

# Integracja parsera zależnościowego z parserem kategoryalnym

Wojciech Jaworski

Instytut Podstaw Informatyki  
Polskiej Akademii Nauk

2 marca 2017

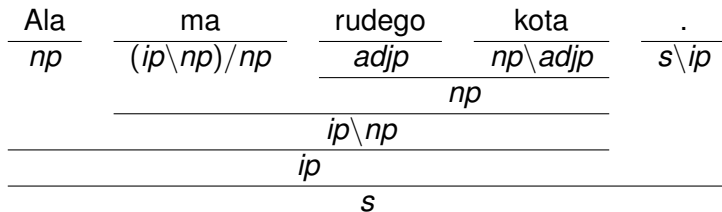
- Wymagania względem reprezentacji podziału na zdania:
  - ▶ reprezentacja podziału na
    - ★ tokeny,
    - ★ zdania,
    - ★ akapity
  - i całego tekstu,
  - ▶ reprezentacja zdania w różnych stadiach przetwarzania:
    - ★ sekwencja znaków,
    - ★ graf tokenów,
    - ★ skompresowany las drzew zależnościowych, ...
  - ▶ reprezentacja mowy niezależnej i zagnieżdżonych zdań,
  - ▶ reprezentacja niejednoznacznego podziału na zdania,
  - ▶ reprezentacja formatów akceptowanych/generowanych przez różne parsery.

- Korzystam z tego, że parsowanie w gramatykach kategoriálních jest dowodzeniem twierdzeń w logice.
- Logika intuicjonistyczna
  - ▶ zbiór przesłanek
  - ▶ przesłanki są łączne, przemienne, można je wielokrotnie używać w dowodzie
- Liniowa logika intuicjonistyczna
  - ▶ multizbiór przesłanek
  - ▶ przesłanki są łączne, przemienne, każdą z nich trzeba użyć dokładnie jeden raz w dowodzie
  - ▶ *glue semantics* w LFG

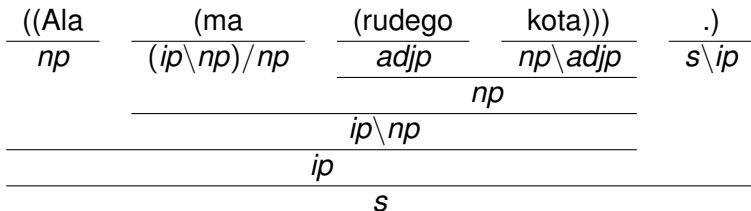
- Niekomutatywna liniowa logika intuicjonistyczna
  - ▶ lista przesłanek
  - ▶ przesłanki są łączne, nieprzemienne, każdą z nich trzeba użyć dokładnie jeden raz w dowodzie
  - ▶ Logical Categorical Grammar
- Logika **NL**, M. Moortgat „Categorical Type Logics”
  - ▶ drzewo przesłanek
  - ▶ przesłanki nie są łączne, nie są przemienne, każdą z nich trzeba użyć dokładnie jeden raz w dowodzie
  - ▶ konwerter drzew zależnościowych

# Przykładowy wywód gramatyczny

- Wywód w LCG:



- Drzewo zależnościowe reprezentują jako nawiasowanie tokenów:  
((Ala (ma (rudego kota))) .)
- Konwersja drzewa zależnościowego:



# Implementacja

- Niewielkie zmiany w parserze:
  - ▶ algorytm CYK zostaje zastąpiony spacerem po drzewie,
  - ▶ reguły formalizmu kategoryjnego pozostają niezmienione,
  - ▶ leksykon gramatyki pozostaje niezmieniony.
- Zaimplementowane nawiasowanie wskazuje funktor oraz listy lewych i prawych argumentów:

$$\frac{\frac{[[\phi, \text{Ala}, \phi], \quad \text{ma} \quad \frac{[[\phi, \text{rudego}, \phi], \quad \text{kota}, \phi]], \quad \text{.}, \phi]}{np \quad (ip \backslash np) / np \quad adjp \quad np \backslash adjp \quad s \backslash ip}}{ip \backslash np}}{ip}}{s}$$

- $\phi$  oznacza listę pustą

- Nawiasowanie odzwierciedla strukturę drzewa zależnościowego,
  - ▶ w którym podrzędniki podzielone są na lewe i prawe
  - ▶ oraz posortowane względem kolejności wystąpienia w tekście.
- W przypadku konstrukcji nieciągłych występujących np. w zdaniu:

*Rudego Ala ma kota.*

nawiasowanie zmieni kolejność tokenów w zdaniu eliminując nieciągłość

$[[\phi, \text{Ala}, \phi], \text{ma} [[\phi, \text{rudego}, \phi], \text{kota}, \phi]], \text{.}, \phi$

# Konstrukcje nieuwzględnione w leksykonie

- Aby parsować konstrukcje nieuwzględnione w leksykonie gramatyki kategoryjnej
  - ▶ do każdego typu dodawane są wielokrotne, opcjonalne argumenty
  - ▶ pasujące do dowolnego symbolu nie będącego funktorem.

$$\frac{\text{ma}}{(ip\{/?X, \?X\} \setminus np) / np}$$
$$\frac{[[\phi, \text{blabla}, \phi], \text{unk}\{/?X, \?X\}]}{\text{unk}\{\setminus?X\}}$$
$$\frac{\text{kota}, \phi]}{np\{/?X, \?X\}}$$
$$\frac{\text{unk}\{\setminus?X\}}{np\{\setminus?X\}}$$
$$\frac{\text{np}\{\setminus?X\}}{np}$$
$$\frac{\text{ip}\{/?X, \?X\} \setminus np}{ip\{/?X, \?X\} \setminus np}$$



# Zmiana powiązań między wierzchołkami w drzewie zależnościowym

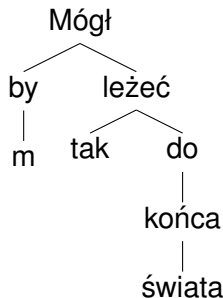
- Aby parsowanie powiodło się, drzewo zależnościowe na wejściu musi mieć strukturę zgodną z kolejnością aplikacji wierzchołków w gramatyce kategoryalnej.
- W *Składnicy Zależnościowej* i w *Krzakach* struktura ta dla pewnych konstrukcji składniowych jest inna od oczekiwanej przez parser kategoryalny.
- Dlatego przed parsowaniem drzewo zależnościowe poddajemy szeregowi lokalnych modyfikacji.

# Tryb przypuszczający

- Czasownik w trybie przypuszczającym opisany jest za pomocą wpisu:

*ip*{ | *by*, | *aglt*, ... }

- Zaś drzewo zależnościowe ma postać:



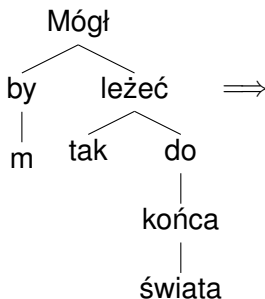
- Aby sparsować drzewo trzeba uczynić „m” bezpośrednim podrzędnikiem „Mógł”

# Tryb przypuszczający

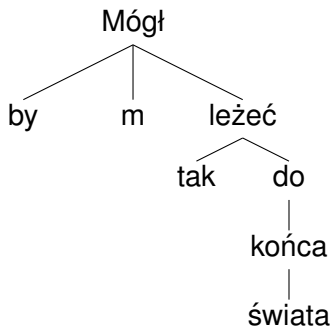
- Czasownik w trybie przypuszczającym opisany jest za pomocą wpisu:

$ip\{\mid by, \mid aglt, \dots\}$

- Zaś drzewo zależnościowe ma postać:



$\Rightarrow$



- Aby sparsować drzewo trzeba uczynić „m” bezpośrednim podrzędnikiem „Mógł”

- usunięcie kropek i przecinków
- zmiana dowiązań myślników przy mowie niezależnej
- rozwiązywanie skrótów
- zmiana struktury wieloczłonowej koordynacji
- przesunięcie spójników podrzędnych, zaimków pytajnych i względnych powyżej czasownika.

# Koncepcja przetwarzania skompresowanych lasów generowanych przez parser *Świga*

- Skompresowany las jest to skierowany graf acykliczny mający wierzchołki:
  - ▶ multiplikatywne — wiążące frazę z jej składowymi,
  - ▶ addytywne — wyrażające możliwe alternatywne rozbiory składniowe fragmentu zdania.
- Wierzchołki multiplikatywne będą przetwarzane tak jak wierzchołki w drzewach zależnościowych.
- Wierzchołki addytywne będą przetwarzane za pomocą aktualnego algorytmu kompresującego niejednoznaczność dla parsera kategoryjnego.
- Przeprowadzenie lokalnych zmian układu wierzchołków będzie wymagać rozkompresowania fragmentu lasu przed parsowaniem.

- Opis działania:

- ▶ Wypełnia tablicę zgodnie z działaniem algorytmu CYK trzymając w każdym polu do 10 alternatywnych poddrzew
- ▶ Przy wypełnianiu kolejnego pola rozpatrywane są osobno przypadki, gdy lewy węzeł jest nadrzędnikiem, i gdy jest nim prawy węzeł.
- ▶ Algorytm rozpatruje każdą parę kandydatów na ojca i syna i wybiera 10 najlepiej ocenionych zgodnie z wagą wyliczoną na podstawie atrybutów.
- ▶ Obsługuje nieciągłości za pomocą przegrupowywania wierzchołków z ograniczeniem poziomu nieciągłości w drzewie.

- Cele reimplementacji:

- ▶ dezambiguacja lasów generowanych przez parser kategorialny.
- ▶ łączenie sparsowanych fragmentów zdań, dla których parser kategorialny nie generuje pełnego rozbioru

# Metoda realizacji

- Kluczowym elementem parsera zależnościowego jest model, który określa wagę powiązania między parą wierzchołków.
- Wagę tą można wykorzystać w różnych algorytmach i do różnych celów.
- Dezambiguację można wykonać uzupełniając algorytm MateParser o uwzględnianie więzów wprowadzanych przez parser kategorialny.
  - ▶ Jest to rozwiązanie analogiczne do proponowanego przy integracji ze Świgrą.
- Uzupełnianie drzew rozbioru można zrealizować podobnie,
  - ▶ z założeniem, że korzenie sparsowanych fragmentów mogą być synami wierzchołków należących do innych fragmentów.

# Dziękuję za uwagę!