

# Towards the Annotation of Named Entities in the National Corpus of Polish

Agata Savary<sup>\*†</sup>  
Jakub Waszczuk<sup>◇†</sup>  
Adam Przepiórkowski<sup>†◇</sup>

\*Université François Rabelais Tours

†Institute of Computer Science, Polish Academy of Sciences

◇University of Warsaw, Poland

LREC'10, May 19-21, 2010

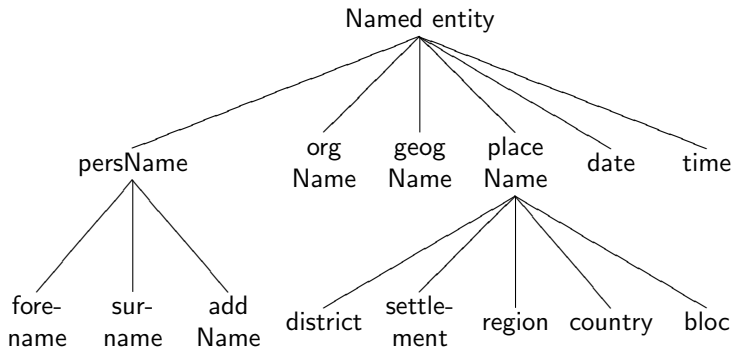
## The project

- consortium: creators of big annotated corpora of Polish
- financed by the Polish Ministry of Science and Higher Education
- period: 2007-2010

## The aim: large national corpus of Polish

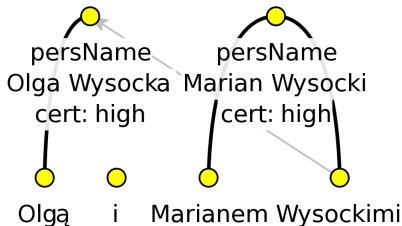
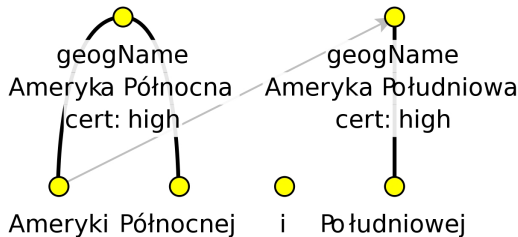
- 1 million words manually annotated, 1 billion words automatically annotated (*Przepiórkowski et al. LREC'2010*)
- representative
- balanced wrt. different genres (*Przepiórkowski et al. 2009*)
- associated to linguistic annotation tools

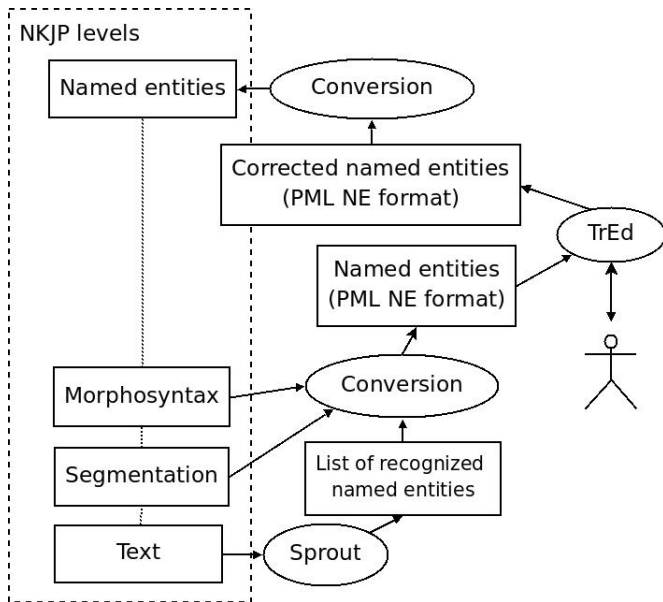
- stand-off
- TEI P5-conformant (*Przepiórkowski & Bański 2009*)
- multi-level
  - \* segmentation
  - \* morphosyntax (*Przepiórkowski & Murzynowski 2009*)
  - \* syntactic words (e.g. bał się)
  - \* syntactic groups (*Głowińska and Przepiórkowski 2010*)
  - \* **named entities**
  - \* word senses
- quality-ensured (double annotation + super-annotation)
- see other presentations in LREC'10: **W4** i **W20**



- vertical hierarchy of **related names**
  - relational adjectives *poznański, ONZ-owski*
  - names of inhabitants and members *poznanianak, Grek, AK-owiec*
- not annotated: quantities, products, periods, events, titles, ...

- Gramatically motivated lemma (*Piskorski et al. 2009*)  
*Stanów Zjednoczonych* → *Stany Zjednoczone*
- Semantically motivated derivation base  
*amerykański* → *Stany Zjednoczone*
- Embedded names annotated (*Galicia-Haro and Gelbukh 2009; Finkel and Manning 2009; Kravalová and Žabokrtský 2009*)  
[[*Tadeusz*]<sub>forename</sub>[*Kościuszko*]<sub>surname</sub>]<sub>persName</sub>
- Discontinuous names  
*Wydział Matematyczny ówczesnej Akademii Krakowskiej*
- Coordinations separated (*Mazur and Dale 2009*)  
*Ameryka Północna i Południowa*





## NLP platform

- fast gazetteer lookup (*Budisak et al. 2009*)
- cascaded unification-based FST grammar parser
- output: feature-structures with user-defined types
- previous NE grammar for Polish (*Piskorski 2005*)

## NKJP resources and rules (*Savary & Piskorski 2010*)

- gazetteer with 300,000 inflected forms
- 120 grammar rules
- precision: 0.88
- recall: 0.61



## Tree Editor (*Pajas & Štěpánek 2008*)

- manipulates tree structures (necessary for embedded, coordinated & discontinuous NEs)
- interoperable
- allows for stand-off multi-level annotations
- PML - open customizable XML abstract data format
- customizable GUI
- easy comparing two annotation
- reliable and well documented

TrEd ver. 1.4295 Default(1/1): /home/agata/Recherche/Moje Publikacje/LREC 20

File Node Tree View Macros Setup Help Mode: NKJP\_names

Współpracował z Radio France Nationale i Rozgłośnią Polską Radia Wolna Europa . 3/99

orgName  
Radio France Nationale  
cert: high

placeName->country  
France  
cert: medium

orgName  
Rozgłośnią Polska  
cert: high

placeName->country  
polski (relAdj, Polska)  
cert: high

orgName  
Radio Wolna Europa  
cert: high

geogName  
Europa  
cert: high

Współpracował z Radio France Nationale i Rozgłośnią Polską Radia Wolna Europa .

Scale: 100%

## What has been done

- State-of-the-art corpus methodology
- Advanced annotation strategies
- Annotator's and super-annotator's platform
- 18,000 annotated sentences (out of 75,000) until mid-May

## To do

- 75% corpus to be annotated
- super-annotation
- machine-learning for 1 billion corpus

## Further perspectives

- extending annotation to new categories (periods, events, . . .)
- corpus studies