

Bank częściowo ujednoznaczionych struktur LFG

Katarzyna Krasnowska¹
Witold Kieraś^{1,2}

¹IPI PAN ²IJP UW

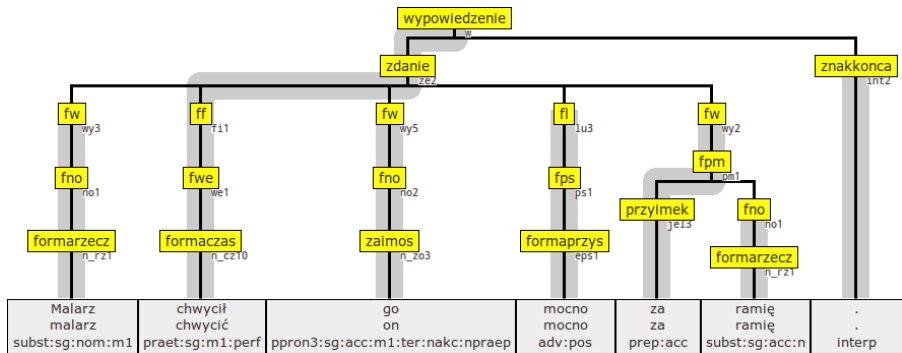
7 października 2013

- Składnica (Woliński *et al.*, 2011)
 - bank drzew składnikowych oparty na GFJP,
 - około 8000 drzew,
 - zasób wciąż rozwijany, podlegający zmianom, rozszerzeniom i poprawkom.
- POLFIE (Patejuk i Przepiórkowski, 2012)
 - polska gramatyka LFG,
 - pierwotnie oparta na regułach GFJP, później modyfikowanych i rozszerzanych.
- INESS
 - system zarządzania treebankami LFG,
 - interfejs webowy do ręcznego ujednoznaczniania struktur.

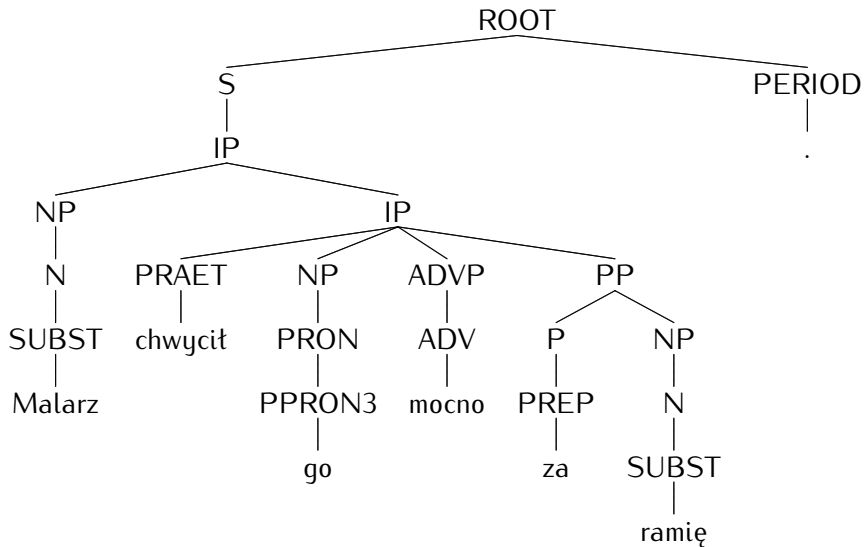
Pomysł na treebank

- Ograniczenie rozbiorów LFG do zgodnych ze Składnicą.
- Ręczne ujednoznacznianie znacznie mniejszego zbioru struktur.
- W efekcie: mniej czasu i pracy potrzebnych do stworzenia nowego zasobu.

Przykład: drzewo ze Składnicy



Przykład: c-struktura w POLFIE



Przykład: f-struktura w POLFIE

PRED		'chwycić<1,2,3>'										
SUBJ	1	<table><tr><td>PRED</td><td>'malarz'</td></tr><tr><td>CASE</td><td>nom</td></tr><tr><td>GEND</td><td>m1</td></tr><tr><td>NUM</td><td>sg</td></tr></table>	PRED	'malarz'	CASE	nom	GEND	m1	NUM	sg		
PRED	'malarz'											
CASE	nom											
GEND	m1											
NUM	sg											
OBJ	2	<table><tr><td>PRED</td><td>'on'</td></tr><tr><td>CASE</td><td>acc</td></tr><tr><td>GEND</td><td>m1</td></tr><tr><td>NUM</td><td>sg</td></tr></table>	PRED	'on'	CASE	acc	GEND	m1	NUM	sg		
PRED	'on'											
CASE	acc											
GEND	m1											
NUM	sg											
OBL	3	<table><tr><td>PRED</td><td>'ramię'</td></tr><tr><td>CASE</td><td>acc</td></tr><tr><td>GEND</td><td>n</td></tr><tr><td>NUM</td><td>sg</td></tr><tr><td>PFORM</td><td>za</td></tr></table>	PRED	'ramię'	CASE	acc	GEND	n	NUM	sg	PFORM	za
PRED	'ramię'											
CASE	acc											
GEND	n											
NUM	sg											
PFORM	za											
ADJUNCT		{ <table><tr><td>PRED</td><td>'mocno'</td></tr><tr><td>DEGREE</td><td>positive</td></tr></table> }	PRED	'mocno'	DEGREE	positive						
PRED	'mocno'											
DEGREE	positive											

Jak porównywać drzewa?

C-struktura:

- jest drzewem składnikowym (podobnie jak drzewa w Składnicy),
- podlega większym zmianom podczas rozwoju gramatyki,
- pełni drugorzędną funkcję w LFG.

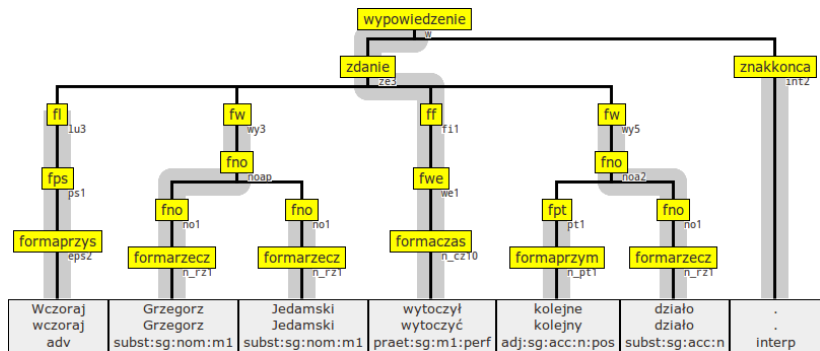
F-struktura:

- odwierciedla zależności funkcyjne pomiędzy predykatami,
- pomimo odmiennej struktury okazała się wygodniejsza do porównywania.

Schemat f-struktury

- Z drzew w Składnicy można odczytać, jak powinna wyglądać struktura predykatów.
- Informacja o elementach głównych poszczególnych fraz przekłada się na strukturę zależności między predykatami.

Przykład



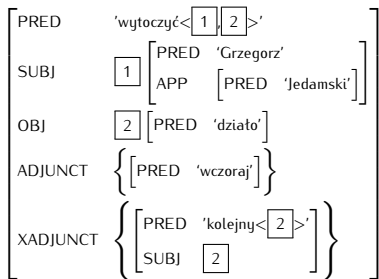
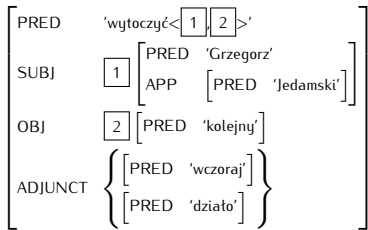
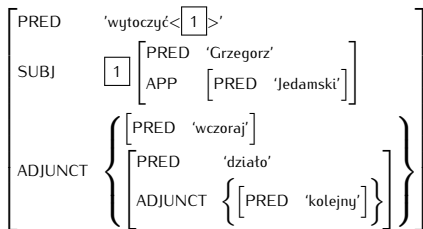
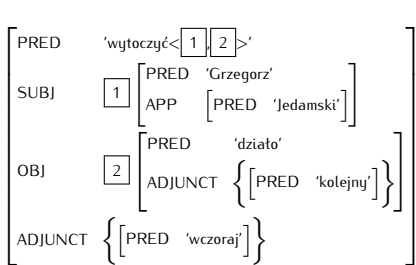
predykat f-struktury

'WYTOCZYĆ'
'GRZEGORZ'
'DZIAŁO'

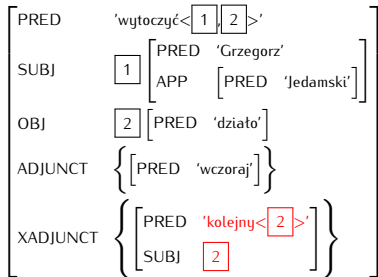
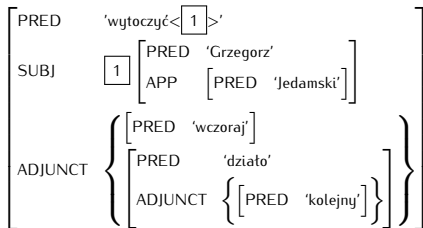
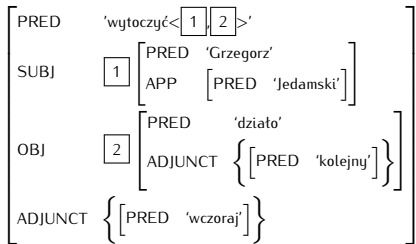
predykaty atrybutów

'WCZORAJ', 'GRZEGORZ', 'DZIAŁO'
'JEDAMSKI'
'KOLEJNY'

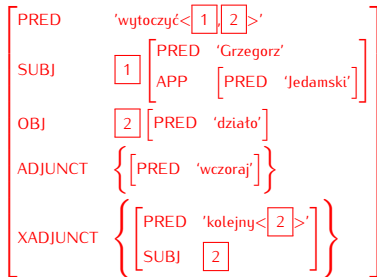
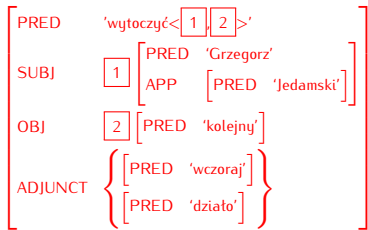
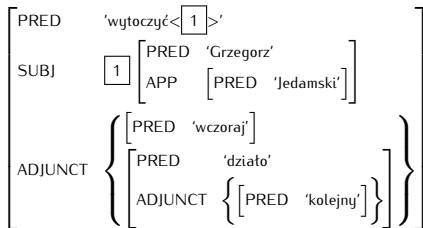
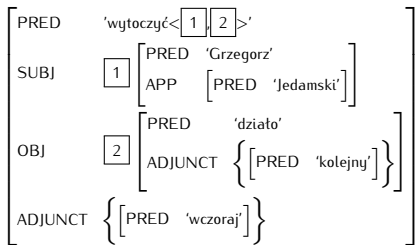
Przykład



Przykład



Przykład



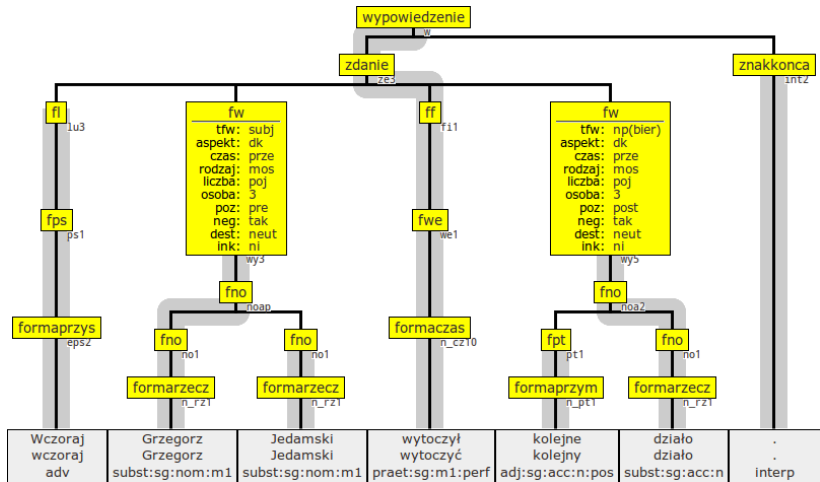
Odwzorowanie tfw na atrybuty LFG

- Typ frazy wymaganej:
 - wymaganie walencyjne centrum finitywnego realizowane przez daną frazę,
 - informacja powierzchniowo-składniowa,
 - np. subj (podmiot składniowy!), np(bier), sentp(że).
- Nazwy atrybutów w LFG:
 - funkcje gramatyczne,
 - np. SUBJ, OBJ, XCOMP-PRED, ADJUNCT.
- Brak jednoznacznej odpowiedniości:
 - *Mruczek jest kotem* – np(narz) / XCOMP-PRED,
 - *Ala macha kotem* – np(narz) / OBL-INST.

Odwzorowanie tfw na atrybuty LFG

typ frazy wymaganej	nazwy atrybutów
subj	SUBJ
np(mian)	XCOMP-PRED, OBL-STR
np(dop)	OBJ, OBL-GEN
np(cel)	OBJ, OBJ-TH
np(bier)	OBJ, OBL-STR
np(narz)	OBL-INST, XCOMP-PRED, OBJ
adjp(mian)	XCOMP-PRED
adjp(narz)	XCOMP-PRED
advp	OBL (lub modyfikator)
sentp(_)	SUBJ, COMP
infp(_)	SUBJ, XCOMP
prepnp(_, _)	OBL, OBL ₂ , OBL ₃ , OBL-AG
prepadjp(_, _)	XCOMP-PRED

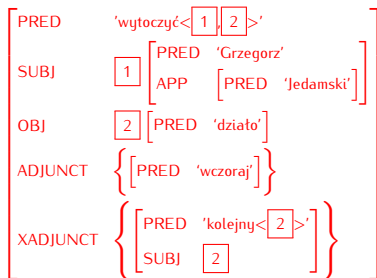
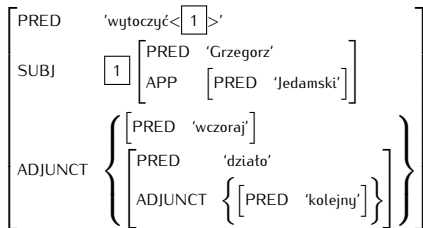
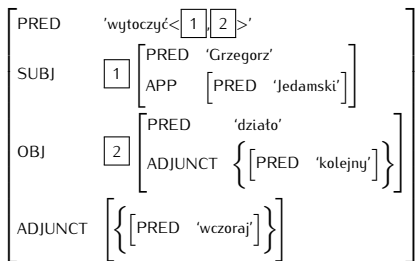
Przykład



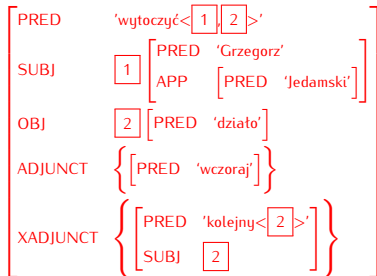
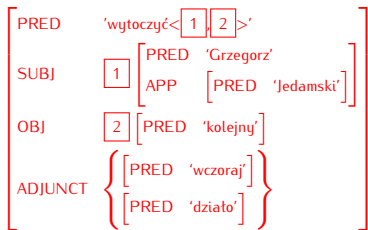
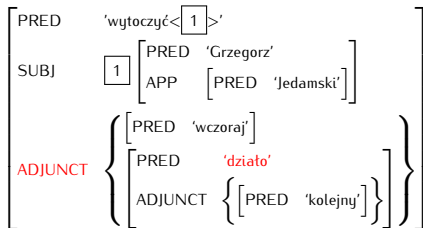
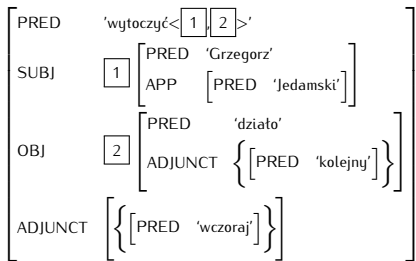
Przykład

predykat f-struktury	nazwa funkcji	predykat atrybutu
'WYTOCZYĆ'	SUBJ OBJ / OBL-STR nieargument	'GRZEGORZ' 'DZIAŁO' 'WCZORAJ'
'GRZEGORZ'	—	'JEDAMSKI'
'DZIAŁO'	—	'KOLEJNY'

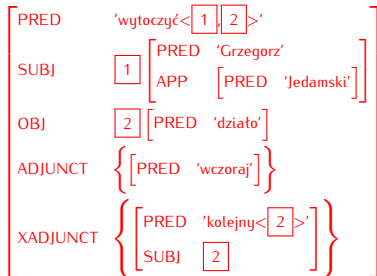
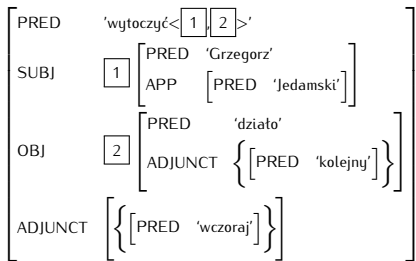
Przykład



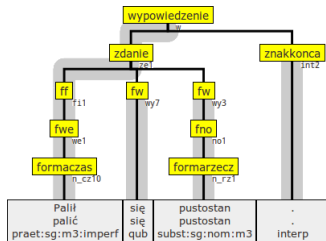
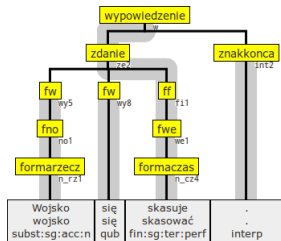
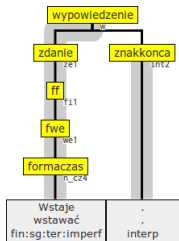
Przykład



Przykład



Podmiot

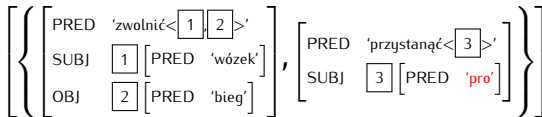
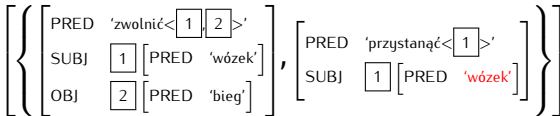
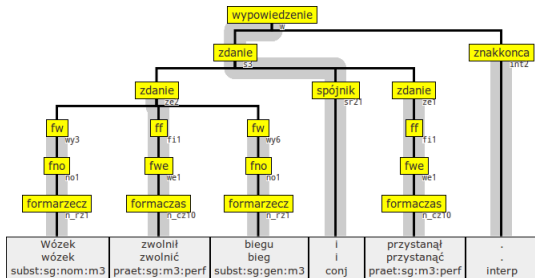


$$\left[\begin{array}{l} \text{PRED } \text{'wstawać} \langle 1 \rangle \\ \text{SUBJ } 1 \left[\text{PRED } \text{'pro'} \right] \end{array} \right]$$

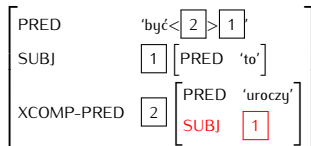
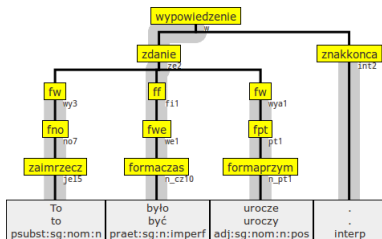
$$\left[\begin{array}{l} \text{PRED } \text{'skasować} \langle 1, 2 \rangle \\ \text{SUBJ } 1 \left[\text{PRED } \text{'pro'} \right] \\ \text{OBJ } 2 \left[\text{PRED } \text{'wojsko'} \right] \end{array} \right]$$

$$\left[\begin{array}{l} \text{PRED } \text{'palić się} \langle 1 \rangle \\ \text{SUBJ } 1 \left[\text{PRED } \text{'pustostan'} \right] \end{array} \right]$$

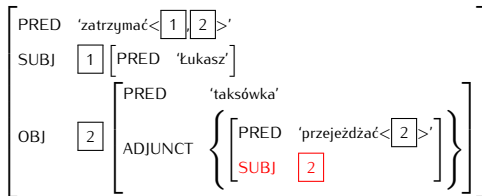
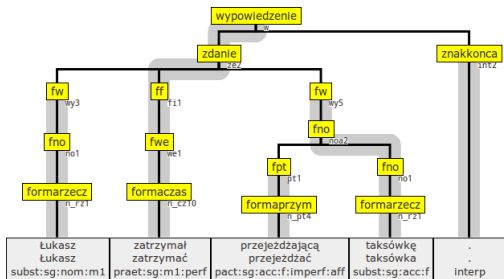
Podmiot — koordynacja



Podmiot — XCOMP-PRED

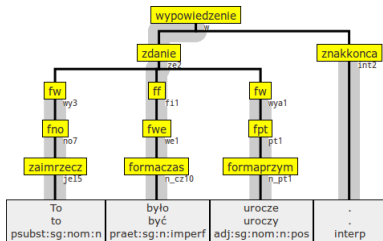


Podmiot — imiestów



- W f-strukturach POLFIE może pojawiać się podmiot nieobecny w Składnicy.
- Nie tylko centrum finitywne może mieć podmiot.
- Proste, ale skuteczne rozwiązanie:
 - jeśli w Składnicy pojawia się podmiot, to w f-strukturze również musi się pojawić i musi być z nim identyczny;
 - jeśli w Składnicy brak podmiotu, to dopuszczamy dowolny podmiot.
- Nie powoduje to dopuszczenia zbyt wielu struktur — niezgodność zostanie wykryta w innym miejscu.

Podmiot

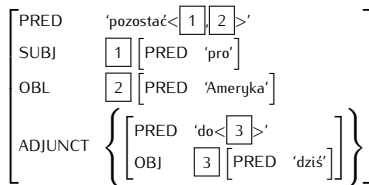
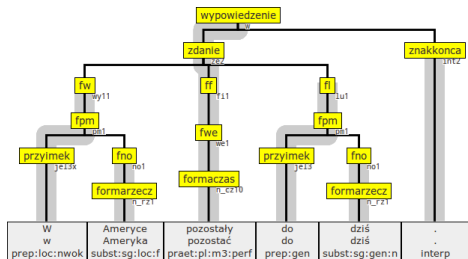
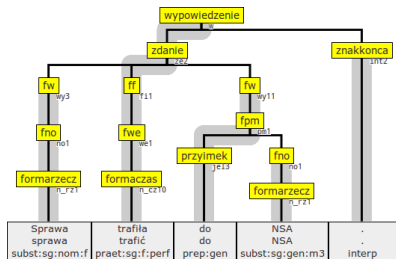


[PRED 'być<1>'
 SUBJ 1 [PRED 'to']
 XADJUNCT { [PRED 'uroczy<1>'] }]

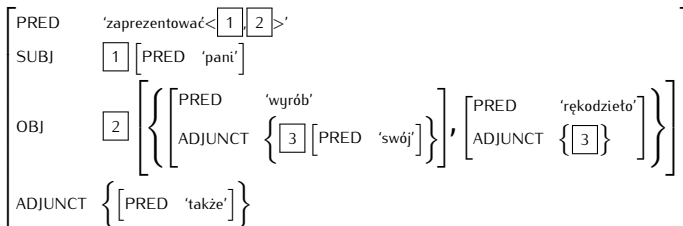
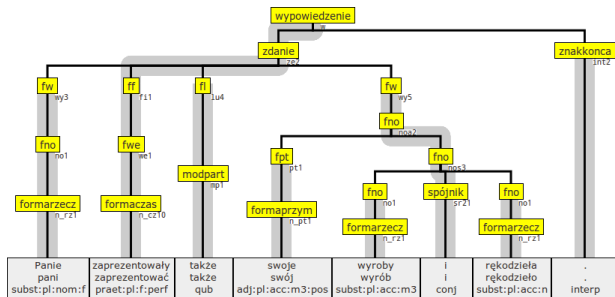
[PRED 'być<2>1'
 SUBJ 1 [PRED 'pro']
 XCOMP-PRED 2 [PRED 'to']
 SUBJ 1
 XADJUNCT { [PRED 'uroczy<1>'] }]

[PRED 'być<2>1'
 SUBJ 1 [PRED 'uroczy']
 XCOMP-PRED 2 [PRED 'to']
 SUBJ 1]

Frazy przyimkowe

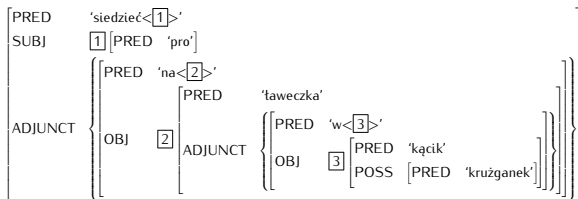
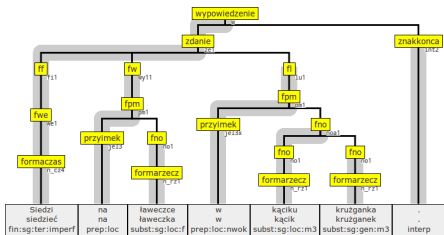


Koordinacja

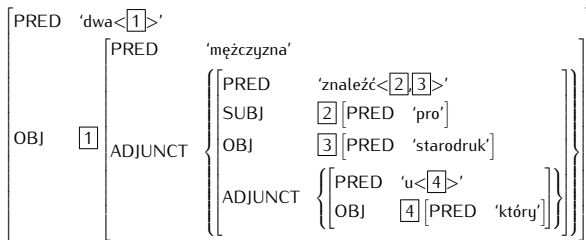
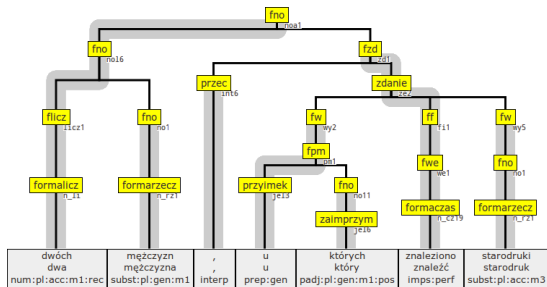


Porównanie z ręczną anotacją

- 206 ręcznie ujednoznacznionych struktur LFG.
- Zgodność dla 97 zdań (92%).
- Przykład niezgodności:



Przykład — brak pełnej zgodności



Szeregowanie rozbiorów LFG

- Ostatnie przykłady pokazują, że ograniczenie rozbiorów do zgodnych ze Składnicą nie zawsze jest najlepszą strategią.
 - Możemy odrzucić rozbiór dla niejednoznacznego zdania, który anotator treebanku LFG uznałby za lepszy.
 - Możemy być zmuszeni odrzucić wszystkie rozbiory.
- Zamiast tego szeregujemy rozbiory LFG według ich zgodności ze Składnicą.
- Porównanie z danymi ręcznie anotowanymi wskazuje, że dla zdecydowanej większości zdań anotator znajdzie właściwy rozbiór już w pierwszym zbiorze.

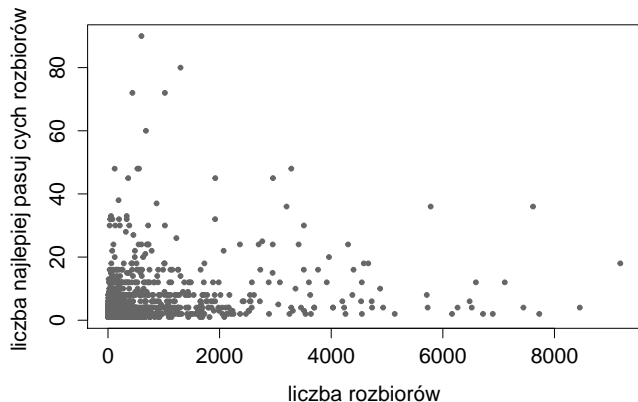
Ograniczenie liczby rozbiorów

Wyniki dla najnowszych wersji Składnicy (25.09.2013) i POLFIE (17.09.2013):

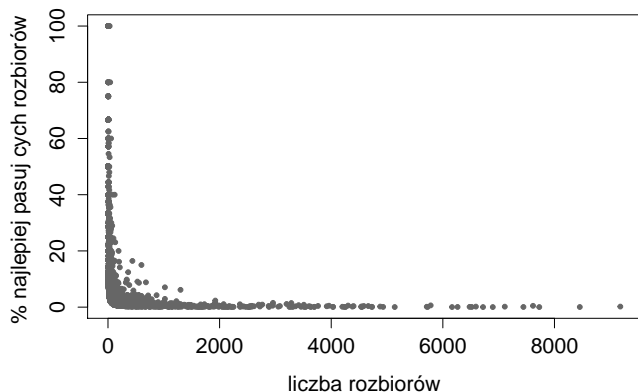
- 6910 zdań, w tym:
 - 770 zdań z 1 rozbiorem,
 - 36 zdań z >10000 rozbiorów.
- Dla pozostałych 6104 zdań:

	rozbiory	
	wszystkie	najlepiej pasujące
min	2	1
max	9180	90
średnia	144	2.6
mediana	12	1.5

Ograniczenie liczby rozbiorów



Ograniczenie liczby rozbiorów



Ograniczenie liczby rozbiorów

Zdania z dużą liczbą najlepiej pasujących rozbiorów:

- >20 dla 51 zdań (0.74%).
- >40 dla 12 zdań (0.17%).
- Koordynacja kilku fraz: f-struktury odpowiadające różnym nawiasowaniom mają tę samą strukturę predykatów.
 - *Towar z ciężarówek pochodził z Belgii, Francji, Hiszpanii, Holandii i Niemiec.* – 90 rozbiorów
- Inicjały: struktury dla M1, M2, M3, F i N.
 - *46-letni Marek Ch. doznał jedynie złamania prawej nogi.*
– 80 rozbiorów.

„Skutki uboczne”

- Dodatkowy test gramatyki POLFIE:
 - analiza przypadków braku pełnej zgodności wskazała, co można poprawić.
- Wykrycie kilkudziesięciu błędów anotacji w Składnicy:
 - podmiot predykatywnego *to*,
 - `fno` → `flicz`,
 - ...

Co dalej?

- W przypadku ewentualnych zmian w GFJP i POLFIE procedura może wymagać pewnych poprawek, ale ponieważ jest ona bardzo ogólna, to nie powinny być one znaczące.

Co dalej?

- W przypadku ewentualnych zmian w GFJP i POLFIE procedura może wymagać pewnych poprawek, ale ponieważ jest ona bardzo ogólna, to nie powinny być one znaczące.
- W przyszłości można ją zmodyfikować w taki sposób, by działała również w drugą stronę, tzn. by na podstawie wyników ręcznej anotacji zdań w formalizmie LFG uzyskać automatycznie częściowo ujednoznacznione rozbiory zdań GFJP.

Co dalej?

- W przypadku ewentualnych zmian w GFJP i POLFIE procedura może wymagać pewnych poprawek, ale ponieważ jest ona bardzo ogólna, to nie powinny być one znaczące.
- W przyszłości można ją zmodyfikować w taki sposób, by działała również w drugą stronę, tzn. by na podstawie wyników ręcznej anotacji zdań w formalizmie LFG uzyskać automatycznie częściowo ujednoznacznione rozbiory zdań GFJP.
- Zaprezentowana procedura jest na tyle ogólna, że po niezbędnych modyfikacjach mogłaby potencjalnie służyć do utworzenia banków częściowo ujednoznacznionych rozbiorów dla różnych innych formalizmów gramatycznych.

Bibliografia

- Patejuk, A. i Przepiórkowski, A. (2012). Towards an LFG parser for Polish: An exercise in parasitic grammar development. W: *Proceedings of the Eighth International Conference on Language Resources and Evaluation, LREC 2012*, str. 3849–3852, Istanbul, Turkey. ELRA.
- Woliński, M., Głowińska, K. i Świdziński, M. (2011). A preliminary version of Składnica—a treebank of Polish. W: Z. Vetulani, red., *Proceedings of the 5th Language & Technology Conference*, str. 299–303, Poznań.