

# Generating sets of synonyms between languages

Ján Genčí<sup>1</sup>, Ondrej Dzurjov<sup>1</sup>,  
Radovan Garabík<sup>2</sup>

<sup>1</sup>Department of Computers and Informatics,  
Technical University of Košice

<sup>2</sup>Ľ. Štúr Institute of Linguistics, Slovak Academy of  
Sciences, Bratislava

# Introduction

- Besides classical bi-/multi-lingual dictionaries on-line, there are EuroWordNet or Global WordNet exist
- Nor EuroWordNet neither Global WordNet cover the Slovak language
- There were several works dealing with automatic creation of synsets in the past
- A requirement to create computer application for building Slovak synsets

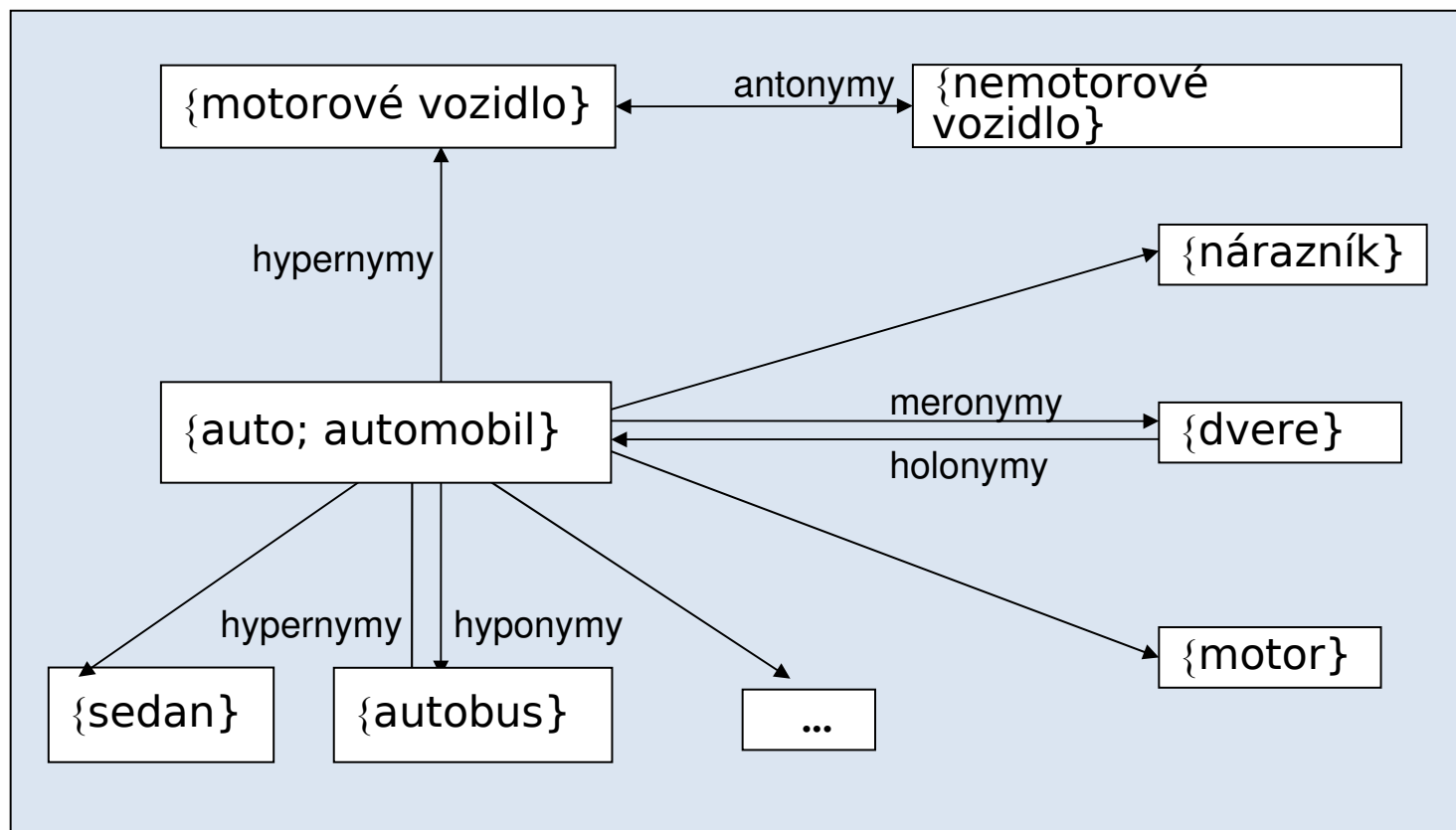
# Our goals

- Design new methods of synsets generation
- Evaluate the quality of synset generation
- Design and implementation of computer application for assisted synset editing

# WordNet

- Project WordNet started in 1990 at Princeton University
- Lexical database of English language; nouns, pronouns, verbs and adverbs
- Two base properties of WordNet
  - Words are organized into groups – sets of synonyms - synsets
  - Synsets are interconnected

# Relationships in WordNet



# Method A

## □ Synonymy among words in synsets

$P_{S1}^{J_1}, P_{S2}^{J_1}, P_{S3}^{J_1}, \dots, P_{Sn}^{J_1}$  — Set of words from synset

$t_{J_1 \rightarrow J_2}(P_{Si}^{J_1}) = \{P_{S1}^{J_2}, P_{S2}^{J_2}, \dots, P_{Sp}^{J_2}\}, i = 1, \dots, n$  — translation

$S := S \cup (t_{J_1 \rightarrow J_2}(P_{Si}^{J_1}) \cap t_{J_1 \rightarrow J_2}(P_{Sj}^{J_1}))$  — result  
*pre*  $i, j = 1, \dots, p; i \neq j$

{kind; sort; form; variety} ->

kind -> druh, rod, kategória

sort -> druh, akosť, trieda, typ, forma, chlap

form -> forma, tvar, podoba, formulár, blanketa,  
formula

variety -> rozmanitosť, odroda, výber, druh, rad,  
množstvo, mnohotvárnosť, rôznosť

∩

{druh, forma}

# Method B

## □ Translation of univocal words

$P_{S1}^{J_1}, P_{S2}^{J_1}, P_{S3}^{J_1}, \dots, P_{Sn}^{J_1}$  - Set of words from synset

$P'_{S1}^{J_1}, P'_{S2}^{J_1}, \dots, P'_{Sm}^{J_1}$  - Subset of univocal words

$t_{J_1 \rightarrow J_2}(P'_{Si}^{J_1}) = \{P'_{S1}^{J_2}, P'_{S2}^{J_2}, \dots, P'_{Sk}^{J_2}\}, i = 1, \dots, m$  - translation

$S := S \cup t_{J_1 \rightarrow J_2}(P'_{Si}^{J_1}), i = 1, \dots, m$  - result



kind – 1 sense -> druh, rod, kategória

sort – 4 senses

form – 16 senses

variety – 6 senses

{druh, rod, kategória}

# Method C

- It uses hyponymy and hyperonymy relationship for synset generation

$P_{S1}^{J_1}, P_{S2}^{J_1}, P_{S3}^{J_1}, \dots, P_{Sn}^{J_1}$  - Words from synset

$P_{H1}^{J_1}, P_{H2}^{J_1}, P_{H3}^{J_1}, \dots, P_{Hm}^{J_1}$  - Set of hyponym and hyperonym related to source synset

$t_{J_1 \rightarrow J_2}(P_{Si}^{J_1}) = \{P_{S1}^{J_2}, P_{S2}^{J_2}, \dots, P_{Sp}^{J_2}\}, i = 1, \dots, n$  - translation

$t_{J_1 \rightarrow J_2}(P_{Hi}^{J_1}) = \{P_{H1}^{J_2}, P_{H2}^{J_2}, \dots, P_{Hq}^{J_2}\}, i = 1, \dots, m$

$S := t_{J_1 \rightarrow J_2}(P_{Si}^{J_1}) \cap t_{J_1 \rightarrow J_2}(P_{Hj}^{J_1}) = \{P_{S1}^{J_2}, P_{S2}^{J_2}, \dots, P_{Sp}^{J_2}\} \cap \{P_{H1}^{J_2}, P_{H2}^{J_2}, \dots, P_{Hq}^{J_2}\}$

*pre*  $i = 1, \dots, n$  a  $j = 1, \dots, m$  - result

{kind; sort; form; variety} -> {hyper,hypo}nyms: {category, type, brand, genus, species}

{kind; sort; form; variety} -> druh, rod, kategória, akosť, trieda, typ, forma, chlap, tvar, podoba, formulár, blanketá, formula, rozmanitosť, odroda, výber, rad, množstvo, mnohotvárnosť, rôznosť

{category, type, brand, genus, species} -> kategória, skupina, trieda, typ, symbol, litera, druh, odroda, značka, označenie, známka, kvalita, akosť, ohorok, rod, forma, tvar

∩

{druh; rod; kategória; akosť; trieda; typ; forma; tvar; odroda}

# Method D

- It uses hyponymy and hyperonymy relationship for synset generation as source synsets

$P_{H1}^{J_1}, P_{H2}^{J_1}, P_{H3}^{J_1}, \dots, P_{Hk}^{J_1}$  - source synset hypernyms

$P_{H'1}^{J_1}, P_{H'2}^{J_1}, P_{H'3}^{J_1}, \dots, P_{H'l}^{J_1}$  - source synset hyponyms

$t_{J_1 \rightarrow J_2}(P_{Hi}^{J_1}) = \{P_{H1}^{J_2}, P_{H2}^{J_2}, \dots, P_{Hq}^{J_2}\}, i = 1, \dots, k$  - translation

$t_{J_1 \rightarrow J_2}(P_{H'i}^{J_1}) = \{P_{H'1}^{J_2}, P_{H'2}^{J_2}, \dots, P_{H'r}^{J_2}\}, i = 1, \dots, l$

$S := t_{J_1 \rightarrow J_2}(P_{Hi}^{J_1}) \cap t_{J_1 \rightarrow J_2}(P_{H'j}^{J_1}) = \{P_{H1}^{J_2}, P_{H2}^{J_2}, \dots, P_{Hq}^{J_2}\} \cap \{P_{H'1}^{J_2}, P_{H'2}^{J_2}, \dots, P_{H'r}^{J_2}\}$

pre  $i = 1, \dots, k$  a  $j = 1, \dots, l$  - target synset

{kind; sort; form; variety}:

hypernym: {category} -> kategória, skupina, trieda

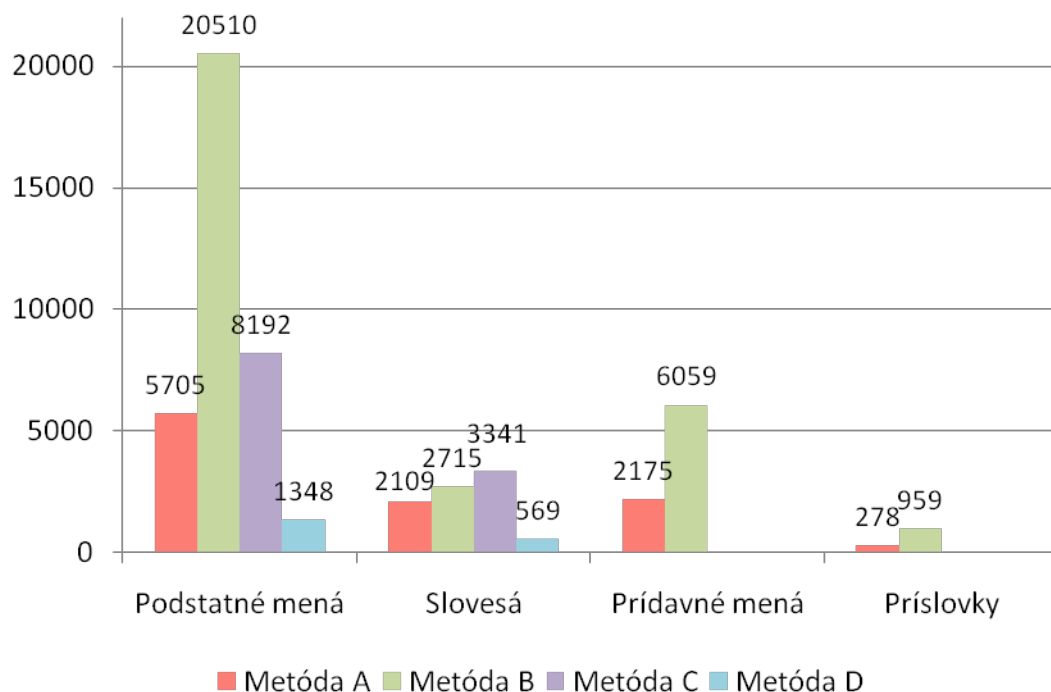
hyponym: {type, brand, genus, species} -> typ,  
symbol, litera, druh, odroda, značka, označenie,  
známka, kvalita, akosť, ohorok, druh, rod, skupina,  
trieda, forma, tvar

∩

{skupina, trieda}

# Results

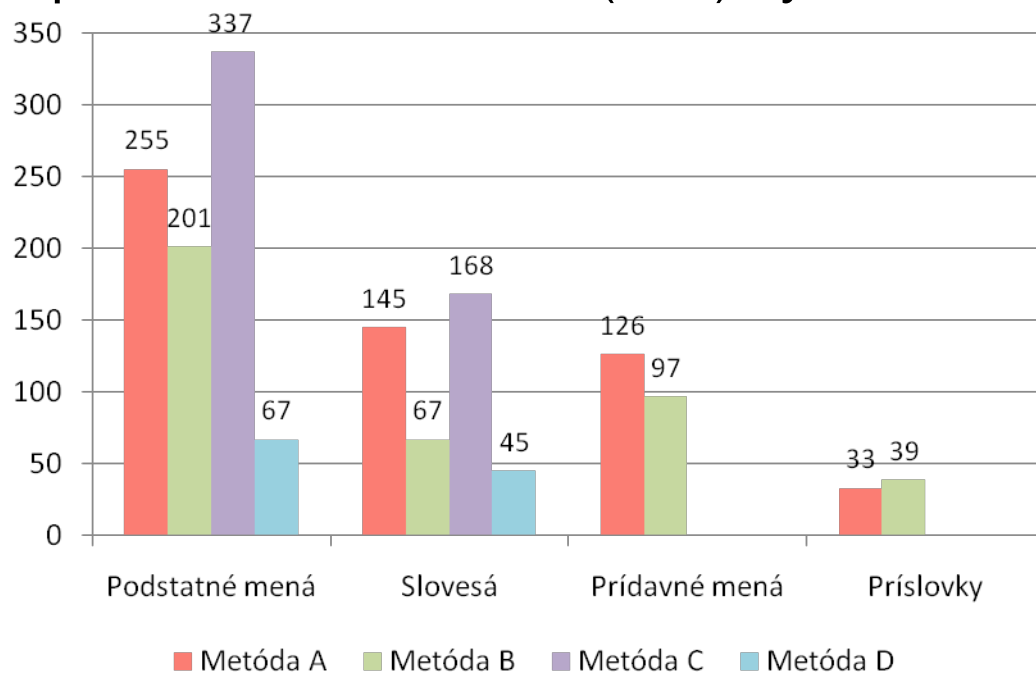
The WordNet contains 117659 synsets, Slovak equivalents were generated for 34.4 % (40521) of English synsets



Method A: 8.7 %  
Method B: 25.7 %  
Method C: 12 %  
Method D: 1.4 %

# Results for frequently used words

For set of 300 more frequently used English words there are 1709 synsets in WordNet, Slovak equivalents were produced for 55.4 % (946) synsets



Method A: 32.7 %

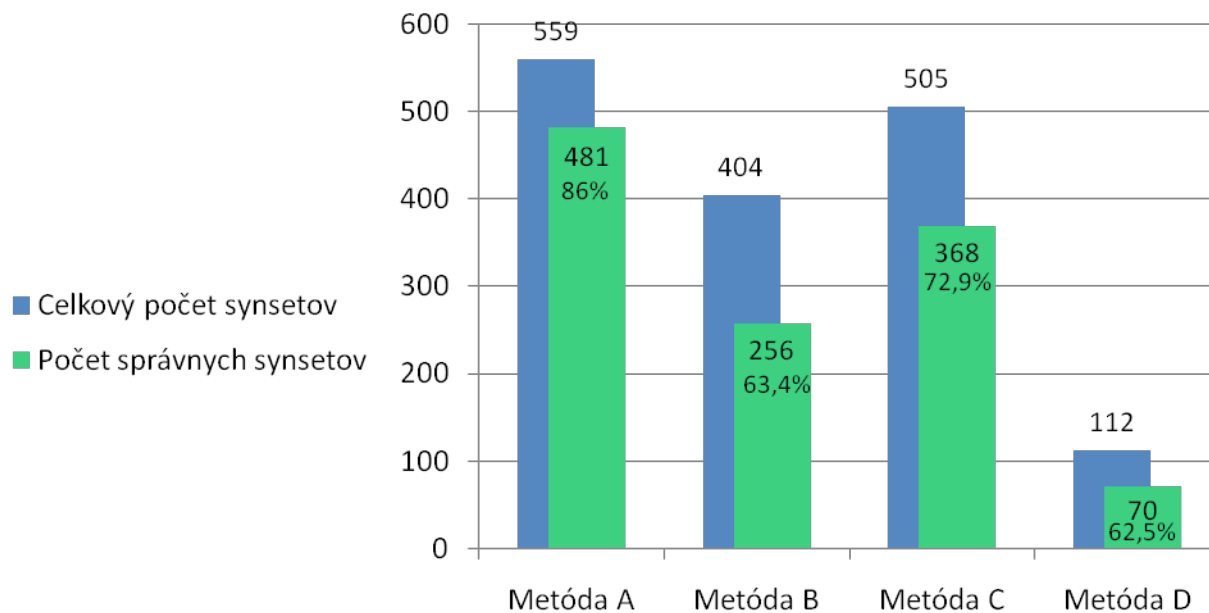
Method B: 23.6 %

Method C: 29.6 %

Method D: 6.6 %

# Inspection of Slovak synsets

Part of speech and semantics of generated Slovak synsets were inspected.





# Results

- New computer application for manual synset translation. Application was “fed” by generated data
- About 60% of synsets were generated
- About 75% of generated synsets were evaluated as “not incorrect“

# Conclusion

- We were able to produce Slovak synsets for 34.4% of English synsets. Manual inspection was required.
- Process of synset generation is influenced by many factors – translation dictionaries used, frequency of words, number of words in synset.
- Application – used to create en-sk-pl-de-it dictionary (with ontology relations)

Filter

Identifikátor anglického synsetu:  (identifikátor synsetu je zložený z 8 číslic)Prechádzať:  ▾Zobraziť prislúchajúce  slovenské,  nemecké,  poľské,  litovské synsety

## Zdrojový anglický synset

Anglické synsety	Slovenské synsety	Nemecké synsety	Poľské synsety	Litovské synsety
-EN synset: <a href="#">06410904</a> { <a href="#">book</a> } ----- a written work or composition that has been published (printed on pages bound together); "I am reading a good book on economics" ----- <a href="#">+SYN</a>   <a href="#">-SYN</a>	<input checked="" type="checkbox"/> SK synset: 3538 { kniha; knižka } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input checked="" type="checkbox"/>   (i)	-GER synset: 1737 { s Buch; s Büchlein } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input type="checkbox"/>   (i)	-PL synset: 1686 { książka; książeczka } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input type="checkbox"/>   (i)	-LAT synset: 1404 { knyga; knygelė } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input type="checkbox"/>   (i)

Antonym - počet anglických synsetov: 0

Hypernym - počet anglických synsetov: 1

Anglické synsety	Slovenské synsety	Nemecké synsety	Poľské synsety	Litovské synsety
-EN synset: <a href="#">06589574</a> { <a href="#">publication</a> } ----- a copy of a printed work offered for distribution ----- <a href="#">+SYN</a>   <a href="#">-SYN</a>	<input checked="" type="checkbox"/> SK synset: 835 { publikácia } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input checked="" type="checkbox"/>   (i)	-GER synset: 1795 { e Publikation; e Veröffentlichung } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input type="checkbox"/>   (i)	-PL synset: 1744 { publikacja } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input type="checkbox"/>   (i)	-LAT synset: 4510 { leidinys } ----- \   <a href="#">+EN</a>   <a href="#">-EN</a>   <a href="#">×</a>   <input type="checkbox"/>   (i)

<b>Part of speech</b>	<b>Nouns</b>	<b>Adjectives</b>	<b>Verbs</b>	<b>Adverbs</b>	<b>Totals</b>
<b>Unique strings</b>	12941	3321	1150	982	18394
<b>Strings with one sense</b>	10239	2305	953	702	14199
<b>Word-sense pairs</b>	18740	5551	1400	1505	27196
<b>Synsets</b>	9317	2329	830	549	13025
<b>Synsets of one word</b>	3916	773	426	141	5256

Thank you for your attention