

## Definicja języka

- Podstawą każdego języka jest **słownik**. W językach formalnych elementy słownika (słowa) nazywamy **symbolami** (podstawowymi).
- Ciąg słów zbudowany według określonych reguł nazywamy **zdaniem**.
- Zbiór reguł lub formuł, które określają zbiór formalnie poprawnych zdań nazywamy **gramatyką** lub **składnią**. Składnia lub struktura języka pozwala na określenie, czy dany ciąg słów jest zdaniem.
- Składnia i **semantyka** (znaczenie) języka są ze sobą blisko powiązane.

## Język formalny

---

- **Język formalny** to podzbiór zbioru łańcuchów skończonych złożonych z elementów zbioru skończonego nazywanego **alfabetem**.
- **Język formalny** jest generowany przez **gramatykę formalną**, przy czym ten sam język może być generowany przez wiele gramatyk.

### PRZYKŁAD 1

---

```
<zdanie> ::= <podmiot> <orzeczenie>  
<podmiot> ::= kwiaty | gwiazdy  
<orzeczenie> ::= kwitną | świecą
```

#### Zdania wyprowadzone z tej gramatyki:

- kwiaty kwitną
- gwiazdy świecą
- gwiazdy kwitną
- kwiaty świecą

## Notacja BNF (Backus-Naur-Form)

- **Symbol początkowy**: <zdanie> ,
- **symbole pomocnicze** (nieterminalne): <zdanie>, <podmiot>, <orzeczenie>;
- **symbole końcowe** (terminalne): kwiaty, gwiazdy, kwitną, świecą;
- **metasymbole** notacji BNF: <, >, ::=, |
- **produkcje** – reguły, poprzez które można zdefiniować język.

### PRZYKŁAD 1

$S ::= A B$   
 $A ::= x \mid y$   
 $B ::= z \mid w$

Zdania wyprowadzone z gramatyki:

xz, yz, xw, yw.

Stosując **regułę kolejnych wyprowadzeń** można zdanie języka wyprowadzić z **symbolu początkowego**:

$S \rightarrow AB \rightarrow xB \rightarrow xw$

$S \rightarrow AB \rightarrow yB \rightarrow yz$

$S \xrightarrow{*} yz$

## Gramatyka Chomsky'ego

– matematyczna definicja języka

---

1. Język  $L = L(T, N, P, S)$  jest definiowany przez:
  - $T$  – słownik symboli końcowych;
  - $N$  – zbiór symboli pomocniczych (kategorii gramatycznych);
  - $P$  – zbiór produkcji (reguł syntaktycznych);
  - $S$  – symbol początkowym, należący do  $N$ , nazywany też głową języka.

## Gramatyka Chomsky'ego, cd.

---

2. Język  $L(T, N, P, S)$  jest zbiorem ciągów symboli końcowych  $\xi$ , które mogą być wyprowadzone z  $S$  zgodnie z podaną poniżej regułą 3.

$$L = \{ \xi \mid S \xrightarrow{*} \xi \text{ i } \xi \in T^* \}$$

gdzie: litery greckie oznaczają ciągi symboli,

$T^*$  oznacza zbiór wszystkich ciągów symboli z  $T$ .

## Gramatyka Chomsky'ego, cd.

---

3. Ciąg  $\delta_n$  może być wyprowadzony z ciągu  $\delta_0$  wtedy i tylko wtedy, gdy istnieją ciągi  $\delta_1, \delta_2, \dots, \delta_{n-1}$  takie, że każdy ciąg  $\delta_i$  może być bezpośrednio wyprowadzony z  $\delta_{i-1}$  zgodnie z podaną poniżej regułą 4.

$$(\delta_0 \xrightarrow{*} \delta_n) \leftrightarrow ((\delta_{i-1} \rightarrow \delta_i) \text{ dla } i = 1, \dots, n)$$

## Gramatyka Chomsky'ego, cd.

---

4. Ciąg  $\eta$  może być bezpośrednio wyprowadzony z ciągu  $\xi$  wtedy i tylko wtedy gdy istnieją ciągi  $\alpha, \beta, \xi', \eta'$  takie, że:
- $\xi = \alpha \xi' \beta$
  - $\eta = \alpha \eta' \beta$
  - P** zawiera produkcję  $\xi' ::= \eta'$

## PRZYKŁAD 2

Gramatyka, w której za pomocą rekursji można wyprowadzić nieskończenie wiele zdań za pomocą skończonego zbioru produkcji.

$$\begin{aligned} S &::= xA \\ A &::= z \mid yA \end{aligned}$$

Zbiór zdań jakie mogą być wyprowadzone z symbolu początkowego S:

xz  
xyz  
xyyz  
xyyyz  
.....

## PRZYKŁAD

Gramatyka, za pomocą której można zdefiniować liczby całkowite:

$$\begin{aligned} \langle \text{liczba całkowita} \rangle &::= \langle \text{cyfra} \rangle \\ &\quad | \langle \text{liczba całkowita} \rangle \langle \text{cyfra} \rangle \\ \langle \text{cyfra} \rangle &::= 0|1|2|3|4|5|6|7|8|9 \end{aligned}$$

## Hierarchia Chomsky'ego

---

Hierarchia Chomsky'ego dzieli języki na cztery typy.

typ 0 – języki rekursywnie przeliczalne

języki rekursywne

typ 1 – języki kontekstowe

typ 2 – języki bezkontekstowe

typ 3 – języki regularne

## Języki regularne

---

**Języki regularne** – języki formalne generowane przez wyrażenia regularne, tzn. definiowane za pomocą produkcji o postaci:

$A ::= a$

$A ::= aB$

$A ::= \varepsilon$

gdzie:  $A \in N, B \in N, a \in T$



## Języki bezkontekstowe i kontekstowe

---

**Język bezkontekstowy** – język formalny zdefiniowany przez zbiór produkcji bezkontekstowych, tzn. produkcji o postaci:

$$A ::= \xi \quad \text{dla } (A \in N, \xi \in (N \cup T)^*)$$

**Język kontekstowy** – język formalny zdefiniowany przez zbiór produkcji o postaci:

$$\alpha A \beta ::= \alpha \xi \beta \quad \text{dla nie pustego } \xi$$

## Języki rekursywne i rekursywnie przeliczalne

---

**Język rekursywny** – język formalny, dla którego istnieje algorytm rozstrzygający, czy dany łańcuch należy do tego języka, czy nie.

**Język rekursywnie przeliczalny** – język formalny, dla którego istnieje algorytm decydujący, czy dany łańcuch należy do tego języka, czy nie; przy tym algorytm ten musi zakończyć działanie z odpowiedzią pozytywną dla łańcuchów należących do języka, zaś dla łańcuchów spoza języka może dać odpowiedź negatywną lub w ogóle nie dać odpowiedzi.

## Analiza zdań

---

- Analiza polega na rozbiórze struktur zdaniowych i zdań.
- Zadaniem teorii **analizy składniowej** jest opracowanie algorytmów rozbioru języka o skomplikowanych regułach strukturalnych.

## Rozbiór zstępujący „Top down”

---

Rozbiór zstępujący lub generacyjny (ang. *Top down*) polega na wyznaczeniu kroków generacyjnych od symbolu początkowego do zdania końcowego.

## Przykład 1

$\langle \text{zdanie} \rangle ::= \langle \text{podmiot} \rangle \langle \text{orzeczenie} \rangle$   
 $\langle \text{podmiot} \rangle ::= \text{kwiaty} \mid \text{gwiazdy}$   
 $\langle \text{orzeczenie} \rangle ::= \text{kwitną} \mid \text{świecą}$

Czy zdanie "gwiazdy świecą" należy do języka?

$\langle \text{zdanie} \rangle$		gwiazdy świecą
$\langle \text{podmiot} \rangle \langle \text{orzeczenie} \rangle$		gwiazdy świecą
gwiazdy	$\langle \text{orzeczenie} \rangle$	gwiazdy świecą
	$\langle \text{orzeczenie} \rangle$	świecą
	świecą	świecą
	---	---

## Przykład 2

$S ::= xA$   
 $A ::= z \mid yA$

Czy zdanie "xyyz" należy do języka?

S	xyyz
xA	xyyz
A	yyz
yA	yyz
A	yz
yA	yz
A	z
z	z
--	--

### Przykład 3

$S ::= A \mid B$   
 $A ::= xA \mid y$   
 $B ::= xB \mid z$

Dokonaj rozbioru zdania "xxxz"

S	xxxz
A	xxxz
xA	xxxz
A	xxz
xA	xxz
A	xz
xA	xz
A	z

### Gramatyki klasy LL(1)

#### REGUŁA 1:

Dla zadanej gramatyki:

$$A ::= \xi_1 \mid \xi_2 \mid \dots \mid \xi_n$$

zbiory pierwszych symboli w zdaniach, które mogą być wyprowadzone z  $\xi_i$  muszą być rozłączne, tzn.

$$\text{pierw}(\xi_i) \cap \text{pierw}(\xi_j) = \emptyset \text{ dla wszystkich } i \neq j.$$

## Gramatyki klasy LL(1)

Zbiór  $\text{pierw}(\xi)$  jest zbiorem wszystkich symboli końcowych, które mogą wystąpić na pierwszej pozycji w zdaniach wyprowadzonych z  $\xi$ . Zbiór ten może być wyznaczony według następujących reguł:

1. jeśli pierwszy symbol argumentu jest symbolem końcowym, to

$$\text{pierw}(a\xi) = \{a\}$$

2. jeśli pierwszy symbol jest symbolem pomocniczym i istnieje produkcja

$$A ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

to

$$\text{pierw}(A\xi) = \text{pierw}(\alpha_1) \cup \text{pierw}(\alpha_2) \cup \dots \cup \text{pierw}(\alpha_n)$$

### Przykład 3

$$\begin{aligned} S &::= A \mid B \\ A &::= xA \mid y \\ B &::= xB \mid z \end{aligned}$$



$$\begin{aligned} S &::= C \mid xS \\ C &::= y \mid z \end{aligned}$$

Dokonaj rozbioru zdania "xxxz"

S	xxxz
xS	xxxz
S	xxz
xS	xxz
S	xz
xS	xz
S	z
C	z
z	z
--	--

## Lewostronna faktoryzacja

Produkcję postaci:

$$A ::= \alpha\xi_1 \mid \alpha\xi_2 \mid \dots \mid \alpha\xi_n$$

należy zastąpić produkcjami:

$$A ::= \alpha A'$$

$$A' ::= \xi_1 \mid \xi_2 \mid \dots \mid \xi_n$$

## Przykład 4

*Dana jest gramatyka:*

$$\begin{array}{l} S ::= Ax \\ A ::= x \mid \varepsilon \end{array}$$

*gdzie:  $\varepsilon$  - pusty ciąg symboli.*

*Dokonaj rozbioru zdania "x"*

S		x
Ax		x
xx		x
x		--

## Gramatyki klasy LL(1)

---

### REGUŁA 2:

Dla każdego symbolu  $A \in N$ , z którego można wyprowadzić pusty ciąg symboli ( $A \xrightarrow{*} \varepsilon$ ), zbiór jego pierwszych symboli musi być rozłączny ze zbiorem symboli, które mogą następować po dowolnym ciągu wyprowadzonym z  $A$ , tzn.

$$\text{pierw}(A) \cap \text{nast}(A) = \emptyset$$

## Gramatyki klasy LL(1)

---

Zbiór  $\text{nast}(A)$  wyznacza się następująco:  
dla każdej produkcji  $P_i$  postaci:

$$X ::= \xi A \eta$$

przez  $S_i$  oznaczamy  $\text{pierw}(\eta_i)$ ; suma wszystkich zbiorów  $S_i$  tworzy zbiór  $\text{nast}(A)$ . Jeśli z co najmniej jednego  $\eta_i$  możemy wyprowadzić pusty ciąg symboli, to zbiór  $\text{nast}(X)$  musi być również włączony do  $\text{nast}(A)$ .

## Rekurencja

---

Produkcja:

$$A ::= B \mid AB$$

opisuje zbiór ciągów: B, BB, BBB, ...

Jej użycie zgodnie z regułą 1 jest niedozwolone,  
ponieważ:

$$\text{pierw}(B) \cap \text{pierw}(AB) = \text{pierw}(B) \neq \emptyset$$

## Rekurencja

---

Produkcja:

$$A ::= \varepsilon \mid AB$$

opisuje zbiór ciągów:  $\varepsilon$ , B, BB, BBB, ...

Jej użycie zgodnie z regułą 2 jest niedozwolone,  
ponieważ:

$$\text{pierw}(A) = \text{pierw}(B)$$

a zatem

$$\text{pierw}(A) \cap \text{nast}(A) \neq \emptyset.$$



## Eliminacja lewostronnej rekurencji

I i II Reguła gramatyczna nie pozwalają na stosowanie definicji lewostronnie rekurencyjnych.

Problem ten można rozwiązać poprzez:

- zamianę rekurencji lewostronnej na prawostronną  
 $A ::= \varepsilon \mid BA$
- zastąpienie rekursji iteracją.

Zapis  $\{B\}$  w notacji EBNF oznacza iterację czyli powtórzenie symbolu B zero, jeden, dwa, ... lub nieskończenie wiele razy. Produkcja:

$$A ::= \{B\}$$

opisuje zbiór ciągów:  $\varepsilon, B, BB, BBB, \dots$

## Przykład 5

Dana jest gramatyka:

$$\begin{array}{l} S ::= A \mid S - A \\ A ::= a \mid b \mid c \end{array} \quad \Rightarrow \quad a - b - c = ((a - b) - c)$$

gdzie:  $\{a, b, c, -\} \in T$

$$\begin{array}{l} S ::= A \mid A - S \\ A ::= a \mid b \mid c \end{array} \quad \Rightarrow \quad (a - (b - c))$$

Te dwie gramatyki nie są równoważne semantycznie

## Eliminacja lewostronnej rekurencji

Produkcję postaci:

$$A ::= A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$

można zastąpić produkcjami:

$$A ::= \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A'$$

$$A' ::= \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \varepsilon$$

## Przykład 5

Dana jest gramatyka:

$$\begin{array}{l} S ::= A \mid S - A \\ A ::= a \mid b \mid c \end{array}$$

$$\Rightarrow a - b - c = ((a - b) - c)$$

$$\begin{array}{l} S ::= A \{- A\} \\ A ::= a \mid b \mid c \end{array}$$

$$\Rightarrow a - b - c$$

$$\begin{array}{l} S ::= A S' \\ S' ::= -AS' \mid \varepsilon \\ A ::= a \mid b \mid c \end{array}$$

$$\Rightarrow a - b - c$$

Te gramatyki są równoważne semantycznie

## Drzewo wyprowadzeń

---

- **Drzewo wyprowadzeń** jest to hierarchiczna struktura, która reprezentuje w sposób graficzny wyprowadzenie zdania z gramatyki.
- **Korzeń** drzewa ma etykietę głowy języka. **Gałęzie** są symbolami pomocniczymi, a **liście** symbolami końcowymi lub pustymi. **Struktura gałęzi** reprezentuje produkcje. Gałąź główna ma etykietę lewej strony produkcji, a gałęzie z niej wyprowadzone etykiety prawej strony produkcji w kolejności od strony lewej do prawej.

## Przykład 6

---

Należy wyprowadzić liczbę 123.456 z gramatyki:

```
<liczba rzeczywista> ::= <część całkowita> . <ułamek>  
<część całkowita> ::= <cyfra> | <część całkowita> <cyfra>  
<ułamek> ::= <cyfra> | <cyfra> <ułamek>  
<cyfra> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
```

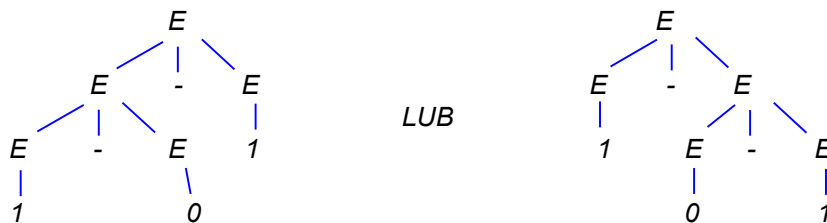
## Niejednoznaczność składni.

Problem **niejednoznaczności gramatyki**, występuje wtedy, gdy generuje ona zdania, które mają więcej niż jedno drzewo wyprowadzeń.

### Przykład 7

Należy wyprowadzić zdanie 1 - 0 - 1 dla gramatyki:

$$E ::= E - E \mid 0 \mid 1$$



Dla zdania 1 - 0 - 1 można narysować dwa różne drzewa wyprowadzeń.

## Przykład 8

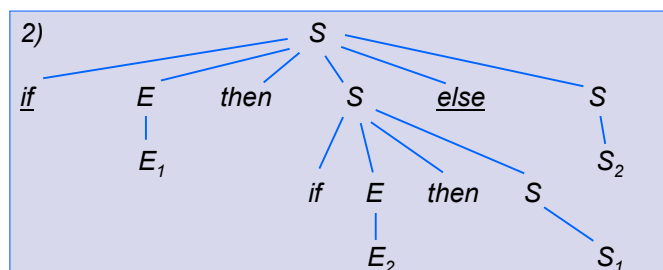
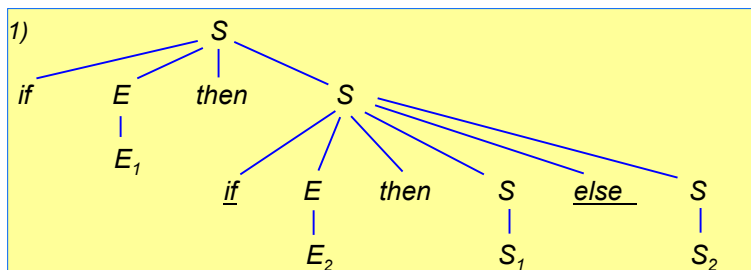
Kolejna niejednoznaczność występuje w języku C i Pascal dla instrukcji warunkowej *IF*.

Jest ona spowodowana wystąpieniem produkcji:

$S :: = \text{if } E \text{ then } S$   
 $S :: \text{if } E \text{ then } S \text{ else } S$       *Pascal*

Dla zdania *if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub>* można wprowadzić dwa różne drzewa.

*if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub>*



## Eliminacja niejednoznaczności if-else

---

Problemem niejednoznaczności *if-else* w językach *Pascal* i *C* został rozwiązany w ten sposób, że słowo kluczowe *else* jest kojarzone z najbliższym występującym słowem kluczowym *if*.

W tym przypadku rozpatrywane zdanie będzie interpretowane tak jak przedstawia to drzewo 1.

## Eliminacja niejednoznaczności if-else w języku Modula-2

---

```
S ::= <empty>
    | stmt
    | if E then SL end
    | if E then SL else SL end
SL ::= SL; S | S
```

*Pascal:*     *if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub>*

*Modula-2:*   1°. *if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub> end end*  
              2°. *if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> end else S<sub>2</sub> end*

## Pascal i Modula-2

```
if E1 then S1
    else if E2 then S2
        else if E3 then S3
            else S4
```

Pascal:

```
if E1 then S1
    else if E2 then S2
        else if E3 then S3
            else S4 end
    end
end
```

Modula-2:

## Modula-2

```
if E1 then S1
    elsif E2 then S2
        elsif E3 then S3
            else S4
    end
```

```
S ::= if E then SL { elsif E then SL } [else SL] end
SL ::= SL; S | S
```

## Zestawienie oznaczeń stosowanych w notacji **BNF** oraz w jej zmodyfikowanej wersji - **MBNF**

	ZAPIS W MBNF	ZNACZENIE	ZAPIS W BNF (EBNF)
1	=	jest zdefiniowane jako	::=
2		lub	
3	.	zakończenie formuły	brak jawnego oznaczenia
4	[x]	0 lub jednokrotne powtórzenie napisu x	nie używa się metasympoli [ ]
5	{x}	0 lub wielokrotne powtórzenie napisu x	{x} w EBNF
6	(x   y ...   z)	dowolny z napisów: x, y, ..., z	nie używa się metasympoli ( )
7	"a"	symbol terminalny (z alfabetu języka)	nie ujmuje się w cudzysłów
8	ciąg-małych-liter	symbol nieterminalny (zmienna metajęzykowa)	zmienne ujmuje się w nawiasy ostre, nie ogranicza się do małych liter i niekoniecznie używa się łącznika między wyrazami

## Separatory i terminatory

*Separatory* – oddzielają elementy

*Terminatory* – pojawiają się za każdym elementem, np. zdaniem

```

S ::= <empty>
    | stmt
    | begin SL end
SL ::= SL; S | S
    
```

*Pascal:*

```

Pascal:
begin S1; S2; S3 end
begin S1; S2; S3; end
begin ;S1;; S2; S3;; end
    
```



## Instrukcja pusta

**Pascal:**  
S ::= <empty>  
| stmt  
| if E then S  
| if E then S else S  
| begin SL end  
| while E do S  
SL ::= SL; S | S

**Modula-2:**  
S ::= <empty>  
| stmt  
| if E then SL end  
| if E then SL else SL end  
| begin SL end  
| while E do SL end  
SL ::= SL; S | S

**Pascal:** if E then ; S<sub>1</sub>; - błąd semantyczny  
if E then S<sub>1</sub>; else S<sub>2</sub>; - błąd syntaktyczny  
**Modula-2:** if E then S<sub>1</sub>; S<sub>2</sub> end  
if E then S<sub>1</sub>; S<sub>2</sub> else S<sub>3</sub>; S<sub>4</sub> end