

Języki formalne i kompilatory

(15 wyk., 15 ćw. audyt., 15 ćw. lab.)

dr inż. **Jolanta Cybulka**

Instytut Automatyki, Robotyki i Inżynierii Informatycznej, pok. 313A Piotrowo bud. A1

e-mail: Jolanta.Cybulka@put.poznan.pl

www: <http://users.man.poznan.pl/jolac>

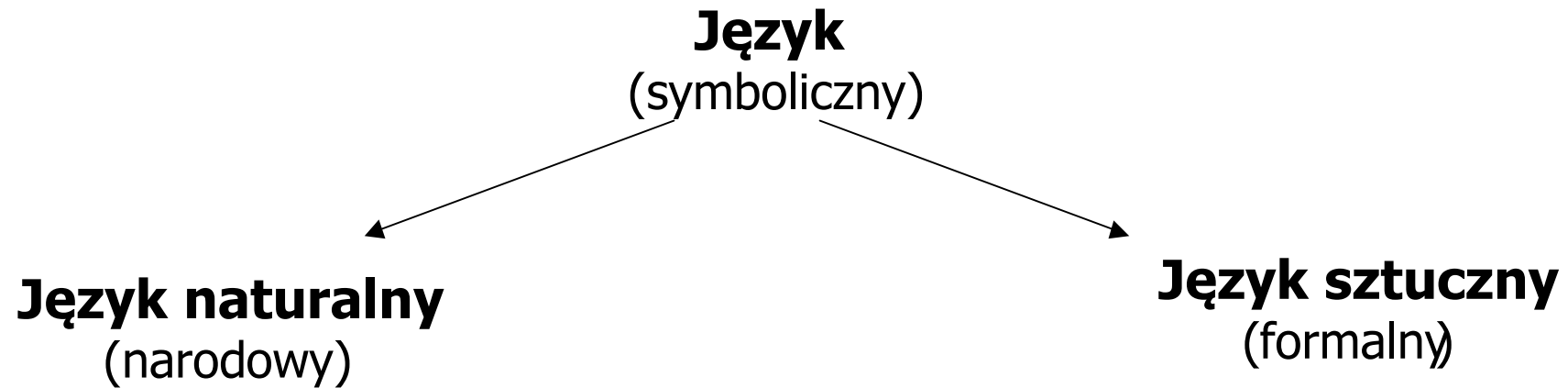
Literatura

1. Aho A V., Sethi R., Ullman J.: [Kompilatory. Reguły metody i narzędzia.](#) WNT, Warszawa 2002.
2. Cybulka J., Jankowska B., Nawrocki J. R.: [Automatyczne przetwarzanie tekstów. AWK, Lex i YACC,](#) Wyd. NAKOM, Poznań, 2002.
3. Dembiński P., Małuszyński J.: [Matematyczne metody definiowania języków programowania,](#) WNT, Warszawa 1981.
4. Hopcroft J.E., Ullman J.D.: [Wprowadzenie do teorii automatów, języków i obliczeń,](#) PWN, Warszawa, 1994.
5. Révész G.E.: [Introduction to Formal languages,](#) McGraw-Hill, New York, 1983.

Spis treści

LITERATURA	2
1. JĘZYKI FORMALNE – PODSTAWOWE POJĘCIA.....	4
2. JĘZYKI REGULARNE (KLASA "3").....	11
WYRAŻENIA REGULARNE.....	11
2.2 JEZYK LEX	14
2.2.1 Wzorce w programach.....	16
2.2.2 Zasady działania przetwornika (reguły dopasowywania wzorców)	21
2.2.3 Akcje w regułach.....	22
2.2.4 Sekcja definicji pomocniczych.....	26
2.2.5 Sekcja podprogramów pomocniczych	27
2.2.6 Uwagi o uruchamianiu generatora Lex i przetwornika tekstu.....	30
3. ELEMENTY TEORII JĘZYKÓW FORMALNYCH	31
METODY DEFINIOWANIA SKŁADNI JĘZYKA $\langle \Sigma, \text{SYN} \rangle$	31
3.2 HIERARCHIA JĘZYKÓW WG N. CHOMSKY'EGO:.....	32
3.3 FORMALNY SYSTEM GENERACJI SKŁADNI (PODEJŚCIE GENERACYJNE)	33
3.4 FORMALNA DEFINICJA AUTOMATU (PODEJŚCIE AKCEPTOROWE)	36
4. TRANSLACJA I KOMPILACJA – PODSTAWOWE POJĘCIA	37
4.1 ZAGADNIENIE TRANSLACJI.....	37
4.2 SKŁADNIA PROSTEGO JĘZYKA PROGRAMOWANIA (PJP).....	43

1. Języki formalne – podstawowe pojęcia



- rozważamy *języki symboliczne*, skonstruowane z symboli należących do zbioru (niepustego, skończonego), zwanego *alfabetem*, np.:

$D = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$

$B = \{0, 1\}$

$H = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F\}$

$Lat = \{a, b, c, \dots, z, A, B, C, \dots, Z\}$

$Pol = \{a, \acute{a}, b, c, \acute{c}, \dots, \acute{z}, dz, d\acute{z}, d\acute{z}, A, A\acute{a}, B, C, \acute{C}, \dots, Z, Dz, D\acute{z}, D\acute{z}\}$

$S = \{\neg, \wedge, \vee, \Leftrightarrow, \Rightarrow, p, f, (,)\}$

i wiele innych (alfabet Morse'a, Braille'a, symbole nutowe, symbole chemiczne, znaki migowe,...);

- z symboli alfabetu, zgodnie z zadanymi *regułami składania* symboli powstają *napisy języka* – zbiór wszystkich napisów należących do języka (i tylko takich) nazywamy jego *składnią* (syntaktyką, syntaksą);

0 123 09 55000 nad alfabetem **D**

0001 11001 0 1 nad alfabetem **B**

0 123F 09A FFAA nad alfabetem **H**

Fundamenta Informaticae Formal Language nad alfabetem **Lat**

Ala mieć ma As Asa nad alfabetem **Pol**

$p \wedge f$ p $(f \vee p) \wedge p \Leftrightarrow p$ nad alfabetem **S**

- składnia może mieć strukturę warstwową, tzn. z ciągu symboli powstają napisy elementarne, których zbiór staje się nowym alfabetem, z którego powstają nowe napisy:

nad alfabetem **L = {Ala, mieć, ma, As, Asa, \ , \n, \t...}**
 można zbudować napis **Ala ma Asa**

nad alfabetem **C = {if, else, (,), {, }, =, ;, ++, >, num, id, \ , \n, \t...}**
 można zbudować napis **if (xsr >12) {y++;} else {y=xsr;}**

- w wypadku języków formalnych, na składni można wykonywać operacje mnogościowe, jak: *suma*, *różnica*, *przecięcie*, *dopełnienie* do uniwersum, *złożenie*, *potęgowanie*, *domknięcie*;

$A = \{a, b, c\}$ alfabet

– potęgowanie: $A^0, A^1, A^2, A^3, \dots$

$A^0 = \{\varepsilon\}$ słowo długości 0, słowo „puste”

$A^1 = \{a, b, c\}$ słowa długości 1

$A^2 = \{aa, ab, ac, ba, bb, bc, ca, cb, cc\}$ słowa długości 2

$A^3 = \{aaa, aab, aac, aba, abb, abc, aca, acb, acc, baa, bab, bac, bba, bbb, bbc, bca, bcb, bcc, caa, cab, cac, cba, cbb, cbc, cca, ccb, ccc\}$ słowa długości 3

...

– domknięcie zwrotne i przechodnie operacji potęgowania: A^*

$$A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \cup \dots = \bigcup_{(i=0, \infty)} A^i$$

domknięcie przechodnie operacji potęgowania: A^+

$$A^+ = A^1 \cup A^2 \cup A^3 \cup \dots = \bigcup_{(i=1, \infty)} A^i$$

- napisy (składnia) służą do notowania znaczeń, zadaniem funkcji semantyki języka formalnego jest przypisanie każdemu elementowi składniowemu (napisowi) znaczenia z pewnej dziedziny znaczeń:

znaczeniem napisu $(2+3)*4$ nad alfabetem $A_n = \{0,1,\dots,9\} \cup A_o = \{+, * \} \cup \{ (,) \}$ jest 20 o ile

$D = \{ \underline{1}, \underline{2}, \underline{3}, \dots \}$ jest zbiorem wartości (liczb) naturalnych,

$F = \{ \mathbf{plus}, \mathbf{times} \}$ jest zbiorem funkcji (dodawania i mnożenia),

$\mathbf{val}_n : A_n \rightarrow D$

$\mathbf{val}_o : A_o \rightarrow F, \quad \mathbf{val}_o(+)= \mathbf{plus}, \quad \mathbf{val}_o(*)= \mathbf{times}$

$\mathbf{value} : (A_n \cup A_o \cup \{ (,) \})^* \rightarrow D$ (dla napisów poprawnych składniowo)

- język formalny L definiujemy matematycznie jako:

$$\mathbf{L} = \langle \Sigma, \mathbf{Syn}, \mathbf{DSem}, \mathbf{Sem} \rangle$$

Σ – *alfabet*, skończony niepusty zbiór symboli (znaków),

\mathbf{Syn} – *syntaktyka* (składnia, syntaksa), zbiór napisów języka, zbudowanych z symboli alfabetu

\mathbf{DSem} – *dziedzina znaczeń*, zbiór bytów semantycznych, przypisywanych napisom języka

\mathbf{Sem} – *relacja semantyczna* ($\mathbf{Sem} \subseteq \mathbf{Syn} \times \mathbf{DSem}$), definiująca związki między napisami języka a elementami dziedziny znaczeń (relacja ta w zastosowaniach praktycznych jest *funkcją semantyki*);

- dodatkowo, język zastosowany w pewnej dziedzinie może wytwarzać odmiany *pragmatyczne*:

$$\mathbf{L}_1 = \langle \Sigma, \mathbf{Syn}, \mathbf{DSem}_1, \mathbf{Sem}_1 \rangle \quad \mathbf{L}_2 = \langle \Sigma, \mathbf{Syn}, \mathbf{DSem}_2, \mathbf{Sem}_2 \rangle$$

Na czym może polegać *definiowanie i przetwarzanie języka formalnego*?

- a) określenie języka (wszystkich jego elementów) w sposób ścisły, za pomocą odpowiednich metod i narzędzi
- b) badanie przynależności napisu do języka (sprawdzenie poprawności napisu)
- c) translacja napisu z jednego języka na drugi (przy zachowaniu znaczenia)

2. Języki regularne (klasa "3")

Wyrażenia regularne

- *wyrażeniem regularnym* (ang. *regular expression*) nad alfabetem Σ nazywamy napis zbudowany za pomocą następującej rekurencyjnej definicji:
 - 1) \emptyset jest w. r. oznaczającym język **pusty** (pusty zbiór napisów)
 - 2) ε jest w. r. oznaczającym język $\{\varepsilon\}$
 - 3) jeśli $a \in \Sigma$, to napis a jest w. r. oznaczającym język $\{a\}$
 - 4) jeśli r_1 i r_2 są w. r., to napis r_1r_2 jest w. r. oznaczającym język $L(r_1)L(r_2)$
 - 5) jeśli r_1 i r_2 są w. r., to napis $r_1|r_2$ jest w. r. oznaczającym język $L(r_1) \cup L(r_2)$
 - 6) jeśli r jest w. r., to napis r^* jest w. r. oznaczającym język $\cup_{(i=0,\infty)} L(r)^i$
 - 7) jeśli r jest w. r., to (r) jest w. r. oznaczającym $L(r)$.

- właściwości wyrażeń regularnych; niech p , q i r oznaczają dowolne wyrażenia regularne, wtedy:

$p (q r) = (p r) q$	$p(qr) = (pq)r$
$p q = q p$	$p^* = \varepsilon pp^*$
$(p q)r = pr qr$	$r(p q) = rp rq$
$\varepsilon p = p\varepsilon = p$	$p^* = (\varepsilon p)^*$
$\emptyset p = p\emptyset = \emptyset$	if $p = r pq \wedge q \neq \varepsilon$ then $p = rq^*$

- przykłady**

Niech $\Sigma = \{0, 1, 2\}$

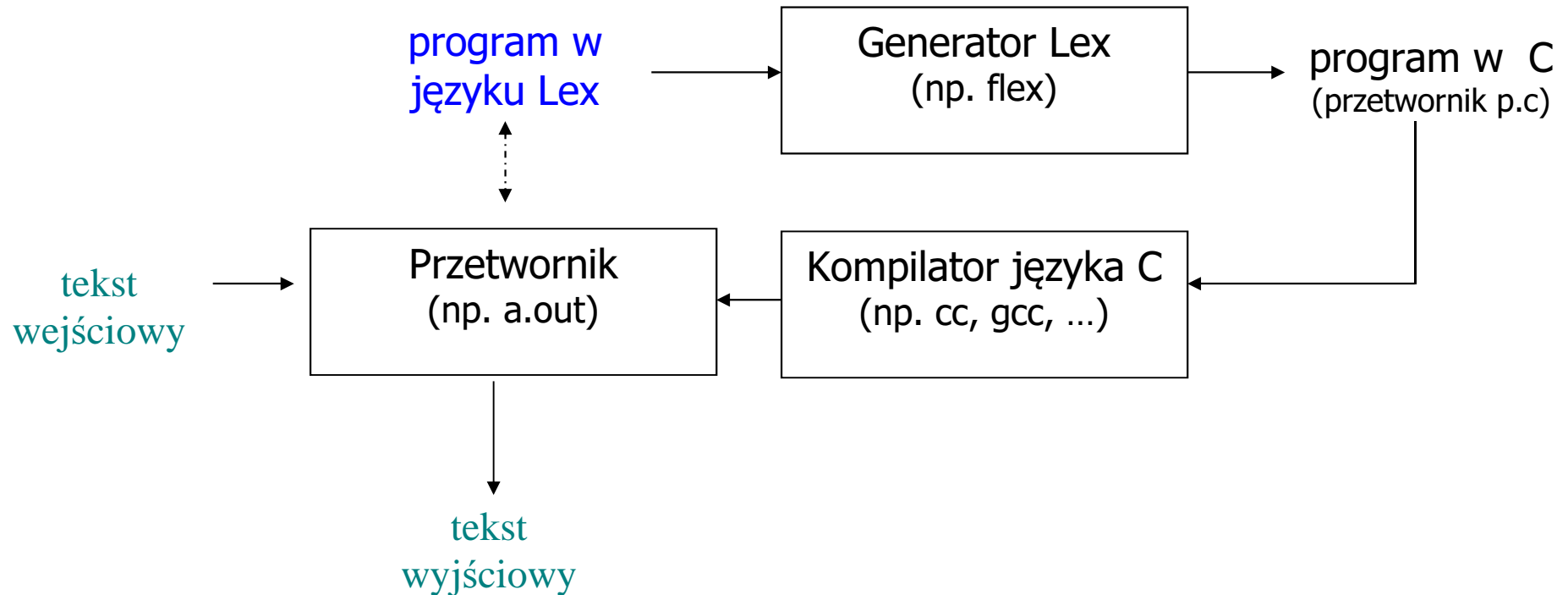
- 11 reprezentuje język $\{11\}$
- $(0 | 1)^*$ reprezentuje język wszystkich napisów zbudowanych z zer i jedynek
- $(0 | 1)^*11(0 | 1)^*$ reprezentuje język wszystkich napisów zbudowanych z zer i jedynek, w których co najmniej raz występują dwie kolejne jedynek
- $0^*1^*2^*$ reprezentuje język składający się z napisów zbudowanych z dowolnej liczby zer, po której następuje dowolna liczba jedynek, a następnie dowolna liczba dwójek;

Niech $\Sigma = \{0, 1, 2, \dots, 9, a, b, \dots, z, A, B, \dots, Z, _, -, +\}$

- 5) $(a|b|\dots|z|A|B|\dots|Z)(a|b|\dots|z|A|B|\dots|Z|0|1|\dots|9)^*$ reprezentuje język niepustych ciągów liter i cyfr, rozpoczynających się od litery
- 6) $(-|+|\epsilon)(0|1|\dots|9)(0|1|\dots|9)^*$ reprezentuje język stałych całkowitych dziesiętnych

2.2 Język Lex

- *programy* w języku Lex **opisują** działanie *przetworników tekstu*
- konkretny system Lex (np. lex, flex) jest generatorem przetworników zdefiniowanych w języku Lex:



- *struktura programu* w języku Lex:

definicje pomocnicze

%%

wzorzec 1 akcja1

wzorzec 2 akcja2

...

wzorzec n akcja n

%%

podprogramy pomocnicze

opcjonalne (dyrektywy, nazwy wyrażeń, deklaracje, ...)

obligatoryjne (reguły przetwarzania)

opcjonalne

2.2.1 Wzorce w programach

- mają postać (poszerzonych o nowe operatory) *definicji regularnych* (w cudzysłowie lub bez)
- definicje regularne są zbudowane nad pewnym *alfabetem* oraz zbiorem *nazw* (nazw wyrażeń regularnych)

np. wyrażenie regularne: `[a-zA-Z][a-zA-Z0-9]*`, nad alfabetem $A = \{a, b, \dots, z, A, B, \dots, Z, 0, 1, \dots, 9\}$

można zastąpić *definicją regularną* postaci:

`litera [a-zA-Z]`

`cyfra [0-9]`

`{litera}({litera}|{cyfra})*`

nad alfabetem

$\{a, b, \dots, z, A, B, \dots, Z, 0, 1, \dots, 9\} \cup \{litera, cyfra\}$

A =

- znaki specjalne:

`. % \ / " * () ? + | ^ $ [] - < > { }`

- priorytety operatorów tworzenia wyrażeń regularnych w porządku malejącym:
 - a) domknięcie zwrotne `*`, domknięcie dodatnie `+`, opcjonalny wybór `?`, zbiór znaków `[]`, dopełnienie zbioru znaków `[^]`, konkatenacja wielokrotna `{ }`
 - b) konkatenacja
 - c) alternatywa `|`
 - d) prawy kontekst `/`, początek wiersza `^`, koniec wiersza `$`
 - e) lewy kontekst `<>`
 - f) koniec pliku `<<EOF>>`

- **przykłady** programów z różnymi typami wzorców:

<code>%%</code>	<code>%%</code>
<code>nie ;</code>	<code>"nie" ;</code>

Wejście: Kto się nie nauczył, ten nie zda.

Wyjście: Kto się nauczył, ten zda.

`%%`

```

^[a-zA-Z][a-zA-Z_0-9]*$    {printf("nazwa: %s", yytext);}
^[0-9]+(\.[0-9]+)?$      {printf("liczba: %s", yytext);}
^[^a-zA-Z0-9]+$          {printf("inne: %s", yytext);}
.                          {;}

```

zastosowanie konkatencji wielokrotnej:

```

%%
[a-zA-Z]{1,8}(\.[a-zA-Z0-9]{0,3})?    ECHO;
.                                       {;}

```

wykorzystanie znaków specjalnych:

```

%%
\ |\t|\n    {printf("biala spacja");}
.           {;}

```

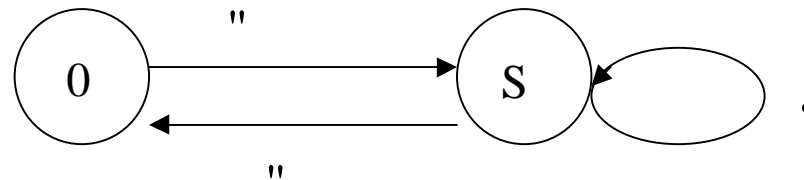
zastosowanie operatora lewego kontekstu:

`%start nazwa stanu` w części definicji pomocniczych (może być wiele)

`<nazwa stanu>` odwołanie do stanu we wzorcach reguł przetwarzania

np.:

```
%start string          definicja stanu S (pomijanie znaków łańcucha)
%%
<string>\ "          {BEGIN 0;}
<string>.            ;
\"                   {BEGIN string;} albo <0>\ "          {BEGIN string;}
```



w stanie 0 następuje przepisywanie wejścia na wyjście, w stanie string (S) – pomijanie znaków

przykład zastosowania definicji regularnych i operatora prawego kontekstu:

wyrażenie regularne / podgląd

wyodrębnienie części pasującej do *wyrażenie regularne*,
o ile występuje po niej fragment pasujący do wyrażenia regularnego *podgląd*

np.:

numb [0-9]+

nazwa wyrażenia

%%

{numb}

{printf("Card ");}

{numb} " . " {numb}? / [^ \ .]

{printf("Real ");}

" . . "

{printf("Range ");}

" . "

{printf("Dot ");}

" "

;

dla wejścia:

1..5 37.5

otrzymamy na wyjściu:

Card Range Card Real

2.2.2 Zasady działania przetwornika (reguły dopasowywania wzorców)

- *cykliczne analizowanie* pliku wejściowego, którego *przedrostki* są *dopasowywane* do *wzorców* zawartych w regułach przetwarzania; wykonywanie akcji związanych z wzorcami
- reguła *najdłuższego dopasowania* (gdy przedrostek pasuje do kilku wzorców; bierze się pod uwagę także długość ciągu podglądanego, który jest z powrotem przesyłany na wejście)
- reguła *pierwszego dopasowania* (wybór wzorca występującego wcześniej, gdy ta sama długość dopasowania)
- ustawienie wartości zmiennych globalnych: *yytext* (dopasowany ciąg znaków, *przedrostek* pliku wejściowego, łącznie z częścią *podgląd* w sytuacji zastosowania prawego kontekstu), *yytext* (długość *przedrostka* bez części dopasowanej do *podgląd*) i *yytext* (długość całości dopasowanej, łącznie z częścią dopasowaną do *podgląd*), wykonanie akcji związanej z wzorcem dopasowania i usunięcie dopasowanego *przedrostka* z pliku wejściowego (część dopasowana do *podgląd* jest zwrótnie przekazywana na wejście)
- reguła *zastępcza* (przepisanie znaku wejściowego na wyjście), w wypadku braku dopasowania.

2.2.3 Akcje w regułach

- akcja ma postać instrukcji języka C (w tym akcja pusta `;` – pominięcie dopasowanego ciągu znaków albo kreska `|` – powtórzenie akcji z następnej reguły przetwarzania)
- w akcjach można wykorzystywać *zmiennne standardowe*, *funkcje standardowe* i *makrowywołania*
- zestawienie *makrowywołań*:

Makrowywołanie	Opis
ECHO	– ECHO – skopiowanie <code>yytext</code> początkowych znaków zmiennej <code>yytext</code> do pliku wyjściowego
BEGIN <i>stan</i>	– BEGIN <i>stan</i> – zmiana bieżącego stanu początkowego przetwornika tekstu na stan <i>stan</i>
REJECT	– REJECT – przekazanie sterowania do kolejnej spośród aktualnie pasujących reguł przetwarzania
<code>yymore()</code>	– <code>yymore()</code> – modyfikacja standardowego algorytmu wyznaczania wartości <code>yytext</code> ; w kolejnym kroku wartość zmiennej <code>yytext</code> powstaje poprzez sklejenie przedrostka ciągu wejściowego dopasowanego do ustalonego w tym kroku wzorca z bieżącą wartością zmiennej <code>yytext</code>
<code>yyterminate()</code>	– odpowiada instrukcji <code>return 0;</code>

<code>yyless(n)</code>	– <code>yyless(n)</code> – zwrotne przesłanie na początek pliku wejściowego wszystkich, z wyjątkiem n początkowych, znaków dopasowanych do aktualnego wzorca
------------------------	--

- zestawienie *funkcji standardowych*:

Funkcja	Opis
<code>input()</code>	– <code>input()</code> – wartością jest jeden wczytany, początkowy znak z nieprzetworzonej części pliku wejściowego
<code>unput(c)</code>	– <code>unput(c)</code> – dołącza wskazany znak c na początek nieprzetworzonej części pliku wejściowego
<code>yywrap()</code>	– <code>yywrap()</code> – obsługuje sytuację: <ul style="list-style-type: none"> – a) zakończeniu działania, w sytuacji napotkania końca aktualnie przetwarzanego pliku i braku dalszych plików do przetwarzania; – b) kontynuacji działania w odniesieniu do nowego pliku wejściowego, w sytuacji napotkania końca aktualnie przetwarzanego pliku wejściowego i wystąpienia żądania przetwarzania kolejnego pliku;
<code>yyerror(msg,...)</code>	– <code>yyerror(msg,...)</code> – wysyła do standardowego pliku z błędami komunikat o treści msg ; po wysłaniu tego komunikatu przetwornik tekstu kontynuuje działanie na pozostałej części pliku wejściowego.

– `yyerror(msg,...)` jest równocześnie funkcją standardową systemu YACC.

• zestawienie *zmiennych standardowych*:

Zmienna	Opis
<code>yytext</code>	<code>char yytext[YYLMAX]</code> (ewentualnie <code>char *yytext</code>) – przechowuje przedrostek dopasowany do wzorca w bieżącej regule przetwarzania; jeśli wzorzec jest postaci <i>wregularne/podglad</i> , to w <code>yytext</code> znajdzie się sklejenie ciągów znaków dopasowanych do <i>rwregularne</i> i <i>podglad</i> ; pojemność <code>YYLMAX</code> tablicy <code>yytext</code> jest ustalana za pomocą dyrektywy <code>#define</code>
<code>yy_end</code>	<code>int yy_end</code> – pamięta długość ciągu znaków aktualnie przechowywanego w <code>yytext</code> ; jeśli <code>yytext</code> przechowuje przedrostek dopasowany do <i>wregularne/podglad</i> , to <code>yy_end</code> pamięta łączną długość tego przedrostka
<code>yyleng</code>	<code>int yyleng</code> – pamięta: a) długość całego ciągu znaków przechowywanego w <code>yytext</code> , jeśli jest nim przedrostek dopasowany do wyrażenia bez operatora prawego kontekstu / b) długość tego fragmentu ciągu znaków przechowywanego w <code>yytext</code> , który pasuje do <i>wregularne</i> , jeśli zmienna <code>yytext</code> przechowuje przedrostek dopasowany do wzorca z operatorem prawego kontekstu /, o postaci: <i>rwregularne/podglad</i>
<code>yylastc</code>	<code>char yylastc</code> – przechowuje ostatni znak przedrostka dopasowanego do wzorca poprzedniej

	aktywnej reguły przetwarzania; jeśli wzorzec ten był postaci <i>wregularne/podglad</i> , to jest to ostatni znak dopasowany do <i>wregularne</i>
<code>yyin</code>	<code>FILE *yyin</code> – wskazuje na aktualnie przetwarzany plik wejściowy; domyślnie jest nim <code>stdin</code> ; redefinicja przyjmuje postać <code>yyin = plikwe</code> , gdzie <i>plikwe</i> to wskazanie na nowy plik wejściowy
<code>yyout</code>	<code>FILE *yyout</code> – wskazuje na plik wyjściowy, domyślnie na <code>stdout</code> ; redefinicja za pomocą instrukcji <code>yyout = plikwy</code> , gdzie <i>plikwy</i> oznacza wskazanie na nowy plik wyjściowy
<code>yylineno</code>	<code>int yylineno</code> – przechowuje numer porządkowy aktualnie przetwarzanego wiersza pliku wejściowego
<code>yystart</code>	<code>int yystart</code> – przechowuje numer tego stanu automatu skończonego, w którym rozpoczyna się realizacja bieżącej reguły przetwarzania (automat skończony opisujący działanie przetwornika, jest budowany równoległe z opisem w języku C, standardowo umieszczanym w pliku <code>lex.yy.c</code>)

2.2.4 Sekcja definicji pomocniczych

- *deklaracje* `extern` zmiennych zewnętrznych programu
- *definicje zmiennych statycznych*, widocznych w pewnym fragmencie lub całym programie (definicje zmiennych automatycznych, których zakres widoczności ma być ograniczony do pojedynczej akcji programu należy umieścić w treści tej akcji)
- *definicje stanów początkowych* wykorzystywanych we wzorcach z operatorem lewego kontekstu; definicje te mają postać:

`%start stan`

gdzie *stan* oznacza nazwę stanu początkowego używanego w programie

- *dyrektywy* `#include` włączania innych plików do programu
 - *dyrektywy* `#define`
 - *definicje nazw wyrażeń regularnych*, postaci:
 - *nazwa* *znaczenie*
- gdzie *nazwa* oznacza kod mnemoniczny wyrażenia regularnego, a *znaczenie* definiuje jego treść

2.2.5 Sekcja podprogramów pomocniczych

- zawierają opisy działań w postaci funkcji języka C
- mogą redefiniować funkcje standardowe lub definiować specyficzne działania niestandardowe
- **przykłady** programów:

```
%{ sekcja definicji pomocniczych
```

```
int num_count = 0;
```

```
%}
```

```
%% sekcja reguł przetwarzania
```

```
[0-9]+ {printf("+"); REJECT;} 
```

```
[\+\-]?[0-9]+ {ECHO; ++num_count;} 
```

```
[\ \t\n] {ECHO;} 
```

```
%% sekcja podprogramów (redefinicja funkcji standardowej)
```

```
int yywrap(void)
```

```
{
```

```
printf("FILE CONSISTS OF %d NUMBERS", num_count);
```

```
return 1;}
```

```
%%
```

```

a          |
ab         |
abc        |
abcd       | ECHO; REJECT;
.|\n      | ;

```

powtórzenie akcji z następnej reguły

```

%{
int Top=0, Stos[100];
%}
%%
[0-9]+    {push1(atoi(ytext));}
"+"       {push1(pop()+pop());}
"*"       {push1(pop()*pop());}
[\\ \t]   ;
\n        {printf("%d\n",pop());}
.         {yyterminate();}
%%
int push1(int x)
    {Stos[++Top]=x; return 1;}
int pop(void)
    {return Stos[Top--];}

```

kalkulator wyrażeń w ONP

dane (wyrażenia arytmetyczne w ONP) i wyniki (wartości wyrażeń):

2	2
3 5 +	8
2 8 *	16
2 3 5 + *	16
2 3 5 * +	17

```
%{
#include <stdlib.h>
#include <string.h>
typedef union {char *txt; float fval; int ival;} u_type;
u_type attrib;
}%
%%
\+          {printf("+");}
\*          {printf("*");}
\(          {printf("(");}
\)          {printf(")");}
[0-9]*\.[0-9]+ {attrib.fval = atof(yytext);ECHO;}
[0-9]+\.[0-9]+ {attrib.fval = atof(yytext);ECHO;}
[0-9]+       {attrib.ival = atoi(yytext);ECHO;}
[A-Za-z][A-Za-z0-9_]* {attrib.txt = strdup(yytext);ECHO;}
[\ \t\n]    {;}
.           {yyerror("Unexpected symbol");}
```

2.2.6 Uwagi o uruchamianiu generatora Lex i przetwornika tekstu

- jeśli plik z programem w języku Lex, opisującym działanie przetwornika tekstu nazywa się *przetw.l*, to polecenie:

```
lex -o przetw.c przetw.l
```

spowoduje utworzenie przetwornika (automatu) w języku C i umieszczenie go w pliku o nazwie *przetw.c*

- przetwornik należy skompilować, korzystając z dostępnego kompilatora języka C (np. `cc` lub `gcc`), wskazując plik biblioteczny systemu Lex o nazwie `l` (opcja ma postać `-ll`):

```
gcc przetw.c -ll
```

- powstanie plik wykonywalny o standardowej nazwie `a.out`, który można następnie stosować do przetwarzania tekstów, np.:

```
a.out < test.in > test.out
```

```
a.out < test.in
```

```
a.out > test.out
```

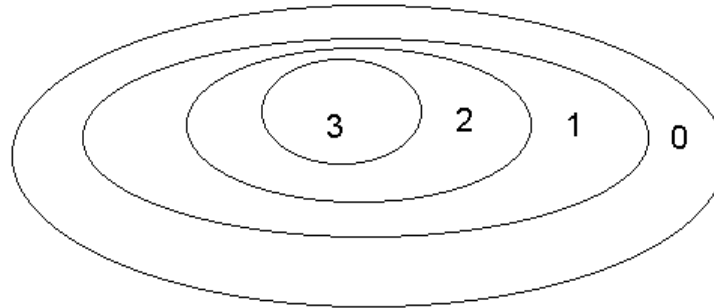
```
a.out
```

3. Elementy teorii języków formalnych

Metody definiowania składni języka $\langle \Sigma, \text{Syn} \rangle$

- metody definiowania składni dzielimy na: *generacyjne* i *akceptorowe*
- w metodach *generacyjnych* podaje się *zbiór reguł generacji* (formalny system generacji składni) poprawnych napisów nad podanym alfabetem (np. gramatyki kombinatoryczne N. Chomsky'ego)
- w metodach *akceptorowych* definiuje się nad alfabetem automat/akceptor (formalna definicja automatu), który potrafi zbadać przynależność napisu do języka (wszystkie napisy zaakceptowane przez automat definiują język)
- każdemu typowi gramatyki odpowiada typ automatu (podejścia są zatem komplementarne)
- *translator* jest przetwornikiem napisów, od akceptora różni się tym, że ma zdefiniowany alfabet wyjściowy i potrafi generować napisy w języku nad tym alfabetem (formalnie, jest automatem „z wyjściem”).

3.2 Hierarchia języków wg N. Chomsky'ego:



- klasy języków uporządkowane relacją zawierania się zbiorów;
- języki klasy "0" nazywamy *obliczalnymi*, są klasą najszerszą, matematycznie są przykładem zbiorów rekurencyjnie przeliczalnych (częściowo rozstrzygalnych);
- języki klasy "1" nazywamy *monotonicznymi*, matematycznie są przykładem zbiorów rekurencyjnych (rozstrzygalnych);
- języki klasy "2" nazywamy *bezkontekstowymi*; wyróżnia się istotne z praktycznego punktu widzenia podklasy **LL** i **LR** (generator YACC);
- języki klasy "3" nazywamy regularnymi; biorą swą nazwę od jednego z formalizmów ich definiowania zwanego *wyrażeniami regularnymi*;

3.3 Formalny system generacji składni (podejście **generacyjne**)

$F = \langle \Sigma, \mathbf{Syn}, \mathbf{P}, \mathbf{s} \rangle$

Σ – alfabet

Syn – syntaktyka

P – zbiór reguł generacji napisów ze zbioru **Syn**

s – wyróżniony zbiór napisów początkowych;

- przykładem systemu generacji jest *formalny system kombinatoryczny*:

$F = \langle \Sigma, \mathbf{Syn}, \mathbf{P}, \mathbf{s} \rangle$, gdzie $\mathbf{P} \subseteq \Sigma^* \times \Sigma^*$, $\mathbf{s} \in \Sigma^*$;

- przykładem systemu kombinatorycznego jest *gramatyka* kombinatoryczna (*klasa 0*), zdefiniowana przez Noama Chomsky'ego:

$\mathbf{G} = \langle \mathbf{N}, \Sigma, \mathbf{P}, \mathbf{S} \rangle \quad \mathbf{N} \cap \Sigma = \emptyset$

N – skończony niepusty alfabet symboli pomocniczych (*nieterminalnych*),

Σ – skończony, niepusty alfabet definiowanego języka (symbole *terminalne*),

P $\subseteq (\mathbf{N} \cup \Sigma)^+ \times (\mathbf{N} \cup \Sigma)^*$, zbiór *produkcji* gramatyki (reguł generacji napisów),

S $\in \mathbf{N}$, wyróżniony symbol pomocniczy, zwany *aksjomatem* gramatyki;

- produkcję gramatyki, opisującą związek pomiędzy dwoma napisami (np. napisem α i napisem β), zapisujemy zwyczajowo $\alpha \rightarrow \beta$, gdzie strzałka oznacza *znak zastąpienia*, a cały napis poleca zastąpić napis α napisem β ;
- relacja *wyvodu bezpośredniego* \Rightarrow (oraz jej domknięcie \Rightarrow^*):

$xabcy \Rightarrow xdv$ o ile w zbiorze P istnieje produkcja $abc \rightarrow dv$;

- język $J(G)$ generowany gramatyką G :

$$J(G) = \{x \mid S \Rightarrow^* x \wedge x \in \Sigma^*\}$$

- *zdanie* (lub słowo) to napis zbudowany wyłącznie z symboli terminalnych, uzyskany za pomocą relacji wyvodu \Rightarrow^* ; napis zawierający symbole terminalne i nieterminalne, również uzyskany za pomocą relacji wyvodu \Rightarrow^* , nazywamy *formą zdaniową*

• **operacje na językach** (zbiorach słów J, J_1, J_2):

suma: $J_1 \cup J_2 = \{x \mid x \in J_1 \vee x \in J_2\}$

przekrój: $J_1 \cap J_2 = \{x \mid x \in J_1 \wedge x \in J_2\}$

różnica: $J_1 - J_2 = \{x \mid x \in J_1 \wedge x \notin J_2\}$

dopełnienie: $J' = \Sigma^* - J$

złożenie: $J_1 J_2 = \{xy \mid x \in J_1 \wedge y \in J_2\}$

potęgowanie: $J^n, J^0 = \{\varepsilon\}, J^k = J J^{k-1}$

domknięcie: $J^* = \bigcup_{i=0}^{\infty} J^i$

domknięcie dodatnie: $J^+ = \bigcup_{i=1}^{\infty} J^i, J^+ = J J^* = J^* J, J^* = J^+ \cup \{\varepsilon\}$

3.4 Formalna definicja automatu (podejście **akceptorowe**)

- w podejściu *akceptorowym*: automat M *akceptuje* słowo (zdanie) S wtw. startując od konfiguracji początkowej ze słowem S na wejściu może wykonać skończony ciąg kroków kończących się konfiguracją końcową
- automat M definiuje język $L(M)$, który jest zbiorem wszystkich słów S akceptowanych przez M

4. Translacja i kompilacja – podstawowe pojęcia

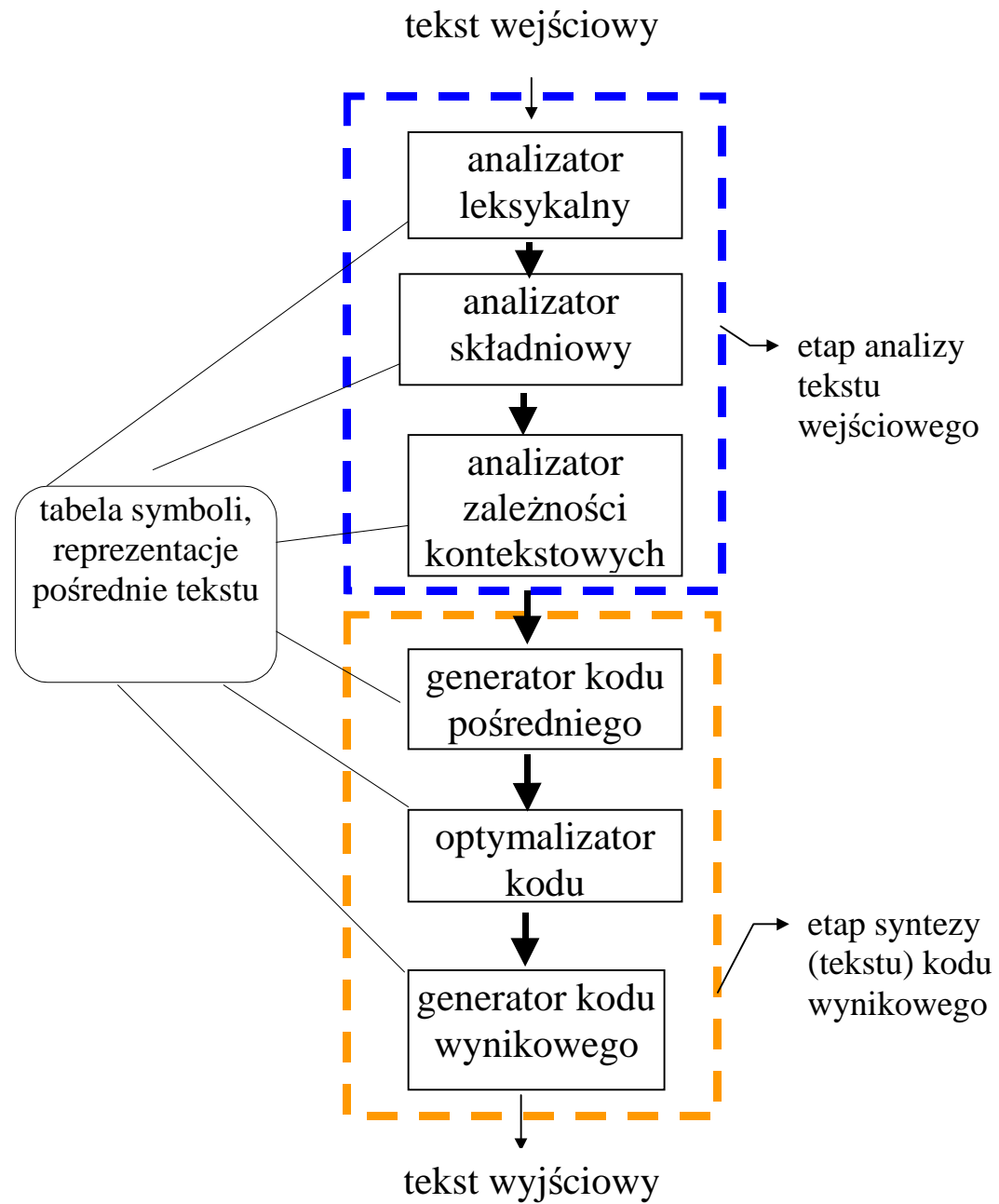
4.1 Zagadnienie translacji

tekst wejściowy -----> **TRANSLACJA** -----> *wynik translacji*

- zajmiemy się tylko wybranymi aspektami translacji: dotyczącymi tekstów napisanych w *językach formalnych* (sztucznych, np. językach programowania, dostępu do baz danych, publikacji informacji w Internecie – HTML, SGML, XML itd.) w modelu *translacji sterowanej składnią* ($L = \langle \Sigma, \text{Syn} \rangle$)
- translacja sterowana składnią – proces, w którym po przeanalizowaniu struktury składniowej tekstu wejściowego jest on *interpretowany* lub na jego podstawie *generowany* jest *tekst wyjściowy*
- *narzędzia programistyczne* realizujące translacje sterowane składnią to np.: AWK, *Lex*, *YACC*, sed, PERL i wiele innych, w tym języki programowania

- wyróżniamy *dwie techniki* translacji: *interpretację* i *kompilację*
- proces kompilacji realizuje program zwany *kompilatorem*:
tekst wejściowy --> **kompilator** --> *tekst wyjściowy*
- kompilacja przebiega w dwóch *etapach*:
 - a) *analizy* tekstu wejściowego i
 - b) *syntezy* tekstu wyjściowego,

które rozpadają się na *fazy* (każda faza przetwarza tekst wejściowy lub jego reprezentację):



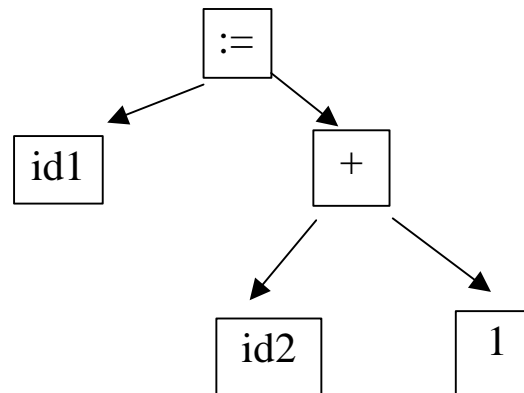
- kompilacja może być *jedno-* lub *wieloprzebiegowa*
- przykład kompilacji:

```
program p1;  
var x: real; i: integer;  
begin  
x:=i+1;  
end.
```

Analizator leksykalny wyodrębni następujący zbiór *jednostek leksykalnych*:

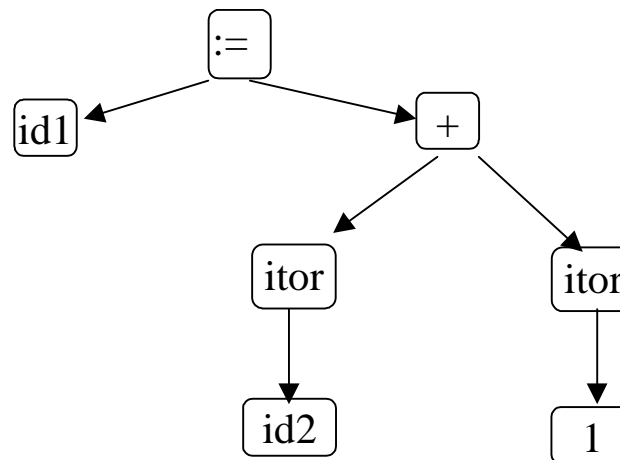
program p1 var x : real ; i : integer begin := + 1 end .

Analizator składniowy utworzy *drzewo rozbioru składniowego*, a informacje o identyfikatorach zmiennych umieści w *tabeli symboli*:



Nazwa	Leksem	Typ	...
id1	x	real	...
id2	i	integer	...

Analizator zależności kontekstowych stwierdzi, że zmienne wykorzystane w instrukcji przypisania zostały zadeklarowane (informacje o tym zawarto w tabeli symboli) oraz zaznaczy w kodzie pośrednim, że należy przeprowadzić konwersję typu stałej **1** i zmiennej **i** z typu `integer` do typu `real` (operator *itor*):



Generator kodu trójadresowego, wygeneruje dwie zmienne okresowe *temp1* i *temp2* i na podstawie powyższego drzewa utworzy zapis:

```
temp1 := itor(1)
temp2 := itor(id2)
id1 := temp1 + temp2
```

Optymalizator kodu dokona redukcji:

```
temp1 := itor(id2)
id1 := temp1 + 1.0
```

Na podstawie kodu pośredniego, **generator kodu wynikowego** utworzy zapis w kodzie asemblera, który następnie poddany zostanie **asemlacji** w celu uzyskania wersji wykonywalnej:

```
MOV id2, R1
ADD R1, #1.0
MOV R1, id1
```

4.2 Składnia Prostej Języka Programowania (PJP)

- notacja BNF (sposób zapisywania gramatyki języka)
- alfabet PJP (A)

`<litera> ::= a|b|c|... z|A|B|C ... Z`

`<cyfra> ::= 0|1|2|3|4|5|6|7|8|9`

`<operator jednoargumentowy> ::= +|-|nie`

`<operator dwuargumentowy> ::= +|-|/|*|:=|=|<|<=|>|>=|<>|i|lub`

`<słowo kluczowe> ::= program|początek|koniec|tablica|jesli|to|`

`przeciwnie|ilsej|dopoki|wykonuj|ikopod|`

`czytaj|drukuj|całkowite|rzeczywiste|boolowskie`

`A = <litera>∪<cyfra>∪<operator jednoargumentowy>∪`

`<operator dwuargumentowy>∪<słowo kluczowe>∪ {:=|(|)[|].}`

- stałe

`<stała> ::= <stała numeryczna>|<stała typu boolowskiego>`

$\langle \text{stała numeryczna} \rangle ::= +\langle \text{liczba} \rangle | -\langle \text{liczba} \rangle | \langle \text{liczba} \rangle$

$\langle \text{liczba} \rangle$ (zdefiniować za pomocą wyrażenia regularnego, jako ćwiczenie)

$\langle \text{stała typu boolowskiego} \rangle ::= \text{prawda} | \text{falsz}$

$\langle \text{identyfikator} \rangle ::= \langle \text{litera} \rangle | \langle \text{identyfikator} \rangle \langle \text{litera} \rangle |$
 $\langle \text{identyfikator} \rangle \langle \text{cyfra} \rangle$

- zmienne

$\langle \text{zmienna} \rangle ::= \langle \text{identyfikator} \rangle | \langle \text{identyfikator} \rangle [\langle \text{stała} \rangle]$

- wyrażenia

$\langle \text{wyrażenie} \rangle ::= \langle \text{stała} \rangle | \langle \text{zmienna} \rangle |$
 $(\langle \text{wyrażenie} \rangle \langle \text{operator dwuargumentowy} \rangle \langle \text{wyrażenie} \rangle) |$
 $(\langle \text{operator jednoargumentowy} \rangle \langle \text{wyrażenie} \rangle)$

- instrukcje

$\langle \text{instrukcje} \rangle ::= \langle \text{instrukcja} \rangle |$
 $\langle \text{instrukcja} \rangle ; \langle \text{instrukcje} \rangle$

- instrukcja

```
<instrukcja> ::= <zmienna> := <wyrażenie> |  
                jeśli <wyrażenie> to <instrukcja> przeciwnie  
                <instrukcja> ilsej |  
                dopoki <wyrażenie> wykonuj <instrukcja> ikopod |  
                początek <instrukcje> koniec |  
                czytaj <zmienna> |  
                drukuj <wyrażenie>
```

- typy

```
<typ> ::= całkowite | rzeczywiste | boolowskie
```

- deklaracje

```
<deklaracje> ::= <deklaracja> |  
                <deklaracja> ; <deklaracje>
```

- deklaracja

```
<deklaracja> ::= <zmienna>:<typ> |  
                <zmienna>:<typ> tablica [<stała>]
```

- programy

```
<program> ::= program <identyfikator> <deklaracje>; <instrukcje>.
```