


Wprowadzenie do informatyki

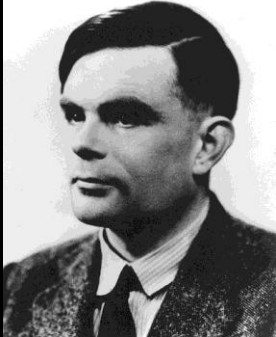
Jerzy Nawrocki
Wydział Informatyki
Politechnika Poznańska
jerzy.nawrocki@put.poznan.pl

**Sztuczna
inteligencja
i język naturalny**

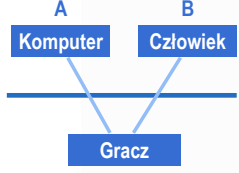


Wprowadzenie do informatyki

Test Turinga (1950)



Alan Turing



Przetwarzanie języka (2)

Wprowadzenie do informatyki

Program ELIZA (1966)



Joseph Weizenbaum

User: Men are all alike.
Eliza: In what way?
User: They're always bugging us about something or other.
Eliza: Can you think of a specific example?
User: My boyfriend made me come here.
Eliza: Your boyfriend made you come here?

Przetwarzanie języka (3)

Wprowadzenie do informatyki

Cel wykładu




- **Wprowadzić w problematykę automatycznego przetwarzania języka naturalnego.**
- **Przedstawić pojęcie gramatyki**

Przetwarzanie języka (4)

Wprowadzenie do informatyki

Plan wykładu




- Wprowadzenie
- Części mowy
- Gramatyka i wywód
- Język formalny
- Gramatyki bezkontekstowe
- Notacja Backusa-Naura

Przetwarzanie języka (5)

Wprowadzenie do informatyki

Dionizjusz Trak (170 – 90 pne)



8 części mowy (ang. *part of speech, POS*):

- Rzeczownik
- Czasownik
- Zaimek
- Przyimek
- Przysłówek
- Spójnik
- Imiesłów
- Przedimek **Przymiotnik?**

<http://www.uni-koeln.de/phil-fak/ifa/NRWakademie/papyrologie/PKoeln/PK5128v.jpg>

Przetwarzanie języka (6)

Wprowadzenie do informatyki

Współczesna systematyka części mowy – język ang.

- Rzeczownik
- Czasownik
- Zaimek
- Przyimek
- Przysłówek
- Spójnik
- Imiesłów
- Przedimek

Część mowy ↔ znacznik

Penn Treebank: 45 znaczników
Znaczniki C5: 61 znaczników
Brown corpus: 87 znaczników
Znaczniki C7: 146 znaczników

(znacznik = ang. tag)

Przetwarzanie języka (7)

Wprowadzenie do informatyki

Znaczniki Penn Treebank

CC Coordinating conjunction *and*
CD Cardinal number *one, two*
DT Determiner *the*
EX Existential *there there are*
FW Foreign word *mea culpa*
IN Preposition or subordinating conjunction *of, in, by*
JJ Adjective *yellow*
JJR Adjective, comparative *bigger*
NN Noun, singular or mass *tiger*
NNS Noun, plural *tigers*
VB Verb, base form *eat*
VBD Verb, past tense *ate*

Przetwarzanie języka (8)

Wprowadzenie do informatyki

Oznaczenie części mowy

The grand jury commented on a number of other topics.

The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

(oznaczenie cz. mowy = ang. tagging)

Przetwarzanie języka (9)

Wprowadzenie do informatyki

Oznaczenie części mowy

Diagram illustrating the process of part-of-speech tagging:

Diagram: Ciąg wyrazów → [Oznaczenie części mowy] ← Wyrazy ze znacznikami

Diagram: Zbiór części mowy → [Oznaczenie części mowy]

Przetwarzanie języka (10)

Wprowadzenie do informatyki

Oznaczenie części mowy – Główny problem

Book that flight.

Book/VB that/DT flight/NN ./.

Inne możliwości:
Book/NN Buy me that book.
that/CC I thought that you knew.

Niejednoznaczność

Przetwarzanie języka (11)

Wprowadzenie do informatyki

Rozstrzygnięcie niejednoznaczności

Book that flight.

Book/VB/NN that/DT/CC flight/NN ./.

Przetwarzanie języka (12)

Wprowadzenie do informatyki

Rozstrzyganie niejednoznaczności

Book that flight.

Book/VB/NN that/DT/CC flight/NN ./.

Książka ten/DT lot.
Książka, że/CC lot.

Reguła:
Jeśli jest konflikt między czasownikiem (VB, VBD, ..) a inną częścią mowy i w zdaniu nie ma innego czasownika, to należy uznać, że to jest czasownik.

Przetwarzanie języka (13)

Wprowadzenie do informatyki

Rozstrzyganie niejednoznaczności

Book that flight.

Book/VB/NN that/DT/CC flight/NN ./.

Zarezerwuj ten/DT lot.
Zarezerwuj, że/CC lot.

Reguła:
Jeśli jest konflikt między „że/CC” a inną częścią mowy i w zdaniu jest tylko jeden czasownik, to wariant „że/CC” należy odrzucić.

Przetwarzanie języka (14)

Wprowadzenie do informatyki

Regułowe oznaczanie części mowy

System ENGTWOL: około 1100 reguł

Przetwarzanie języka (15)

Wprowadzenie do informatyki

Eliminowanie niejednoznaczności

- Podejście regułowe
- Podejście stochastyczne (ukryte modele Markowa – HMM; pewne części mowy występują częściej niż inne i widać to jeszcze bardziej, gdy uwzględnimy kontekst, w którym występuje dany wyraz)

Przetwarzanie języka (16)

Wprowadzenie do informatyki

Plan wykładu

- Wprowadzenie
- Części mowy
- Gramatyka i wywód
- Język formalny
- Gramatyki bezkontekstowe
- Notacja Backusa-Naura

Przetwarzanie języka (17)

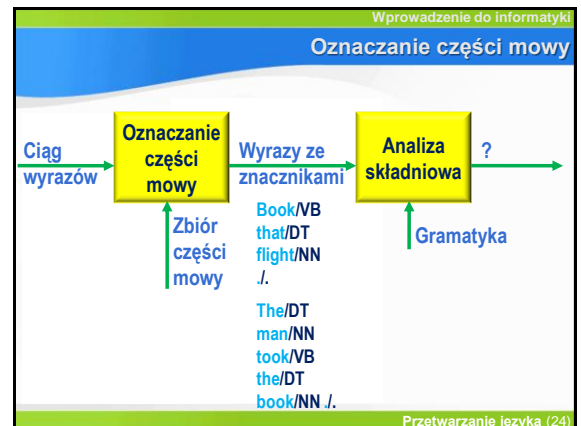
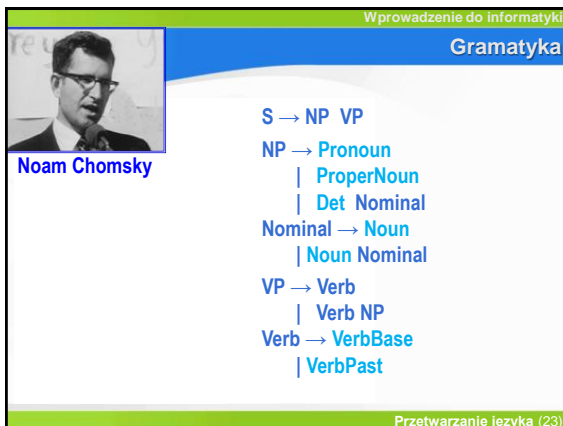
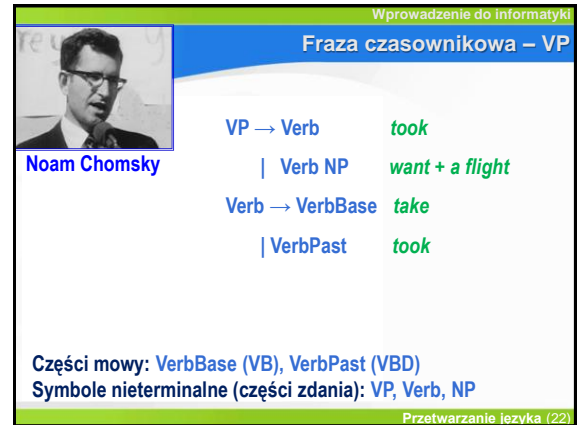
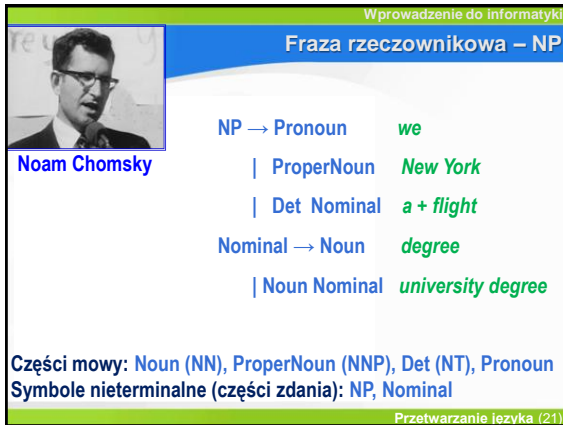
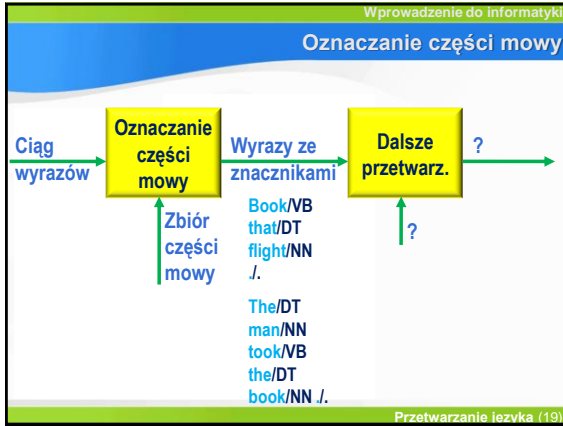
Wprowadzenie do informatyki

Oznaczanie części mowy

Book/VB that/DT flight/NN ./.

The/DT man/NN took/VB the/DT book/NN ./.

Przetwarzanie języka (18)



Wprowadzenie do informatyki

Gramatyka

$S \rightarrow NP VP$

$NP \rightarrow$ Pronoun
| ProperNoun
| Det Nominal

Nominal \rightarrow Noun
| Noun Nominal

$VP \rightarrow$ Verb
| Verb NP

Verb \rightarrow VerbBase
| VerbPast

The/DT
man/NN
took/VB
the/DT
book/NN
./.

Przetwarzanie języka (25)

Wprowadzenie do informatyki

Drzewo rozbioru (parse tree)

```

graph TD
    S[S] --- NP1[NP]
    S --- VP[VP]
    NP1 --- The[The]
    NP1 --- man[man]
    VP --- took[took]
    VP --- NP2[NP]
    NP2 --- the[the]
    NP2 --- book[book]
    
```

Przetwarzanie języka (26)

Wprowadzenie do informatyki

Oznaczenie części mowy

Ciąg wyrazów

→

Oznaczenie części mowy

Wyrazy ze znacznikami

→

Analiza składniowa

Drzewo rozbioru

→

Zbiór części mowy

↑

Book/VB
that/DT
flight/NN
./.

↑

Gramatyka

↑

The/DT
man/NN
took/VB
the/DT
book/NN ./.

Przetwarzanie języka (27)

Wprowadzenie do informatyki

Produkcje i wywód

$1^+ = \{1, 11, 111, \dots\}$

Symbol początkowy: S

Produkcje: 1) $S \rightarrow 1$
(reguły zastępowania) 2) $S \rightarrow S1$

Wywód:

1: $S \xrightarrow{1} 1$

111: $S \xrightarrow{2} S1 \xrightarrow{2} S11 \xrightarrow{1} 111$

Przetwarzanie języka (28)

Wprowadzenie do informatyki

Inne produkcje

- 1) $S \rightarrow AB$
- 2) $A \rightarrow 1$
- 3) $A \rightarrow A1$
- 4) $B \rightarrow 0$
- 5) $B \rightarrow B0$

Wywód:

10: $S \xrightarrow{1} AB \xrightarrow{2} 1B \xrightarrow{4} 10$

100:

Przetwarzanie języka (29)

Wprowadzenie do informatyki

Gramatyka

$S \rightarrow AB$

$A \rightarrow 1$

$A \rightarrow A1$

$B \rightarrow 0$

$B \rightarrow B0$

• Symbol początkowy

Przetwarzanie języka (30)

Wprowadzenie do informatyki

Gramatyka

$S \rightarrow AB$
 $A \rightarrow 1$
 $A \rightarrow A1$
 $B \rightarrow 0$
 $B \rightarrow B0$

- Symbol początkowy
- Symbole nieterminalne $N = \{S, A, B\}$

Przetwarzanie języka (31)

Wprowadzenie do informatyki

Gramatyka

$S \rightarrow AB$
 $A \rightarrow 1$
 $A \rightarrow A1$
 $B \rightarrow 0$
 $B \rightarrow B0$

- Symbol początkowy
- Symbole nieterminalne $N = \{S, A, B\}$
- Symbole terminalne $T = \{0, 1\}$

Przetwarzanie języka (32)

Wprowadzenie do informatyki

Gramatyka


$S \rightarrow AB$
 $A \rightarrow 1$
 $A \rightarrow A1$
 $B \rightarrow 0$
 $B \rightarrow B0$

- Symbol początkowy
- Symbole nieterminalne $N = \{S, A, B\}$
- Symbole terminalne $T = \{0, 1\}$
- Produkcje

Przetwarzanie języka (33)

Wprowadzenie do informatyki

Plan wykładu



- Wprowadzenie
- Części mowy
- Gramatyka i wywód
- Język formalny
- Gramatyki bezkontekstowe
- Notacja Backusa-Naura

Przetwarzanie języka (34)

Wprowadzenie do informatyki

Domknięcie relacji wywodu

- 1) $S \rightarrow AB$
- 2) $A \rightarrow 1$
- 3) $A \rightarrow A1$
- 4) $B \rightarrow 0$
- 5) $B \rightarrow B0$

Wywód:

$S \xrightarrow{1} AB \xrightarrow{2} 1B \xrightarrow{4} 10$
 $S \xrightarrow{3} A1$

$S \xrightarrow{*} 10$ Z S można wywieść 10 stosując 1 lub więcej produkcji

Przetwarzanie języka (35)

Wprowadzenie do informatyki

Zbiór ciągów nad alfabetem

$S \rightarrow AB$
 $A \rightarrow 1$
 $A \rightarrow A1$
 $B \rightarrow 0$
 $B \rightarrow B0$

Alfabet = Zbiór symboli terminalnych
 $T = \{0, 1\}$

Zbiór ciągów nad alfabetem T^* :

Zbiór wszystkich ciągów skończonych zbudowanych z elementów zbioru T .

Jeśli $T = \{0, 1\}$ to $T^* = \{\epsilon, 0, 1, 00, 01, 10, 11, 000, \dots\}$

Jeśli $T = \{a, b, c\}$ to $T^* = \{\epsilon, a, b, c, aa, ab, ac, ba, bb, bc, \dots\}$

Przetwarzanie języka (36)

Wprowadzenie do informatyki

Język formalny

Gramatyka $G = \langle S, N, T, P \rangle$

S – Symbol początkowy
 N – Zbiór symboli nieterminalnych
 T – Zbiór symboli terminalnych
 P – Zbiór produkcji

Język formalny L zdefiniowany przez gramatykę G:

$$L(G) = \{x \in T^* : S \xRightarrow{+} x\}$$

Przetwarzanie języka (37)

Wprowadzenie do informatyki

Język formalny

- 1) $S \rightarrow AB$
- 2) $A \rightarrow 1$
- 3) $A \rightarrow A1$
- 4) $B \rightarrow 0$
- 5) $B \rightarrow B0$

$$L(G) = \{x \in T^* : S \xRightarrow{+} x\}$$

$$S \xRightarrow{1} AB \xRightarrow{2} 1B$$

$$S \xRightarrow{+} 1B$$

Czy $1B$ należy do $L(G)$?

$11 \in T^*$

Czy 11 należy do $L(G)$?

Przetwarzanie języka (38)

Wprowadzenie do informatyki

Równoważność gramatyk

Gramatyki G_1 i G_2 są równoważne wtedy i tylko wtedy, gdy


$$L(G_1) = L(G_2)$$

G1	G2	G3
$S \rightarrow AB$	$S \rightarrow S0$	$S \rightarrow 1S$
$A \rightarrow 1$	$S \rightarrow A0$	$S \rightarrow 1A$
$A \rightarrow A1$	$A \rightarrow 1$	$A \rightarrow 0$
$B \rightarrow 0$	$A \rightarrow A1$	$A \rightarrow 0A$
$B \rightarrow B0$		

Przetwarzanie języka (39)

Wprowadzenie do informatyki

Plan wykładu




- Wprowadzenie
- Części mowy
- Gramatyka i wywód
- Język formalny
- Gramatyki bezkontekstowe
- Notacja Backusa-Naura


Przetwarzanie języka (40)

Wprowadzenie do informatyki

Klasyfikacja Chomsky'ego



Noam Chomsky



Gramatyki klasy 0

Gramatyki kontekstowe


Gramatyki bezkontekstowe

Gramatyki liniowe

Przetwarzanie języka (41)

Wprowadzenie do informatyki

Klasyfikacja Chomsky'ego



Gramatyki liniowe

Przetwarzanie języka (42)

Wprowadzenie do informatyki

Gramatyki liniowe

a+ b+

1. $S \rightarrow S b$	1. $S \rightarrow a S$
2. $S \rightarrow A b$	2. $S \rightarrow a B$
3. $A \rightarrow a$	3. $B \rightarrow b B$
4. $A \rightarrow A a$	4. $B \rightarrow b$

Lewoliniowa Prawoliniowa

Twierdzenie.
Dla każdego wyrażenia regularnego istnieje gramatyka lewoliniowa (prawoliniowa) opisująca ten sam język.

Przetwarzanie języka (43)

Wprowadzenie do informatyki

Klasyfikacja Chomsky'ego

Przetwarzanie języka (44)

Wprowadzenie do informatyki

Gramatyka bezkontekstowa

1. $W \rightarrow (W)$
2. $W \rightarrow 1$

Jeden nieterminal

Przetwarzanie języka (45)

Wprowadzenie do informatyki

Gramatyka bezkontekstowa

1. $W \rightarrow S$
2. $W \rightarrow W + S$
3. $S \rightarrow C$
4. $S \rightarrow S * C$
5. $C \rightarrow L$
6. $C \rightarrow (W)$
7. $L \rightarrow 1$
8. $L \rightarrow 2$
9. $L \rightarrow 3$

Jeden nieterminal

Przetwarzanie języka (46)

Wprowadzenie do informatyki

Gramatyki bezkontekstowe potrafią więcej

Gramatyki bezkontekstowe
Gramatyki liniowe

Przetwarzanie języka (47)

Wprowadzenie do informatyki

Gramatyki bezkontekstowe potrafią więcej

Język $0^n 1^n$, gdzie $n \geq 1$.

0011 OK.
0001 Error
 $S \rightarrow 0 S 1$
 $S \rightarrow 0 1$

Język $0^n 1^k$, gdzie $n, k \geq 1$.

0001 OK.
1000 Error
 $S \rightarrow 0 S$
 $S \rightarrow 0 J$
 $J \rightarrow 1 J$
 $J \rightarrow 1$

Gramatyki bezkontekstowe
Gramatyki liniowe

Przetwarzanie języka (48)

Wprowadzenie do informatyki

Gramatyki bezkontekstowe potrafią więcej

Język $0^n 1^n$, gdzie $n \geq 1$.

Gramatyki bezkontekstowe

Gramatyki liniowe

Przetwarzanie języka (49)

Wprowadzenie do informatyki

Klasyfikacja Chomsky'ego

Gramatyki kontekstowe

Gramatyki bezkontekstowe

Gramatyki liniowe

Przetwarzanie języka (50)

Wprowadzenie do informatyki

Gramatyka kontekstowa

1. $S \rightarrow aXY$
2. $S \rightarrow aSXY$
3. $aX \rightarrow ab$
4. $bX \rightarrow bb$
5. $cX \rightarrow cc$
6. $bY \rightarrow bc$
7. $cY \rightarrow cc$

Przetwarzanie języka (51)

Wprowadzenie do informatyki

Plan wykładu

- Wprowadzenie
- Części mowy
- Gramatyka i wywód
- Język formalny
- Gramatyki bezkontekstowe
- Notacja Backusa-Naura

Przetwarzanie języka (52)

Wprowadzenie do informatyki

Rozszerzona notacja Backusa-Naura

John Backus

Produkcje + wyrażenia regularne

$\langle C \rangle ::= '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'$

$\langle L \rangle ::= \langle C \rangle^+$

$\langle L \rangle ::= \langle C \rangle^* \langle C \rangle$

$\langle S \rangle ::= (\langle L \rangle^*)^* \langle L \rangle$

$\langle W \rangle ::= (\langle S \rangle^+)^* \langle S \rangle$

Przetwarzanie języka (53)

Wprowadzenie do informatyki

Przejście z EBNF na gramatyki

$\langle J \rangle ::= \langle A \rangle^* \langle B \rangle$

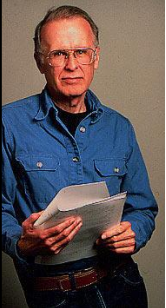
$J \rightarrow B$

$J \rightarrow A J$

Przetwarzanie języka (54)

Wprowadzenie do informatyki

Przejęcie z EBNF na gramatyki



```

<C> ::= '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' |
      '8' | '9'
<L> ::= <C>* <C>
<S> ::= (<L> (*)) * <L>
<W> ::= (<S> '+' ) * <S>
    
```

```

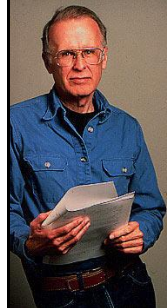
C → '0'
C → '1'
C → '2'
...
C → '9'
    
```

John Backus

Przetwarzanie języka (55)

Wprowadzenie do informatyki

Przejęcie z EBNF na gramatyki



```

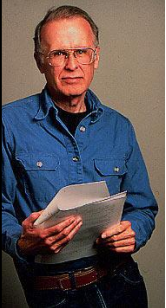
<C> ::= '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' |
      '8' | '9'
<L> ::= <C>* <C>
<S> ::= (<L> (*)) * <L>
<W> ::= (<S> '+' ) * <S>
<J> ::= <A>* <B>
J → B
J → A J
L → C
L → C L
    
```

John Backus

Przetwarzanie języka (56)

Wprowadzenie do informatyki

Przejęcie z EBNF na gramatyki



```

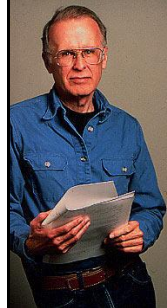
<C> ::= '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' |
      '8' | '9'
<L> ::= <C>* <C>
<S> ::= (<L> (*)) * <L>
<W> ::= (<S> '+' ) * <S>
<J> ::= <A>* <B>
J → B
J → A J
S → L
S → L (*) S
    
```

John Backus

Przetwarzanie języka (57)

Wprowadzenie do informatyki

Przejęcie z EBNF na gramatyki




```

<C> ::= '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' |
      '8' | '9'
<L> ::= <C>* <C>
<S> ::= (<L> (*)) * <L>
<W> ::= (<S> '+' ) * <S>
<J> ::= <A>* <B>
J → B
J → A J
W → S
W → S '+' W
    
```

John Backus

Przetwarzanie języka (58)

Wprowadzenie do informatyki



Podsumowanie

Przetwarzanie języka (59)

Wprowadzenie do informatyki

Przykładowa budowa narzędzia NLP

```

Tekst → Podział tekstu na zdania
        ↓ Ciąg zdań
        Podział zdania na wyrazy
        ↓ Ciąg wyrazów
        Oznaczenie części mowy
        ↓ Wyrazy ze znacznikami
        Lematyzacja
        ↓ Dochodzi postać podst. wyrazów
        Analiza składniowa
        ↓ Drzewo rozbioru składniowego
    
```

Przetwarzanie języka (60)

Wprowadzenie do informatyki

Oznaczanie części mowy

The grand jury commented on a number of other topics.

↓

The/DT grand/JJ jury/NN commented/VBD on/IN
a/DT number/NN of/IN other/JJ topics/NNS ./.

Przetwarzanie języka (61)

Wprowadzenie do informatyki

Lematyzer

The/DT man/NN took/VB the/DT book/NN ./.

↓ ↓ ↓ ↓ ↓ ↓

the man take the book .

Przetwarzanie języka (62)

Wprowadzenie do informatyki

Drzewo rozbioru (parse tree)

```

    graph TD
      S[S] --- NP1[NP]
      S --- VP[VP]
      NP1 --- The[The]
      NP1 --- man[man]
      VP --- took[took]
      VP --- NP2[NP]
      NP2 --- the[the]
      NP2 --- book[book]
    
```

Przetwarzanie języka (63)

Wprowadzenie do informatyki

Podsumowanie

- Gramatyka formalna
- Wywód zdania
- Język formalny
- Gramatyki bezkontekstowe
- Notacja EBNF

Przetwarzanie języka (64)

Wprowadzenie do informatyki

Literatura

Daniel Jurafsky, James Martin:
Speech and Language
Processing, Prentice-Hall, 2008.

Przetwarzanie języka (65)

Wprowadzenie do informatyki

Dziękuję za uwagę!

Przetwarzanie języka (66)