

2 1976

prace



P.22.29/76

Instytutu  
Maszyn  
Matematycznych

rok XVIII



Rok XVIII

2 1976



P.2229/76

prace

Instytutu

Maszyn

Matematycznych

Roman KRZEMIENI  
KONSTRUKCJA ANALIZATORA LEKSYKALNEGO  
W METATRANSLATORZE

Praca doktorska napisana  
pod kierunkiem prof. dr hab. W. M. Turskiego



Zjednoczenie Przemysłu Automatyki i Aparatury  
Pomiarowej "MERA"

INSTYTUT MASZYN MATEMATYCZNYCH



Copyright © 1976 - by Instytut Maszyn Matematycznych  
Poland

Wszelkie prawa zastrzeżone

#### KOMITET REDAKCYJNY

Jerzy GRADOWSKI, Andrzej JANICKI (z-ca red. nacz.),  
Roman KULESZA (redaktor naczelny),  
Antoni MAZURKIEWICZ, Tomasz PAWLAK, Ryszard PREGIEL,  
Władysław M. TURSKI (z-ca red. nacz.), Zbigniew WIERZBICKI

Sekretarz Redakcji: Romana NITKOWSKA

Redaktor Techniczny: Maria KOZŁOWSKA

Adres Redakcji: Instytut Maszyn Matematycznych  
Branżowy Ośrodek INTE  
ul. Krzywickiego 34, 02-078 WARSZAWA  
tel. 28-37-29

Przedstawiona Czytelnikom praca doktorska  
ROMANA KRZEMIENIA została wyróżniona  
nagrodą Ministra Przemysłu Maszynowego  
w 1976 r.

# S p i s t r e ś c i

|   | str. |
|---|------|
| 1. WSTĘP  | 5    |
| 2. ANALIZA LEKSYKALNA W TRANSLATORZE  | 8    |
| 2.1. Cele analizy leksykalnej   | 8    |
| 2.2. Realizacja analizatora leksykalnego  | 12   |
| 3. PRZEGLĄD METOD AUTOMATYCZNEJ KONSTRUKCJI ANALIZATORÓW LEKSYKALNYCH W ZNANYCH METATRANSLATORACH | 16   |
| 3.1. Najprostsze rozwiązania problemów automatyzacji konstrukcji L-analizatorów                   | 16   |
| 3.2. System AED RWORD   | 19   |
| 3.2.1. Opis klasy znaków  | 20   |
| 3.2.2. Opis atomu leksykalnego  | 21   |
| 3.3. L-analizator Lecarma   | 24   |
| 4. OPIS KONSTRUKCJI MODUŁU ANALIZY LEKSYKALNEJ  | 26   |
| 4.1. Organizacja ogólna   | 27   |
| 4.2. Podmoduł ładowacza   | 27   |
| 4.2.1. Metajęzyk opisu składni języka źródłowego  | 28   |
| 4.2.2. Procedura ŁADOWACZ   | 30   |
| 4.2.3. Postać gramatyki języka źródłowego w pamięci maszyny                                       | 30   |
| 4.3. Podmoduł L-konstruktora  | 36   |
| 4.3.1. Procedura REGULAR  | 36   |
| 4.3.2. Uzasadnienie teoretyczne procedury REGULAR   | 40   |
| 4.3.3. Procedura LKONSTRUKTOR   | 42   |
| 4.4. Podmoduł L-analizatora   | 46   |
| 4.4.1. Procedura LANALIZATOR  | 46   |
| 4.4.2. Algorytm realizowany przez procedurę LANALIZATOR   | 47   |
| 4.4.3. Ograniczenia klasy gramatyk wejściowych  | 52   |
| 5. ZAKOŃCZENIE  | 56   |
| BIBLIOGRAFIA  | 57   |
| Dodatek A. Program realizujący moduł analizy leksykalnej  | 59   |
| Dodatek B. Składnia języka Euler  | 74   |
| Dodatek C. Przykład zastosowania modułu analizy leksykalnej do programu w języku Euler            | 77   |



Leksykograf zajmuje się korelowaniem form językowych w pewien szczególny sposób, mianowicie łączeniem w pary synonimów.

... leksykograf chce wiedzieć jakie formy są synonimiczne, czyli tożsame pod względem znaczenia.

("Z punktu widzenia logiki". Quine W.V.)

## KONSTRUKCJA ANALIZATORA LEKSYKALNEGO W METATRANSLATORZE

Roman KRZEMIEN

Pracę złożono 20.09.1975

W pracy przedstawiono nową metodę konstrukcji analizatorów leksykalnych w metatranslatorach. Proponowana metoda jest całkowicie różna od metod stosowanych dotychczas, ponieważ atomy leksykalne są znajdowane automatycznie jedynie na podstawie opisu składu języka.

### 1. WSTĘP

Analizę leksykalną można traktować jako wstępną analizę syntaktyczną. Program dokonujący analizy leksykalnej (L-analizy) czyta tekst w języku źródłowym, rozumiany jako ciąg znaków, i grupom znaków przyporządkowuje pewne elementarne jednostki syntaktyczne, zwane atomami leksykalnymi.

Program ten, zwany analizatorem leksykalnym lub L-analizatorem korzysta jednak ze znacznie prostszego i efektywniej-

szego aparatu lingwistycznego niż program dokonujący analizy syntaktycznej, zwany S-analizatorem lub analizatorem syntaktycznym.

Składnię atomów leksykalnych można bowiem zapisać za pomocą gramatyki regularnej (np.  $\langle \text{identyfikator} \rangle ::= \text{litera} \mid \langle \text{identyfikator} \rangle \text{litera} \mid \langle \text{identyfikator} \rangle \text{cyfra}$ , gdzie przyjmujemy, że litera i cyfra są symbolami terminalnymi). Do rozbioru atomów leksykalnych można zatem użyć automatu skończonego (zamiast automatu ze stosem), co znacznie przyspiesza analizę. Łatwość zaprogramowania automatu skończonego umożliwi również automatyczną konstrukcję L-analizatora.

Praca niniejsza dotyczy właśnie problemów związanych z automatyczną konstrukcją analizatorów leksykalnych w metatranslatorach.

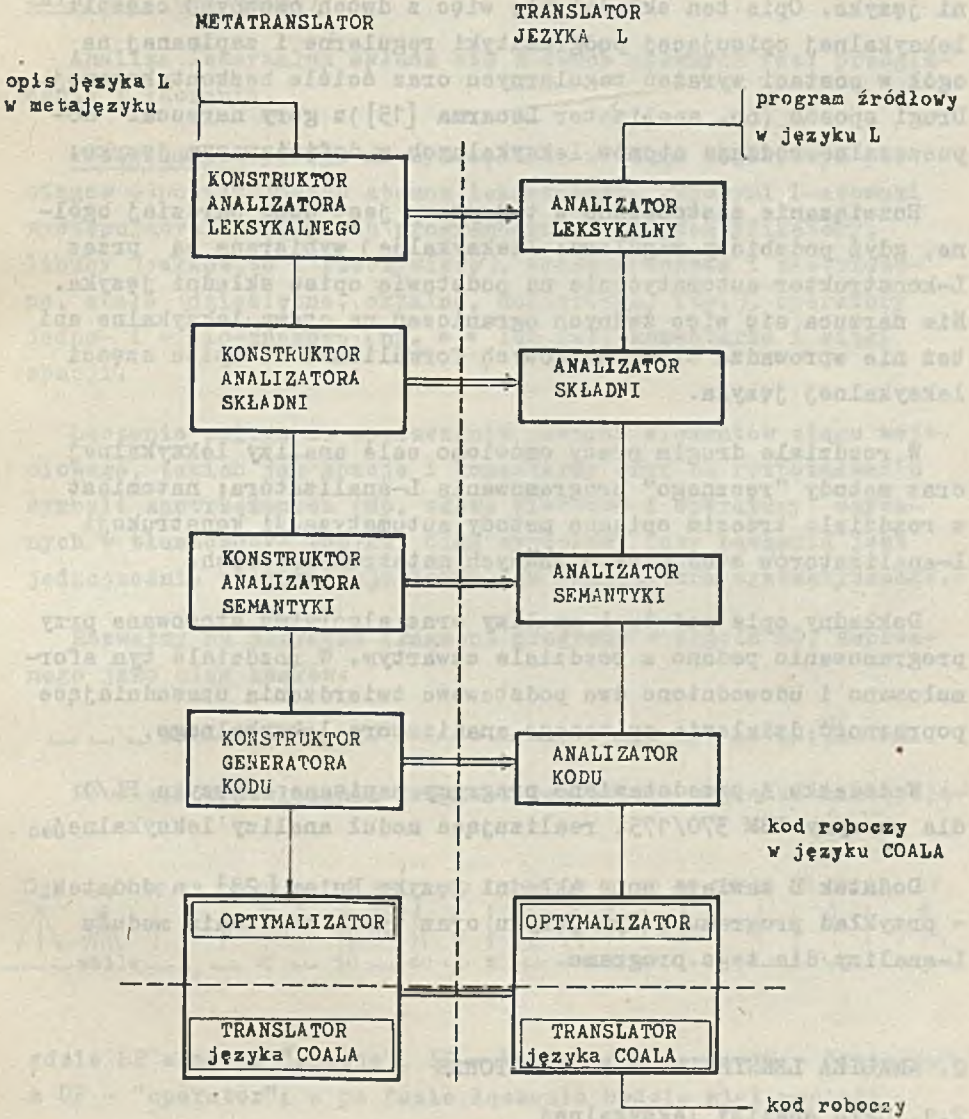
Ogólny schemat metatranslatora (ang. compiler-compiler) realizowanego w Zakładzie Teorii Translatorów Instytutu Maszyn Matematycznych przedstawia rysunek 1.

Schemat ten jest realizowany przez pakiet modułów, zwanych odpowiednio: modułem analizy leksykalnej, analizy syntaktycznej, analizy semantycznej, itd.

Głównym wynikiem pracy jest zaprogramowanie, opisanie i uzasadnienie modułu L-analizy, na który składają się programy ładowacza, L-konstruktora i L-analizy.

Dotychczas stosowane dwa sposoby rozwiązania problemu analizy leksykalnej w metatraslatorach. Pierwszy (np. AED-RWORD





Rys. 1.

[11]) polegał na wydzieleniu części leksykalnej w opisie składni języka. Opis ten składał się więc z dwóch osobnych części: leksykalnej opisującej podgramatyki regularne i zapisanej na ogół w postaci wyrażeń regularnych oraz ściśle bezkontekstowej. Drugi sposób (np. analizator Lecarma [15]) z góry narzucał dopuszczalne rodzaje atomów leksykalnych w definiowanym języku.

Rozwiązanie zastosowane w tej pracy jest dużo bardziej ogólne, gdyż podzbiory regularne (leksykalne) wybierane są przez L-konstruktor automatycznie na podstawie opisu składni języka. Nie narzuca się więc żadnych ograniczeń na atomy leksykalne ani też nie wprowadza się dodatkowych formalizmów do opisu części leksykalnej języka.

W rozdziale drugim pracy omówiono cele analizy leksykalnej oraz metody "ręcznego" programowania L-analizatora; natomiast w rozdziale trzecim opisano metody automatyzacji konstrukcji L-analizatorów stosowane w znanych metatranslatorach.

Dokładny opis modułu L-analizy oraz algorytmy stosowane przy programowaniu podano w rozdziale czwartym. W rozdziale tym sformułowano i udowodniono dwa podstawowe twierdzenia uzasadniające poprawność działania opisanego analizatora leksykalnego.

W dodatku A przedstawiono programy napisane w języku PL/1 dla maszyny IBM 370/175, realizujące moduł analizy leksykalnej.

Dodatek B zawiera opis składni języka Euler [28], a dodatek C - przykład programu w tym języku oraz wynik działania modułu L-analizy dla tego programu.

## 2. ANALIZA LEKSYKALNA W TRANSLATORZE

### 2.1. Cele analizy leksykalnej

Na pierwszym etapie procesu tłumaczenia program źródłowy ma postać ciągu znaków. Analizator leksykalny (będący pierwszym przebiegiem translatora) czyta tekst programu źródłowego rozumiany jako ciąg znaków i grupom znaków przyporządkowuje nazwę



elementarne jednostki syntaktyczne, zwane dalej atomami leksykalnymi (L-atomami).

Analiza leksykalna składa się z dwóch głównych faz: przeglądania i łączenia.

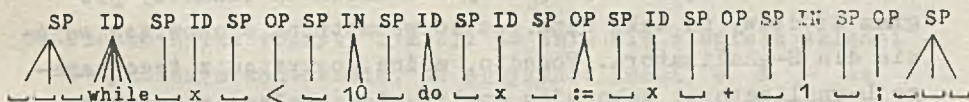
Przeglądanie polega na wydzieleniu w ciągu wejściowym podciągów odpowiadających atomom leksykalnym. Typowymi L-atomami występującymi w językach programowania są: identyfikatory, liczby (całkowite i rzeczywiste), słowa kluczowe i zastrzeżone, stałe (dziesiętne, oktalne, Holleritha, itp.), operatory jedno- i wielo-znakowe (np. \* \* lub :=), komentarze i ciągi spacji.

Łączenie polega na opuszczaniu pewnych elementów ciągu wejściowego, takich jak spacje i komentarze oraz na rozpoznawaniu symboli zastrzeżonych (np. słowa kluczowe i operatory) używanych w tłumaczonym języku. Ciąg wyjściowy fazy łączenia jest jednocześnie ciągiem wejściowym dla analizatora syntaktycznego.

Rozważmy na przykład fragment programu w Algolu 60, zapisanego jako ciąg znaków:

while x < 10 do x := x + 1;

Po fazie przeglądania fragment ten będzie wyglądał następująco:



gdzie SP oznacza "spację", ID - "identyfikator", IN - "integer", a OP - "operator"; a po fazie łączenia będzie miał postać:

|       |    |   |    |    |    |    |    |   |    |   |
|-------|----|---|----|----|----|----|----|---|----|---|
| while | ID | < | IN | do | ID | := | ID | + | IN | ; |
|       |    |   |    |    |    |    |    |   |    |   |
|       | x  |   | 10 |    | x  |    | x  |   | 1  |   |



Analizator syntaktyczny otrzymuje na wejściu tylko ciąg symboli umieszczony w górnym rzędzie na powyższym rysunku, ale dodatkowe informacje (np. jakiemu konkretnemu identyfikatorowi odpowiada dany symbol ID) są przechowywane w tablicach dla przebiegu analizy semantycznej.

Redukcja tekstu wejściowego do ciągu atomów leksykalnych, które są traktowane przez S-analizator jako symbole terminalne, jest głównym celem analizy leksykalnej. Jednakże istnieje jeszcze kilka innych przyczyn wyodrębniania w translatorze przebiegu analizy leksykalnej:

1. Program źródłowy może być zapisywany za pomocą różnych urządzeń zewnętrznych, o rozmaitych zbiorach znaków. Może się też zdarzyć, że część programu znajduje się w pamięci maszyny, a pozostała jego część jest wczytywana z urządzeń wewnętrznych. Bardzo często dopuszczalne są również różne reprezentacje tego samego języka. Na przykład w niektórych reprezentacjach Algolu 60 takie ograniczniki jak begin, else, if, end, itp. występują w programie ujęte w apostrofy, spacje zaś są ignorowane, natomiast w innych reprezentacjach - ograniczniki te są traktowane jako słowa zastrzeżone, a spacje są używane jako separatory. W przypadku wyodrębnienia analizy leksykalnej możemy napisać jeden S-analizator oraz kilka L-analizatorów - znacznie prostszych i łatwiejszych do zaprogramowania - po jednym dla każdej realizacji języka i urządzenia wejściowego. Każdy L-analizator tłumaczy program źródłowy na taką samą postać wewnętrzną stanowiącą wejście dla S-analizatora. Ponadto, można korzystać z tego samego L-analizatora, dokonując analizy leksykalnej różnych języków programowania.
2. Znaczną część czasu maszyny przeznaczonego na tłumaczenie programu zajmuje badanie poszczególnych jego znaków. Wyodrębnienie analizy leksykalnej umożliwi przyspieszenie tego procesu, na przykład dzięki zaprogramowaniu części lub całego L-analizatora w języku wewnętrznym.

3. Skrócenie tekstu programu źródłowego przez L-analizator pozwala zaoszczędzić zarówno miejsca pamięci jak i czas maszyny. Jest to bardzo istotne w przypadku translatorów wieloprzebiegowych dla małych maszyn cyfrowych, gdzie program jest wielokrotnie przesyłany między pamięcią operacyjną i pomocniczą.
4. Pewne zależności kontekstowe mogą być łatwo wykryte na etapie analizy leksykalnej. Na przykład Algol 60 mógłby być zdefiniowany przez gramatykę z pierwszeństwem, gdyby nie dwuznaczne wystąpienia symbolu ":" jako ogranicznika etykiety lub w deklaracjach tablic. Na etapie analizy leksykalnej można łatwo zastąpić dwukropek dwoma różnymi symbolami - zależnie od miejsca występowania w programie, a wobec tego w S-analizatorze Algolu 60 można zastosować metody rozbioru przewidziane dla gramatyki z pierwszeństwem.
5. L-analizator może wyszukiwać (lub usuwać) niektóre trywialne błędy w programie (na przykład błędy w stałych, brak odpowiedniości między lewymi i prawymi nawiasami, itp.). Znaczne przyspieszenie translacji można osiągnąć przez usuwanie na etapie analizy leksykalnej zbędnych spacji, znaków kontynuacji, komentarzy, itp. lub przez klasyfikowanie instrukcji według typów (dla Fortranu IV ANSI Standard ta ostatnia operacja może być wykonywana przez automat skończenie stanowy).
6. Często kilka różnych symboli pojawia się w opisie składni w tym samym kontekście. Na przykład: znaki + i - są używane w ten sam sposób, symbole real, integer mogą pojawić się w tym samym miejscu programu w Algolu 60. L-analizator może przydzielić jeden symbol dla każdej grupy takich syntaktycznie synonimicznych symboli.

Założmy, na przykład, że gramatyka języka źródłowego zawiera między innymi następujące produkcje:

```
<type> ::= real | integer | boolean
<declaration> ::= <type> <identifier list>
<expression> ::= <term> | <expression> + <term> | <expression>
- <term>
```



Gramatyka ta będzie miała po L-analizie prostszą postać, mianowicie:

```
<declaration> ::= TYPE <identifier list>  
<expression> ::= <term> | <expression> AO <term>
```

gdzie TYPE oraz AO będą traktowane przez S-analizator jako symbole terminalne.

## 2.2. Realizacja analizatora leksykalnego

Jak już zaznaczyliśmy, składnia atomów leksykalnych jest bardzo prosta i może być zapisana w postaci gramatyki regularnej. Dlatego też analizator leksykalny będzie na ogół automatem skończenie stanowym akceptującym słowa należące do tej gramatyki. Automat ten będziemy zapisywać w postaci tzw. macierzy przejść, to jest dwuwymiarowej tablicy, której wiersze będą określały stany automatu, kolumny zaś - symbole pojawiające się na wejściu. Elementami tablicy będą wartości funkcji przejścia i funkcji wyjścia automatu.

L-analizator może być zrealizowany albo jako oddzielny przebieg, który wykonuje analizę leksykalną całego programu źródłowego i dostarcza S-analizatorowi tablicę zawierającą ten program w postaci wewnętrznej, albo jako podprogram wywoływany przez S-analizator w celu otrzymania nowego symbolu potrzebnego do analizy. W poniższym przykładzie (w którym zakładamy ten drugi wariant) korzystać będziemy z dwóch zmiennych roboczych: ZNAK, w której będziemy przechowywać aktualnie analizowany znak ciągu wejściowego, oraz ATOM, w której będziemy kompletować atom leksykalny. Przed wywołaniem analizatora leksykalnego wartość zmiennej ATOM będzie zerowana przez S-analizator.

Przykład 2.1. Rozważmy język zdefiniowany następującą gramatyką:



|              |   |
|--------------|---|
| <program>    | ::= .begin.<block>.end.                                     |
| <block>      | ::= <statement>   <block> <statement> ;                     |
| <statement>  | ::= <leftside> <number>                                     |
| <leftside>   | ::= <identifier> :=   |
| <identifier> | ::= <letter>   <identifier> <letter>   <identifier> <digit> |
| <number>     | ::= <integer>   <real>                                      |
| <integer>    | ::= <digit>   <integer> <digit>                             |
| <real>       | ::= <integer> . <integer>                                   |
| <letter>     | ::= a   b   d   e   g   i   n                               |
| <digit>      | ::= 0   1   |

W języku tym występują następujące atomy leksykalne:

1. operator - znak "==" lub " ; "
2. liczba całkowita - dowolny ciąg zer i jedynek (np. 101001)
3. liczba rzeczywista - atomy te są utworzone z części całkowitej i ułamkowej. Część całkowita jest ciągiem zer i jedynek; część ułamkowa ciągiem zer i jedynek poprzedzonym kropką (np. 10.011)
4. identyfikator - dowolny ciąg liter (a, b, d, e, g, i, n) i cyfr (0, 1) zaczynający się od litery (np. a10)
5. słowo kluczowe - identyfikator ograniczony z obu stron kropkami (np. .begin.).

Macierz przejść dla naszego języka jest przedstawiona w poniższej tablicy:

| STAN                    | SYMBOL WEJŚCOWY           |        | 0 1    | kropka | spacja,<br>nowa<br>linia | opera-<br>tor<br>:= lub; | inny<br>znak |
|-------------------------|---------------------------|--------|--------|--------|--------------------------|--------------------------|--------------|
|                         | a b d e g i n<br>f.przej. | f.wyj. |        |        |                          |                          |              |
| 1 (początkowy)          | 2                         | CONT   | 3,CONT | 7,CONT | 1,CONT1                  | 1,END1                   | -,ERR        |
| 2 (identyfika-<br>tor)  | 2                         | CONT   | 2,CONT | 2,END  | 2,END                    | 2,END                    | -,ERR        |
| 3 (l. całkowita)        | 3                         | END    | 3,CONT | 6,CONT | 3,END                    | 3,END                    | -,ERR        |
| 4 (l. rzeczy-<br>wista) | 4                         | END    | 4,CONT | 4,END  | 4,END                    | 4,END                    | -,ERR        |
| 5 (słowo klu-<br>czowe) | 5                         | CONT   | 5,CONT | 5,END1 | -,ERR                    | -,ERR                    | -,ERR        |
| 6                       | -                         | ERR    | 4,CONT | -,ERR  | -,ERR                    | -,ERR                    | -,ERR        |
| 7                       | 5                         | CONT   | -,ERR  | -,ERR  | -,ERR                    | -,ERR                    | -,ERR        |

funkcje wyjścia realizują następujące zadania<sup>\*</sup>):

CONT: ATOM——ATOM || ZNAK;  
Wykonaj funkcję CONT1.

CONT1: ZNAK——kolejny znak ciągu wejściowego;  
analizuj dalej tablicę.

END1: ZNAK——kolejny znak ciągu wejściowego;  
wykonaj funkcję END.

END: Załaduj ATOM do słownika symboli;  
wróć do S-analizatora podając adres ATOMU w słowniku  
symboli i zaznaczając, że ATOM jest:

- operatorem                      - jeżeli stan = 1,
- identyfikatorem                - jeżeli stan = 2,

<sup>\*</sup>) Znak —— oznacza operację przypisania, a znak || operację konkatenacji



- liczbą całkowitą - jeżeli stan = 3,
- liczbą rzeczywistą - jeżeli stan = 4,
- słowem kluczowym - jeżeli stan = 5.

ERR: Wróć do analizatora syntaktycznego sygnalizując błąd.

Przyporządkujmy poszczególnym atomom leksykalnym następujące numery, stanowiące ich umowną reprezentację:

- nie zdefiniowany : 0,
- operator : 1,
- identyfikator : 2,
- liczba całkowita : 3,
- liczba rzeczywista : 4,
- słowo kluczowe : 5.

Przykład 2.2. Rozpatrzmy następujący program w naszym języku:

```
.begin. ab1 := 11.01; ab2 := 11.01; .end.
```

etapy pracy L-analizatora wyglądają następująco:

| etap | wyjście | słownik symboli |     |    |       |   |     |       |
|------|---------|-----------------|-----|----|-------|---|-----|-------|
|      |         | 1               | 2   | 3  | 4     | 5 | 6   | 7     |
| 1    | 5,1     | .begin.         | -   | -  | -     | - | -   | -     |
| 2    | 2,2     | .begin.         | ab1 | -  | -     | - | -   | -     |
| 3    | 1,3     | .begin.         | ab1 | := | -     | - | -   | -     |
| 4    | 4,4     | .begin.         | ab1 | := | 11.01 | - | -   | -     |
| 5    | 1,5     | .begin.         | ab1 | := | 11.01 | ; | -   | -     |
| 6    | 2,6     | .begin.         | ab1 | := | 11.01 | ; | ab2 | -     |
| 7    | 1,3     | .begin.         | ab1 | := | 11,01 | ; | ab2 | -     |
| 8    | 4,4     | .begin.         | ab1 | := | 11.01 | ; | ab2 | -     |
| 9    | 1,5     | .begin.         | ab1 | := | 11.01 | ; | ab2 | -     |
| 10   | 5,7     | .begin.         | ab1 | := | 11.01 | ; | ab2 | .end. |

Na każdym etapie L-analizator generuje na wyjściu dwie wartości, z których pierwsza oznacza umowną reprezentację atomu leksykalnego, a druga adres, pod którym dany atom jest pamiętany w słowniku symboli.



### 3. PRZEGLĄD METOD AUTOMATYCZNEJ KONSTRUKCJI ANALIZATORÓW LEKSYKALNYCH W ZNANYCH METATRANSLATORACH

#### 3.1. Najprostsze rozwiązania problemów automatyzacji konstrukcji L-analizatorów

Większość istniejących metatranslatorów w ogóle nie generuje przebiegu analizy leksykalnej (BMCC [22], Ingerman [10], TMG [18]); w niektórych zaś sprowadza się go do wykonywania kilku zupełnie trywialnych operacji.

Na przykład w Cogencie [21] można podawać opis znaków, tzn. dowolnie definiować kod reprezentantów znaków źródłowych; można też wprowadzać do definicji składni stałe teksty (np. słowa zastrzeżone języka), które przed S-analizatorem przetwarza analizator leksykalny i zamienia na określony układ bitów. Dokonuje się tego za pomocą definicji znaków, które określają wartości poszczególnych reprezentantów, wyrażone liczbami ósemkowymi.

Opis znaków w Cogencie ma następującą postać:

```
<opis znaków> ::= <puste> § CHARDEF <ciąg definicji
                    znaków>
<ciąg definicji znaków> ::= <puste> | <ciąg definicji znaków>
                    <definicja znaku>
<definicja znaku> ::= <reprezentant znaku źródłowego> =
                    <ciąg kodów źródłowych> <ciąg ko-
                    dów nieistotnych> <kod wynikowy>
<ciąg kodów źródłowych> ::= (<ciąg kodów>)
<ciąg kodów nieistotnych> ::= <puste> | (<ciąg kodów>)
<kod wynikowy> ::= <dodatnia całkowita ósemkowa>
<ciąg kodów> ::= <dodatnia całkowita ósemkowa> |
                    <ciąg kodów>, <dodatnia całkowita
                    ósemkowa>
<reprezentant znaku źródłowego> ::= <normalny znak> | (<spe-
                    cjalny znak>) | (§ <tekst nazwowy>)
```

Każda definicja znaku określa wartości reprezentanta występującego po lewej stronie znaku równości. Przede wszystkim zadaje ona jedną lub kilka wartości ósemkowych, które ten repre-

zementant może przyjmować na wejściu. Na przykład, gdyby opis języka mógł wchodzić do maszyny w kodach M2, BCDIC oraz KW8<sup>\*</sup>), wówczas w przypadku definicji litery A, składowa <ciąg kodów źródłowych> miałyby postać: (30, 301, 202).

Składowa <ciąg kodów nieistotnych> określa kodowe wartości znaków, które należy opuścić przy wczytywaniu, jeżeli występują one za zdefiniowanym znakiem. Na przykład, dla kodu KW8 następująca definicja spacji (niepełna - bez określenia składowej <kod wynikowy>)

§ SP = (100) (100,24)

spowoduje, że po wczytaniu jednej spacji jako znaku z wejścia, następujące za nią spacje oraz nowe linie będą ignorowane.

Składowa <kod wynikowy> określa wartość kodową opisanego znaku na wyjściu konstruowanego kompilatora. Zakładając na przykład, że mamy do czynienia z kodem KW8 zarówno na wejściu, jak i na wyjściu, możemy znak ";" zdefiniować następująco:

§ CL = (164 ` 164

Zbiór znaków metajęzyka Cogent ma (w realizacji na CDC 3600) zadane z góry standardowe definicje. Dlatego w opisie języka źródłowego, sformułowanym w metajęzyku Cogent, definiujemy jedynie te znaki, które do tego metajęzyka nie należą oraz np. słowa zastrzeżone. Początek opisu języka może mieć następującą postać (przyjmując kod KW8):

§ CHARDEF § CL = (164)164 § SE = (166)166

§ BEGIN = (302)0 § END = (303)0

§ REAL = (304)0 § GOTO = (307)0

Trochę bardziej rozbudowany aparat do analizy leksykalnej istnieje w języku Trandir [20]. Rolę analizy leksykalnej pełni w tym języku specjalna procedura SCAN. Wywołanie SCAN(n) powoduje wczytanie n słów z ciągu wejściowego i umieszczenie ich

<sup>\*</sup>) Przez kod rozumiemy przypisanie każdemu znakowi pewnej liczby. Na przykład: litera A ma w kodzie M2 wartość: 30, w kodzie BCDIC: 301, a w kodzie KW8: 202



w odpowiednich tablicach. Projektant języka definiuje tablice słów (ich pola i wymiary) za pomocą odpowiednich reguł języka Trandir. Jednocześnie może on definiować klasy słów syntaktycznie synonimicznych oraz równoważne postacie zewnętrzne dla tych samych symboli. Część programu w języku Trandir może wyglądać następująco:

```

DECLARE TABLES          : deklaracja tablic
OPTAB (1.25)            : słownik operacji języka o wymiarze 25 oraz o dwóch polach
WITH FIELDS
  COPLM
  OPCLASS
END
DECLARE ENTRIES
OPTAB FORMAT           : definicja elementów tablicy
% DECLARE..
% INT..
% GO..
% EQ..          EQUAL          RELOP
% NE..          NEQUAL         RELOP
% GT..          GREAT          RELOP
% GE..          GREATEREQ     RELOP
% LT..          LESS           RELOP
% LE..          LESSEQ        RELOP
% END..
+..            PLUS           ADOP
-..            MINUS          ADOP
END
                |               |
                |               |
      (pole COPLM)   (pole OPCLASS)

```

W tablicy OPTAB są zgromadzone operatory języka wejściowego, ich odpowiedniki w języku przejściowym (pole COPLM) oraz klasy równoważności do celów analizy syntaktycznej (pole OPCLASS).

Jeszcze inne możliwości zastosowania analizy leksykalnej zapewnia metatranslator ATWS [24]. W systemie tym projektant języka definiuje atomy leksykalne oraz procedury do rozpozna-

wania tych atomów. Na podstawie tych definicji konstruowany jest L-analizator służący do wydzielania L-atomów w programie źródłowym.

Przykład programu w języku ATWS:

```
begin
.....
class ELEM, LETTER, DIGIT, TOKEN, KEYWORD, ..., SYMBOL: lista
                                atomów leksykalnych
phase <identifier> = begin                                : defini-
                                <identifier ::=      cja procedury rozpozn-
                                <l>{<ld>}              najacej
                                <ld> ::= exclusive {<l>}{<d>}
                                <l> ::= LETTER
                                <d> ::= DIGIT
                                end
phase <keyword> = ...
phase <string> = ...
init   ELEM with (<...>), : definicja L-atomów
        LETTER with (a, b, ..., z),
        .
        .
        ID with <identifier>,
        KEYWORD with <keyword>,
        ...
        SYMBOL with (begin, end, int, +, ...)
```

### 3.2. System AED RWORD

Omówione powyżej systemy albo nie są zrealizowane w praktyce (ATWS), albo nie są metatranslatorami lecz językami do pisania translatorów, w związku z czym L-analizator musi być programowany "ręcznie" przez programistę, nie zaś generowany automatycznie na podstawie zadanej gramatyki. Jedynymi działającymi obecnie rozwiązaniami automatycznej konstrukcji L-analizatora są: system AED RWORD [11] oraz L-analizator opisany w pracy [15], zwany poniżej L-analizatorem Lecarma.



System AED RWORD (Read a WORD) jest pierwszym systemem, w którym L-analizator jest generowany automatycznie. Jest on częścią systemu AED-1, którego celem jest generowanie kompilatorów, interpretatorów, systemów operacyjnych, itd. System AED RWORD jest stosowany nie tylko do celów analizy leksykalnej. Uniwersalność systemu powoduje, że jest on mniej efektywny. L-analizator skonstruowany przez RWORD na maszynie IBM 7094 przetwarzał 3000 znaków/s, podczas gdy L-analizator zakodowany "ręcznie" przetwarzał 8000 znaków/s.

W systemie RWORD zakłada się wydzielenie w opisie składni języka podgramatyk regularnych, na podstawie których konstruowane są automaty skończone. Podgramatyki te stanowią dane wejściowe do systemu i są zapisywane w następującej postaci:<sup>\*</sup>

```
BEGIN { < opis klasy znaków > } END
BEGIN { < opis atomu leksykalnego > } END
FINI
```

### 3.2.1. Opis klasy znaków

Klasa znaków może być zdefiniowana dwoma sposobami: w pierwszym znaki wymienia się *explicite*, w drugim - podaje się ich wewnętrzną reprezentację (w kodzie EBCDIC). Na przykład:

```
LET = /ABC... YZ/
PUNCTUATION = A+-.,; =/A
XX = §21,22,25§
```

Pierwszy znak następujący po znaku "=" i różny od spacji jest ogranicznikiem klasy znaków. Użycie znaku "§" jako ogranicznika oznacza, że definiujemy klasę przez podanie reprezentacji znaków w kodzie EBCDIC. Klasa LET składa się więc z wszystkich liter, a klasa PUNCTUATION - ze znaków + - . , ; = oraz /. Klasa XX składa się ze znaków o reprezentacji wewnętrznej 21, 22, 25.

<sup>\*</sup> Zapis {A} oznacza ciąg (być może pusty) powtórzeń symbolu A

Nazwa klasy musi być różna od BEGIN, END, FINI. Nazwa IGNORE jest używana do zdefiniowania klasy znaków ignorowanych przez L-analizator.

### 3.2.2. Opis atomu leksykalnego

Opis atomu leksykalnego ma następującą postać:

$\langle \text{identyfikator} \rangle \langle \langle \text{integer} \rangle \rangle = \langle \text{wyrażenie regularne} \rangle \text{\$}$

gdzie zmiennymi wyrażenia regularnego mogą być:

- a) dowolne znaki kodu EBCDIC różne od |, (, ) , ' , \text{\\$} i znaku spacji,
- b) puste słowo  $\epsilon$ ,
- c) nazwy klas znaków,
- d) znak apostrofu z następującym po nim dowolnym znakiem kodu EBCDIC różnym od spacji; wartością takiej pary jest znak następujący po apostrofie (jest to jedyna możliwość używania znaków niedopuszczalnych w punkcie a).

Przykłady:

$\text{PUNCT (5)} = \text{PUNCTUATION } \text{\$}$

Klasa leksykalnych atomów PUNCT składa się ze znaków z klasy PUNCTUATION; atomy z tej klasy mają umowną reprezentację 5.

$\text{BEG (3)} = \text{'BEGIN'} \text{\$}$

oznacza, że słowo 'BEGIN' ma umowną reprezentację 3,

$\text{ID (1)} = \text{LET (LET | DIGIT)}^* \text{\$}$

oznacza, że słowo złożone najwyżej z sześciu liter i cyfr, którego pierwszym znakiem jest litera, ma umowną reprezentację 1.

Specjalną klasą jest klasa IGNORE. Każde słowo z tej klasy jest ignorowane. Ponieważ umowna reprezentacja nie jest w tym przypadku potrzebna, definicja tej klasy ma następującą postać:



IGNORE = < wyrażenie regularne > §

Zauważmy różnicę między klasą znaków IGNORE i klasą atomów leksykalnych IGNORE. Na przykład:

```
BEGIN IGNORE = / / LET = /AB/ END
BEGIN ID (1) = LET (LET) § END
FINI
```

Ciąg znaków AB A będzie traktowany jako identyfikator ABA, ponieważ znak spacji jest ignorowany. Natomiast przy następującej definicji:

```
BEGIN SPACE = / / LET = /AB/ END
BEGIN IGNORE = SPACE § ID = LET (LET) § END
FINI
```

ciąg AB A będzie traktowany jako dwa identyfikatory AB i A.

Klasę atomów leksykalnych można również zadeklarować w następujący sposób:

<identyfikator>(<integer>, <sub name>) = <wyrażenie regularne > §

wtedy za każdym razem po skonstruowaniu atomu należącego do klasy <identyfikator > zostaje wywołana procedura <sub name>. Procedura związana z identyfikatorem może umieszczać go w słowniku symboli i sprawdzać, czy dany identyfikator jest słowem zastrzeżonym, czy też nie.

W wielu przypadkach można dokonać rozbioru wejściowego ciągu znaków kilkoma sposobami. System RWORD wykrywa wszystkie wieloznaczności i podejmuje jedną z następujących akcji:

1. Podejmuje decyzję usunięcia wieloznaczności. Na przykład ciąg znaków AB może oznaczać zarówno identyfikator AB jak i dwa identyfikatory A oraz B. System RWORD konstruuje wtedy najdłuższe możliwe słowo i przechodzi do dalszej analizy.

2. Podejmuje decyzję usunięcia wieloznaczności i powiadamia o tym operatora. Na przykład ciąg znaków X123 może być traktowany jako identyfikator X123 lub też jako identyfikator X oraz liczba całkowita 123. W takim przypadku RWORD też konstruuje najdłuższe słowo.
3. Komunikuje, że wieloznaczność jest zbyt skomplikowana i wymaga analizy kontekstu (na przykład instrukcja Fortranu DO10I = 1.20, której pierwsza część (DO10I) mogłaby zostać potraktowana jako identyfikator), po czym kończy pracę.

Poniżej pokażemy, jak można by zapisać jako dane wejściowe systemu RWORD część języka z przykładu 2.1.

```
BEGIN SPACE = / /
LET = /abdegin/
DIG = /01/
OP = /;:=/
DOT = ./

END
BEGIN OPR (1) = OP §
ID (2,LOOKUP) = LET (LET|DIG)* §
INT(3,CONVERTI) = DIG (DIG)* §
REAL (4,CONVERTR) = DIG (DIG)*.DIG (DIG)* §
KEY (5) = (.begin|.end.) §
END FINI
```

Procedura LOOKUP będzie umieszczala identyfikator w słowniku symboli, a procedury CONVERTI oraz CONVERTR obliczały wartość numeryczną liczby.

W pewnych językach programowania można wydzielić kilka podjęzyków, w których mogą być zupełnie inne reguły tworzenia słów (np. w Algolu 60 spacje mają znaczenie wewnątrz tekstów, a poza nimi są ignorowane). W takim przypadku dobrze jest mieć kilka różnych L-analizatorów (po jednym dla każdego podjęzyka) oraz aparat pozwalający przechodzić od jednego L-analizatora do drugiego. Na przykład, gdyby w powyższym języku



założyć, że komentarz zaczyna się i kończy znakiem "/", wówczas można by wprowadzić definicję:

```
SL '6,COMMENT' = '/§
```

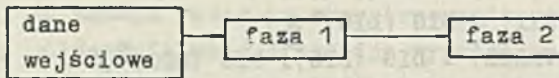
gdzie procedura COMMENT wywołuje inny, podany poniżej L-analizator, który ominie komentarz i powróci do poprzedniego L-analizatora.

L-analizator akceptujący komentarz wygląda następująco:

```
BEGIN IGNCHAR = §00,01,...,77§ 'wszystkie znaki poza "/"
END
BEGIN IGNORE = IGNCHAR §
TERMIN (1,ENDCOMMENT) = '/§
END FINI
```

gdzie procedura ENDCOMMENT przekazuje sterowanie do głównego L-analizatora,

System RWORD składa się z dwóch części. Pierwszą z nich można zilustrować następującym diagramem:



W fazie 1 są tworzone automaty niedeterministyczne dla każdego wyrażenia regularnego. W fazie 2 automaty te są łączone i powstaje jeden automat deterministyczny zapisany w języku AED-0. Składa się on z wywołań procedur i tablic skonstruowanych w obu fazach.

Zadaniem drugiej części systemu RWORD jest optymalizacja wyprodukowanego programu ze względu na czas wykonywania i wykorzystanie pamięci maszyny.

### 3.3. L-analizator Lecarma

W systemie Lecarma przyjęto jako główne założenie fakt, że w większości języków programowania powtarzają się typowe atomy

leksykalne. Wobec tego w systemie tym zdefiniowano następujące L-atomy:

- a) identyfikator - ciąg liter i cyfr zaczynający się od litery; długość maksymalna identyfikatora jest ustalana przez projektanta języka;
- b) słowo zastrzeżone - atom o analogicznej składni jak identyfikator, o ustalonym z góry znaczeniu;
- c) słowo kluczowe - ciąg liter i cyfr zaczynający się od litery i ograniczony z obu stron specjalnym znakiem zdefiniowanym przez projektanta języka;
- d) operator - jeden lub dwa znaki różne od litery, cyfry i spacji;
- e) liczba całkowita - ciąg cyfr dziesiętnych (bez znaku);
- f) liczba rzeczywista - atomy te są utworzone z części całkowitej, części ułamkowej i części wykładniczej, przy czym druga lub trzecia część może nie występować; część całkowita jest ciągiem cyfr dziesiętnych; część ułamkowa składa się z ciągu cyfr dziesiętnych następującego po znaku DELFRACTION - definiowanym przez projektanta (domyślnie przyjmuje się kropkę); część wykładnicza jest oznaczona znakiem DELEXPOSAN - również definiowanym przez projektanta (domyślnie litera E) - po którym następuje ciąg cyfr poprzedzony ewentualnie znakiem + lub - ;



- g) tekst - ciąg dowolnych znaków ograniczony z lewej strony przez znak DEBCHaine, a z prawej strony przez znak FINCHaine;
- h) komentarz - ciąg dowolnych znaków ograniczony z lewej strony przez znak DEBCOMMENT, a z prawej strony przez znak FINCOMMENT.

Projektant języka definiuje tylko pewne uzupełnienia do powyższych definicji, jak np.:

- DEBCHaine, FINCHaine - znaki początku i końca tekstu,  
DEBCOMMENT, FINCOMMENT - znaki początku i końca komentarza,  
DELEXPOSANT, DELFRACTION - znaki reprezentujące odpowiednio wykładnik i kropkę dziesiętną w liczbie rzeczywistej,  
DELMOTCLÉ - znak ograniczający słowo kluczowe (z obu stron),  
MAXIDENT - maksymalna dopuszczalna długość identyfikatora.

W związku z tym automaty akceptujące atomy leksykalne są zaprogramowane "z góry" i wprowadza się do nich - dla konkretnego języka - tylko nieznaczne poprawki, co umożliwia bardzo efektywną pracę L-analizatora.

Wadą tej metody jest fakt ograniczenia liczby atomów leksykalnych, jednakże autorzy twierdzą, że w praktyce bardzo rzadko spotyka się języki mające inne L-atomy.

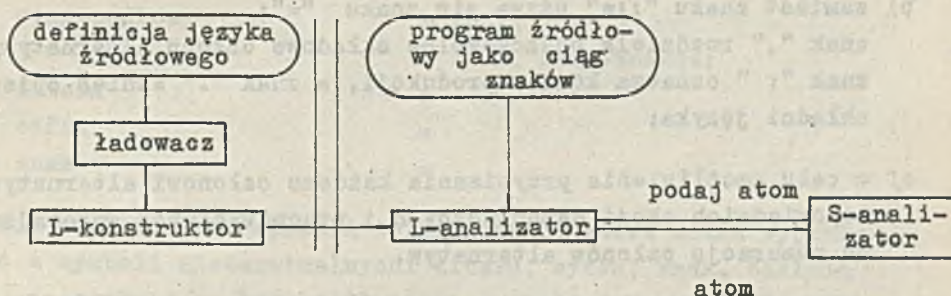
#### 4. OPIS KONSTRUKCJI MODUŁU ANALIZY LEKSYKALNEJ

Przedstawimy w tym rozdziale nowe rozwiązanie koncepcji L-analizatora, zastosowane w metatranslatorze realizowanym w Instytucie Maszyn Matematycznych. Proponowana metoda polega na automatycznym wydzieleniu podzbiorów regularnych języka, dokonywanym na podstawie opisu jego składni. W związku z tym nie ma potrzeby wprowadzania dodatkowego metajęzyka, służącego do opisu części leksykalnej (co jest najpoważniejszą wadą systemu

RWORD), Nie narzuca się również (w odróżnieniu od L-analizatora Lecarma) żadnych ograniczeń na postać atomów leksykalnych.

#### 4.1. Organizacja ogólna

Moduł analizy leksykalnej składa się z trzech podmodułów: ładowacza, L-konstruktora i L-analizatora. Ogólny schemat działania tych podmodułów wygląda następująco:



Wszystkie podmoduły zostały zaprogramowane w języku PL/1 dla maszyny IBM 370/175. Podmoduły ładowacza i L-analizatora są realizowane przez procedury o nazwach ŁADOWACZ I LANALIZATOR, natomiast na podmoduł L-konstruktora składają się dwie procedury - LKONSTRUKTOR oraz REGULAR. Całość programu zawiera około 800 instrukcji w języku PL/1 i zajmuje w pamięci maszyny około 22000 bajtów.

W dodatku A zamieszczono teksty wszystkich programów.

#### 4.2. Podmoduł ładowacza

Zadaniem podmodułu jest wczytanie opisu składni języka źródłowego, umieszczenie jej w pamięci maszyny oraz utworzenie pewnych tablic, z których będą korzystały kolejne podmoduły. Realizuje to procedura zewnętrzna - ŁADOWACZ, zawierająca około 200 instrukcji w języku PL/1.



#### 4.2.1. Metajęzyk opisu składni języka źródłowego

Metajęzyk opisu składni oparty jest w zasadzie na metajęzyku Backusa-Naura, do którego wprowadzono następujące zmiany:

- a) w celu umożliwienia formalnego zdefiniowania metajęzyka opisu składni w tym metajęzyku, wprowadzono apostrofy ograniczające symbole terminalne gramatyki; nie stosuje się natomiast nawiasów kątowych, używanych w notacji BNF do ograniczania symboli nieterminalnych;
- b) zamiast znaku " := " używa się znaku " = "; znak " , " rozdziela poszczególne składowe członu alternatywy; znak " ; " oznacza koniec produkcji, a znak " . " koniec opisu składni języka;
- c) w celu umożliwienia przypisania każdemu członowi alternatywy odpowiednich akcji semantycznych i pragmatycznych wprowadzono numerację członów alternatyw.

Formalna definicja metajęzyka opisu składni zapisana w tym metajęzyku jest następująca:

|                         |   |
|-------------------------|---|
| 1: opis-gramatyki       | = ciąg-produkcji ' . ' ;                                    |
| 2: ciąg-produkcji       | = produkcja   |
| 3:                      | produkcja, ' ; ' , ciąg-produkcji;                          |
| 4: produkcja            | = ident-czł-altern, ' : ' , l-str-prod, ' = ' , p-str-prod; |
| 5: ident-czł-altern     | = cyfra   |
| 6:                      | . cyfra, ident-czł-altern;                                  |
| 7: l-str-prod           | = nieterminal;  |
| 8: nieterminal          | = litera  |
| 9:                      | nieterminal, znak-1;  |
| 10: znak-1              | = litera  |
| 11:                     | cyfra   |
| 12:                     | ' - ' ;   |
| 13: p-str-prod          | = alternatywa;  |
| 14: alternatywa         | = ciąg-członów-altern;                                      |
| 15: ciąg-członów-altern | = człon-alternatywy   |

|     |                            |  |
|-----|----------------------------|--|
| 16: |                            | człon-alternatywy, ' ', ident-<br>czł-altern, ':', ciąg-członów-<br>altern;    |
| 17: | człon-alternatywy          | = ciąg-składowych-czł-altern;  |
| 18: | ciąg-składowych-czł-altern | = skład-czł-altern  <br>skład-czł-altern, ',', ciąg-<br>składowych-czł-altern; |
| 20: | skład-czł-altern           | = symbol;  |
| 21: | symbol                     | = terminal   |
| 22: |                            | nieterminal;   |
| 23: | terminal                   | = '""', ciąg-znaków, '""';   |
| 24: | ciąg-znaków                | = znak   |
| 25: |                            | znak, ciąg-znaków;   |
| 26: | litera                     | =  |
| 27: | cyfra                      | =  |
| 28: | znak                       | =  |

Nie zdefiniowano symboli terminalnych, które można wyprowadzić z symboli nieterminalnych: litera, cyfra, znak. Zakłada się, że symbole te będą definiowane w każdej konkretnej realizacji.

Przykład 4.1. Rozważmy gramatykę z przykładu 2.1 usuwając z niej tylko kropki ograniczające słowa begin i end. W naszym metajęzyku gramatyka ta będzie zapisana następująco:

|     |            |                          |
|-----|------------|--------------------------|
| 1:  | program    | = 'begin', block, 'end'; |
| 2:  | block      | = statement, ';'         |
| 3:  |            | block, statement, ';' ;  |
| 4:  | statement  | = leftside, number;      |
| 5:  | leftside   | = identifier, ':=' ;     |
| 6:  | identifier | = letter                 |
| 7:  |            | identifier, letter       |
| 8:  |            | identifier, digit;       |
| 9:  | number     | = integer                |
| 10: |            | real;                    |
| 11: | integer    | = digit                  |
| 12: |            | integer, digit;          |
| 13: | real       | = integer, '.', integer; |



|     |        |   |     |   |
|-----|--------|---|-----|---|
| 14: | letter | = | 'a' |   |
| 15: |        |   | 'b' |   |
| 16: |        |   | 'd' |   |
| 17: |        |   | 'e' |   |
| 18: |        |   | 'g' |   |
| 19: |        |   | 'i' |   |
| 20: |        |   | 'n' |   |
| 21: | digit  | = | '0' |   |
| 22: |        |   | '1' | . |

#### 4.2.2. Procedura ŁADOWACZ

Procedura ta wczytuje opis gramatyki języka źródłowego i sprawdza jego poprawność. Wykrywane i sygnalizowane są następujące błędy:

- po lewej stronie znaku równości w produkcji występuje symbol terminalny;
- brak znaku "=" kończącego lewą stronę produkcji;
- brak dwukropka po identyfikatorze produkcji;
- symbol nieterminalny nie zaczyna się od litery;
- brak aksjomatu gramatyki.

Po stwierdzeniu poprawności opisu gramatyki, procedura ŁADOWACZ umieszcza go w odpowiednich tablicach w pamięci maszyny.

#### 4.2.3. Postać gramatyki języka źródłowego w pamięci maszyny

Ładowacz umieszcza wczytaną gramatykę w czterech tablicach: SYMBTAB, PRODAT, PROD i ALTER, ponadto wyszukuje aksjomat gramatyki<sup>\*)</sup> i umieszcza go w specjalnej zmiennej - STARTSYMB.

Tablice SYMBTAB (tablica symboli gramatyki) i PRODAT (tablica adresów produkcji) mają rozmiary równe liczbie wszystkich symboli (nieterminalnych i terminalnych) gramatyki języka wejściowego.

\*) Za aksjomat przyjmujemy symbol nieterminalny nie występujący po prawej stronie żadnej produkcji

Każdy wiersz tablicy SYMTAB jest strukturą (w sensie języka PL/1) składającą się ze zmiennej znakowej SYMBOL (w zmiennej tej zapisywany jest odpowiedni symbol gramatyki) oraz zmiennej bitowej TCNT określającej, czy dany symbol jest symbolem terminalnym czy też nieterminalnym.

Tablica PRODAT zawiera dla każdego symbolu adres, pod którym są zapisane wszystkie produkcje mające jako lewą stronę dany symbol (w przypadku symbolu terminalnego adres ten jest równy 0).

Tablica PROD (tablica produkcji) jest strukturą o dwóch polach. Pierwsze z nich (ILOŚĆ) jest zmienną liczbową określającą liczbę członów alternatywy będącej prawą stroną produkcji dla danego symbolu nieterminalnego. Drugie pole (ADRALT) - jest tablicą zawierającą adresy kolejnych członów alternatywy; rozmiar tej tablicy zależy od pola ILOŚĆ.

Tablica ALTER (tablica alternatyw) jest strukturą złożoną z trzech elementów. Dwa pierwsze elementy (DALT i IDALT) są zmiennymi liczbowymi określającymi odpowiednio: długość i identyfikator członu alternatywy. Trzeci element (SYMBOL) jest tablicą o rozmiarze równym długości członu alternatywy; w kolejnych wierszach tablicy są zapisywane adresy w tablicy SYMTAB kolejnych składowych członów alternatywy.



Przykład 4.2. Pokażemy, jak po zakończeniu działania procedury ŁADOWACZ będzie rozmieszczona w pamięci maszyny gramatyka podana w przykładzie 4.1.

|    | SYMBOL     | TCNT |
|----|------------|------|
| 1  | program    | 0    |
| 2  | begin      | 1    |
| 3  | block      | 0    |
| 4  | end        | 1    |
| 5  | statement  | 0    |
| 6  | ;          | 1    |
| 7  | leftside   | 0    |
| 8  | number     | 0    |
| 9  | identifier | 0    |
| 10 | :=         | 1    |
| 11 | letter     | 0    |
| 12 | digit      | 0    |
| 13 | integer    | 0    |
| 14 | real       | 0    |
| 15 | .          | 1    |
| 16 | a          | 1    |
| 17 | b          | 1    |
| 18 | d          | 1    |
| 19 | e          | 1    |
| 20 | g          | 1    |
| 21 | i          | 1    |
| 22 | n          | 1    |
| 23 | 0          | 1    |
| 24 | 1          | 1    |

tablica SYMBTAB

|       |
|-------|
| ADR1  |
| 0     |
| ADR2  |
| 0     |
| ADR3  |
| 0     |
| ADR4  |
| ADR5  |
| ADR6  |
| 0     |
| ADR7  |
| ADR8  |
| ADR9  |
| ADR10 |
| 0     |
| 0     |
| 0     |
| 0     |
| 0     |
| 0     |
| 0     |
| 0     |
| 0     |
| 0     |
| 0     |

← adres produkcji, której lewą stroną jest symbol "program"

tablica PRODAB

|       | liczba członów altern. | adres członu altern. |
|-------|------------------------|----------------------|
| ADR1: | IŁOŚĆ                  | ADRALT               |
|       | 1                      | ADR11                |

tablica PROD

|        | długość członu altern. | identyf. członu altern. | adresy w SYMBTAB kolej.skł. członu alternat. |   |   |
|--------|------------------------|-------------------------|--|---|---|
| ADR11: | DALT                   | IDALT                   | SYMBOL                                       |   |   |
|        | 3                      | 1                       | 2  | 3 | 4 |

tablica ALTER

|       | IŁOŚĆ | ADRALT |       |
|-------|-------|--------|-------|
| ADR2: | 2     | ADR12  | ADR13 |

tablica PROD

|        | DALT | IDALT | SYMBOL |   |
|--------|------|-------|--------|---|
| ADR12: | 2    | 2     | 5      | 6 |

tablica ALTER

|        | DALT | IDALT | SYMBOL |   |   |
|--------|------|-------|--------|---|---|
| ADR13: | 3    | 3     | 3      | 5 | 6 |

tablica ALTER

|       | IŁOŚĆ | ADRALT |
|-------|-------|--------|
| ADR3: | 1     | ADR14  |

tablica PROD

|        | DALT | IDALT | SYMBOL |   |
|--------|------|-------|--------|---|
| ADR14: | 2    | 4     | 7      | 8 |

tablica ALTER

|       | IŁOŚĆ | ADRALT |
|-------|-------|--------|
| ADR4: | 1     | ADR15  |

tablica PROD

|        | DALT | IDALT | SYMBOL |    |
|--------|------|-------|--------|----|
| ADR15: | 2    | 5     | 9      | 10 |

tablica ALTER

\*) rysunek ADRi: 

|  |  |
|--|--|
|  |  |
|--|--|

 oznacza, że pierwszy element tablicy A zapamiętany jest pod adresem ADRi.

tablica A



|       |       |        |       |  |
|-------|-------|--------|-------|--|
| ADR5: | ILOŚĆ | ADRALT |       |  |
|       | 2     | ADR16  | ADR17 |  |

tablica PROD

|        |      |       |        |  |
|--------|------|-------|--------|--|
| ADR16: | DALT | IDALT | SYMBOL |  |
|        | 1    | 9     | 13     |  |

tablica ALTER

|        |      |       |        |  |
|--------|------|-------|--------|--|
| ADR17: | DALT | IDALT | SYMBOL |  |
|        | 1    | 10    | 14     |  |

tablica ALTER

|       |       |        |       |       |
|-------|-------|--------|-------|-------|
| ADR6: | ILOŚĆ | ADRALT |       |       |
|       | 3     | ADR18  | ADR19 | ADR20 |

tablica PROD

|        |      |       |        |  |
|--------|------|-------|--------|--|
| ADR18: | DALT | IDALT | SYMBOL |  |
|        | 1    | 6     | 11     |  |

tablica ALTER

|        |      |       |        |    |
|--------|------|-------|--------|----|
| ADR19: | DALT | IDALT | SYMBOL |    |
|        | 2    | 7     | 9      | 11 |

tablica ALTER

|        |      |       |        |    |
|--------|------|-------|--------|----|
| ADR20: | DALT | IDALT | SYMBOL |    |
|        | 2    | 8     | 9      | 12 |

tablica ALTER

|       |       |        |       |       |
|-------|-------|--------|-------|-------|
| ADR7: | ILOŚĆ | ADRALT |       |       |
|       | 7     | ADR21  | ADR22 | ADR23 |
|       |       | ADR24  | ADR25 | ADR26 |
|       |       | ADR27  |       |       |

tablica PROD

|        |      |       |        |  |
|--------|------|-------|--------|--|
| ADR21: | DALT | IDALT | SYMBOL |  |
|        | 1    | 14    | 16     |  |

tablica ALTER

|        |      |       |        |  |
|--------|------|-------|--------|--|
| ADR22: | DALT | IDALT | SYMBOL |  |
|        | 1    | 15    | 17     |  |

tablica ALTER

ADR23:

| DALT | IDALT | SYMBOL |
|------|-------|--------|
| 1    | 16    | 18     |

tablica ALTER

ADR24:

| DALT | IDALT | SYMBOL |
|------|-------|--------|
| 1    | 17    | 19     |

tablica ALTER

ADR25:

| DALT | IDALT | SYMBOL |
|------|-------|--------|
| 1    | 18    | 20     |

tablica ALTER

ADR26:

| DALT | IDALT | SYMBOL |
|------|-------|--------|
| 1    | 19    | 21     |

tablica ALTER

ADR27:

| DALT | IDALT | SYMBOL |
|------|-------|--------|
| 1    | 20    | 22     |

tablica ALTER

ADR8::

| ILOŚĆ | ADRALT |       |
|-------|--------|-------|
| 2     | ADR28  | ADR29 |

tablica PROD

ADR28:

| DALT | IDALT | SYMBOL |
|------|-------|--------|
| 1    | 21    | 23     |

tablica ALTER

ADR29:

| DALT | IDALT | SYMBOL |
|------|-------|--------|
| 1    | 22    | 24     |

tablica ALTER



ADR9:

| ILOŚĆ | ADRALT |       |
|-------|--------|-------|
| 2     | ADR30  | ADR31 |

tablica PROD

ADR30:

| DALT | IDALT | SYMBOL |  |
|------|-------|--------|--|
| 1    | 11    | 12     |  |

tablica ALTER

ADR31:

| DALT | IDALT | SYMBOL |    |
|------|-------|--------|----|
| 2    | 12    | 13     | 12 |

tablica ALTER

ADR10:

| ILOŚĆ | ADRALT |  |
|-------|--------|--|
| 1     | ADR32  |  |

tablica PROD

ADR32:

| DALT | IDALT | SYMBOL |    |    |
|------|-------|--------|----|----|
| 3    | 13    | 13     | 15 | 13 |

tablica ALTER

#### 4.3. Podmoduł L-konstruktora

Podmoduł ten znajduje podjęzyki regularne zadanego języka źródłowego i tworzy deterministyczny automat skończenie stanowy akceptujący te podjęzyki. Określenie podjęzyków regularnych realizuje procedura REGULAR, a tworzenie automatu - procedura LKONSTRUKTOR. Każda z procedur zawiera około 80 instrukcji w języku PL/1.

##### 4.3.1. Procedura REGULAR

Wejściem dla procedury REGULAR są tablice SYMBTAB, PRODAT, PROD i ALTER; wyjściem - tablica RCNR. Tablica ta ma rozmiar równy rozmiarowi tablicy SYMBTAB. Każdy element tablicy RCNR przyjmuje wartość 1 lub 0 zależnie od tego, czy odpowiadający mu symbol tablicy SYMBTAB, przyjęty jako aksjomat języka, generuje język regularny czy też nie.

Algorytm działania procedury REGULAR jest następujący:

Krok 1: Dla każdego  $j$  od 1 do  $N$  ( $N$  jest rozmiarem tablicy SYMBTAB) wykonaj:

$$RCNR(j) \leftarrow TCNT(j)$$

Krok 2:  $ROB \leftarrow 0$

Krok 3: Dla każdego  $j$  od 1 do  $N$ , takiego że  $RCNR(j) = 0$  wykonaj:

Oznaczmy przez  $A$  wartość pola SYMBOL w  $j$ -tym wierszu tablicy SYMBTAB, a przez  $P_A$  zbiór wszystkich członów alternatywy produkcji, której lewą stroną jest  $A$ . Jeżeli dla każdego członu  $p$  ze zbioru  $P_A$ ,  $p$  jest jednej z następujących postaci:  $B$  lub  $CB$ , lub  $AB$  i przy tym  $RCNR(k) = 1$ ,  $TCNT(l) = 0$  i  $RCNR(l) = 1$ , gdzie  $k, l$  są odpowiednio adresami  $B, C$  w SYMBTAB, to wykonaj:

$$RCNR(j) \leftarrow 1 \quad \text{oraz} \quad ROB \leftarrow 1;$$

Krok 4: Jeżeli  $ROB = 1$ , to przejdź do wykonywania kroku 2, w przypadku przeciwnym zakończ pracę procedury REGULAR.

Przykład 4.3. Zilustrujemy działanie powyższego algorytmu dla gramatyki podanej w przykładzie 4.1.

Po wykonaniu kroków 1 i 2 zmienna  $ROB = 0$ , a tablica  $RCNR$  wygląda następująco:



|    |   |
|----|---|
| 1  | 0 |
| 2  | 1 |
| 3  | 0 |
| 4  | 1 |
| 5  | 0 |
| 6  | 1 |
| 7  | 0 |
| 8  | 0 |
| 9  | 0 |
| 10 | 1 |
| 11 | 0 |
| 12 | 0 |
| 13 | 0 |
| 14 | 0 |
| 15 | 1 |
| 16 | 1 |
| 17 | 1 |
| 18 | 1 |
| 19 | 1 |
| 20 | 1 |
| 21 | 1 |
| 22 | 1 |
| 23 | 1 |
| 24 | 1 |

tablica RCNR

Po wykonaniu kroku 3 otrzymamy trzy nowe zmienne "regularne", mianowicie: digit, letter, integer; tzn. zostaną wykonane operacje przypisania:

RCNR(11) ← 1; RCNR(12) ← 1; RCNR(13) ← 1; oraz ROB ← 1.

Po ponownym wykonaniu kroku 2 oraz kroku 3 dochodzi nowa zmienna "regularna" - identifier, czyli RCNR(9) ← 1. Ponieważ ROB = 1, więc znów zostaną wykonane krok 2 i krok 3; otrzymany wówczas RCNR(7) = 1. Dopiero po czwartej iteracji będzie ROB = 0 (tzn. nie otrzymamy żadnej nowej zmiennej "regularnej"), co kończy pracę algorytmu. Na wyjściu otrzymujemy następującą tablicę RCNR:

|    |   |
|----|---|
| 1  | 0 |
| 2  | 1 |
| 3  | 0 |
| 4  | 1 |
| 5  | 0 |
| 6  | 1 |
| 7  | 1 |
| 8  | 0 |
| 9  | 1 |
| 10 | 1 |
| 11 | 1 |
| 12 | 1 |
| 13 | 1 |
| 14 | 0 |
| 15 | 1 |
| 16 | 1 |
| 17 | 1 |
| 18 | 1 |
| 19 | 1 |
| 20 | 1 |
| 21 | 1 |
| 22 | 1 |
| 23 | 1 |
| 24 | 1 |

TABLICA RCNR



#### 4.3.2. Uzasadnienie teoretyczne procedury REGULAR

Pokażemy, że procedura REGULAR znajduje rzeczywiście podjęzyki regularne zadanego języka źródłowego.

##### Twierdzenie 4.1

Klasa języków regularnych jest zamknięta ze względu na operację podstawienia.

Dowód patrz [6].

Definicja 4.1. Niech  $G = \langle V_N, V_T, P, S \rangle$  będzie gramatyką bezkontekstową. Oznaczmy  $V = V_N \cup V_T$ .

Określmy funkcję częściową  $\rho : V \rightarrow I^0$  następująco<sup>\*</sup>):

1.  $\rho(A) = 0 \iff A \in V_T$
2.  $\rho(A) = n \iff n$  jest najmniejszą z możliwych liczb  $K$  takich, że dla każdego  $\alpha$  takiego, że  $(A, \alpha) \in P$   $\alpha = B$  lub  $\alpha = CB$  lub  $\alpha = AB$ , gdzie  $\rho(B) < K$  i  $0 < \rho(C) < K$ .

Funkcję  $\rho$  nazywamy rzędem regularności gramatyki  $G$ .

Definicja 4.2. Gramatyka  $G = \langle V_N, V_T, P, S \rangle$  jest quasi-regularna  $\iff$  dla każdego  $A \in V$  zachodzi  $\rho(A) \in I^0$ .

Symbol  $A$  należący do  $V$  będziemy nazywać quasi-regularnym  $\iff \rho(A) \in I^0$ .

Symbol  $A \in V$  jest regularny rzędu  $n$   $\iff \rho(A) = n$ .

##### Twierdzenie 4.2

Jeżeli gramatyka  $G = \langle V_N, V_T, P, S \rangle$  jest gramatyką quasi-regularną, to język  $L(G)$  jest regularny.

---

<sup>\*</sup>) Przez  $I^0$  rozumiemy zbiór liczb całkowitych nieujemnych

Dowód:

Zauważmy najpierw, że dla każdego  $A \in V$  zachodzi  $\rho(A) \leq \rho(S)$ .

Dowód twierdzenia przeprowadzimy przez indukcję względem rzędu regularności aksjomatu gramatyki.

Jeżeli  $\rho(S) = 1$ , to jedynymi produkcjami z  $P$  mającymi po lewej stronie symbol  $S$  mogą być produkcje postaci  $S \rightarrow a$  lub  $S \rightarrow Sa$ , gdzie  $a \in V_T$ ; czyli język  $L(G)$  jest regularny<sup>\*</sup>.

Założmy teraz, że twierdzenie jest prawdziwe, jeżeli  $\rho(S) < n$ . Niech  $\rho(S) = n$ , tzn. produkcje z  $P$  mające po lewej stronie  $S$  są jednej z następujących postaci:

$S \rightarrow B$  lub  $S \rightarrow CD$  lub  $S \rightarrow SE$ , gdzie  $\rho(B) < n$ ;  $\rho(D) < n$ ;  
 $\rho(E) < n$  oraz  $0 < \rho(C) < n$ .

Dla ustalenia uwagi założmy, że  $S \rightarrow B$ ,  $S \rightarrow CD$ ,  $S \rightarrow SE$  są jedynymi produkcjami z  $P$  mającymi po lewej stronie  $S$ .

Rozważmy gramatykę  $G' = \langle V'_N, V'_T, P', S' \rangle$ , w której:

$V'_N = \{S, S_1\}$ ;  $V'_T = \{B, C, D, E\}$ ;  $S' = S$ ;  $P' = \{S \rightarrow B, S \rightarrow SE,$   
 $S \rightarrow S_1 D, S_1 \rightarrow C\}$ .

Gramatyka  $G'$  jest regularna; czyli język  $L(G')$  jest regularny.

Oznaczmy:  $G_A = (V_N, V_T, P, A)$  i określmy na zbiorze  $V' = V'_N \cup V'_T$  funkcję  $f$  następująco:

$$f(A) = \begin{cases} L(G_A) & \text{dla } A \in V'_T \\ A & \text{dla } A \in V'_N \end{cases}$$

Języki  $L(G_B)$ ,  $L(G_C)$ ,  $L(G_D)$ ,  $L(G_E)$  są regularne z założenia indukcyjnego.

Ponieważ  $L(G')$  jest regularny, więc z twierdzenia 4.1 otrzymujemy, że  $f(L(G'))$  jest także regularny. Ale  $f(L(G')) = L(G)$ , czyli  $L(G)$  jest językiem regularnym.

c.b.d.o.

<sup>\*</sup>) Parę  $(A, B) \in P$  będziemy również zapisywać w postaci  $A \rightarrow B$



## Wniosek

Z algorytmu konstrukcji tablicy RCNR wynika, że dla dowolnego  $k$ , jeżeli  $RCNR(k) = 1$ , to  $\rho(A) \in I^0$ , gdzie  $k$  jest adresem  $A$  w SYMBTAB. Z tego faktu i z twierdzenia 4.2 wynika, że konstruowana przez procedurę REGULAR tablica RCNR rzeczywiście "wskazuje" symbole gramatyki generujące podjęzyki regularne.

### 4.3.3. Procedura LKONSTRUKTOR

Procedura LKONSTRUKTOR, korzystając z tablic utworzonych przez procedury ŁADOWACZ i REGULAR, tworzy nowe tablice: TABSTAN, TABTERM i AUTOMAT będące wejściem dla L-analizatora.

Tablice TABTERM i TABSTAN (o rozmiarach - odpowiednio ILT oraz ILST) służą do zapamiętywania symboli terminalnych oraz stanów automatu.

Tablica AUTOMAT jest tablicą dwuwymiarową o liczbie wierszy i kolumn - odpowiednio ILST oraz ILT.

Każdy element tablicy AUTOMAT składa się z trzech pól: (1) pola PRZEJSCIE, reprezentowanego zmienną liczbową i określającego funkcję przejścia, (2) pole WYJSCIE (również reprezentowanego zmienną liczbową), w którym jest zapisany identyfikator produkcji odpowiadającej napotkanemu w danym stanie symbolowi terminalnemu oraz (3) jednobitowego pola REDUKCJA, określającego czy należy wczytywać nowy znak ciągu wejściowego, czy też nie.

Chociaż podjęzyki znajdowane przez procedurę REGULAR są regularne (tw. 4.2), to jednak automat konstruowany przez LKONSTRUKTOR nie jest automatem skończonego stanu i program obsługujący go jest dużo bardziej skomplikowany niż program obsługi automatu skończonego. Wynika to z faktu, że gramatyki podjęzyków regularnych nie muszą być regularne. Oczywiście można skonstruować automat skończony akceptujący te podjęzyki; wymagałoby to jednak zmiany niektórych produkcji gramatyki wejściowej, a w związku z tym również związanych z nimi akcji semantycznych, co w przypadku automatycznej konstrukcji translatora jest niemożliwe.

Algorytm realizowany przez procedurę LKONSTRUKTOR jest następujący:

Krok 1: TABTERM (1) ← 'TERM1';  
TABSTAN (1) ← 'STAN1';

Krok 2: dla każdego j od 1 do N (N jest rozmiarem tablicy SYMBTAB) takiego, że RCNR(j) = 1 i TCNT(j) = 0 (tzn. j-ty symbol tablicy SYMBTAB jest quasi-regularny i nie-terminalny),

wykonaj:

Oznaczmy przez A wartość pola SYMBOL w j-tym wierszu tablicy SYMBTAB, a przez  $P_A$  zbiór wszystkich produkcji gramatyki, których lewą stroną jest A. Ponieważ A jest zmienną quasi-regularną więc produkcje należące do  $P_A$  mogą być tylko jednej z czterech postaci:

- (I)  $A \rightarrow a \quad a \in V_T$
- (II)  $A \rightarrow Ba \quad B \in V_N; \rho(B) \in I^0; a \in V_T$
- (III)  $A \rightarrow B \quad B \in V_N; \rho(B) \in I^0$
- (IV)  $A \rightarrow BC \quad B \in V_N; C \in V_N; \rho(C) \in I^0; \rho(B) \in I^0;$   
 $\rho(B) \neq 0$

Dla każdej produkcji postaci (I) wykonaj:

AUTOMAT (1,ADR2) . PRZEJSCIE ← ADR3<sup>\*</sup>

AUTOMAT (1,ADR2) . WYJSCIE ← NRPROD

AUTOMAT (1,ADR2) . REDUKCJA ← 0

gdzie ADR2 oznacza adres "a" w tablicy terminali (TABTERM),

ADR3 - adres "A" w tablicy stanów (TABSTAN),

a NRPROD - identyfikator danej produkcji  $A \rightarrow a$ .

Dla każdej produkcji postaci (II) wykonaj:

AUTOMAT (ADR1, ADR2) . PRZEJSCIE ← ADR3

AUTOMAT (ADR1, ADR2) . WYJSCIE ← NRPROD

AUTOMAT (ADR1, ADR2) . REDUKCJA ← 1

\* ) Zapis AUTOMAT (i, j) . PRZEJSCIE oznacza pole PRZEJSCIE elementu i, j dwuwymiarowej tablicy AUTOMAT



gdzie ADR1 jest adresem "B" w TABSTAN, ADR2 - adresem "a" w TABTERM, ADR3 - adresem "A" w TABSTAN, a NRPROD - identyfikatorem produkcji: A → Ba.

Dla każdej produkcji postaci (III) wykonaj:

AUTOMAT (ADR1,1) . PRZEJŚCIE ← ADR3  
AUTOMAT (ADR1,1) . WYJSCIE ← NRPROD  
AUTOMAT (ADR1,1) . REDUKCJA ← 0

gdzie ADR1, ADR3 są adresami odpowiednio "B" i "A" w TABSTAN.

Dla każdej produkcji postaci (IV) wykonaj:

AUTOMAT (ADR1,ADR2) . PRZEJSCIE ← ADR3  
AUTOMAT (ADR1,ADR2) . WYJSCIE ← NRPROD  
AUTOMAT (ADR1,ADR2) . REDUKCJA ← 1

gdzie ADR1, ADR3 są adresami odpowiednio "B" i "A" w TABSTAN, a ADR2 - adresem "C" w TABTERM.

Przykład 4.4. Podamy tablice TABTERM, TABSTAN oraz AUTOMAT skonstruowane przez procedurę LKONSTRUKTOR dla gramatyki z przykładu 4.1.

|    |        |
|----|--------|
| 1  | TERM1  |
| 2  | a      |
| 3  | b      |
| 4  | d      |
| 5  | e      |
| 6  | g      |
| 7  | i      |
| 8  | n      |
| 9  | o      |
| 10 | 1      |
| 11 | letter |
| 12 | digit  |
| 13 | :=     |

Tablica TABTERM

|   |              |
|---|--------------|
| 1 | STAN 1       |
| 2 | letter       |
| 3 | digit        |
| 4 | identifiiier |
| 5 | integer      |
| 6 | leftside     |

Tablica TABSTAN

pole PRZEJSCIE

pole WYJSCIE

pole REDUKCJA

| TERMI-<br>NAL     | (TERM1) | (a)        | (b)        | (d)        | (e)        | (g)        | (l)        | (n)        | (o)        | (i)        | (let-<br>ter) | (di-<br>git) | (:=)    |
|-------------------|---------|------------|------------|------------|------------|------------|------------|------------|------------|------------|---------------|--------------|---------|
| STAN              | 1       | 2          | 3          | 4          | 5          | 6          | 7          | 8          | 9          | 10         | 11            | 12           | 13      |
| (STAN1)           | 1       | $2_0^{14}$ | $2_0^{15}$ | $2_0^{16}$ | $2_0^{17}$ | $2_0^{18}$ | $2_0^{19}$ | $2_0^{20}$ | $3_0^{21}$ | $3_0^{22}$ |               |              |         |
| (letter)          | 2       | $4_0^6$    |            |            |            |            |            |            |            |            |               |              |         |
| (digit)           | 3       | $5_0^{11}$ |            |            |            |            |            |            |            |            |               |              |         |
| (identi-<br>fier) | 4       |            |            |            |            |            |            |            |            |            | $4_1^7$       | $4_1^7$      | $6_1^5$ |
| (integer)         | 5       |            |            |            |            |            |            |            |            |            |               | $5_1^{12}$   |         |
| (left-<br>side)   | 6       |            |            |            |            |            |            |            |            |            |               |              |         |

TABLICA AUTOMAT



#### 4.4. Podmoduł L-analizatora

Podmoduł ten wykonuje następujące zadania:

- a) czyta tekst programu źródłowego (tekst źródłowy);
- b) wydziela w tekście źródłowym atomy leksykalne;
- c) omija komentarze;
- d) tablicuje teksty;
- e) omija spacje i ciągi spacji.

Realizuje to procedura LANALIZATOR, która zawiera około 400 instrukcji w języku PL/1.

##### 4.4.1. Procedura LANALIZATOR

Procedura ta jest wywoływana przez S-analizator w celu podania nowego symbolu potrzebnego do analizy syntaktycznej. Danymi wejściowymi procedury są tablice TABSTAN, TABTERM i AUTOMAT - skonstruowane przez L-konstruktor, zmienne COMBEG, COMEND, STRINGDEL - zdefiniowane przez projektanta języka oraz tekst programu źródłowego w postaci ciągu znaków. Na podstawie tych tablic i zmiennych LANALIZATOR łączy pojedyncze znaki tekstu źródłowego w atomy leksykalne, usuwa komentarze<sup>\*\*</sup>), tablicuje teksty<sup>\*\*\*</sup>) i ciągi spacji zastępuje jedną spacją<sup>\*\*\*\*</sup>).

Procedura LANALIZATOR korzysta z trójelementowej tablicy STOS (elementy tej tablicy oznaczać będziemy przez STOS1, STOS2, STOS3). Każdy element tablicy STOS składa się z czterech

<sup>\*</sup>) Komentarz jest to ciąg znaków umieszczony między znakami COMBEG i COMEND definiowanymi przez projektanta języka.

<sup>\*\*</sup>) Tekst jest to ciąg znaków ograniczony z obu stron znakiem STRINGDEL, również definiowanym przez projektanta.

<sup>\*\*\*</sup>) W różnych językach programowania spacje są odmiennie traktowane. Na przykład w Algolu 60 spacje są ignorowane poza tekstami, w Fortranie są znaczące (do 6 kolumny), w Pascalu są separatorami, a w Snobolu 4 - operatorami.

Aktualna wersja L-analizatora traktuje spacje jako separatory, tzn. łączy ciągi spacji występujące w programie (poza tekstami) w jedną spację, która zostaje przesłana do analizatora syntaktycznego. W kolejnej wersji zostanie wprowadzona specjalna funkcja, za pomocą której projektant języka będzie mógł definiować sposoby analizowania spacji.

pól: tablicy PRODUKCJE - o zmiennym rozmiarze, dwóch pól liczbowych STAN i TERMINAL oraz pola znakowego TKN. Pola STAN i TERMINAL określają odpowiednio aktualny stan automatu i aktualnie analizowany symbol, pole TKN - ciąg znaków wejściowych już przeanalizowanych, a tablica PRODUKCJE zawiera identyfikatory produkcji użytych w dotychczasowej analizie.

Wyjściem procedury LANALIZATOR jest trójpolowa struktura ZWROT o polach TOKEN, TOKENS i WSK. Pole liczbowe TOKEN określa adres w tablicy SYMBTAB aktualnie skonstruowanego atomu leksykalnego, pole TOKENS zawiera atom leksykalny w postaci ciągu znaków, a pole WSK identyfikatory produkcji użytych do konstrukcji danego atomu (a właściwie adres tablicy zawierającej te identyfikatory).

#### 4.4.2. Algorytm realizowany przez procedurę LANALIZATOR

Ponieważ program realizujący ten algorytm zawiera dużą liczbę zagnieżdżeń instrukcji warunkowych, jego opis w języku potocznym byłby zupełnie nieczytelny. Wobec tego podamy tylko bardzo ogólną ideę algorytmu, a jego działanie przedstawimy na przykładzie analizy programu źródłowego, napisanego w języku zdefiniowanym w przykładzie 4.1.

Procedura LANALIZATOR bada, czy w gramatyce istnieje produkcja, której lewa strona należy do  $V_R$  ( $V_R = \{A \in V : \varphi(A) \in I^0\}$ ), a prawa strona jest identyczna z pierwszymi dwoma (lub tylko z pierwszym) symbolami ciągu wejściowego. Jeżeli taka produkcja istnieje, to procedura LANALIZATOR wykonuje redukcję, tzn. symbole składające się na prawą stronę tej produkcji, zastępuje w ciągu wejściowym symbolem będącym lewą stroną tej produkcji.

Czynności te są powtarzane dopóty, dopóki nie można już dokonać żadnej redukcji; następuje wówczas przekazanie sterowania do programu realizującego moduł S-analizatora.



4.4.2.1. Przykład 4.5.

Analiza programu: begin tab1 := 11.01; end

Krok 1. Wczytujemy trzy pierwsze znaki programu. Otrzymujemy następującą sytuację:

|           | STOS1                 | STOS2 | STOS3 |
|-----------|-----------------------|-------|-------|
| PRODUKCJE | 0                     | 0     | 0     |
| TERMINAL  | 3 (adres b w TABTERM) | 5     | 6     |
| STAN      | 1 (stan początkowy)   | 1     | 1     |
| TKN       | b                     | e     | g     |

Krok 2. Ponieważ

AUTOMAT (1,3). PRZEJSCIE = 2

AUTOMAT (1,3). WYJSCIE = 15

wpisujemy

STOS.STAN ← 2

STOS1.PRODUKCJE ← 15

STOS1.TERMINAL ← 11 (jest to adres nazwy stanu-2 (czyli letter) w TABTERM).

Otrzymujemy następującą sytuację:

|           | STOS1 | STOS2 | STOS3 |
|-----------|-------|-------|-------|
| PRODUKCJE | 14    | 0     | 0     |
| TERMINAL  | 11    | 5     | 6     |
| STAN      | 2     | 1     | 1     |
| TKN       | b     | e     | g     |

Krok 3. Ponieważ

STOS1.STAN ≠ 1

próbujemy analizować symbol zapisany w zmiennej STOS2.TERMINAL. Ale AUTOMAT (2,5). PRZEJSCIE = 0, próbujemy więc przejść do kolejnego stanu, niezależnie od analizowanego symbolu terminalnego. Ponieważ

AUTOMAT (2,1).PRZEJSCIE = 4

AUTOMAT(2,1).WYJSCIE = 6

wpisujemy

STOS1.STAN ← 4

STOS1.PRODUKCJE ← 6

STOS1.TERMINAL ← 0

otrzymujemy

|           | STOS1 | STOS2 | STOS3 |
|-----------|-------|-------|-------|
| PRODUKCJE | 15,6  | 0     | 0     |
| TERMINAL  | 0     | 5     | 6     |
| STAN      | 4     | 1     | 1     |
| TKN       | b     | e     | g     |

Krok 4. Ponieważ

AUTOMAT (4,5). PRZEJSCIE = 0

oraz

AUTOMAT (4,1). PRZEJSCIE = 0

próbujemy wykonać krok pracy automatu "na STOSIE 2".

Krok 5. Ponieważ

AUTOMAT (1,5). PRZEJSCIE = 2

AUTOMAT (1,5). WYJSCIE = 17

otrzymujemy:

|           | STOS1 | STOS2 | STOS3 |
|-----------|-------|-------|-------|
| PRODUKCJE | 15,6  | 17    | 0     |
| TERMINAL  | 0     | 11    | 6     |
| STAN      | 4     | 2     | 1     |
| TKN       | b     | e     | g     |

Krok 6. "Wracamy" do stanu zapamiętanego w zmiennej STOS1.STAN, analizując symbol zapisany w zmiennej STOS2.TERMINAL



AUTOMAT (4,11). PRZEJSCIE = 4

AUTOMAT (4,11). WYJSCIE = 7

otrzymujemy:

|           | STOS1  | STOS2 | STOS3 |
|-----------|--------|-------|-------|
| PRODUKCJE | 15,6,7 | 17    | 0     |
| TERMINAL  | 0      | 11    | 6     |
| STAN      | 4      | 2     | 1     |
| TKN       | b      | e     | g     |

Ponieważ

AUTOMAT (4,11). REDUKCJA = 1

"dołączamy" zawartość zmiennej STOS2 do zawartości zmiennej STOS1, zawartość zmiennej STOS3 przepisujemy do zmiennej STOS2, a do zmiennej STOS3 przesyłamy następnny znak wejściowy. Otrzymujemy:

|           | STOS1     | STOS2 | STOS3 |
|-----------|-----------|-------|-------|
| PRODUKCJE | 15,6,17,7 | 0     | 0     |
| TERMINAL  | 0         | 6     | 7     |
| STAN      | 4         | 1     | 1     |
| TKN       | be        | g     | i     |

Postępując analogicznie jak powyżej, otrzymamy w pewnej chwili następującą sytuację:

|           | STOS1                    | STOS2 | STOS3 |
|-----------|--------------------------|-------|-------|
| PRODUKCJE | 15,6,17,7,18,7,19,7,20,7 | 0     | 0     |
| TERMINAL  | 0                        | 0     | 2     |
| STAN      | 4                        | 1     | 1     |
| TKN       | begin                    | ⌊     | a     |

Ponieważ spacja jest traktowana jako ogranicznik, analiza ciągu wejściowego zostaje przerwana. Następuje sprawdzenie, czy atom zapisany w zmiennej STOS1.TKN jest słowem kluczowym. Następnie wypełniana jest struktura ZWROT następująco:

ZWROT.TOKEN ← 2 (adres "begin" w SYMBTAB)  
ZWROT.WSK ← adres pola STOS1.PRODUKCJE  
ZWROT.TOKENS ← begin

po czym następuje powrót do S-analizatora.

Przy następnym wywołaniu L-analizatora, zostanie przesłana na wyjście spacja i otrzymamy następującą sytuację:

|           | STOS1 | STOS2 | STOS3 |
|-----------|-------|-------|-------|
| PRODUKCJE | 0     | 0     | 0     |
| TERMINAL  | 2     | 3     | 10    |
| STAN      | 1     | 1     | 1     |
| TKN       | a.    | b     | 1     |

Kolejne wywołanie L-analizatora spowoduje następującą sytuację:

|           | STOS1            | STOS2       | STOS3 |
|-----------|------------------|-------------|-------|
| PRODUKCJE | 14,6,15,7,22,8,5 | 22,11,22,12 | 0     |
| TERMINAL  | 0                | 0           | 0     |
| STAN      | 6                | 5           | 1     |
| TKN       | ab1:=            | 11          | .     |

Ponieważ w tej sytuacji nie możemy wykonać kroku automatu na żadnym ze "STOSÓW", L-analizator kończy pracę, przesyłając na wyjście zawartość zmiennej STOS1, czyli

ZWROT.TOKEN = 7 (adres leftside w SYMBTAB)  
ZWROT.WSK = adres pola STOS1.PRODUKCJE  
ZWROT.TOKENS = ab1:=

Przy kolejnych wywołaniach L-analizatora otrzymamy następujące wyniki:



Wywołanie 4:

ZWROT.TOKEN = 13  
ZWROT.WSK = 22,11,22,12 (są to oczywiście numery produkcji)  
ZWROT.TOKENS = 11

Wywołanie 5:

ZWROT.TOKEN = 15  
ZWROT.WSK = 0  
ZWROT.TOKENS = .

Wywołanie 6:

ZWROT.TOKEN = 13  
ZWROT.WSK = 21,11,22,12  
ZWROT.TOKENS = 01

Wywołanie 7:

ZWROT.TOKEN = 6  
ZWROT.WSK = 0  
ZWROT.TOKENS = ;

Wywołanie 8:

ZWROT.TOKEN = 4  
ZWROT.WSK = 0  
ZWROT.TOKENS = end

Wywołanie 9:

ZWROT.TOKEN = 0  
ZWROT.WSK = 0  
ZWROT.TOKENS = 1 (umowny znacznik końca programu).

#### 4.4.3. Ograniczenia klasy gramatyk wejściowych

Wykażemy teraz, że dla pewnych klas gramatyk analiza leksykalna wykonywana według opisanego powyżej algorytmu<sup>\*)</sup> nie będzie powodowała niejednoznaczności analizy syntaktycznej ciągu wejściowego.

\*) W punkcie tym procedurę LANALIZATOR, realizującą powyższy algorytm, będziemy nazywać L-analizatorem (analizatorem leksykalnym).

Niech  $G = \langle V_N, V_T, P, S \rangle$  będzie gramatyką bezkontekstową.  
 Oznaczmy  $V = V_N \cup V_T, V_R = \{A \in V: \varphi(A) \in I^0\}$ .

Definicja 4.3

Niech  $G = \langle V_N, V_T, P, S \rangle$

Zapis  $\alpha \xrightarrow{*} \beta$  (dla  $\alpha, \beta \in V^*$ ) oznacza, że istnieje taki ciąg  $\gamma_1, \dots, \gamma_n$  ( $n \geq 1$ ), że  $\alpha = \gamma_1 \Rightarrow \gamma_2 \Rightarrow \gamma_3 \Rightarrow \dots \Rightarrow \gamma_n = \beta$

Zapis  $\alpha \Rightarrow \beta$  (dla  $\alpha, \beta \in V^*$ ) oznacza, że  $\alpha = \alpha_1 \varphi \alpha_2, \beta = \alpha_1 \psi \alpha_2$   
 i  $(\varphi, \psi) \in P$

Słowo  $\alpha \in V^*$  jest formą zdaniową gramatyki  $G \iff S \xrightarrow{*} \alpha$

Definicja 4.4

Gramatyka  $G = \langle V_N, V_T, P, S \rangle$ . Zdefiniujemy następujące relacje (tzw. relacje pierwszeństwa) pomiędzy symbolami  $S_1$  i  $S_2$  ze zbioru  $V$ :

1.  $S_1 \stackrel{\ominus}{\sim} S_2$  jeżeli do  $P$  należy produkcja postaci:  $U \rightarrow \alpha_1 S_1 S_2 \alpha_2$
2.  $S_1 \stackrel{\circ}{\sim} S_2$  jeżeli
  - (a) do  $P$  należy produkcja postaci  $U \rightarrow \alpha_1 U_1 S_2 \alpha_2$   
 i przy tym  $U_1 \xrightarrow{*} \alpha_3 S_1$
  - lub (b) do  $P$  należy produkcja postaci  $U \rightarrow \alpha_1 U_1 U_2 \alpha_2$   
 i przy tym  $U_1 \xrightarrow{*} \alpha_3 S_1$  oraz  $U_2 \xrightarrow{*} S_2 \alpha_4$
3.  $S_1 \stackrel{\otimes}{\sim} S_2$  jeżeli do  $P$  należy produkcja postaci:  
 $U \rightarrow \alpha_1 S_1 U_1 \alpha_2$  i  $U_1 \xrightarrow{*} S_2 \alpha_3$

Gramatyka  $G$  jest gramatyką z pierwszeństwem jeżeli:

1. Dla każdej pary symboli  $S_1, S_2$  z  $V$  zachodzi najwyżej jedna relacja pierwszeństwa.
2. Wszystkie produkcje ze zbioru  $P$  mają różne prawe strony.



Definicja 4.5

Niech  $\alpha = \alpha_1 \alpha_2 \alpha_3 \dots \alpha_n$  będzie formą zdaniową gramatyki G. Słowo  $\alpha_i \alpha_{i+1} \dots \alpha_{K-1} \alpha_K$  ( $1 \leq i, K \leq n$ ) jest rdzeniem formy  $\alpha$  jeżeli:

- (a)  $\alpha_j \oplus \alpha_{j+1}$  ( $i \leq j < K$ )
- (b)  $\alpha_{i-1} \otimes \alpha_i$
- (c)  $\alpha_K \oplus \alpha_{K+1}$
- (d) nie istnieje ciąg  $\alpha_p \alpha_{p+1} \dots \alpha_{r-1} \alpha_r$  ( $1 \leq p, r < i$ ) taki, że
  - $\alpha_s \oplus \alpha_{s+1}$  ( $p \leq s < r$ )
  - $\alpha_{p-1} \otimes \alpha_p$
  - $\alpha_r \oplus \alpha_{r+1}$

Ponieważ analiza syntaktyczna języka definiowanego przez gramatykę z pierwszeństwem polega na redukcji rdzeni kolejnych form zdaniowych, więc analiza leksykalna nie będzie powodowała niejednoznaczności, jeżeli również będzie redukowałą tylko rdzenie.

Twierdzenie 4.3. Jeżeli gramatyka bezkontekstowa  $G = \langle V_N, V_T, P, S \rangle$  spełnia następujące warunki:

- (\*) G jest gramatyką z pierwszeństwem,
- (\*\*) dla każdych dwóch produkcji  $p_1, p_2 \in P$  takich, że  $p_1 = (D, \gamma)$  ;  $p_2 = (E, \beta_1 \gamma \beta_2)$ ;  $\beta_1, \beta_2, \gamma \in V^*$ ; jeżeli  $D \in V_R$  to  $E \in V_R$

to analizator leksykalny zredukuje rdzeń ciągu wejściowego.

Dowód

Niech  $\alpha = BC \alpha_1$  będzie aktualnie analizowaną przez L-analizator formą zdaniową, gdzie B, C,  $\alpha_1 \in V^*$ . Nie zmniejszając ogólności, możemy założyć, że B jest częścią rdzenia formy  $\alpha$ . Ponieważ L-analizator redukuje najwyżej dwa symbole formy  $\alpha$ , otrzymamy dwa przypadki, które rozważymy kolejno:

(a) L-analizator zredukuje w formie  $\alpha$  symbol BC.

Wobec tego musi należeć do P produkcja:  $A \rightarrow BC$ , gdzie  $A, B, C \in V_R$ . Pokażemy, że rdzeniem formy  $\alpha$  jest rzeczywiście BC. Przypuśćmy, że tak nie jest. Rdzeniem formy  $\alpha$  mogą być zatem: (1)  $BC\alpha_1$  lub (2) B. W przypadku (1) musi należeć do P produkcja:  $F \rightarrow BC\alpha_1$ . Oczywiście  $F \notin V_R$  (ponieważ po prawej stronie produkcji występuje więcej niż dwa symbole), a więc otrzymujemy sprzeczność z założeniem (\*\*\*) ( $A \rightarrow BC, F \rightarrow BC\alpha_1$  należą do P, przy czym  $A \in V_R$  i  $F \notin V_R$ ). W przypadku (2) otrzymujemy natychmiast, że  $B \supseteq C$ . Ponieważ jednocześnie  $B \not\subseteq C$  (ponieważ produkcja  $A \rightarrow BC$  należy do P), więc otrzymujemy sprzeczność z założeniem (\*).

(b) L-analizator zredukuje w formie  $\alpha$  symbol B.

Wobec tego produkcja  $A \rightarrow B$  należy do P;  $A, B \in V_R$ . Pokażemy, że rdzeniem formy  $\alpha$  jest B. Przypuśćmy, że to nie zachodzi. Rdzeniem formy  $\alpha$  mogą być zatem (1)  $BC\alpha_1$  lub (2) BC. Przypadek (1) jest analogiczny jak (1) w punkcie (a). W przypadku (2) do P musi należeć produkcja  $F \rightarrow BC$ . F nie może należeć do  $V_R$ , bo L-analizator konstruuje najdłuższy możliwy atom leksykalny i wobec tego zredukowałby słowa BC na F, a nie B na A. Wobec tego  $F \notin V_R$  i otrzymujemy sprzeczność z założeniem (\*\*).

obdo.

#### Definicja 4.6

Określmy funkcję  $h: V_R \rightarrow I^0$  następująco:

1.  $h(A) = 0 \iff A \in V_T$  lub dla każdej produkcji  $p = (A, \alpha) \in P$  p jest jednej z następujących postaci:  
 $A \rightarrow a$  lub  $A \rightarrow Ba$  lub  $A \rightarrow B$ , gdzie  
 $a \in V_T, B \in V_N$ .

2.  $h(A) = n \iff n$  jest najmniejszą z możliwych liczb K takich, że dla każdej produkcji postaci  $A \rightarrow BC$  (gdzie  $B \in V_N, C \in V_N$ ),  $h(C) < K$ .

Funkcję  $h$  nazywamy głębokością regularności gramatyki G.



#### Definicja 4.7

Gramatyka  $G = \langle V_N, V_T, P, S \rangle$  jest gramatyką o głębokości regularności  $n$   $\iff n$  jest najmniejszą z takich liczb  $K$ , że dla każdego  $A \in V_R$  zachodzi  $h(A) \leq K$ .

Z przyjętej organizacji procedury LANALIZATOR wynika dodatkowe ograniczenie, aby głębokość regularności gramatyki wejściowej była nie większa niż 2 (analiza odbywa się w trójelementowej zmiennej STOS, w związku z czym są dostępne tylko trzy kolejne symbole analizowanej formy zdaniowej). Ograniczenie to przyjęto wyłącznie ze względów technicznych, nie zaś koncepcyjnych.

Praktycznie warunek ten nie ma specjalnego znaczenia, ponieważ większość gramatyk języków programowania spełnia to ograniczenie.

#### 5. ZAKOŃCZENIE

Proponowana metoda pozwala uniknąć wielu wad występujących w stosowanych dotychczas rozwiązaniach, powoduje jednakże pewne zawężenie klasy języków wejściowych. Wydaje się jednak, że przy nieznacznej modyfikacji koncepcji ogólnej L-analizatora można by wykonywać analizę leksykalną dla dowolnego języka bezkontekstowego. Zmiana ta powinna polegać na tym, że S-analizator wywołując L-analizator wskazywałby mu jednocześnie, jaki atom leksykalny powinien pojawić się na wejściu. Pozwoliłoby to uniknąć niejednoznaczności napotykanych podczas analizy leksykalnej, które w opisanej realizacji są usuwane przez ograniczenia wymienione w punkcie 4.4.3.

Taka metoda byłaby optymalna, ponieważ L-analizator byłby generowany całkowicie automatycznie, a jednocześnie można by go było stosować do analizy wszystkich języków bezkontekstowych. Wymagałoby to jednak zmiany koncepcji również S-analizatora.

Chciałbym podziękować prof. dr hab. W.M. Turskiemu za liczne konsultacje i uwagi dotyczące pracy. Dziękuję również doc. dr hab. A.W. Mazurkiewiczowi, którego cenne uwagi pomogły poprawić pierwotną wersję pracy. Ponadto dziękuję p.p. J. Koncewicz-Krzemiń, J. Borowcowi, A. Łukasiewiczowi i J. Witaszkowi za pomoc przy rozwiązywaniu wielu problemów związanych z pracą.

### Bibliografia

- [1] BLIKLE A.: Automaty i gramatyki. Wstęp do lingwistyki matematycznej, PWN, Warszawa 1971.
- [2] COCKE J., SCHWARTZ J.T.: Programmi z Languages and Their Compilers. Preliminary notes. Courant Institute of Mathematical Sciences; New York University, 1970.
- [3] FELDMAN J., GRIES D.: Translator Writing System, CACM, 1968, 2, 77-113.
- [4] GRIES D.: Compiler Construction for Digital Computers, John Wiley and Sons, Inc. 1971.
- [5] HOPCROFT J.: An NlogN Algorithm for Minimising States in a Finite Automaton, STAN-CS-74-190. 1971.
- [6] HOPCROFT J., ULLMAN J.D.: Formal Languages and Their Relation to Automata, Addison-Wesley Publishing Company, Inc. 1969.
- [7] HOPGOOD F.R.A.: Compiling Techniques, American Elsevier Publishing Company, Inc. New York 1969.
- [8] HOPGOOD F.R.A., BELL A.G.: The ATLAS ALGOL Preprocessor for non Standard Dialects. 1967 CJ, 4, 360-364.
- [9] IGLEWSKI M., KRZEMIŃ R.: Analiza leksykalna, ETO Nowości 1974, 3.
- [10] INGERMAN P.Z.: A Syntax-Oriented Translator, Academic Press, 1966.
- [11] JOHNSON W.L., PORTER J.H., ACKLEY S.I., ROSS D.T.: Automatic Generation of Efficient Lexical Processors Using Finite State Techniques, CACM, 1968, 12, 805-813.
- [12] JOHANSEN P.: Construction of Recognition Devices for Regular Languages from Their Backus Normal Form Definition, BIT, 1966, 6, 294-309.
- [13] KNUTH D.E.: Sorting and Searching: the Art of Computer Programming, Addison-Wesley Publishing Co., Inc. 1970, t. 3.
- [14] KONCEWICZ J.: Preliminary Analysis of ALGOL-60 Source Programs for ZAM Computers, Algorytmy, 1969, 10, 77-88.
- [15] LECARME O.: Un générateur d'analyseurs lexicaux. Doc. de travail: Université de Montreal, 1973, nr 40.
- [16] LYNCH W.C., PIERSON H.L.: A Finite State Transducer Model for Compiler Lexical Scanners, IFIP 68, North-Holland Publishing Co., Amsterdam 1969.



- [17] MALUSZYŃSKI J.: Podgramatyki elementarne składni Algolu 60, Algotytm, 1970, 13, 17-34.
- [18] Mc CLURE R.M.: TMS - a Syntax Directed Compiler, Proceedings ACM National Conf., 1965.
- [19] NAUR P. ed.: Revised Report on the Algorithmic Language ALGOL-60, CACM, 1962, 1, 1-17.
- [20] PLASKOV J., SCHUMAN S.: The TRANGEN System on the M460 Computer, Computer Associates Inc. Report, 1966.
- [21] REYNOLDS J.C.: COGENT Programming Manual, ANL-7022 Report, 1965.
- [22] ROSEN S.: A Compiler Building System Developed by Brooker and Morris, CACM, 1964, 7, 403-414.
- [23] ROSS D.T.: The AED Approach to Generalized Computer-Aided Design, Report ESL-R-305, MIT, 1967.
- [24] SCHEIDIG H.: ATWS: A Translator Writing System, Bericht nr 7219, TUM, 1971.
- [25] SQUIRES B.E.: Lexical Analysis by a Precedence Grammar, University Illinois Dept. Computer Science Report, 1966.
- [26] TURSKI W.M.: Struktury danych, WNT, Warszawa 1971.
- [27] WALIGÓRSKI S.: Algebraiczna teoria automatów, Algotytm, 1969, 11.
- [28] WIRTH N., WEBER H.: EULER - a Generalization of ALGOL, and its Formal Definition, CACM, 1966, 1, 13-25; CACM, 1966, 2, 89-99.
- [29] Zakład Teorii Translatorów IMM: Projekt koncepcyjny metatranslatora, Raport wewnętrzny, 1974.

DODATEK A. Program realizujący moduł analizy leksykalnej

RCMAN: PROCEDURE OPTIONS(MAIN);

NEST

RCMAN: PROCEDURE OPTIONS(MAIN);

DCL M1 FIXED BINARY;  
DCL M2 FIXED BINARY;  
DCL M5 FIXED BINARY;  
DCL N FIXED BINARY;  
DCL I FIXED BINARY;  
DCL ILST FIXED BINARY;  
DCL ILT FIXED BINARY;  
DCL A(1:80) CHAR(1);  
DCL 1 SYMTAB(1:N) CONTROLLED,  
2 SYMCL CHAR(20) VARYING,  
3 TCNT BIT(1);  
DCL PRCCAT(1:N) PCINTER CONTROLLED;  
DCL FCAR(1:N) BIT(1) CONTROLLED;  
DCL STARTSYME FIXED BINARY;  
DCL 1 PRCC BASEC(F),  
2 ILCSC FIXED BINARY,  
3 ACPALT(1:M1 REFER(ILCSC)) PCINTER;  
DCL 1 ALTER BASEC(R),  
2 CALT FIXED BINARY,  
3 ICALT FIXED BINARY,  
4 SYMCL(1:M2 REFER(CALT)) FIXED BINARY;  
DCL SYSIN RECCRE INPLT;  
DCL TARTERM(1:ILT) CHAR(20) VARYING CONTROLLED;  
DCL TARBSTAN(1:ILST) CHAR(20) VARYING CONTROLLED;  
DCL 1 AUTCMAT(1:ILST,1:ILT) CONTROLLED,  
2 PRZEJSCIE FIXED BINARY, /\* NR PRODUKCJI \*/  
3 WYJSCIE FIXED BINARY, /\* NR PRODUKCJI \*/  
4 REDUKCJA BIT(1);  
DCL 1 ZWACT, /\*CLA ANALIZATORA SYNT \*/  
2 TCKEN FIXED BINARY,  
3 WSK PCINTER,  
4 TCKENS CHAR(20) VARYING;  
DCL C PCINTER;  
DCL 1 STCS(1:21),  
2 WSKAZNIK PCINTER, /\*PCINTER DO PRODUKCJI\*/  
3 TERMINAL FIXED BINARY,  
4 STAN FIXED BINARY,  
5 TKA CHAR(20) VARYING;  
DCL 1 PRCCUKCJE BASEC(F),  
2 LRPRC FIXED BINARY,  
3 NRPRCC(1:M5 REFER(LRPRC)) FIXED BINARY;

M1=1;

M2=1;

M5=30;

I = 0;

N=160;

ILT=45;

ILST=6;

ALLOCATE PRCCAT;

ALLOCATE SYMTAB;

ALLOCATE AUTCMAT;

ALLOCATE TARBSTAN;

ALLOCATE TARTERM;

ALLOCATE RCMP;

CC J=1 TO N;



RCFAN: PROCEDURE OPTIONS(MAIN);

NCST

1  
1  
1  
1  
1

```
FRCCAT(J)=NULL;  
SYMETAB.SYMBCL(J)='';  
SYMETAB.TCNT(J)='1'B;  
RCNR(J)='0'B;  
ENC;  
SYMETAB.SYMBCL(1)='SPACJA';  
OPEN FILE(SYSIN);  
ON ENDFILE(SYSIN) GO TO KONIEC;  
CALL LACCHACZ;  
PUT DATA (SYMETAB);  
RCNP=SYMETAB.TCNT;  
PUT DATA (RCNR);  
CALL REGULAR;  
CALL LKCNSTRUKTOR;  
PUT DATA (SYMETAB);  
PUT DATA (STARTSYMB);  
PUT DATA (RCNR);  
PUT DATA (ILT);  
PUT DATA (TABTERM);  
PUT DATA (ILST);  
PUT DATA (TAESTAN);  
PUT DATA (AUTCMAT);  
CALL LANALINIT;  
PUT DATA (ZWRCT.TCKENS,ZWRGT.TCKEN);  
FU1: CALL LANALIZATOR;  
PUT DATA (ZWRCT.TCKENS,ZWRGT.TCKEN);  
GO TO FU1;  
LACCHACZ: PROCEDURE;  
DCL NRFRCD FIXED BINARY;  
DCL 1 WYNIK,  
      2 NAZWA CHAR(20) VARYING,  
      2 TNT PIT(1);  
DCL ACRES FIXED BINARY;  
DCL FRCC(1:20) PCINTER;  
DCL LFS FIXED BINARY;  
DCL RFS(1:20) FIXED BINARY;  
PUT LIST ('LACCHACZ');  
READ FILE(SYSIN) INTO(A);  
PUT DATA (A);  
L1: I=I+1;  
CALL NUMER;  
CALL TCKEN;  
IF WYNIK.TNT='1'E THEN GO TO BLAD1;  
CALL SFACJA;  
IF A(I)='=' THEN GO TO PLAC2;  
I=I+1;  
CALL FASF;  
LFS=ACRES;  
L2: CALL TCKEN;  
CALL FASF;  
RFS(M2)=ACRES;  
RCNP(ACRES)='1'B;  
CALL SPACJA;  
IF A(I)=',' THEN  
CC;
```

RCMAN: PROCEDURE OPTICKS(PAIN);

NEST

```

1          PUT LIST('PRZECINFK');
1          M2=M2+1;
1          I=I+1;
1          CC TC L2;
1          ENC;
          IF A(I)=';' THEN
1          CC;
1          CALL SLASH;
1          CALL SRECNIK;
1          CC TC L1;
1          ENC;
          IF A(I)='|' THEN
1          CC;
1          CALL SLASH;
1          I=I+1;
1          CALL NUMER;
1          M1=M1+1;
1          CC TC L3;
1          ENC;
          IF A(I)='.' THEN
1          CC;
1          PUT LIST('KRCPKA');
1          CALL SLASH;
1          CALL SRECNIK;
1          CALL AKSJCMAT;
1          RETURN;
1          ENC;
AKSJCMAT: PROCEDURE;
          CC J=1 TC A;
1          IF (FCAR(J)='O'B & TCNT(J)='C'B) THEN
1          CC;
2          STARTSYMB=J;
2          RETURN;
2          ENC;
1          ENC;
          CC TC ELACC;
          END AKSJCMAT;
          PUT LIST ('SRECNIK');
SRECNIK: PROCEDURE;
          ALLOCATE FRCC;
          PRCCAT(LFS)=F;
          CC J=1 TC M1;
1          FRCC.ACFALT(J)=FRCC1(J);
1          ENC;
          M1=1;
          END SRECNIK;
          SLASH: PROCEDURE;
          PUT LIST ('SLASH');
          ALLOCATE ALTER;
          FRCC1(M1)=F;
          ALTER.ICALT=ARPRCC;
          CC J=1 TC M2;
1          ALTER.SYMBOL(J)=RHS(J);
1          ENC;
          M2=1;

```



PROGRAM: PROCEDURE OPTICS(PAIR);

NEST

```

      ENC SLASH;
FASH: PROCEDURE;
      DCL F1 FIXED BINARY;
      DCL F2 CHAR(26) VARYING;
      DCL F3 CHAR(26) VARYING;
      PUT LIST ('FASH');
      F1=1;
      F2=KYNIK.MAWA;
      F4: F3=SYMTAB.SYMBOL(F1);
      IF F2=F3 THEN GO TO F5;
      IF F3=' ' THEN
        CC;
        SYMTAB(F1)=KYNIK;
        GO TO F5;
      ENC;
      IF F1<N THEN
        CC;
        F1=F1+1;
        GO TO F4;
      ENC;
      ELSE GO TO ELAC2;
      H5: ACRES=F1;
      ENC FASH;
SFACJA: PROCEDURE;
      PUT LIST ('SFACJA');
      S1: IF A(I)=' ' THEN
        CC;
        IF I<EC THEN
          CC;
          I=I+1;
          GO TO S1;
        ENC;
        ELSE
          CC;
          READ FILE(SYSIN) INTO(A);
          PUT DATA(A);
          I=1;
          GO TO S1;
        ENC;
      ENC SFACJA;
NUMER: PROCEDURE;
      DCL CYFFA CHAR(10) INIT('1234567890');
      DCL N1 FIXED BINARY;
      DCL N2 FIXED BINARY;
      PUT LIST ('NUMER');
      N1=0;
      CALL SFACJA;
      N3: N2=INDEX(CYFFA,A(I));
      IF N2=0 THEN
        CC;
        IF N2=10 THEN N2=0;
        N1=N1+10+N2;
        I=I+1;
        GO TO N3;
      ENC;

```

PCMAN: PROCEDURE OPTICNS (MAIN);

NEST

```

ELSE
CC;
CALL SFACJA;
IF A(I)~=':' THEN GO TO BLAC4;
ELSE
CC;
NRPRCC=N1;
I=I+1;
PUT DATA(NRPRCC);
END NUPER;
TOKEN: PROCEDURE;
DCL LITERA CHAR(26) INIT('ABCDEFGHIJKLMNQRSTUWXYZ');
DCL IDENT CHAR(27) INIT('ABCDEFGHIJKLMNQRSTUWXYZ0123456789-');
DCL T1,
2 T2 CHAR(20) VARYING,
2 T2 BIT(1);
PUT LIST ('TOKEN');
T1,T2='';
CALL SFACJA;
IF A(I)~='' THEN
T5: BEGIN;
I=I+1;
IF A(I)~=''' THEN
CC;
T1,T2=T1,T2||A(I);
GO TO T5;
END;
ELSE
CC;
T1,T2=' 'E;
WYNIK=T1;
I=I+1;
RETURN;
END;
END;
ELSE
CC;
IF VERIFY (A(I),LITERA)~=0 THEN GO TO BLAD5;
T1,T2=T1,T2||A(I);
T6: I=I+1;
IF VERIFY (A(I),IDENT)=0 THEN
CC;
T1,T2=T1,T2||A(I);
CC TO T6;
END;
ELSE
CC;
T1,T2=' 'E;
WYNIK=T1;
T7: END TOKEN;
LKNIEC: END LACCWAC2;
REGULAR: PROCEDURE;
DCL IFRACC FIXED BINARY;
DCL FRCSK BIT(1);
DCL RECALC BIT(1);

```

1  
1  
1  
1  
2  
2  
2  
2

1  
1  
1

1  
1  
1  
1  
1

1  
1  
1  
1  
2  
2  
2  
2  
1  
1  
2  
2  
2



RCNR: PROCEDURE OPTICS(MAIN);

NEST

```

REGWSK='0'B;
R7: CC J=1 TC N;
  IF RCNR(J)='0'B THEN
    CC;
    F=FRCCAT(J);
    ILFRCC=FRCC.ILCSC;
    CALL REG1;
    RCNR(J)=REGALT;
    REGWSK=(REGALT|REGWSK);
  END;
  IF RCNR(J)='1'B THEN
    CC;
    PUT DATA(RCNR);
    REGWSK='0'B;
    GC TC R7;
  END;
  REGWSK='1'B;
R8: CC J=1 TC N;
  IF(RCNR(J)='1'B & TCNT(J)='0'B) THEN
    CC;
    F=FRCCAT(J);
    ILFRCC=FRCC.ILCSC;
    CALL REG2;
    RCNR(J)=REGALT;
    REGWSK=(REGALT & REGWSK);
  END;
  IF REGWSK='0'B THEN
    CC;
    PUT DATA(RCNR);
    REGWSK='1'B;
    CC TC R8;
  END;
REG1: PROCEDURE;
  CC K=1 TC ILFRCC;
  R=FRCC.ACFALT(K);
  CALL REG2;
  IF REGALT='0'B THEN RETURN;
  END;
REG2: PROCEDURE;
  DCL PCB1 FIXED BINARY;
  DCL PCB2 FIXED BINARY;
  REGALT='0'B;
  IF ALTER.CALT>2 THEN RETURN;
  IF ALTER.CALT=1 THEN
    CC;
    PCB1=ALTER.SYMPCL(1);
    REGALT=RCNR(PCB1);
    RETURN;
  END;
  ELSE
    PCB2=ALTER.SYMPCL(2);
    PCB1=ALTER.SYMPCL(1);

```

PCNAT: PROCEDURE CFTYCS(PAIN);

AFST

```

REGALT=(-TCNT(RCB1)*PCNR(PCB2));
END REC2;
REC3: PROCEDURE;
  CC K=1 TO ILFRCD;
  R=PRCC.ACFALT(K);
  CALL REC4;
  IF REGALT='0'B THEN RETURN;
END;
END REC3;
REC4: PROCEDURE;
  DCL RCB1 FIXED BINARY;
  RCB1=ALTER.SYMBOL(1);
  REGALT=PCNR(RCB1);
  END REC4;
  END REGULAR;
LKNSTRUKTUR: PROCEDURE;
  DCL LFS FIXED BINARY;
  DCL RHS1 FIXED BINARY;
  DCL RHS2 FIXED BINARY;
  DCL ILFRCD FIXED BINARY;
  DCL LK2 CHAR(20) VARYING;
  DCL NRFRCD FIXED BINARY;
  DCL LK3 FIXED BINARY;
  PUT LIST ('LKNSTRUKTUR');
  CC J=1 TO ILT;
  TABTERM(J)='';
  END;
  CC J=1 TO ILST;
  TABSTAN(J)='';
  CC K=1 TO ILT;
  AUTMAT(J,K).PRZEJSCIE=C;
  END;
  END;
  TABTERM(1)='TERM1';
  TABSTAN(1)='STAN1';
  CC J=1 TO N;
  IF (PCNR(J)='1'B & TCNT(J)='0'B) THEN
  CC;
  P=PRCCAT(J);
  ILFRCD=PRCC.ILFRCD;
  LK2=SYMTAB.SYMBOL(J);
  CALL SEARCH(TABSTAN,LFS,ILST);
  CC K=1 TO ILFRCD;
  P=PRCC.ACFALT(K);
  NRFRCD=ALTER.ICALT;
  LK3=ALTER.SYMBOL(1);
  LK2=SYMTAB.SYMBOL(LK3);
  IF ALTER.CALT=2 THEN
  CC;
  CALL SEARCH(TABSTAN,RHS1,ILST);
  LK3=ALTER.SYMBOL(2);
  LK2=SYMTAB.SYMBOL(LK3);
  CALL SEARCH(TABTERM,RHS2,ILT);
  AUTMAT(RFS1,RFS2).PRZEJSCIE=IFHS;
  AUTMAT(RFS1,RFS2).WYJSCIE=NRFRCD;

```



RCHAN: PROCEDURE OPTICNS(PAIN);

NEST

```

4     AUTOMAT(RFS1,RFS2).REDUKCJA='1'B;
4     ENC;
3     ELSE /* TZN.CAL1=1*/
3     CC;
4     IF SYMETAB.TCN1(LK2)='1'B THEN
4     CC;
5     CALL SEARCH(TABTERM,RHS1,ILT);
5     AUTOMAT(1,RFS1).PRZEJSCIE=LHS;
5     AUTOMAT(1,RFS1).WYJSCIE=NRPRCD;
5     AUTOMAT(1,RFS1).REDUKCJA='0'B;
4     ENC;
4     ELSE
5     CC;
5     CALL SEARCH(TABSTAN,RHS1,ILST);
5     AUTOMAT(RFS1,1).PRZEJSCIE=LHS;
5     AUTOMAT(RFS1,1).WYJSCIE=NRPRCD;
5     AUTOMAT(RFS1,1).REDUKCJA='0'B;
4     ENC;
3     END;
2     ENC;
1     ENC;

```

```

SEARCH: PROCEDURE(TAB,ACR,WYM);
CCL TAB(*) CHAR(20) VARYING;
CCL ACR FIXED BINARY;
CCL WYM FIXED BINARY;
CCL RCE1 CHAR(20) VARYING;
PUT LIST('SEARCH');
CC 1=1 TO WYM;
RCE1=TAB(1);
IF RCE1=LK2 THEN
CC;
ACR=1;
FETUFA;
ENC;
IF RCE1='' THEN
CC;
ACR=1;
TAB(1)=LK2;
RETURN;

```

1  
1  
1  
2  
2  
2  
1  
1  
2  
2  
2  
2  
2  
1

```

ENC;
ENC;
ENC SEARCH;
ENC LKCNSTRUKTOR;
LANALINIT: PROCEDURE;
CCL B FIXED BINARY;
CCL REDUKUJ BIT(1);
CCL SZT CHAR(30) VARYING;
CCL LA4 FIXED BINARY;
CCL ZWFCTNICA BIT(1);
CCL LA6 FIXED BINARY;
/* PC DCLACZENIU CC LKCNSTR DCCAC DEKLARACJE A,I*/
/* SYSIN, SYMETAB, ILT, ILST, TABTERM ITD */
CH ENDFILE(SYSIN) GO TO KONIECLA;
READ FILE(SYSIN) INTO(A);

```

RCMAN: PROCEDURE OPTICS(MAIN);

NEST

```

1      I=1;
1      CC J=1 TO 3;
1      ALLCCATE PRCEUKCJE;
1      WSKAZNIK(J)=F;
1      LRFRCC=1;
1      NRFRCC(LRFRCC)=0;
1      STCS.STAN(J)=1;
1      TKN(J)=A(I);
1      SZT=A(I);
1      CALL SZUKTERM;
1      TERMINAL(J)=LA4;
1      I=I+1;
1      ENC;
1      ALLCCATE PRCEUKCJE;
1      WSK=F;
LANALIZATOR: ENTRY;
      LA01: BEGIN;
1      ECL LA02 BIT(1);
1      LA02='0'B;
1      LA0: CC WFILE (TKN(1)=' ');
1      LA02='1'B;
1      SP1: STCS(1)=STCS(2);
1      STCS(2)=STCS(3);
1      F=WSKAZNIK(3);
1      LRFRCC=1;
1      NRFRCC(LRFRCC)=0;
1      STAN(2)=1;
1      IF I>60 THEN
2          CC;
2          REAC FILE(SYSIN) INTC(A);
2          I=1;
2          ENC;
1      SZT=A(I);
1      CALL SZUKTERM;
1      TERMINAL(2)=LA4;
1      TKN(2)=A(I);
1      I=I+1;
1      ENC;
1      IF LA02='1'B THEN CC;
1      ZWRCT.TCKEN=1;
1      ZWRCT.WSK=NULL;
1      ZWRCT.TCKENS=' ';
1      RETURN;
1      ENC;
1      ENC;
LA1: B=1;
      CALL ANAL;
      IF ZWRCT.NICA='1'B THEN GO TO LA1;
LA2: B=2;
      CALL ANAL;
      IF ZWRCT.NICA='1'B THEN GO TO LA1;
      B=3;
      CALL ANAL;
      IF ZWRCT.NICA='1'B THEN GO TO LA2;
      IF (STCS.STAN(1) & STCS.STAN(2)=1) THEN
```



RCFAN: PROCEDURE OPTIONS(MAIN);

NEST

```
1 CC;
1 TCKENS=TKN(1)||TKN(2);
1 CALL KEYWCRC;
1 IF TCKEN =0 THEN
1 CC;
2 F=WSK;
2 LFRCC=1;
2 NFRCC(LFRCC)=C;
2 STCS(1)=STCS(3);
2 CC J=2 TC 3;
3 F=WSKAZNIK(J);
3 LFRCC=1;
3 NFRCC(LFRCC)=0;
3 STAN(J)=1;
3 IF I >EO THEN
3 CC;
4 REAC FILE (SYSIN) INTO(A);
4 I=1;
4 ENC;
3 TKN=A(I);
3 SZT=A(I);
3 CALL SZUKTERP;
3 TERMINAL(J)=LA4;
3 J=J+1;
3 ENC;
2 RETURN;
2 ENC;
1 ENC;
IF -(STCS.STAN(2)=1 & STCS.STAN(3)=1) THEN
1 CC;
1 TCKENS=TKN(2)||TKN(3);
1 CALL KEYWCRC;
1 IF TCKEN =0 THEN
1 CC;
2 TKN(2)=TCKENS;
2 SZT=TCKENS;
2 CALL SZUKTERP;
2 TERMINAL(2)=LA4;
2 F=WSKAZNIK(2);
2 LFRCC=1;
2 NFRCC(LFRCC)=C;
2 IF I >EO THEN
2 CC;
3 REAC FILE(SYSIN) INTO(A);
3 I=1;
3 ENC;
2 TKN(3)=A(I);
2 SZT=A(I);
2 CALL SZUKTERP;
2 TERMINAL(3)=LA4;
2 J=J+1;
2 F=WSKAZNIK(3);
2 LFRCC=1;
2 NFRCC(LFRCC)=C;
2 CC TC LA1;
```

RCFAN: PRCECUPE OPTICS(MAIN);

NEST

```
2          END;
1          END;
          TCKENS=TKN(1);
          CALL KEYWCRC;
          IF ((TCKEN=C) && (SYMTAB.TCNT(TOKEN)='C'B)) THEN
1          CC;
1          F=WSK;
1          LRFRCC=1;
1          NRFRCC(LRFRCC)=C;
1          END;
          ELSE
1          CC;
1          TCKENS=TABSTAN(STAN(1));
1          CALL KEYWCRC;
1          TCKENS=TKN(1);
1          F=STCS.WSKAZNIK(1);
1          WSK->PRCEUKCJE =PRCEUKCJE;
1          LRFRCC=LRFRCC-1;
1          END;
          STCS(1)=STCS(2);
          STCS(2)=STCS(3);
          F=WSKAZNIK(3);
          LRFRCC=1;
          NRFRCC(LRFRCC)=0;
          STAN(2)=1;
          IF I>80 THEN
1          CC;
1          REAC FILE (SYSIN) INTO(A);
1          I=1;
1          END;
          SZT=A(1);
          CALL SZUKTERM;
          TERMINAL(2)=LA4;
          TKN(2)=A(1);
          I=I+1;
          RETURN;
ANAL: PRCECUPE;
          IF STCS.STAN(B)=1 THEN
1          CC;
1          IF STCS.TERMINAL(B)=C THEN /*DWLZNAKCY SYMPL*/
2          CC;
2          IF ((B=3) || (TERMINAL(B+1)=C) || (STOS.STAN(B+1)=1)) THEN
3          CC;
3          ZHACTRICA='C'B;
3          RETURN;
3          END;
2          TCKENS=TKN(B) || TKN(B+1);
2          CALL KEYWCRC;
2          IF TCKEN=C THEN
3          CC;
3          ZHACTRICA='C'B;
3          RETURN;
2          END;
2          TKN(B)=TCKENS;
2          SZT=TCKENS;
```



ROZKAZ: PROCEDURA OPTICAS (MAIN):

END

```

2      CALL SZUKTERM;
2      TERMINAL(P)=LA4;
2          IF P=1 THEN STOS(2)=STOS(3);
2          IF I>80 THEN
3              CC;
3              READ FILE(SYSIN) INTO(A);
3              I=1;
2              END;
2      SZT=A(1);
2      CALL SZUKTERM;
2      TERMINAL(2)=LA4;
2          TKN(2)=A(1);
2          STAN(2)=1;
2          F=WSKAZNIK(2);
2          LRFRCC=1;
2          NRFRCC(LRFRCC)=C;
2          ZWRCTNICA='1'B;
2          RETURN;
2          END;
1      IF AUTMAT(1,TERMINAL(B)).PRZEJSCIE=0 THEN
1          CC;
2          ZWRCTNICA='0'B;
2          RETURN;
2          END;
1          ELSE
1          CC;
2      STOS.STAN(P)=AUTMAT(1,TERMINAL(B)).PRZEJSCIE;
2          F=WSKAZNIK(E);
2          CC L=2 TO LRFRCC;
3          NRFRCC(L)=NRFRCC(L-1);
3          END;
2          NRFRCC(1)=AUTMAT(1,TERMINAL(B)).WYJSCIE;
2          LRFRCC=LRFRCC+1;
2          TCKENS=TARSTAN(STOS.STAN(B));
2          SZT=TCKENS;
2          CALL SZUKTERM;
2          TERMINAL(E)=LA4;
2          ZWRCTNICA='1'B;
2          RETURN;
2          END;
1      END;
1          ELSE /*TKN STAN --=1 */
1          CC;
1          IF P=3 THEN
2              CC;
2              ZWRCTNICA='0'B;
2              RETURN;
2              END;
1          IF AUTMAT(STOS.STAN(B),TERMINAL(B+1)).PRZEJSCIE=0 THEN
1          IF AUTMAT(STOS.STAN(B),1).PRZEJSCIE=0 THEN
1              CC;
2              ZWRCTNICA='0'B;
2              RETURN;
2              END;
1          ELSE

```

PCFAN: PROCEDURE OPTICS(PAIN);

AE5T

```

1.      CC;
2      RECUKUJ=AUTCMAT(STAN(P),1).RECLKCJA;
2      STCS.STAN(E)=AUTCMAT(STCS.STAN(P),1).PRZEJSCIE;
2      TOKENS=TABSTAN(STCS.STAN(P));
2      SZT=TOKENS;
2      CALL SZUKTERM;
2      TERMINAL(E)=LA4;
2      F=STCS.WSKAZNIK(P);
2      DO L=2 TO LRPRCC;
3      NRPRCC(L)=NRPRCC(L-1);
3      ENC;
2      NRPRCC(1)=AUTCMAT(STCS.STAN(P),1).WYJSCIE;
2      IF RECUKUJ='C'B THEN
2      DO TO LA6;
2      ELSE
2      DO TO LA9;
1      ENC;
1      ELSE
2      CC; /* ANALIZA STOS(1),STOS(2)*/
2      RECUKUJ=AUTCMAT(STAN(P),TERMINAL(P+1)).RECLKCJA;
2      STCS.STAN(E)=AUTCMAT(STCS.STAN(P),TERMINAL(P+1)).PRZEJSCIE;
2      TOKENS=TABSTAN(STCS.STAN(P));
2      SZT=TOKENS;
2      CALL SZUKTERM;
2      TERMINAL(E)=LA4;
2      F=STCS.WSKAZNIK(P);
2      DO L=2 TO LRPRCC;
3      NRPRCC(L)=NRPRCC(L-1);
3      ENC;
2      NRPRCC(1)=AUTCMAT(STCS.STAN(P),TERMINAL(P+1)).WYJSCIE;
2      IF RECUKUJ='C'E THEN
2      DO TO LA6;
2      ELSE DO TO LA9;
1      ENC;
1      LA9: TKN(E)=TKN(P)||TKN(P+1);
1      IF E=1 THEN
1      CC;
2      F=WSKAZNIK(1);
2      C=WSKAZNIK(2);
2      LA6=C->LRPRCC-1;
2      DO K=1 TO LA6;
3      NRPRCC(LRPRCC)=C->NRPRCC(K);
3      LRPRCC=LRPRCC+1;
3      ENC;
2      STCS(2)=STCS(3);
2      C=WSKAZNIK(3);
1      ENC;
1      ELSE
1      CC;
2      F=WSKAZNIK(2);
2      C=WSKAZNIK(3);
2      LA6=C->LRPRCC-1;
2      DO K=1 TO LA6;
3      NRPRCC(LRPRCC)=C->NRPRCC(K);
3      LRPRCC=LRPRCC+1;

```



PCMAN: PROCEDURE CPTICNS(MAIN);

NEST

```

2 .          ENC;
2          ENC;
1          C->LRFRCC=1;
1          ST#N(2)=1;
1          IF J>EO THEN
1            CC;
2            REAC FILE(SYSIN) INTC(A);
2            I=1;
2            ENC;
1            F=WSKAZNIK(2);
1            TKN(2)=A(1);
1            SZT=A(1);
1            CALL SZUKTERM;
1            TERMINAL(2)=LA4;
1            I=I+1;
1            LRFRCC=1;
1            NRFRCC(LRFRCC)=C;
1            CC TC LA10;
1            LRFRCC=LRFRCC+1;
1            LAB:  ZWRCTAICA='1'B;
1            RETURN;
1            ENC;
1          ENC ANAL;
1        KCNIECLA: TCKEN=0;
1          F=WSK;
1          LRFRCC=1;
1          NRFRCC(LRFRCC)=0;
1          RETURN;
1        SZUKTERM: PROCEDURE;
1          CC P=1 TC ILT;
1          IF TAETERM(M)=SZT THEN CC;
1          LA4=M;
1          RETURN;
1          ENC;
1          LA4=0;
1          ENC SZUKTERM;
1        KEYWGPC: PROCEDURE;
1          CC P=1 TC N;
1          IF SYMETAB,SYMBOL(M)=TCKENS THEN
1            CC;
1            TCKEN=M;
1            RETURN;
1            ENC;
1          ENC;
1          TCKEN=C;
1          ENC KEYWGPC;
1          ENC LANALINIT;
1        KCNIEC: PUT LIST ('CRUK SYMETAB I PRODUKCJI');
1          PUT DATA (STARTSYMB);
1          CC TC FCM;
1        ELAD1: PUT LIST ('TERMINAL PO LEWEJ STRONIE PRODUKCJI') SKIP;
1          CC TC FCM;
1        BLAD2: PUT LIST ('BRAK ZNAKU = W PRODUKCJI') SKIP;
1          CC TC FCM;

```

RCFAN: PROCEDURE OPTICNS(MAIN):

NEST

```
ELAD3: PUT LIST ('ZA KRÓTKA SYMBŁÓW') SKIP;  
      GC TC FCM;  
ELAD4: PUT LIST ('PRAK CŁUKROPKA PC NLMERZE PRODUKCJI') SKIP;  
      GC TC FCM;  
ELAD5: PUT LIST ('ZŁA PRAWA STRONA PRODUKCJI') SKIP;  
      GC TC FCM;  
ELAD6: PUT LIST ('EFAN SYMBŁU STARTOWEGO') SKIP;  
      GC TC FCM;  
      /* ZAMIĄST FCM DĄC KONIEC FC WYRZUCENIU WYDRUKOW*/  
ROM:END RCFAN;
```



DODATEK B. Składnia języka Euler

FORM: PROCEDURE OPTIONS(MAIN);

NEŚI

FORM: PROCEDURE OPTIONS(MAIN);

```
1:      PROCFAH = '3', ELCK, '4';
2:      VARCECI = 'NFW', IDENTIFIER;
3:      FCFCECI = 'FORMAL', IDENTIFIER;
4:      LABLECI = 'LFEEL', IDENTIFIER;
5:      VAFI = IDENTIFIER;
6:      VAFI, ' ', EXPR, '1';
7:      VAFI, ' ':
8:      VAF = VAFI;
9:      LCCV/L = 'TRCE';
10:     'FALSE';
11:     IDENTIFIER = LETTER;
12:     IDENTIFIER, LETTER;
13:     IDENTIFIER, DIGIT;
14:     LETTER = 'A';
15:     'E';
16:     'C';
17:     'E';
18:     'E';
19:     'F';
20:     'C';
21:     'F';
22:     'J';
23:     'J';
24:     'K';
25:     'L';
26:     'M';
27:     'N';
28:     'C';
29:     'P';
30:     'R';
31:     'S';
32:     'T';
33:     'C';
34:     'L';
35:     'V';
36:     'W';
37:     'X';
38:     'Y';
39:     'Z';
40:     DIGIT = '0';
41:     '1';
42:     '2';
43:     '3';
44:     '4';
45:     '5';
46:     '6';
47:     '7';
48:     '8';
49:     '9';
50:     INTEGER1 = DIGIT;
51:     INTEGER1, DIGIT;
52:     INTEGER2 = INTEGER1;
53:     PRZL! = INTEGER2, ' ', INTEGER2;
54:     INTEGER2;
```

RCMAN: PROCEDURE OPTICS(MAIN);

NEST

```

55:      NUMBER = REAL1
56:      REAL1, 'EXP', INTEGER2
58:      'EXP', INTEGER2;
59:      REFERENCE = '%', VAR;
60:      LISTHEAD = LISTHEAD, EXPR, ',';
61:      '(';
62:      LIST1 = LISTHEAD, EXPR, '),';
63:      LISTHEAD, '),';
64:      PROCHEAD = PROCHEAD, PROCDECL, ',';
65:      'CFEN';
66:      PROCDEF = PROCHEAD, EXPR, 'CLOSE';
67:      PRIMARY = VAR|
68:      VAR, LIST1|
69:      LCCVAL|
70:      NUMBER|
71:      REFERENCE|
72:      LIST1|
73:      'TAIL', PRIMARY|
74:      PROCDEF|
75:      'CMECA';
76:      '_, EXPR, ' ';
77:      'IN';
78:      'ISA', VAR|
79:      'ISA', VAR|
80:      'ISR', VAR|
81:      'ISL', VAR|
82:      'ISL1', VAR|
83:      'ISY', VAR|
84:      'ISP', VAR|
85:      'ISV', VAR|
86:      'AES', PRIMARY|
87:      'LENGTH', VAR|
88:      'INTEGER', PRIMARY|
89:      'REAL', PRIMARY|
90:      'LCCJCAL', PRIMARY|
91:      'LIST', PRIMARY;
92:      FACTCF1 = PRIMARY|
93:      FACTCF1, '**', PRIMARY;
94:      FACTCF = FACTCF1;
95:      TERM1 = FACTCF|
96:      TERM1, '**', FACTCF|
97:      TERM1, '/', FACTCF|
98:      TERM1, '^', FACTCF|
99:      TERM1, 'MCC', FACTCF;
100:      TERM = TERM1;
101:      SUM1 = TERM|
102:      '+', TERM|
103:      '-', TERM|
104:      SUM1, '+', TERM|
105:      SUM1, '-', TERM;
106:      SLP = SUM1;
107:      CHCISE1 = SLP|
108:      CHCISE1, 'MIN', SLP|
109:      CHCISE1, 'MAX', SLP;
110:      CHCISE = CHCISE1;

```



RCMPL: PROCELFEE CTICES (MAIN):

NEST

```

111:      RELATION      =      CFCISE |
112:      CFCISE, '=' , CFCISE |
113:      CFCISE, '<=' , CFCISE |
114:      CFCISE, '<' , CFCISE |
115:      CFCISE, '<=' , CFCISE |
116:      CFCISE, '>' , CFCISE |
117:      CFCISE, '>' , CFCISE ;
118:      NEGATION      =      RELATION |
119:      '~' , RELATION ;
120:      CCNJFFZL      =      NEGATION, '&' ;
121:      CCNJI         =      CCNJHEAD, CCNJI |
122:      NEGATION ;
123:      CCNJ          =      CCNJI ;
124:      CJSJFFZC      =      CCNJ, '!' ;
125:      CJSJ          =      CJSJHEAD, CJSJ |
126:      CCNJ ;
127:      CATENA        =      CATENA, '||' , PRIMARY |
128:      CJSJ ;
129:      TRLEFFAT      =      EXPR, 'ELSE' ;
130:      IFCLAUSE      =      'IF', EXPR, 'THEN' ;
131:      EXFF          =      ELCK |
132:      IFCLAUSE, TRLEPART, EXPR |
133:      VAR, '=' , EXPR |
134:      'CCTC', PRIMARY |
135:      'CLT', EXPR |
136:      CATENA ;
137:      EXFF          =      EXPR ;
138:      STAT          =      LAECEF, STAT |
139:      EXPR ;
140:      STAT          =      STAT ;
141:      LAECEF        =      IDENTIFIER, ':' ;
142:      ELCKHEAD      =      'BEGIN' |
143:      ELCKHEAD, VARDECL, ':' ;
144:      ELCKHEAD, LAEDECL, ':' ;
145:      ELCKECCY      =      ELCKHEAD |
146:      ELCKECCY, STAT, ':' ;
147:      ELCK          =      ELCKECCY, STAT, 'END' .

```

47

















## КОНСТРУКЦИЯ ЛЕКСИЧЕСКИХ АНАЛИЗАТОРОВ В МЕТАТРАНСЛЯТОРАХ

### Резюме

В работе представлен новый метод конструкции лексических анализаторов в метатрансляторах. Предлагаемый метод на много отличается от использованных до сих пор методов поскольку лексические атомы обнаруживаются только лишь на основе описания синтаксиса языка.

## CONSTRUCTION OF LEXICAL ANALYSERS IN COMPILER-COMPILERS

### Summary

In this paper we present a new method for the construction of lexical analysers in compiler-compilers. The method presented is essentially different from those applied up to now since the lexical atoms are found automatically solely on the basis of the syntax description of the language.

BIBLIOTEKA GŁÓWNA  
Politechniki Śląskiej

P 2229/76