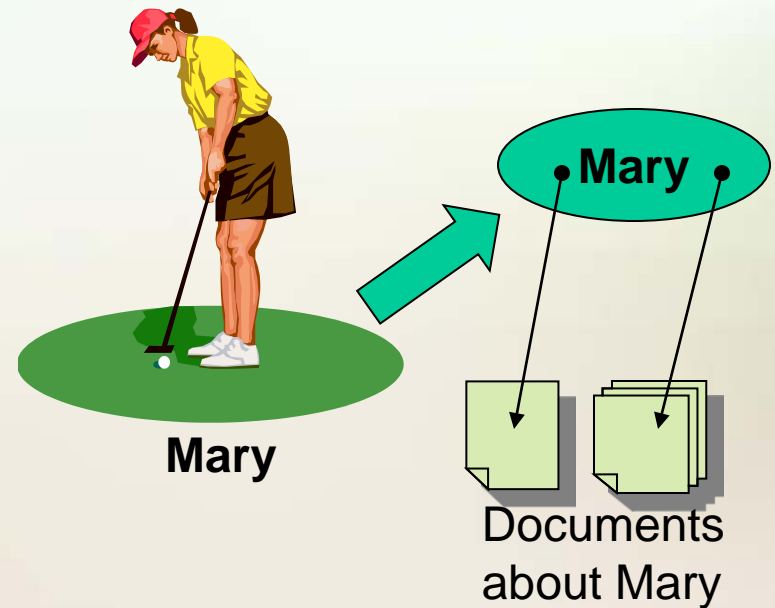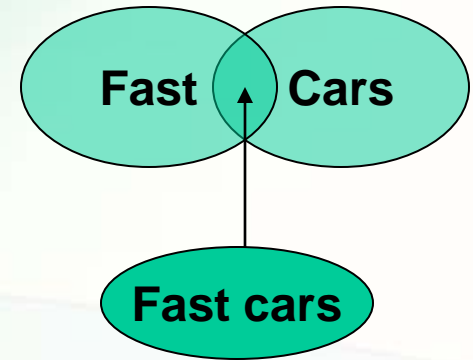  - **Designed (also) for automated reasoning**

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- **Disambiguating Labels**
  - **Named Entity Locating**
  - **Part-of-Speech Tagging**
  - **Word Sense Disambiguation**
- **Disambiguating Edges**
- **Applications**
- **Conclusions**

# Our approach

- **Build a bridge from classifications to (lightweight) ontologies in order to automate operations on classifications**

- **Operationally, our proposal is twofold:**
  - **(Automatically) extract semantics from the classifications' labels and structure, thus converting the classifications into (lightweight) ontologies**
  - **Encode operations on classifications as reasoning problems on lightweight ontologies**

# Our approach

- As from [1], WordNet senses of adjectives and common nouns as well as proper nouns become atomic concepts

- Extension of a common noun concept is the set of documents about objects of the class, denoted by the noun

- Extension of an adjective concept is the set of documents about objects, which possess the qualities, denoted by the adjective

- Extension of a proper name concept is the set of documents about the individual referenced by the proper name
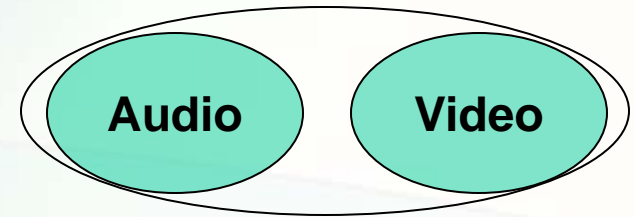
[1] F. Giunchiglia, M. Marchese, and I. Zaihrayeu: **Encoding classifications into lightweight ontologies.** In JoDS VIII, LNCS

Fast | Cars

Fast cars

Mary

Mary

Documents about Mary

# Our approach

•**Syntactic relations between words in the label, coordinating conjunctions, prepositions are translated into logical connectives to build complex formulas, e.g.:**

  •**Coordinating conjunctions "and" and "or" are translated into the logical disjunction**

  •**Prepositions are converted into the logical conjunction**

  •**Words denoting exceptions are translated into the logical negation (almost no such words in classification labels!)**

**Audio and video**
audio ⊔ video

**Audio**   **Video**

**Life in Trento**
life ⊓ trento

**Life** **Trento**

**Runners except sprinters**
runner ⊓ ¬ sprinter

**Runners** **Sprinters**

# Our approach

- **E.g., label "*Bank and personal details of George Bush*" should be translated in DL as:**

$$(\text{bank-noun-1} \sqcup \text{personal-adj-1}) \sqcap \text{detail-noun-1} \sqcap \text{george\_bush}_{NNP}$$

- **A concept, whose extension is the intersection of three sets of documents:**
  - **(i) documents about the President George W. Bush,**
  - **(ii) documents containing isolated facts about something (i.e., details), and**
  - **(iii) the union of documents about bank institutions and documents concerning a particular person or his/her private life**

# Our approach

label: "*Bank and personal details of George Bush*"

- **Despite its seeming simplicity, the translation process is subject to various mistakes originating from inaccurate NLP**
- **Due to a mistake in POS tagging, the word "personal" might be recognized as a noun defined as "*a short newspaper article about a particular person or group*"**
- **Due to a mistake in WSD, the word "bank" might be identified as "*sloping land (especially the slope beside a body of water)*"**
- **Due to a mistake in NE locating, the proper name "George Bush" might not be located and might then be considered as two distinct nouns**
- **Due to a mistake in (syntax) parsing, the input label might be translated into:**

bank-noun-1 ⊔ personal-adj-1 ⊓ detail-noun-1 ⊓ george_bush$_{NNP}$

- **Thus, proper NLP tools are crucial for correct translation**
- **How much of the std. NLP technology can be reused?**

# Our approach

•**The NLP framework, which enables the conversion of classification NL labels into DL formulas is depicted below:**

NL Label

WordNet

Tokenizer → POS Tagger → WSD → Parser

NE Locator

DL Concept

☐ Use existing approaches  ☐ Addressed in this paper  ☐ Future work

# Our approach

- We have analysed the DMoz web directory as our study case
- Standart NLP technology is primarily used on full-fledged sentences, however:
- Web directory labels are short phrases, which provide very limited context for NLP
- Most of the words in a Web directory are nouns, adjectives, articles, conjunctions and prepositions. The verbs and pronouns are very rare in a Web directory while being common in full-fledged sentences
- Named Entities occur densely in a Web directory
- The capital rule is different in a Web directory. In full-fledged sentences, the first words of sentences and the words in proper names are initialized with capital letters. In a Web directory, however, most often every word begins with a capital letter except for prepositions and conjunctions
- The intended sense of a word may depend on the meaning of a word appearing in a label located higher in the classification tree. For instance, noun "Java" means an island if it appears under a node with label "Geography"
- New approaches are required!

# Our approach

- The dataset we used is the English part of DMoz:
  - # of labels: 474,389
  - Avg. length: 1.91 tokens
  - Avg. depth: 7.01
- We randomly selected 12,365 labels (2.61%) for analysis and manual annotation:
  - 8177 non-NE labels (66.13%)
  - 4188 NE labels (33.87%)
  - Nearly all NEs take entire labels: only 7 exceptional labels (0.06%)
- Statistics of POS occurrences in the non-NE labels in the data set is shown below:

| POS | NN | NNS | CC | JJ | NNP | IN | , | TO | CD |
|---|---|---|---|---|---|---|---|---|---|
| Occurrence | 7714 | 3619 | 2893 | 1020 | 239 | 235 | 72 | 18 | 11 |
| Percentage | 48.68 | 22.84 | 18.26 | 6.44 | 1.51 | 1.48 | 0.45 | 0.11 | 0.07 |

| POS | : | VBN | DT | RB | POS | VB | NPS | JJR |
|---|---|---|---|---|---|---|---|---|
| Occurrence | 9 | 6 | 4 | 2 | 2 | 1 | 1 | 1 |
| Percentage | 0.06 | 0.04 | 0.03 | 0.01 | 0.01 | < 0.01 | < 0.01 | < 0.01 |

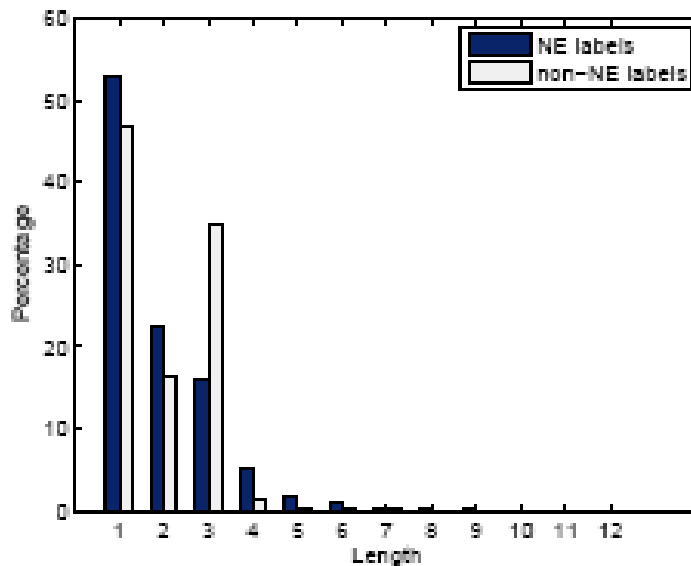Table 1. Statistics of POS occurrences in the data set.

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- **Disambiguating labels**
  - **Named Entity Locating**
  - **Part-of-Speech Tagging**
  - **Word Sense Disambiguation**
- **Disambiguating edges**
- **Applications**
- **Conclusions**

# Named Entity Locating: Approach

- By analyzing the data set, we noticed the following characteristics of NEs:

- Rare labels tend to be NEs. A general label (such as "Arts and Entertainment") can occur thousands of times, while NE labels occur much more rarely. Most of NE labels, such as "Schindler's List" (a movie name), occur only once

- Labels in which most of the tokens are rare words tend to be NEs, e.g., "Agios Dometios" is a geography name and each of its tokens occurs only once

- Letter bars such as single letter "A", "B", ..., "Z" and also double letters "Aa", "Ab", ..., "Zz" are created for the convenience of navigation. They are good indicators of NEs, as nearly all children of these labels are NEs

- In an NE label, initial articles, such as "the", "a" and "an", are usually put at the end after a comma. For example, "The Magic School Bus" is written as "Magic School Bus, The"

# Named Entity Locating: Approach

- The NE and non-NE labels distribute differently on their lengths (see Fig. 2(a))
- The NE and non-NE labels distribute differently on their depths (see Fig. 2(b))



Fig. 2. (a) Label length distribution; (b) Label depth distribution.

# Named Entity Locating: Approach

- **Taking these characteristics into account, we implemented the NE locator using Conditional Maximum Entropy Model (CMEM) with Gaussian smoothing**

- **The feature classes for CMEM have been chosen according to the characteristics described above:**
  - *WordsInLabel*: **The first two and the last two tokens in the label**
  - *WordsInPath*: **The first and the last tokens in the label's parent, grandparent, the farthest ancestor (excluding the root "Top") and the second farthest ancestor**
  - *LengthOfLabel*: **The number of tokens in the label**
  - *DepthOfLabel*: **Depth of the label (distance from the root node)**
  - *FrequencyOfLabel*: **Count how many times the label occurs in the whole directory**
  - *AveFrequencyOfTokens*: **Count how many times each token in the label occurs in the whole directory, and calculate the average**

# Named Entity Locating: Evaluation

- **First, we trained the NE locator by using each feature class to compare their contributions**
- **Then, we trained the NE locator again with some combinations of feature classes to see the best performance we can reach**
- **The results are reported below:**

| Feature Class | PNE | RNE | FNE |
|---|---|---|---|
| 1. *WordsInLabel* | 81.49 | 94.33 | 87.45 |
| 2. *WordsInPath* | 89.48 | 79.36 | 84.12 |
| 3. *FrequencyOfLabel* | 75.04 | 91.30 | 82.37 |
| 4. *AveFrequencyOfTokens* | 76.05 | 82.95 | 79.35 |
| 5. *DepthOfLabel* | 53.13 | 78.76 | 63.45 |
| 6. *LengthOfLabel* | 64.64 | 8.05 | 14.32 |
| 1+2 | 92.08 | 94.20 | 93.13 |
| 1+2+3+4+5+6 | 93.45 | 94.04 | 93.75 |

Table 3. Performance results of the NE locator.

- **One state-of-the-art system [4] of NE locating in the Web environment on full-fledged sentences has the performance of 59% in precision, 66% in recall and 38% in F-score**
- **NE locating on web directories is an easier task, as we only need to tell whether a label is an NE or not**

[4] D. Downey, M. Broadhead, and O. Etzioni: **Locating complex named entities in web text**. In Proc. of IJCAI, 2007, 2007.

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- <span style="color:red">**Disambiguating labels**</span>
  - **Named Entity Locating**
  - <span style="color:red">**Part-of-Speech Tagging**</span>
  - **Word Sense Disambiguation**
- **Disambiguating edges**
- **Applications**
- **Conclusions**

# Part-of-Speech Tagging: Approach

- **In our experiments, we employed two POS taggers:**
    - **FudanNLP POS tagger (based on the Conditional Random Field model) [5]**
    - **OpenNLP POS tagger (based on the Conditional Maximum Entropy Model) [6]**
- **We retrained these tools on our data set and checked if we gain an improvement in accuracy w.r.t. the case when the tools are trained on full-fledged sentences**
- **To avoid a negative influence of NE labels on the training of a POS tagger, both POS taggers were trained and tested only on the non-NE labels in the data set**

[5] X. Qian: **A CRF-based pos tagger**. Technical Report FDUCSE 07302, Fudan University, 2007.

[6] The OpenNLP project. See http://opennlp.sourceforge.net/.

# Part-of-Speech Tagging: Evaluation

- **The following 2 measures were used to evaluate the performance of the POS taggers:**
  - *Precision of POS tagger by Tokens* **(PPT): count tokens which are tagged with the correct tag, and calculate the percentage**
  - *Precision of POS tagger by Labels* **(PPL): count labels whose tokens are all correctly tagged, and calculate the percentage**

| | $PPT_0$ | $PPT_1$ | Gain | $PPL_0$ | $PPL_1$ | Gain |
|---|---|---|---|---|---|---|
| OpenNLP | 91.27 | 97.23 | +6.16 | 84.68 | 96.00 | +11.52 |
| FudanNLP | 96.12 | 97.33 | +1.21 | 92.72 | 96.02 | +3.30 |

Table 4. Performance results of the OpenNLP and FudanNLP POS taggers before and after retraining.

- **The performance of a state-of-the-art POS tagger on full-fledged sentences is 97.24% [7] in token precision (PPT) which is very close to ours**

[7] Toutanova, K., Klein, D., Manning, C.D., Singer, Y.: **Feature-rich part-of-speech tagging with a cyclic dependency network**. In: Proc. of Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology, vol. 1, pp. 173–180 (2003)

# NE Locating & POS Tagging: Evaluation

- To check whether our data set is properly sized, we performed incremental training, namely, keeping the testing set unchanged, we checked how performance varied with the growing size of the training set



- As it can be observed, performance measures increase significantly when the size of the training set grows from 1000 to 5000-7000 samples, and when the number of samples becomes greater, the performance measures change only slightly

- Empirically, we conclude that our NE locating and POS tagging models are effective and stable enough to be used on web directories such as DMoz

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- <span style="color:red">**Disambiguating labels**</span>
  - **Named Entity Locating**
  - **Part-of-Speech Tagging**
  - <span style="color:red">**Word Sense Disambiguation**</span>
- **Disambiguating edges**
- **Applications**
- **Conclusions**

# Word Sense Disambiguation: Approach

- **The proposed WSD algorithm traverses the nodes of the classification tree in the BFS or DFS order**

- **Then, at each node, it first finds concept tokens, i.e., tokens which are present in WordNet as adjectives and/or as nouns**

- **Next, it identifies ambiguous concept tokens, i.e., concept tokens which have more than one sense**

- **Ambiguous concept tokens of each node are processed by the WSD algorithm**

- **The ultimate goal of the algorithm is to select only one sense for each ambiguous concept token**

# Word Sense Disambiguation: Approach

1. Identify the POS of the token and, if the token has senses of this POS, then preserve these senses and discard senses belonging to the other POS, if any
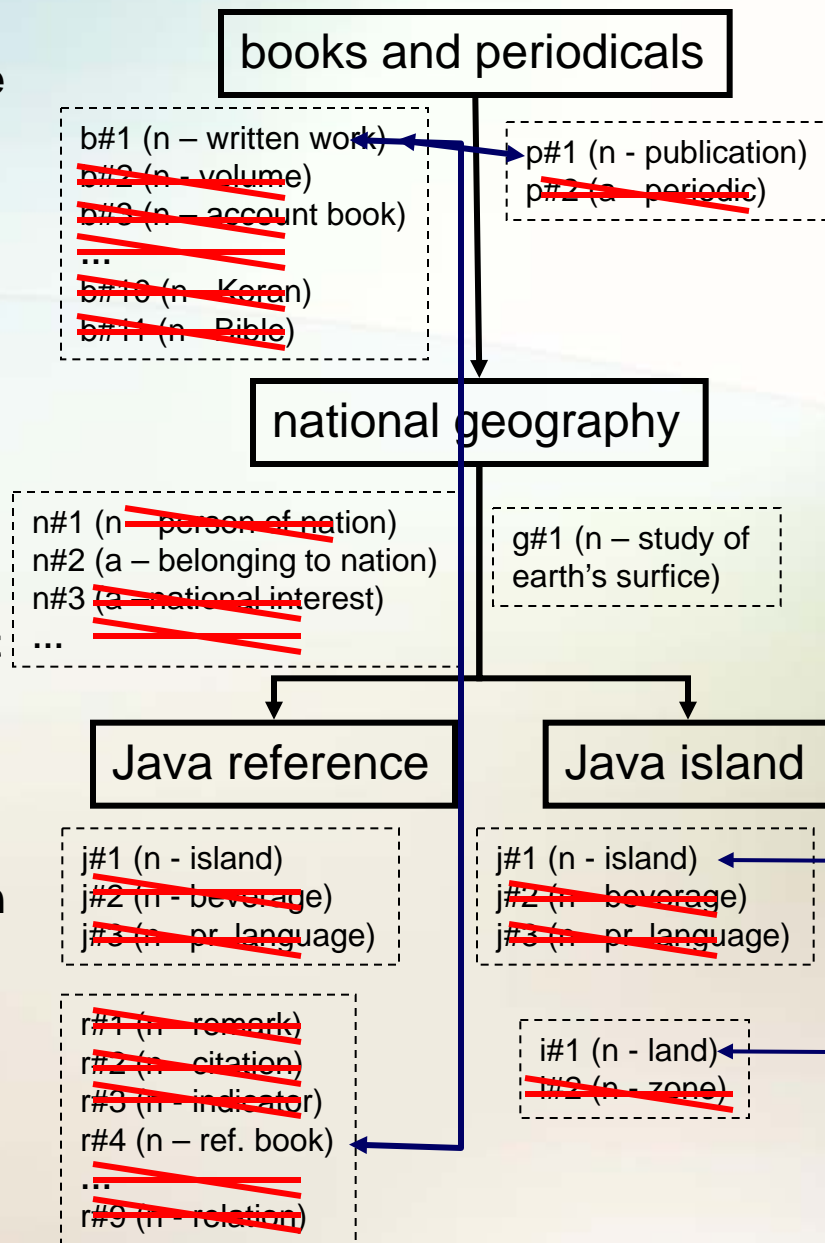
2. Preserve noun token senses if they are hypernyms or hyponyms of active noun senses of other concept tokens in the label, and discard the other senses

3. Preserve noun token senses if they are located within a certain distance in the WordNet hypernymy hierarchy from active noun senses of other concept tokens in the label

4. Preserve noun token senses if they are hyponyms of active noun senses of concept tokens appearing in the label of an ancestor node, and discard the other senses

5. Preserve noun token senses if they are located within a certain distance in the WordNet hypernymy hierarchy from active noun senses of concept tokens appearing in the labels of ancestor nodes, and discard the other senses

6. Preserve the first active noun sense (in WordNet) and discard the other active senses. If there is no active noun sense, then preserve the first active adjective sense and discard the other active senses

**books and periodicals**

b#1 (n – written work)
b#2 (n – volume)
b#3 (n – account book)
...
b#10 (n – Koran)
b#11 (n – Bible)

p#1 (n - publication)
p#2 (a – periodic)

**national geography**

n#1 (n – person of nation)
n#2 (a – belonging to nation)
n#3 (a – national interest)
...

g#1 (n – study of earth's surfice)

**Java reference**

j#1 (n - island)
j#2 (n - beverage)
j#3 (n – pr. language)

r#1 (n – remark)
r#2 (n – citation)
r#3 (n - indicator)
r#4 (n – ref. book)
...
r#9 (n - relation)

**Java island**

j#1 (n - island)
j#2 (n – beverage)
j#3 (n – pr. language)

i#1 (n - land)
i#2 (n - zone)

# Word Sense Disambiguation: Evaluation

- **To evaluate the performance of our WSD algorithm, we have selected a DMoz subtree rooted at Top/Business/Consumer_Goods_and_Services**
- **Subtree characteristics:**
  - **781 nodes, which have 1368 tokens in total**
  - **1107 concept tokens, out of which 845 are ambiguous**
  - **4.05 is the average polysemy of ambiguous concept tokens**
  - **6 is the maximal depth**
  - **4.22 is the average branching factor**

| | # | Step 1 | | Step 2 | | Step 3 | | | Step 4 | | Step 5 | | | Step 6 | | Accur. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | tok. | acc. | tok. | acc. | thr. | tok. | acc. | tok. | acc. | thr. | tok. | acc. | tok. | acc. | |
| **Baseline** | 1 | na | na | na | na | na | na | na | na | na | na | na | na | 845 | 63.90 | 63.90 |
| | 2 | 84 | 98.81 | na | na | na | na | na | na | na | na | na | na | 761 | 60.84 | 64.62 |
| | 3 | 84 | 98.81 | 11 | 100 | na | na | na | na | na | na | na | na | 750 | 61.60 | 65.80 |
| | 4 | 84 | 98.81 | 11 | 100 | 10 | 250 | 50.80 | na | na | na | na | na | 500 | 61.20 | 62.37 |
| | 5 | 84 | 98.81 | 11 | 100 | 2 | 24 | 87.50 | na | na | na | na | na | 726 | 60.88 | 65.92 |
| **Best result** | 6 | 84 | 98.81 | 11 | 100 | 2 | 24 | 87.50 | 8 | 87.50 | na | na | na | 718 | 61.28 | 66.51 |
| | 7 | 84 | 98.81 | 11 | 100 | 2 | 24 | 87.50 | 8 | 87.50 | 10 | 379 | 33.24 | 339 | 42.77 | 46.51 |
| | 8 | 84 | 98.81 | 11 | 100 | 2 | 24 | 87.50 | 8 | 87.50 | 2 | 43 | 41.86 | 675 | 60.00 | 64.49 |

Table 5. Performance results of the WSD algorithm.

# Word Sense Disambiguation: Evaluation

- **The best accuracy of the WSD algorithm presented in [8] is 47.3% for polysemous nouns**

- **Similar to our case, in [8] the best accuracy is only slightly higher than the baseline**

- **A more recent work, [9], uses a web search engine (together with WordNet) for WSD and reaches 76.55% in accuracy for polysemous nouns in the best case**

- **While the average polysemy of nouns is close to ours (4.08), the size of the context window varied from 3 to 7 words that are known to WordNet, what is not possible to have in our case**

- **Empirically, we conclude that the result of our WSD algorithm is comparable to the state-of-the-art in this field of NLP, however, it is a (very) hard problem to solve**

[8] E. Agirre and G. Rigau: **A proposal for word sense disambiguation using conceptual distance**. In the First International Conference on Recent Advances in NLP, Tzigov Chark, Bulgaria, September 1995.
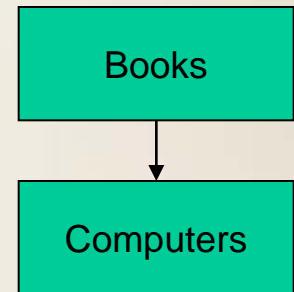
[9] C. Yang and J. C. Hung: **Word sense determination using wordnet and sense co-occurrence**. In Proceedings of AINA 2006 - Volume 1, pages 779–784, Washington, DC, USA, 2006. IEEE Computer Society.

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- **Disambiguating labels**
  - **Named Entity Locating**
  - **Part-of-Speech Tagging**
  - **Word Sense Disambiguation**
- **Disambiguating edges**
- **Applications**
- **Conclusions**

# Disambiguating edges

- **In terms of the extension, the meaning of a node is the set of documents which are about the node's label** and **about the parent node**

- **We represent this property as logical conjunction and we encode it as** *concept at node* $l_i^N$**:**

  $l_i^N = l_i^F$**, if $n_i$ is the root; otherwise:**

  $l_i^N = l_i^F \sqcap l_i^F$**, where $n_i$ is the parent of $n_i$**

- **E.g., the concept of the child node on the right is** [computer*] $\sqcap$ [books*]**, and its extension is the set of documents, which are (about) computer books**

Books

Computers

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- **Disambiguating labels**
  - **Named Entity Locating**
  - **Part-of-Speech Tagging**
  - **Word Sense Disambiguation**
- **Disambiguating edges**
- <span style="color:red">**Applications**</span>
- **Conclusions**

# Applications

- **Once NL labels are translated in DL formulas, the artifact (e.g., web directory) becomes a sort of lightweight ontology suitable for automated reasoning, e.g., for:**

- **Semantic matching (used, e.g., in data integration) [2]**

- **Document classification [3]**

- **Semantic search [1]**

[2] F. Giunchiglia, M. Yatskevich, P. Shvaiko: **Semantic Matching: Algorithms and Implementation**.  In JoDS IX, LNCS 4601, pp. 1-38, 2007

[3] Fausto Giunchiglia, Ilya Zaihrayeu, and Uladzimir Kharkevich: **Formalizing the get-speciffic document classification algorithm**. In  ECDL2007, Budapest, Hungary, September 2007

[1] F. Giunchiglia, M. Marchese, and I. Zaihrayeu: **Encoding classifications into lightweight ontologies.** In JoDS VIII, LNCS 4830, winter 2006

# Reasoning

- **A problem expressed in propositional DL can be translated into an equivalent propositional satisfiability problem which can be solved using a (sounds and complete) SAT decider**

- **If we need to check if relation *rel* holds between two concepts A and B, we check for validity:**

$$KB \rightarrow rel(A, B)$$

  **where KB is a set of axioms, which represents our a priori knowledge**

- **We use a lexical data base (WordNet) to build the core of the KB:**

  - **A is a hypernym of B becomes [A ← B]**
  - **A is holonym  of B becomes [A→B]**
  - **A is a synonym of B becomes [A ←→ B]**

# Index

- **Motivation**
- **Classifications vs. Ontologies**
- **Our Approach**
- **Disambiguating labels**
  - **Named Entity Locating**
  - **Part-of-Speech Tagging**
  - **Word Sense Disambiguation**
- **Disambiguating edges**
- **Applications**
- **Conclusions**

# Conclusions

- **We have presented an approach to converting classifications into lightweight ontologies**
- **In principle, this approach allows users to create (lightweight) ontologies as a a *by-product* of normal computer use, which lowers down the barrier of entering the SW and can help solve the chicken-and-egg problem**
- **The approach aims at the long tail and not at the head of ontology-based knowledge and data organization systems, which can potentially help the SW to scale on the large**
- **However, the quality of these ontologies crucially depends on  NLP tasks, related to the conversion process**
- **The NLP analysis reported in this paper, to the best of our knowledge, is the first investigation of how NLP technology can be applied on web directory labels and, more generally, on short natural language (noun) phrases**
- **However, we still need to understand whether NLP for short phrases is a new domain which requires new methodologies and tools or those used in standard NLP can be suitably adopted**

# Thank you for the warm welcome

# Thank you

# Questions?

# References

[1] F. Giunchiglia, M. Marchese, and I. Zaihrayeu: **Encoding classifications into lightweight ontologies.** In JoDS VIII: Special Issue on Extended Papers from 2005 Conferences, LNCS 4830, winter 2006

[2] F. Giunchiglia, M. Yatskevich, P. Shvaiko: **Semantic Matching: Algorithms and Implementation**. Journal on Data Semantics (JoDS), IX, LNCS 4601, pp. 1-38, 2007

[3] Fausto Giunchiglia, Ilya Zaihrayeu, and Uladzimir Kharkevich: **Formalizing the get-speciffic document classification algorithm**. In 11th European Conference on Research and Advanced Technology for Digital Libraries (ECDL2007), Budapest, Hungary, September 2007

[4] D. Downey, M. Broadhead, and O. Etzioni: **Locating complex named entities in web text**. In Proc. of IJCAI, 2007, 2007.

[5] X. Qian: **A CRF-based pos tagger**. Technical Report FDUCSE 07302, Fudan University, 2007.

[6] The OpenNLP project. See http://opennlp.sourceforge.net/.

[7] Toutanova, K., Klein, D., Manning, C.D., Singer, Y.: **Feature-rich part-of-speech tagging with a cyclic dependency network**. In: Proc. of Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology, vol. 1, pp. 173–180

[8] E. Agirre and G. Rigau: **A proposal for word sense disambiguation using conceptual distance**. In the First International Conference on Recent Advances in NLP, Tzigov Chark, Bulgaria, September 1995.

[9] C. Yang and J. C. Hung: **Word sense determination using wordnet and sense co-occurrence**. In Proceedings of the 20th International Conference on Advanced Information Networking and Applications (AINA) - Volume 1, pages 779–784, Washington, DC, USA, 2006. IEEE Computer Society.