Goals & Motivation
000

Methodology
00000000

Results
0000

# The Polish Cyc lexicon as a bridge between Polish language and the Semantic Web

Aleksander Pohl

Computational Linguistics Departament
Jagiellonian University

18-20, September 2010

Goals & Motivation
○○○

Methodology
○○○○○○○○

Results
○○○○

# Agenda

Goals & Motivation

Methodology

Results

Goals & Motivation
000

Methodology
00000000

Results
0000

# Agenda

## Goals & Motivation

Methodology

Results

Goals & Motivation
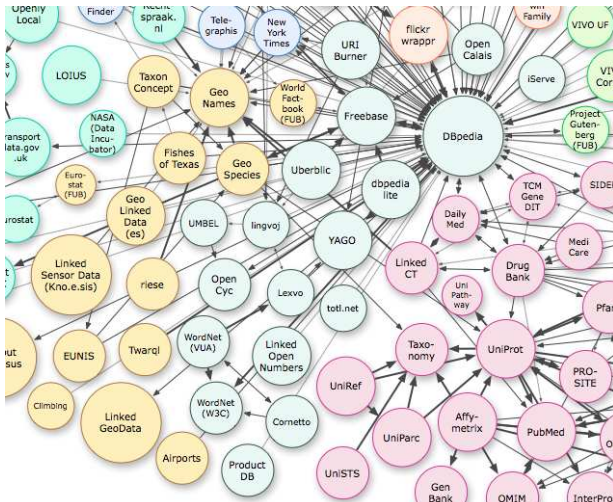●○○

Methodology
○○○○○○○○

Results
○○○○

## Strategic goals

▶ knowledge based **information extraction**:
„Google Inc. is road-testing cars that steer, stop and start
without a human driver, the company says."
```
(#$and
(#$isa #$Test-123 #$PerformenceTesting)
(#$performedBy #$Test-123 #$GoogleInc)
(#$objectOfEvaluation #$Test-123
(#$InstanceFn
#$TransportationDevice-Unmanned)))
```

▶ knowledge based **text generation**:
„Google testuje pojazdy autonomiczne."

Aleksander Pohl

## Tactical goals

- ▶ Polish mappings for Cyc concepts which form semantic restrictions of Cyc relations arguments
    - ▶ (#$arg1Isa #$hasMembers #$Organization)
    - ▶ (#$arg2Isa #$hasMembers #$Agent-Generic)
    - ▶ (#$interArgIsa1-2 #$hasMembers #$SoccerTeam #$SoccerPlayer)
- ▶ Mapping of (Polish) Wikipedia semantic categories to Cyc concepts
    - ▶ polski piłkarz → #$SoccerPlayer, #$PolishPerson
        - ▶ Andrzej Salach
        - ▶ Józef Kałuża
        - ▶ ...772 more
    - ▶ polski klub piłkarski → #$SoccerTeam
        - ▶ Wisła Kraków
        - ▶ Jarota Jarocin
        - ▶ ...210 more

Goals & Motivation
○○●

Methodology
○○○○○○○○

Results
○○○○

## Motivation – **Data + Knowledge**

Goals & Motivation
○○○

Methodology
○○○○○○○○

Results
○○○○

# Agenda

Goals & Motivation

Methodology

Results

Goals & Motivation
000

Methodology
●○○○○○○○

Results
○○○○

## Lexicon creation

Previous results (IIS 2009):

- ▶ overall accuracy: 54%
- ▶ concepts translated: approx. 20K

Current attempt:

- ▶ computer-aided translation
- ▶ transfer-base approach
- ▶ multi-word expressions
- ▶ preservation of full lexical information
- ▶ several thousands of concepts

Aleksander Pohl

Goals & Motivation
000

Methodology
0●000000

Results
0000

## The Algorithm

- ▶ *translate* the English mapping of the concept into Polish (many results might be produced)
- ▶ *map* the words of each translation to the entries of Polish inflectional dictionary
- ▶ *transform* the translations to match syntax constraints
- ▶ *rank* the translations
- ▶ *present* the results to the human operator
- ▶ *store* the selected result in the database
- ▶ *search* for semantic categories extracted from the Polish Wikipedia, corresponding to the translation
- ▶ *merge* or *link* the selected categories with the Cyc concept

# Translation, Mapping & Transformation

- ▶ Translation
  - ▶ English-Polish large transfer dictionary Oxford/PWN
  - ▶ ignoring additional information (gender, categorization, domain, etc.)
  - ▶ each word of a multi-word expression translated separately
- ▶ Mapping to inflectional dictionary
  - ▶ lexemes indexed by base form and inflectional paradigm, e.g. <uzależniający,CAA>
- ▶ Transformation
  1. noun + adjective
  2. noun + noun
  3. noun + verb
  4. noun + other
  5. other

Aleksander Pohl

Goals & Motivation
000

Methodology
0000●0000

Results
0000

# Ranking & Selection

- ► Ranking
  - ► IPI PAN corpus (250 mil. segments), not balanced
  - ► Google considered, but too slow for on-line translation
- ► Selection

Goals & Motivation
○○○

Methodology
○○○○○●○○○

Results
○○○○

# Interpretation

Goals & Motivation
○○○

Methodology
○○○○○●○○

Results
○○○○

# Categories selection

$$R_i = \frac{cm_{i,j}}{cl_i} * \frac{cm_{i,j}}{cl_j} * children_i \tag{1}$$

Goals & Motivation
○○○

Methodology
○○○○○○●○

Results
○○○○

# Verification



▲ 1  ▼ 24  ● 2     **Organism-Whole**                                    organizm

Organizm jest rodzajem materii organicznej.
Organizm jest rodzajem naturalnego obiektu materialnego.
Organizm jest rodzajem struktury.
Organizm jest rodzajem rzeczy wielowymiarowej.
Organizm jest rodzajem żywego obiektu biologicznego.
Drobnoustrój jest rodzajem organizmu.
Zwierzę jest rodzajem organizmu.

Goals & Motivation
○○○

Methodology
○○○○○○○●

Results
○○○○

## Example of a translation

- ► `#$AddictiveSubstance` → **addictive substance**

- ► → `[(uzależniający,`
  `wciągający)`$_1$`,(substancja, istota, ciężar,`
  `waga, podstawa, treść, realność, majątek)`$_2$`]`

- ► → `[(<uzależniać,BDA>, <uzależniać się,BDA>,`
  `<uzależniający,CAA>, ...)`$_1$`,`
  `(<substancja,ADACBAA>, <istota,ADAAA>,`
  `<ciężar,ACAAAAA>, <Ciężar,AAAAD>,...)`$_2$`]`

- ► → `[uzależniająca substancja, uzależniająca`
  `się substancja, wciągająca substancja,`
  `wciągająca się substancja,...]`

- ► → „substancja uzależniająca"

Goals & Motivation
○○○

Methodology
○○○○○○○○

Results
○○○○

# Agenda

Goals & Motivation

Methodology

Results

Goals & Motivation
ooo

Methodology
oooooooo

Results
●ooo

# Results

- ▶ Precision: 37% (560 translations)
- ▶ Baseline: 19% (Google Translate)
- ▶ Recall: 88%
- ▶ Inter-translator agreement: 56%
- ▶ Concepts translated : **1128** (target approx. 4000)
- ▶ Concepts mapped : 493
- ▶ Total number of concepts covered: **218942**

Goals & Motivation
○○○

Methodology
○○○○○○○○

Results
○●○○

Thank You!

Goals & Motivation
000

Methodology
00000000

Results
00●0

## Why not WordNet?

WordNet and OpenCyc contents is overlapping:

- ▶ dog *direct hypernym* canine
- ▶ (#$genls #$Dog #$CanineAnimal)

but:

- ▶ Cyc relations have formal semantics
- ▶ CycL expressiveness is higher (rules, functions, microtheories, arbitrary arity relations):
  - ▶ (#$relationAllExists #$bodyPartsUsed #$AnimalWalkingProcess #$Leg)
  - ▶ (#$relationAllExists #$properPhysicalParts #$CanineAnimal #$Leg)
- ▶ Cyc is shipped with sophisticated *inferencing engine*

Aleksander Pohl

Goals & Motivation
000

Methodology
00000000

Results
00●0

## Why not WordNet?

WordNet and OpenCyc contents is overlapping:

- ▶ dog *direct hypernym* canine
- ▶ (#$genls #$Dog #$CanineAnimal)

but:

- ▶ Cyc relations have formal semantics
- ▶ CycL expressiveness is higher (rules, functions, microtheories, arbitrary arity relations):
  - ▶ (#$relationAllExists #$bodyPartsUsed #$AnimalWalkingProcess #$Leg)
  - ▶ (#$relationAllExists #$properPhysicalParts #$CanineAnimal #$Leg)
- ▶ Cyc is shipped with sophisticated *inferencing engine*

Aleksander Pohl

Goals & Motivation
○○○

Methodology
○○○○○○○○

Results
○○○●

# Why not DBpedia (YAGO, SUMO)?

- ▶ DBpedia
  - ▶ 259 classes
  - ▶ 1200 relations
  - ▶ no rules
  - ▶ SPARQL end-point
  - ▶ no inflectional data for Polish labels
- ▶ Cyc is the biggest ontology
  - ▶ 500K symbols (70K collections)
  - ▶ 17K (26K) relations
  - ▶ 5M assertions (mostly *defining* the terms)
  - ▶ rules
  - ▶ sophisticated inferencing engine

Goals & Motivation
○○○

Methodology
○○○○○○○○

Results
○○○●

# Why not DBpedia (YAGO, SUMO)?

- ▶ DBpedia
    - ▶ 259 classes
    - ▶ 1200 relations
    - ▶ no rules
    - ▶ SPARQL end-point
    - ▶ no inflectional data for Polish labels
- ▶ Cyc is the biggest ontology
    - ▶ 500K symbols (70K collections)
    - ▶ 17K (26K) relations
    - ▶ 5M assertions (mostly *defining* the terms)
    - ▶ rules
    - ▶ sophisticated inferencing engine