

# Semantic relation recognition within Polish noun phrase: A rule-based approach

**Paweł Kędzia**

Institute of Informatics  
Wrocław University of Technology  
pawel.kedzia@pwr.wroc.pl

**Marek Maziarz**

Institute of Informatics  
Wrocław University of Technology  
marek.maziarz@pwr.wroc.pl

## Abstract

The paper<sup>1</sup> presents a rule-based approach to semantic relation recognition within the Polish noun phrase. A set of semantic relations, including some thematic relations, has been determined for the need of experiments. The method consists in two steps: first the system recognizes word pairs and triples, and then it classifies the relations. Evaluation was performed on random samples from two balanced Polish corpora.

## 1 Introduction

Semantic relation recognition is a well-known task in natural language processing. Although the relation recognition within noun phrase and between nominals was studied intensely, the task is still challenge for semantic analysis of Polish. We are aware of few papers and projects dealing with Semantic Role Labelling between predicates and their arguments, cf. (Gołuchowski and Przepiórkowski, 2012) or (Lun, 2009), but of none concerning semantic relation recognition inside Polish noun phrase.

## 2 Related work

In (Nastase et al., 2006) authors classify semantic relations between a head and a modifier of a noun phrase. Number of all relation types was equal to 30. These relations were grouped into 5 more general groups. The authors experimented with decision trees, instance-based learning and Support Vector Machines. For each relation they learnt the binary classifier; as the baseline for F-measure they used the model with all of examples classified as positive and recall being equal to 100%. With

<sup>1</sup>Work financed by The National Centre for Research and Development project SP/I/1/77065/10.

regard to the semantic relation the baseline ranged between 17.78% and 60.35%.

Identifying the semantic relations inside compound nouns was presented in (Uchiyama et al., 2008). The authors used SVM classifier and in the best configuration of features, they achieved accuracy of about 84%.

In (Rosario and Hearst, 2001) authors used neural networks to determine 20 semantic relations – similarly to (Nastase et al., 2006) – between a head and a modifier of noun phrase. They used a domain-specific lexical hierarchy of medicine. The authors achieved accuracy of about 60%.

The workshop SemEval-2010 (task 8) concerned the recognition of semantic relations between nominals. In (Tratz and Hovy, 2010) the authors developed a system based on the Maximum Entropy classifier, able to detect 10 bidirectional semantic relations Achieved F-measures depended on the system configuration and lay between 66, 68% and 77, 75%. The same set of semantic relations was used in (Rink and Harabagiu, 2010). The authors used Support Vector Machines classifier and a very rich set of features (i.e., part of speech for all constituents of a semantic relation pair, number of words between the nominals, features based on paths in the dependency tree from Stanford dependency parser). F-measure of this approach was 82.19%.

Authors in (Tymoshenko and Giuliano, 2010) used shallow syntactic parsing and semantic information from ResearchCyc (Lenat, 1995) in the same task of recognizing semantic relations. They used liner combination of kernels (semantic and syntactic) using Support Vector Machines classifier. For the best combination of kernels, they obtained F-measure equal to 77.62%.

There are some works, where rule-based approaches were used. In (Huang, 2009) there has been proposed an approach for automatic construction of rules identifying ten types of seman-

tic relations, using five types of input informations. The relation instances were extracted from Modern Chinese Standard Dictionary. The authors achieved very high precision (range from 0,81 to 0.99), but recall was low - about 0,2. In (Hearst, 1992) authors used set of manually written rules for identification of hyperonymy relations. (Ben Abacha and Zweigenbaum, 2011) used linguistic patterns (built semi-automatically from corpora) to identify semantic relations in medical texts. In this domain-specific task they achieved 75.72% precision and 60,46% recall.

### 3 Recognized semantic relation types

We seek for semantic relations within nominal phrases. The relation set consists of 12 semantic relations, of which 5 are thematic (semantic) roles<sup>2</sup>. Definitions of our semantic relations are based on works of (Kearns, 2011), (Palmer et al., 2010), (Van Valin, 2004), (Larson, 1996), (Dowty, 1991), (Jędrzejko, 1993), (Laskowski and Wróbel, 1997). We tried to select relations that are very frequent or frequent in Polish texts.<sup>3</sup> The relation set is following (thematic roles are marked with *theta*, other relations – with *rho*):

**Proto-Agent**<sub>*θ*</sub> – it is an instigator of an action or an entity that is in a particular state, it may undergo change of state not caused by another participant; for predicates denoting relations – it is the first element of the relation: (człowiek) wykształcony przez Jana<sub>*θ*</sub> ‘(man) educated by John<sub>*θ*</sub>’, wyjący wilk<sub>*θ*</sub> ‘howling wolf<sub>*θ*</sub>’. The Proto-Agent macrorole covers subroles of Agent, Causer and non-agentive non-causative Actor (cf. *Actor* macrorole in (Kearns, 2011)).

**Proto-Patient**<sub>*θ*</sub> is the second macrorole – it is an entity undergoing action, event or change of state caused by another participant; for predicates denoting relations – it is the second element of a given relation: wykształcenie kogoś<sub>*θ*</sub> ‘educating someone<sub>*θ*</sub>’, (Jan) posiadający majątek<sub>*θ*</sub> ‘(John) possessing an estate<sub>*θ*</sub>’. According to (Dowty, 1991)

<sup>2</sup>In Polish, as in other Indo-European languages, verbs could be nominalized during a process of syntactic transformation (Jędrzejko, 1993), (Kolln, 1990). Such nominalized predicates could be linked with nouns by thematic relations.

<sup>3</sup>Rationale for selection of the presented semantic relation types was their frequencies in a four-text sample taken from a Polish corpus *KPW*. Together chosen relations account for ca 80% of all semantic relation occurrences in these texts. Most of our relation types could be found on the list of the most frequent relation types in the English noun phrase (Moldovan et al., 2004, Tab. 1).

many thematic roles come down to the macroroles of Proto-Patient and Proto-Agent.

**Instrument**<sub>*θ*</sub> is a tool, a device or means used by someone in order to cause something, it is sometimes regarded as a secondary cause of situation or change of state: przesyty włócznią ‘speared with a spear’, lina<sub>*θ*</sub> cumownicza<sub>*adjective*</sub> ‘a hawser, lit. mooring rope<sub>*θ*</sub>’.

**Material**<sub>*θ*</sub> is an entity that is used by someone to produce something from it, material undergoes change of state resulting in its disappearance and emerging of a result: zrobiony z miedzi<sub>*θ*</sub> ‘made out of brass<sub>*θ*</sub>’, miedziana<sub>*θ*</sub> figurka ‘brass<sub>*θ*</sub> statuette’.

**Purpose**<sub>*θ*</sub> – an entity or a situation toward which the event is directed or an individual which benefits from the event (purpose combines goal, beneficiary and recipient roles): wręczenie (medali) olimpijczykom ‘giving (medals) to Olympians<sub>*θ*</sub>, sala koncertowa<sub>*θ*</sub> ‘a concert<sub>*θ*</sub> hall’.

**Location** is a physical place at which a given event is localised, a place being destination of an event, a path or a source of motion, or simply a place at which a particular individual is situated: wręczenie (medali) w auli<sub>*θ*</sub> ‘giving (medals) at the lecture theatre<sub>*θ*</sub>’, przedzieranie się przez moczary<sub>*θ*</sub> ‘struggling through the swamp<sub>*θ*</sub>’.

**Time** is a particular moment or a duration of an event – it localises a situation within the flow of events or gives its duration: przedzieranie się przez godzinę/w środę<sub>*θ*</sub> ‘struggling for an hour<sub>*θ*</sub>/on Wednesday<sub>*θ*</sub>’.

**Temporal/spatial meronymy** – these relations point onto a spatial or temporal part of a place/location/time/period): poniedziałkowy poranek<sub>*θ*</sub> ‘Monday morning<sub>*θ*</sub>’, środek<sub>*θ*</sub> zimy ‘middle<sub>*θ*</sub> of the winter’, koniec<sub>*θ*</sub> drogi ‘end<sub>*θ*</sub> of the road’, stolica<sub>*θ*</sub> kraju ‘capital<sub>*θ*</sub> of the country’.

**Attribute** is a property of an individual or an event, such as colour, size, weight, intensity, duration etc., which might be expressed with a qualitative adjective: czerwony<sub>*θ*</sub> samochód ‘red<sub>*θ*</sub> car’, głośna<sub>*θ*</sub> muzyka ‘loud<sub>*θ*</sub> music’.

**Family (member)** is a relative or an in-law to someone, the relation is bidirectional and reflexive: syn<sub>*θ*</sub> króla<sub>*θ*</sub> ‘king’s<sub>*θ*</sub> son<sub>*θ*</sub>’, moja<sub>*θ*</sub> żona<sub>*θ*</sub> ‘my wife<sub>*θ*</sub>’ (I am a relative to my wife).

**Order** gives a position of an entity or an event in an ordered sequence/chain: druga<sub>*θ*</sub> odpowiedź ‘2nd answer’, lata 80<sub>*θ*</sub>. ‘eighties, lit. eightieth<sub>*θ*</sub> years’.

**Quantity** is an amount of something or a cardinality of a given set: pięciu<sub>*θ*</sub> panów ‘five<sub>*θ*</sub> men’,

kieliszek<sub>Q</sub> wina ‘glass<sub>Q</sub> of wine’.

## 4 Semantic relation recognition rule-based algorithm

Our rule-based system proceeds in two steps<sup>4</sup>: first it recognizes word pairs and triples, then operators classifying relations enter.

### 4.1 Recognizing word pairs and triples

Since we consider relations within noun phrases, we must identify them correctly. We made use of a CRF shallow parser (Radziszewski and Pawlaczek, 2012) trained on an annotated corpus of Polish (KPWr) (Broda et al., 2012) which comprises shallow syntactic annotation level (Radziszewski et al., 2012).

KPWr contains 326 annotated text samples representing different genres and styles: blogs, press articles, official and legal texts and Polish Wikipedia articles, it comprises 106358 annotations (phrases and phrase heads, and predicate-argument relations).

Noun and preposition phrases (NPs/PPs) from the corpus correspond to arguments of predicate-argument structure. Each such NP/PP consists of one or several smaller phrases based on agreement (*AgPs*, for details, please look at cited works). Here is an example NP from the corpus (a head of the phrase is boldfaced, *AgP* heads are underlined):

[[**samolot** wyprodukowany]<sub>AgP</sub> [przez PZL]<sub>AgP</sub> [w roku 1938]<sub>AgP</sub> [w Łodzi]<sub>NP</sub>

‘aircraft made by PZL in (year) 1938 in Łódź (city)’

There is no reliable deep parser for Polish (Gołuchowski and Przepiórkowski, 2012), thus we decided to construct a simple rule-based algorithm for deepened shallow parsing of Polish NPs/PPs. The algorithm works on tagged texts – we used (Radziszewski, 2013) tagger. Parsing rules make use of an output from the CRF shallow parser (Radziszewski and Pawlaczek, 2012), in particular: borders of whole NPs/PPs, and of their constituents (i.e., phrases based on agreement, *AgPs*). Found pairs and triples are directly connected within a syntactic structure.

Hand-written rules act like a partial dependency parser. The pairs consist of one subordinate and

<sup>4</sup>Similarly to system presented in (Gamallo et al., 2002).

one superordinate token, the triples comprise one superordinate token and a subordinate preposition phrase (preposition + governed nominal head of a subordinate noun phrase).

The whole algorithm runs in a main loop which iterates *AgP<sub>i</sub>* heads. We start from the first *AgP<sub>0</sub>* head to the left, then we proceed to the right, jumping from *AgP<sub>i</sub>* head to the closest *AgP<sub>i+1</sub>* head to the right. For every *AgP<sub>i</sub>* head we run a cascade-like chain of rules (numbered from 1 to 7) for genitives, nominatives, small preposition phrases (being a part of larger NPs or PPs), coordination, other known to the tagger tokens, other unknown to the tagger tokens and for modifiers. The algorithm in pseudocode was shown in Algorithm 1

The algorithm gives following description for just analysed phrase, “R + number” denotes the number of a rule in the Algorithm 1 activated on the word pair or triple (for instance, *R3* means that the rule number 3 was activated): *R7*: samolot ← wyprodukowany ‘plane made’, *R3*: wyprodukowany ← przez PZL ‘by PZL’, *R3*: wyprodukowany ← w roku ‘in year’, *R3*: wyprodukowany ← w Łodzi ‘in Łódź’.

Such simple shallow parsing algorithm operates quite well on an annotated part of KPWr with F-measure equal to 84%, P = 88%, R = 80%.<sup>5</sup>

### 4.2 Applying WCCL operators

Having identified pairs and triples we run on them operators written in a constraint language WCCL (Radziszewski et al., 2011). The operators are language-specific and utilize morphosyntactic features (POS, case, number and gender), domains of Polish WordNet lexical units (word-sense pairs (Maziarz et al., 2012)), thousands of derivational relation instances between nouns, adjectives and verbs from the wordnet<sup>6</sup> and information about syntactic frames of nominalized predicates, taken from Polish valence dictionary (Dębowski and Woliński, 2007).

Each of written operators refers to one semantic relation. In other words, each semantic relation is described by one or by many WCCL operators. If an operator is successfully applied to a pair (or a

<sup>5</sup>Random sample of 200 NPs/PPs taken from KPWr, 331 relation instances, bootstrap confidence intervals are following P = 83 ÷ 91%, R = 76 ÷ 84%, F = 79 ÷ 87%,  $\alpha = 0.05$ . The corpus was divided by us into two parts: one working set for testing and preparing parsing rules and semantic operators - consisting of 300 texts, and a smaller evaluation part of 26 texts.

<sup>6</sup>Since we do not use any word sense disambiguation system, we simply take the first sense of every given word.

---

**Algorithm 1** Rule-based algorithm for the recognition of word pairs and triples

---

1. **genitive attachment** – link  $AgP_i$  head in genitive to the closest  $AgP_{i-1}$  head to the left or to the closest nominalized predicate to the left:
  - if there is none - link it to the closest predicate to the right;
  - if there is none - link the considered  $AgP_i$  head to the head of the whole NP/PP;
2. **nominative attachment** – link  $AgP_i$  head in nominative to the closest  $AgP_{i-1}$  head to the left or to the closest nominalized predicate to the left:
  - if there is none - link it to the closest  $AgP_{i+1}$  head to the right or to the closest nominalized predicate to the right;
  - if there is none - link the considered  $AgP_i$  head to the head of the whole NP/PP;
3. **small PP attachment** – link a head of  $AgP_i$  containing a small PP to the closest nominalized predicate to the left:
  - if there is none - to the closest nominalized predicate to the right such that it is not an element of  $AgP_{j>i}$  containing a preposition;
  - if there is none - to the closest  $AgP_{i-1}$  head to the left;
  - if there is none - link  $AgP_i$  with our whole NP/PP head;
4. **coordinated syntactic groups** – look for such  $AgP_i$  that is preceded by a coordination conjunction (i.e., *i* ‘and’, *oraz* ‘and’, *lub* ‘or’) or by coordinating comma (‘coordinating comma’ is such a comma that is placed between two AgPs whose heads are agreed on case), such coordination marker cannot be an element of any AgP:
  - if there is such a marker, look to the left in order to find such  $AgP_{j<i}$  head which is agreed on case with our  $AgP_i$  head – then create a new relation instance by copying the link  $AgP_j \rightarrow X$  and replacing  $AgP_j$  head by the  $AgP_i$  head in that copied linkage, i.e., create the relation instance  $AgP_i \rightarrow X$ ;
  - if it is not possible – do not introduce any relation;

5. head token provided with **POS known to the CRF tagger** – link the  $AgP_i$  head to the closest nominalized predicate to the left:
    - if there is none - to the closest nominalized predicate to the right such that it is not an element of  $AgP_{j>i}$  containing a preposition;
    - if there is none - link  $AgP_i$  head to the closest  $AgP_{j<i}$  head to the left such that  $AgP_{j<i}$  does not contain any preposition;
    - if there is none such  $AgP_{j<i}$  – connect  $AgP_i$  to the whole NP/PP head;
  6. **other cases** (the  $AgP_i$  head was not provided any known POS by the CRF tagger) – in such cases link  $AgP_i$  head to the closest  $AgP_{j<i}$  head to the left; if there is none – do not make any decision;
  7. **relations inside AgPs** – link adjectival and participial modifiers to the head of  $AgP_i$ .
- 

triple), then we know what semantic relation between the pair (or triple) occurs. Otherwise, we assume that the semantic relation does not occur.

For example, our `Proto-Patient` relation was described by the 6 `WCCL` operators. One of them is presented in Listing 1. This operator uses two dictionaries with valence frames (`acc` - a list of verbs possessing any accusative frame, `frames` - a list of verbs described in the Polish valence dictionary (Dębowski, 2013)) and morphosyntactic information about part of speech (`class`) and `case`.

This operator `PROTO-PATIENT-acc` captures pairs like *dręczący<sub>pact</sub> Janka<sub>noun.acc-θ</sub>* ‘tormenting John<sub>θ</sub>’ with a noun playing a Proto-Patient role of the predicate *dręczący*. The operator first checks whether a predicate (active participle) has an accusative frame or is outside the dictionary of Dębowski (“frames”). Since *dręczyć* ‘to torment’ is in `acc` dictionary and since *Janek* ‘John’ has `subst` class and `acc` case - the boolean operator returns ‘true’.

Let us present another example: the Proto-Agent macrorole is recognized by 5 operators, in Listing 2 was shown one of them. The `PROTO-AGENT-ger-przez-acc` operator is written for triples, i.e., for a triple *wydanie<sub>pact</sub> przez<sub>pron</sub> wydawcę<sub>noun.acc-θ</sub>* ‘publishing by the

publisher<sub>θ</sub>'. The first element in the triple is a gerund form of verb *wydać* 'to publish'. The operator checks whether the verb *wydać* has in its frame accusative/genetive or whether it cannot be found in Dębowski's dictionary (position 0 in the triple, frames).

Listing 1: One of the WCCL operators describing Proto-Patient relation. Language details has been described in (Radziszewski et al., 2011), abbreviations for grammatical categories has been explained in (Przepiórkowski et al., 2012)

```
@b:"PROTO-PATIENT-acc" (
  and(
    // 0 - accusative frame
    equal(class[0], pact),
    or(
      equal(lex(base[0], "acc"), ["1"]),
      not(equal(
        lex(base[0], "frames"), ["1"]))
    ),
    // 1 - noun or adj. & accusative
    in(class[1], {subst, depr, ger, adj}),
    equal(cas[1], acc)
  )
)
```

Next the operator seeks for the preposition *przez* 'by' at position 1. Then it tests if the first meaning of the lemma *wydawca* 'publisher' does not belong to the domain 'time' (= Polish *czas*) in Polish WordNet (position 2). Indeed, the first meaning of *wydawca* is in the domain 'person' (that information is available in the dictionary *noun\_domain*). At the end, we check whether the last token of our triple is in accusative. Because all of these conditions are fulfilled, the operator returns 'true', and we may assume that the last token takes the role of Proto-Agent.

Listing 2: A WCCL operator for the Proto-Agent relation

```
@b:"PROTO-AGENT-ger-przez-acc" (
  and(
    // 0 - gerund
    equal(class[0], {ger}),
    or(
      equal(lex(base[0], "acc"), ["1"]),
      equal(lex(base[0], "gen"), ["1"]),
      not(equal(
        lex(base[0], "frames"), ["1"]))
    ),
    // 1 - preposition "przez"
    equal(orth[1], "przez"),
    // 2 - not 'time' & accusative
    equal(cas[2], acc),
    not(
      equal(lex(if(
        equal(class[2], {ger}),
        lex(base[2], "ger_base"), base[2]),
        "noun_domain"), ["czas"]))
    )
  )
)
```

In Listing 3 one operator for family relation was shown. *FAMILY-agpp* used to recognize this relation for word pairs. The operator, inter alia, uses semantic dictionary of kinship names built on the basis of Polish WordNet (the dictionary *kinship*), lammas of possessive pronouns (e.g., *mój* 'my', *twój* 'yours').

Listing 3: Two WCCL operators describing Family relation

```
@b:"FAMILY-agpp" (
  and(
    // agreement
    agrpp(0,1, {nmb, gen, cas}),
    // position 0
    in(base[0], ["mój", "twój",
      "swoj", "nasz", "wasz"])
    // position 1
    equal(lex(base[1], "kinship"), ["1"]),
    equal(lex(
      base[1], "noun_domain"), ["os"]),
    in(class[1], {ger, subst, depr}),
  )
)
```

## 5 Results and conclusions

Evaluation of the presented semantic relation recognition algorithm was performed in three steps. First experiment (labelled *kpwr*) was performed on a random sample of the *KPWr* corpus (26 out of 326 texts, approximately one thirteenth of the corpus). In this experiment we made use of syntactic annotations from *KPWr* (cf. Tab. 1). Second experiment was performed on a random sample of 100 texts taken from yet another Polish corpus, called *NKJP* (Przepiórkowski et al., 2012, *nkjp*, approximately one tenth of the corpus)<sup>7</sup>. Since *NKJP* lacked syntactic annotations of *KPWr* style, we were forced to run on it the CRF shallow parser (described in Sec. 4.1). This experiment gave us information about performance of our algorithm on a 'bare' text (see Tab. 2). Evaluation in the experiments was done by a professional linguist.

At last, four baseline models were constructed and evaluated on the two corpora (Tab. 3). We created baselines similar to that presented in (Uchiyama et al., 2008), which was majority model. We chose the most frequent relation, which in the sample from *KPWr* was Proto-Patient (with the number of 113 instances out of 268 relation instances), this relation type was also the most frequent in the sample of *NKJP* (411 out of 1950 relation instances). For each corpora two baselines

<sup>7</sup>We focused on one-million balanced version of the much bigger corpus.

Relation	TP/FP/FN	P [%]	R [%]	F1 [%]
Proto-Agent	7/5/16	58.3	30.4	40.0
Proto-Patient	45/8/68	84.9	39.8	54.2
Instrument	0/0/7	—	0.0	—
Material	0/0/3	—	0.0	—
Purpose	1/7/30	12.5	3.2	5.1
location	3/9/25	25.0	10.7	15.0
sp. meronymy	0/3/2	0.0	0.0	—
time	2/2/3	50.0	40.0	44.4
t. meronymy	1/0/1	—	—	—
attribute	14/18/10	43.8	58.3	50.0
family	0/0/2	—	—	—
order	5/0/5	100.0	50.0	66.7
quantity	10/2/8	83.3	55.6	66.7
<b>All</b>	<b>88/54/186</b>	<b>53.8-70.1</b>	<b>*26.8-38.0</b>	<b>*36.2-48.3</b>

Table 1: Results of the algorithm on a sample from KPWr: P = Precision, R = recall, F1 = F-measure, TP = true positives, FP = false positives, FN = false negatives, sp. = spatial, t. = temporal. Percentile bootstrap confidence intervals are calculated at  $\alpha = 0.05$ . Asterisks denote significant differences between  $k_{pwr}$  and  $n_{kjp}$  in one-tailed tests,  $\alpha = 0.05$

Relation	TP/FP/FN	P [%]	R [%]	F1 [%]
Proto-Agent	75/7/143	91.5	34.4	50.0
Proto-Patient	181/17/230	91.4	44.0	59.4
Instrument	2/1/8	66.7	20.0	30.8
Material	3/4/36	42.9	7.7	13.0
Purpose	13/7/94	65.0	12.2	20.5
location	90/75/202	54.6	30.8	39.4
sp. meronymy	12/11/25	52.2	32.4	40.0
time	25/16/75	61.0	25.0	35.5
t. meronymy	2/0/66	100	2.9	57.1
attribute	200/248/64	44.6	75.8	56.2
family	18/0/6	100.0	60.0	85.7
order	33/0/100	100.0	24.8	39.8
quantity	113/68/146	62.4	43.6	51.4
<b>All</b>	<b>767/454/1195</b>	<b>60.1-65.6</b>	<b>*36.9-41.2</b>	<b>*46.0-50.3</b>

Table 2: Results of the algorithm on a sample from NKJP, labels as in the previous table. Percentile bootstrap confidence intervals are calculated at  $\alpha = 0.05$ . Asterisks denote significant differences between  $k_{pwr}$  and  $n_{kjp}$  in one-tailed tests,  $\alpha = 0.05$

were calculated: in Baseline #1 we assumed that we had perfectly recognized all occurrences of semantic relations (of any type), in Baseline #2 we simply signed with ‘Proto-Patient’ label every recognized by our system semantic relation instance. Baseline #2 is realistic, while #1 is idealistic, since to obtain #1 we should be able to recognize every single relation instance within a corpus. Baselines #1 are upper limits for all majority models (including #2). Our two idealistic baselines are higher than the realistic baselines (see Tab. 3).

Percentile bootstrap methods (DiCiccio and Efron, 1996), (DiCiccio and Romano, 1988) were applied to statistical significance and confidence interval (CI) analysis of the data.<sup>8</sup> We took 10000

<sup>8</sup>Our data for NKJP were merged, so cross-validation was

$k_{pwr}$	P	R	F1
Baseline #1	*42.2%	*42.2%	42.2%
Baseline #2	*26.2%	*20.0%	*22.7%
Experiment	<b>62.0%</b>	32.8%	<b>42.9%</b>
$n_{kjp}$	P	R	F1
Baseline #1	*21.1%	*21.0%	*21.0%
Baseline #2	*14.9%	*9.2%	*11.4%
Experiment	<b>62.5%</b>	<b>38.9%</b>	<b>47.9%</b>

Table 3: Precision, recall and F1 for baselines (#1 & #2) and experiments ( $k_{pwr}$ ,  $n_{kjp}$ ). Asterisks denote significant differences between an experiment and a baseline in one-tailed test at  $\alpha = 0.05$

bootstrap resamplings for each measure (P, R, F1),  $\alpha$  was equal to 0.05 for each one-tailed test and CI (a percentile CI need not be symmetrical).

In  $n_{kjp}$  we have beaten both idealistic and realistic baselines. Precision, recall and F1 for  $k_{pwr}$  are higher than Baseline #2. Only idealistic Baseline #1 for the KPWr corpus has overtaken our rule-based algorithm with regard to recall (42.2% vs. 32.8%), while its precision is lower and F1’s are statistically indistinguishable.

Results are promising, precisions go above 50% (the lower endpoint for the  $k_{pwr}$  confidence interval), for  $n_{kjp}$  we may assess it even more precisely as 60%-65%. Some semantic relations are recognized with higher precision: Proto-Agent ( $n_{kjp}$ : 89-100%,  $k_{pwr}$ : 90-100%,  $\alpha = 0.05$ ), Proto-Patient ( $n_{kjp}$ : 88-95%,  $k_{pwr}$ : 83%-98%), family ( $n_{kjp}$ : 90-100%) and order ( $n_{kjp}$ : 91-100%). Our system is thus comparable in this aspect to the systems described in Sec. 2.<sup>9</sup>

Overall recall is low, but higher than realistic baselines. In  $k_{pwr}$  we obtained  $R = 27-38\%$ , while for  $n_{kjp}$  we got statistically higher interval of 37-41%. It seems that recall was not affected by lack of marked NP/PP borders in the corpus (these should have been brought out by the CRF shallow parser). F-measures calculated on our both corpora are also much higher than realistic baselines #2.

We can already conclude that our preliminary experiments turned successful. Now we are aiming at improving our operators to raise their recall and at expanding the semantic role set (e.g., for Agent, Causer, Experiencer, Possessor or Result). Parallel, we start work on construction of automatic algorithms for relation recognition.

not available.

<sup>9</sup>Not directly, of course.

## References

- Asma Ben Abacha and Pierre Zweigenbaum. 2011. Automatic extraction of semantic relations between medical entities: a rule based approach. *Journal of Biomedical Semantics*, 2(5):1–11.
- Bartosz Broda, Michał Marcińczuk, Marek Maziarz, Adam Radziszewski, and Adam Wardyński. 2012. KPWr: Towards a Free Corpus of Polish. In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of LREC'12*, Istanbul, Turkey. ELRA.
- Thomas J. DiCiccio and Bradley Efron. 1996. Bootstrap confidence intervals. *Statistical Science*, 11(3):189–212.
- Thomas J. DiCiccio and Joseph P. Romano. 1988. A review of bootstrap confidence intervals. *Journal of the Royal Statistical Society. Series B (Methodological)*, 50(3):338–354.
- David Dowty. 1991. Thematic proto-roles and argument selection. *Language*, 67(3):547–619.
- Łukasz Dębowski and Marcin Woliński. 2007. Argument co-occurrence matrix as a description of verb valence. In *Proc. of the 3rd Language & Technology Conference*, volume 3, pages 260–264, Poznań, Poland.
- Łukasz Dębowski. 2013. Polish valence dictionary (<http://nlp.ipipan.waw.pl/ppjp/slownik/swigra/koncowy.txt>).
- Pablo Gamallo, Marco Gonzalez, Alexandre Agustini, Gabriel Lopes, and Vera S. De Lima. 2002. Mapping syntactic dependencies onto semantic relations. In *ECAI Workshop on Machine Learning and Natural Language Processing for Ontology Engineering*.
- Konrad Gołuchowski and Adam Przepiórkowski. 2012. Semantic role labelling without deep syntactic parsing. In Hitoshi Isahara and Kyoko Kanzaki, editors, *Advances in Natural Language Processing*, volume 7614 of *Lecture Notes in Computer Science*, pages 192–197. Springer Berlin Heidelberg.
- Marti A. Hearst. 1992. Automatic acquisition of hyponyms from large text corpora. In *Proceedings of the 14th conference on Computational linguistics - Volume 2*, COLING '92, pages 539–545, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Rongfeng Huang. 2009. Semantic relation extraction by automatically constructed rules. In *Proceedings of the International Conference on Artificial Intelligence and Computational Intelligence*, AICI '09, pages 425–434, Berlin, Heidelberg. Springer-Verlag.
- Ewa Jędrzejko. 1993. *Nominalizacje w systemie i tekstach współczesnej polszczyzny*. Uniwersytet Śląski, Katowice.
- Kate Kearns. 2011. *Semantics*. Palgrave.
- Martha J. Kolln. 1990. *WordNet: Understanding English Grammar*. Macmillan.
- Gabriel Larson, Richard & Segal. 1996. *Knowledge of Meaning: An Introduction to Semantic Theory*. The MIT Press.
- Roman Laskowski and Henryk Wróbel, editors. 1997. *Gramatyka współczesnego języka polskiego. Morfologia [Grammar of contemporary Polish: Morphology]*. PWN.
- Douglas B. Lenat. 1995. Cyc: a large-scale investment in knowledge infrastructure. *Commun. ACM*, 38(11):33–38, November.
2009. Luna project: Spoken language understanding in multilingual communication systems (<http://zil.ipipan.waw.pl/luna>).
- Marek Maziarz, Adam Radziszewski, and Jan Wiczorek. 2011. Chunking of Polish: guidelines, discussion and experiments with Machine Learning. In *Proc. of the LTC*.
- Marek Maziarz, Maciej Piasecki, and Stanisław Szpakowicz. 2012. Approaching plWordNet 2.0. In *Proceedings of the 6th Global Wordnet Conference*, Matsue, Japan, January.
- Dan Moldovan, Adriana Badulescu, Marta Tatu, Daniel Antohe, and Roxana Girju. 2004. Models for the semantic classification of noun phrases. In *Proceedings of the HLT-NAACL Workshop on Computational Lexical Semantics*, CLS '04, pages 60–67, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Vivi Nastase, Jelber Sayyad Shirabad, Marina Sokolova, and Stan Szpakowicz. 2006. Learning noun-modifier semantic relations with corpus-based and wordnet-based features. In *Proceedings of the 21st National Conference on Artificial Intelligence*.
- Palmer, Martha & Gildea, Daniel & Xue, and Nianwen. 2010. *Semantic Role Labeling*. Morgan & Claypool Publishers.
- Adam Przepiórkowski, Mirosław Banko, Rafał Górski, and Barbara Lewandowska-Tomaszczyk, editors. 2012. *Narodowy Korpus Języka Polskiego*. Wydawnictwo Naukowe PWN.
- Adam Radziszewski and Adam Pawlaczek. 2012. Large-scale experiments with NP chunking of Polish. In *Proceedings of TSD 2012*, Brno, Czech Republic. Springer.
- Adam Radziszewski, Adam Wardyński, and Tomasz Śniatowski. 2011. WCCL: A Morpho-syntactic Feature Toolkit. In *Proceedings of the Balto-Slavonic Natural Language Processing Workshop*. Springer.

- Adam Radziszewski, Marek Maziarz, and Jan Wiczorek. 2012. Shallow syntactic annotation in the Corpus of Wrocław University of Technolog. *Cognitive Studies*, 12.
- Adam Radziszewski. 2013. A tiered CRF tagger for Polish. In R. Bembenik, Ł. Skonieczny, H. Rybiński, M. Kryszkiewicz, and M. Niezgódka, editors, *Intelligent Tools for Building a Scientific Information Platform: Advanced Architectures and Solutions*.
- Bryan Rink and Sanda Harabagiu. 2010. Utd: Classifying semantic relations by combining lexical and semantic resources. In *Proceedings of the 5th International Workshop on Semantic Evaluation, SemEval '10*, pages 256–259, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Barbara Rosario and Marti Hearst. 2001. Classifying the semantic relations in noun compounds via a domain-specific lexical hierarchy. In *Proceedings of the 2001 Conference on Empirical Methods in Natural Language Processing (EMNLP-01)*, pages 82–90.
- Stephen Tratz and Eduard Hovy. 2010. Isi: Automatic classification of relations between nominals using a maximum entropy classifier. In *Proc. of the 5th Intern. Workshop on Semantic Evaluation*, Sweden. Association for Computational Linguistics.
- Kateryna Tymoshenko and Claudio Giuliano. 2010. Fbk-irst: Semantic relation extraction using cyc. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 214–217, Uppsala, Sweden, July. Association for Computational Linguistics.
- Kiyoko Uchiyama, Shunsuke Aihara, and Shun Ishizaki. 2008. Identifying semantic relations in japanese compound nouns for patent documents analysis. In *Proceedings of the 3rd international conference on Large-scale knowledge resources: construction and application, LKR'08*, pages 75–81, Berlin, Heidelberg. Springer-Verlag.
- Robert D. Van Valin. 2004. Semantic macroroles in role and reference grammar. In Rolf Kailuweit and Martin Hummel, editors, *Semantische Rollen*, pages 62–82. Tuebingen Narr.