# Worth its Weight in Gold or Yet Another Resource

## A Comparative Study of Wiktionary, OpenThesaurus and GermaNet

TECHNISCHE UNIVERSITÄT DARMSTADT

**Christian M. Meyer and Iryna Gurevych**

First Workshop on Automated Knowledge Base Construction (AKBC), Grenoble, France, May 2010.

Ubiquitous Knowledge Processing

# Motivation
## *NLP Tasks and Lexical Semantic Knowledge*



*Applications*

**Expert-built**

*Lexical Semantic Resources*

**GermaNet**

*GermaNet*

*WordNet*

*OpenCyc*

# Motivation
### *Expert-built Lexical Semantic Resources*

**Applications**

YAHOO!® BABEL FISH

Google translate

Semantic Search
on its way!

WORD SENSE
THE FAST-PACED WORD RACE!

**Expert-built**

*Lexical Semantic Resources*

**GermaNet**

*GermaNet*

*WordNet*

*OpenCyc*

✓ used for many years

✓ well studied

✗ high construction cost

✗ limited size

✗ hard to keep up-to-date

Ubiquitous
Knowledge
Processing

# Motivation
*Collaboratively-built Lexical Semantic Resources*



Applications

- ✓ emerging
- ✓ freely available
- ✓ constantly updated
- ✓ competitive to expert-built
- ✗ **structure and content related properties are largely unknown**

**Collaboratively-built**

*Lexical Semantic Resources*

**Wiktionary**
[ˈwɪkʃənrɪ] *n.,*
a wiki-based Open
Content dictionary

# Motivation
## *Collaboratively-built Lexical Semantic Resources*

**Structure and content related properties of collaborative resources are largely unknown**

- **How are the resources organized?**
- **Which kind of semantic knowledge is encoded?**
- **What are their strengths and drawbacks?**

**Collaboratively-built**

*Lexical Semantic Resources*

**Wiktionary**
[ˈwɪkʃənɹɪ] *n.*,
a wiki-based Open
Content dictionary

**openthesaurus**.de

**WIKIPEDIA**
*The Free Encyclopedia*

## Motivation
### *Collaboratively-built Lexical Semantic Resources*

**TECHNISCHE UNIVERSITÄT DARMSTADT**

---

**Structure and content related properties of collaborative resources are largely unknown**

**→ Perform a comparative study of resources**

### Expert-built

*Lexical Semantic Resources*

**GermaNet**

*GermaNet*

*WordNet*

*OpenCyc*

### Collaboratively-built

*Lexical Semantic Resources*

**Wiktionary**
[ˈwɪkʃənrɪ] *n.*,
a wiki-based Open
Content dictionary

**openthesaurus**.de

**WIKIPEDIA**
*The Free Encyclopedia*

---

Ubiquitous Knowledge Processing

# Lexical Semantic Resources
## *Wiktionary*

**Collaboratively created online dictionary**

- Language
- Etymology
- Pronunciation
- Part-of-speech
- Word senses
- Synonyms
- Derived Terms
- Translations
- …

### boat

#### English

Most common English words: due « Henry « society « #797: boat » heaven » v. » difficult

#### Pronunciation

- (*RP*) enPR: bōt, IPA: /bəʊt/, SAMPA: /b@Ut/
- (*GenAm*) enPR: bōt, IPA: /boʊt/, SAMPA: /boUt/
- 🔊 Audio (US)^help, file
- Rhymes: -əʊt

#### Etymology

From Old English *bāt* < Proto-Germanic *baitaz. Related to Old Norse *bátr*, *beit* (Icelandic: *bátur*). Related to German Boot and Dutch boot.

#### Noun

*Word Senses*

boat (*plural* boats)

[1] A craft used for transportation of goods, fishing, racing, recreational cruising, or military use on or in the water, propelled by oars or outboard motor or inboard motor or by wind.

[2] (*poker slang*) A full house.

[3] (*chemistry*) One of two possible conformers of cyclohexane rings (the other being chair), shaped roughly like a boat.

A boat kept on land

#### Synonyms

*Semantic Relations*

[1] craft, ship, vessel

#### Hyponyms

[1] ark, bangca, barge, canoe, catamaran, caravel, carrack, coracle, cruiser, cutter, dhow, dinghy, dory, dragon boat, Dutch barge, East Indiaman, felucca, ferry, ferryboat, fishing boat, folding boat, galley, galleon, gig, go-fast boat, houseboat, hovercraft, hydrofoil, hydroplane, inflatable boat, inflatable raft, jetboat, jetski, junk, kayaaki,

# Lexical Semantic Resources
## *GermaNet and OpenThesaurus*

**GermaNet**

- Semantic Network for the German Language
- Created by lexicographers
- WordNet-like structure
- [Kunze and Lemnitzer, 2002]



```
<con_rel name="hyperonymy" dir="one"  xmlns:xlink="http:/.
  <locator xlink:type="locator" xlink:href="nomen.Pflanze
  <locator xlink:type="locator" xlink:href="nomen.Pflanze
  <arc xlink:type="arc" xlink:from="nPflanze.
</con_rel>
```

**OpenThesaurus**

- Collaborative (but moderated) collection of synonyms
- Used in OpenOffice
- [Naber, 2005]



openthesaurus.de

bank | Synonym finden

Home · Impressum · Login · twitter

**bank – Synonyme bei OpenThesaurus**
- Bank · Sitzbank – [ändern]
- Bank · Bankhaus · Finanzinstitut · Geldhaus · Geldinstitut · Geschäftsbank · Kreditanstalt · Kreditinstitut · Sparkasse – [ändern]

# A Uniform Representation of Resources
*Splitting of Synsets*

*{vessel, watercraft} is hypernym of {boat}*
*{vessel} is hypernym of {tank, storage tank}*

Insert synonymy relations within a synset

vessel 1

boat 1

*{vessel, watercraft}*

watercraft 1

Insert semantic relations between each individual word sense

tank 1

*{tank, storage tank}*

vessel 2

storage tank 1

# Word Sense Disambiguation in Wiktionary
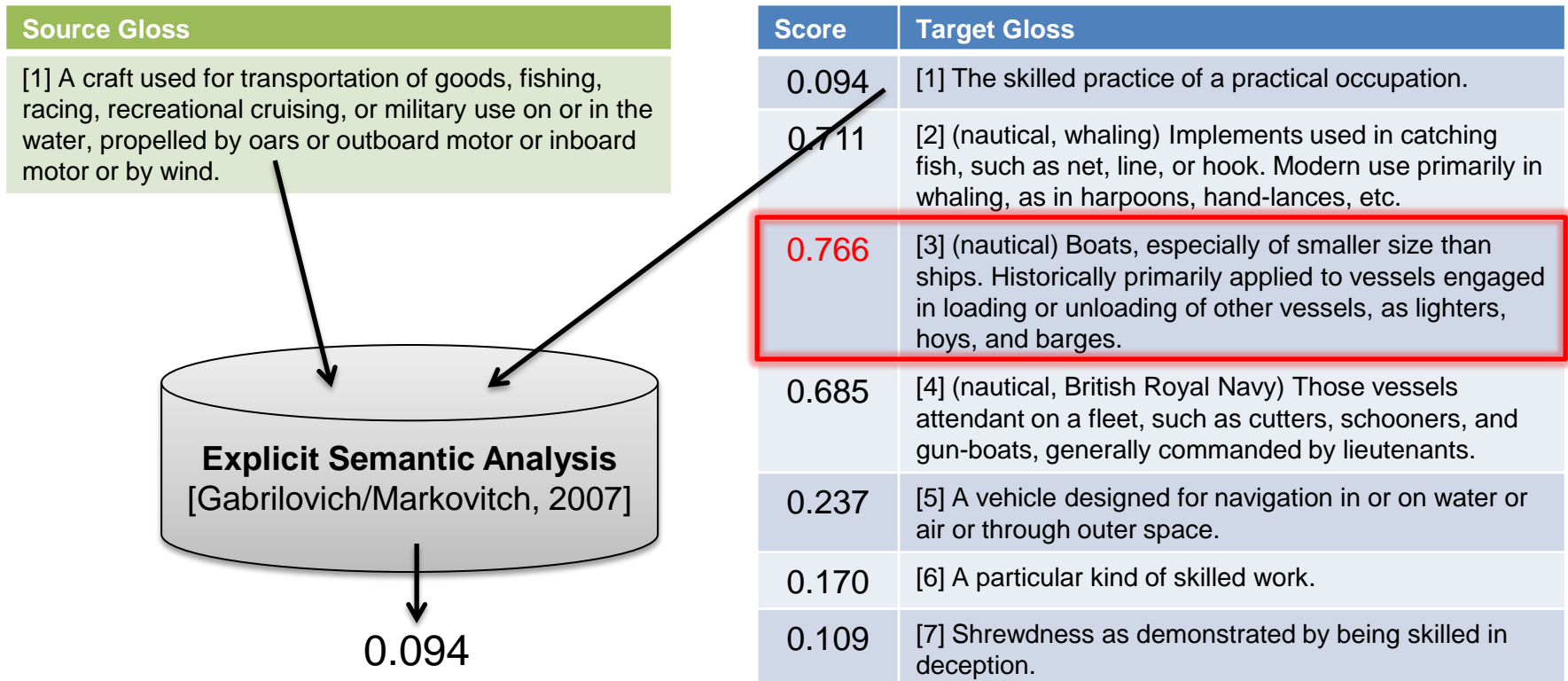## *Finding the Correct Target Word Sense*



Sense [1] encodes a synonymy relation to "craft".
**But which sense?**

# Word Sense Disambiguation in Wiktionary
## *Finding the Correct Target Word Sense – Approach*

- We apply a method based on semantic relatedness here.
- Finding the best method for this task is a subject of our current studies.

| Source Gloss |
|---|
| [1] A craft used for transportation of goods, fishing, racing, recreational cruising, or military use on or in the water, propelled by oars or outboard motor or inboard motor or by wind. |

**Explicit Semantic Analysis**
[Gabrilovich/Markovitch, 2007]

0.094

| Score | Target Gloss |
|---|---|
| 0.094 | [1] The skilled practice of a practical occupation. |
| 0.711 | [2] (nautical, whaling) Implements used in catching fish, such as net, line, or hook. Modern use primarily in whaling, as in harpoons, hand-lances, etc. |
| 0.766 | [3] (nautical) Boats, especially of smaller size than ships. Historically primarily applied to vessels engaged in loading or unloading of other vessels, as lighters, hoys, and barges. |
| 0.685 | [4] (nautical, British Royal Navy) Those vessels attendant on a fleet, such as cutters, schooners, and gun-boats, generally commanded by lieutenants. |
| 0.237 | [5] A vehicle designed for navigation in or on water or air or through outer space. |
| 0.170 | [6] A particular kind of skilled work. |
| 0.109 | [7] Shrewdness as demonstrated by being skilled in deception. |

# Word Sense Disambiguation in Wiktionary
*Finding the Correct Target Word Sense – Evaluation*

- **Upper bound:** 2 human annotators A and B judged 250 randomly sampled relations (= 920 pairs of source and target candidate)

- **Lower bound:** always choose the first word sense (this is usually the most frequent one)

| | Lower Bound | | Our Algorithm | | Upper Bound |
|---|---|---|---|---|---|
| | *0–A* | *0–B* | *M–A* | *M–B* | *A–B* |
| $A_O$ | .791 | .780 | .820 | .791 | .886 |
| $\kappa$ | .498 | .452 | .567 | .480 | .728 |
| $\alpha$ | .679 | .620 | .726 | .649 | .866 |

$A_O$: Percentage of agreement
$\kappa$:  Cohen's Kappa (chance-corrected)
$\alpha$:  Krippendorff's Alpha with set-valued distance function (MASI)

# Structural Analysis
## *Topological Results*

## Analysis:

- Connectivity
- Degree distribution
- Network organization
- Cluster analysis



*log(#nodes)*

*log(degree)*

## Results:

- The **largest connected component** contains the bulk of semantic knowledge
- The Wiktionary graph is **scale-free** and allows to predict analysis results to future (larger) versions
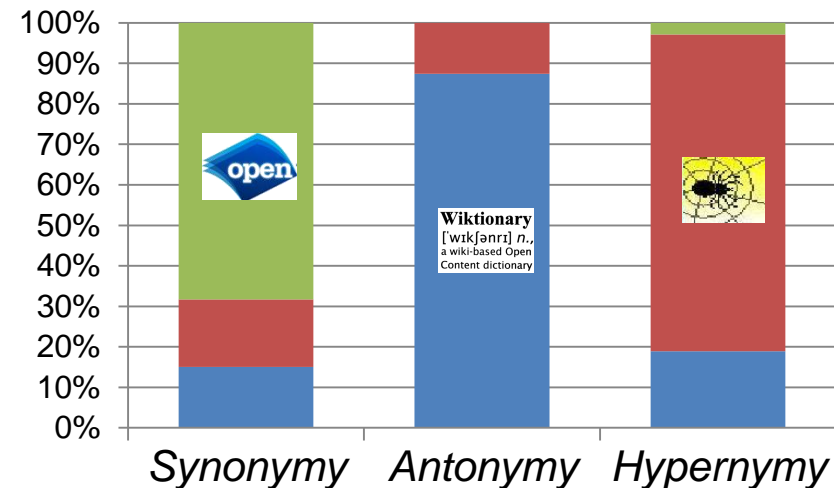- All graphs are **small world graphs**; they show organizational patterns that significantly differ from random graphs

# Content Analysis
## *Encoded Lexemes, Word Senses and Semantic Relations*

## Analysis:

- Resource Size
- Polysemy
- Relation Type
- Unidirectional Relations



## Results:

- Wiktionary has most **lexemes** and **word senses**
- GermaNet has most **semantic relations**
- In general, more **polysemous lexemes** in Wiktionary
- Many **"dangling" articles** in Wiktionary
- Predominant **type of relation** for each resource
- Most Wiktionary relations are **unidirectional**

# Conclusions
## *Take-home Message*

- **How are the resources organized?**
  - Uniform representation needed
  - Largest connected component is sufficient
  - Small world property

- **Which kind of semantic knowledge is encoded?**
  - More polysemous lexemes in Wiktionary
  - Many dangling lexemes

- **What are their strengths and drawbacks?**
  - Predominant type of relation for each resource
  - Number of semantic relations can be increased in Wiktionary

# Conclusions
## *Future Work*

- Study English resources
- Improve word sense disambiguation in Wiktionary
- How large is the information overlap of the resources?
- **Combine the resources at the word sense level**

# Thank you for your attention!

**Ubiquitous Knowledge Processing**



**Additional Online Material:**

`http://www.ukp.tu-darmstadt.de/data/lexical-resources/`

# Thank you for your attention!



## Additional Online Material:
**`http://www.ukp.tu-darmstadt.de/data/lexical-resources/`**

# Kontakt / Contact

**Christian M. Meyer**
Technische Universität Darmstadt
Ubiquitous Knowledge Processing Lab

✉ Hochschulstr. 10, 64289 Darmstadt, Germany
☏ +49 (0)6151 16–7477
🖷 +49 (0)6151 16–5455
✉ meyer (at) ukp.informatik.tu-darmstadt.de