



AKADEMIA GÓRNICZO-HUTNICZA
IM. STANISŁAWA STASZICA W KRAKOWIE

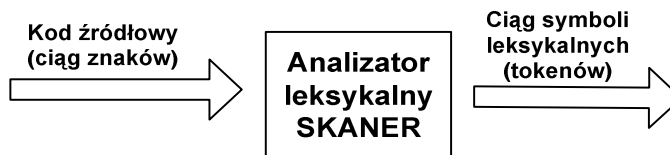
Analiza leksykalna – 1

Języki formalne i automaty

Dr inż. Janusz Majewski
Katedra Informatyki



Zadanie analizy leksykalnej



Przykład:

We: $COST := (PRICE + TAX) * 0.98$

Wy: $id_1 := (id_2 + id_3) * num_4$

Tablica symboli:

Adres	Nazwa/wartość	Charakter	Dodatkowa informacja
1	COST	zmienna
2	PRICE	zmienna
3	TAX	zmienna
4	0.98	stała



Zadanie analizy leksykalnej

Przykład:

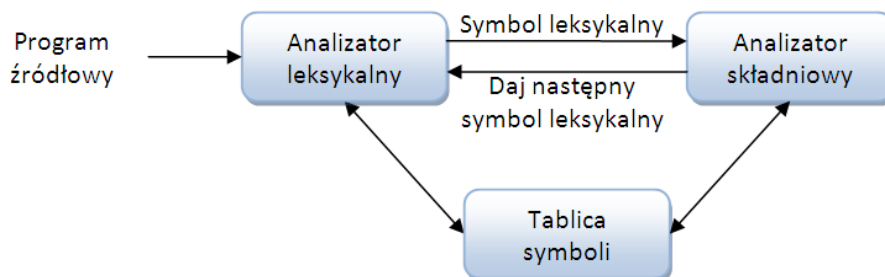
We: $COST := (PRICE + TAX) * 0.98$

Wy: $id_1 := (id_2 + id_3) * num_4$

Wyjście skanera:	<=	Wejście skanera:
(<u>id</u> , 1)		COST
(<u>equ</u> ,)		:=
(<u>left-par</u> ,)		(
(<u>id</u> , 2)		PRICE
(<u>plus</u> ,)		+
(<u>id</u> , 3)		TAX
(<u>right-par</u> ,))
(<u>mult</u> ,)		*
(<u>num</u> , 4)		0.98



Współpraca z parserem





Zadania analizatora leksykalnego

Zadania analizatora leksykalnego:

- wyodrębnianie symboli leksykalnych (tokenów)
- ignorowanie komentarzy
- ignorowanie białych znaków (spacji, tabulacji, znaków nowej linii...)
- korelowanie błędów zgłaszanych przez kompilator z numerami linii
- tworzenie kopii zbioru wejściowego (źródłowego) łącznie z komunikatami o błędach
- czasami realizacja funkcji preprocessingu, rozwijanie makrodefinicji

Rozdzielenie etapu analizy na dwie odrębne funkcje: analizę leksykalną i analizę syntaktyczną sprawia, że jedna i druga mogą być wykonywane przy użyciu bardziej efektywnych algorytmów, gdyż algorytmy te istotnie się różnią, wykorzystując inne pryncypia formalne i realizacyjne.



Podstawowe pojęcia: token, leksem, wzorzec

Przykład:

```
const pi2 = 6.2832;
```

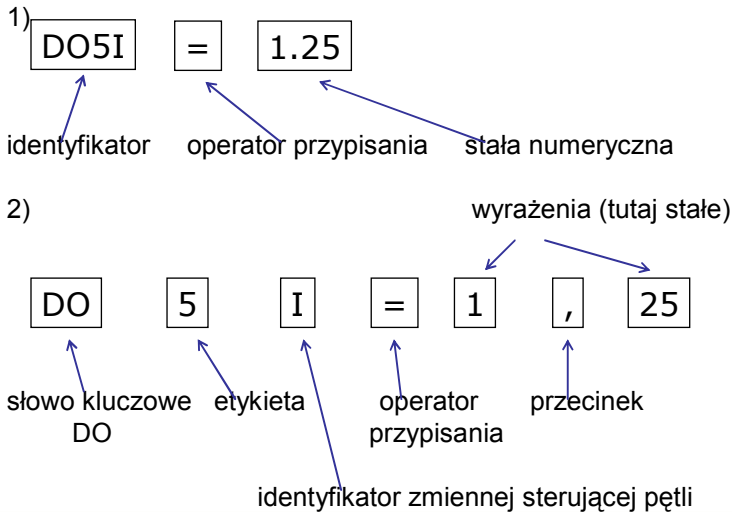
token	leksem	wzorzec (pattern)
<u>const</u> (słowo kluczowe)	const	const
<u>id</u>	pi2	<i>litera (litera cyfra)*</i>
<u>relop</u>	=	< > <= >= = <>
<u>num</u>	6.2832	<i>cyfra⁺ (. cyfra⁺)? ((E e) (+ -)? cyfra⁺)?</i>

źródło: const pi2 = 6.2832
leksemy: const pi2 = 6.2832
tokeny: const id relop num



Trudności w budowaniu analizatora leksykalnego

Przykład c.d.:



Trudności w budowaniu analizatora leksykalnego

Przykład c.d.:

DO 5 I = 1 {
 :
 n
 :
}

Po przeczytaniu znaków DO nie można dokonać uzgodnienia tokenu "Słowo kluczowe DO" dopóki nie zbada się prawego kontekstu i nie znajdzie się przecinka (wtedy rzeczywiście uzgadnia się "DO") lub kropki bądź znaku nowej linii (wtedy mamy instrukcje podstawienia).



Definicje regularne

Do opisu wzorców dla skanera stosujemy definicje regularne:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

.....

$$d_n \rightarrow r_n$$

gdzie:

d_i - unikalna nazwa

r_i - wyrażenie regularne nad symbolami alfabetu

$$\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$$



Przykład definicji regularnych (1)

Stałe bez znaku w Pascal'u:

$$\underline{cyfra} \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$$

$$\underline{cyfry} \rightarrow \underline{cyfra} \underline{cyfra}^*$$

$$\underline{część-ułamkowa} \rightarrow . \underline{cyfry} | \varepsilon$$

$$\underline{wykładnik} \rightarrow (E | e) (+ | - | \varepsilon) \underline{cyfry} | \varepsilon$$

$$\underline{num} \rightarrow \underline{cyfry} \underline{część-ułamkowa} \underline{wykładnik}$$



Rozszerzenie zbioru operatorów

Dla ułatwienia wprowadza się nowe operatory w wyrażeniach regularnych, np.:

w^+ - oznacza jedno lub więcej wystąpień wzorca w

$$w^+ = w w^*$$

$w?$ - oznacza zero lub jedno wystąpienie wzorca w

$$w? = w \mid \varepsilon$$



Przykład definicji regularnych (2)

Stałe bez znaku w Pascal'u zapisane po rozszerzeniu zbioru operatorów w wyrażeniach regularnych:

cyfra $\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

cyfry \rightarrow cyfra⁺

część-ułamkowa $\rightarrow (\cdot \text{cyfry}) ?$

wykładnik $\rightarrow ((E \mid e) (+ \mid -) ? \text{cyfry}) ?$

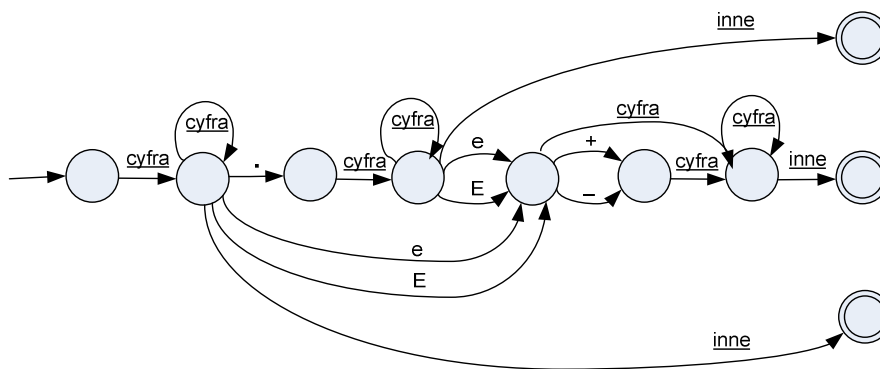
num \rightarrow cyfry część-ułamkowa wykładnik



Diagramy przejść (1)

Liczba bez znaku w Pascalu

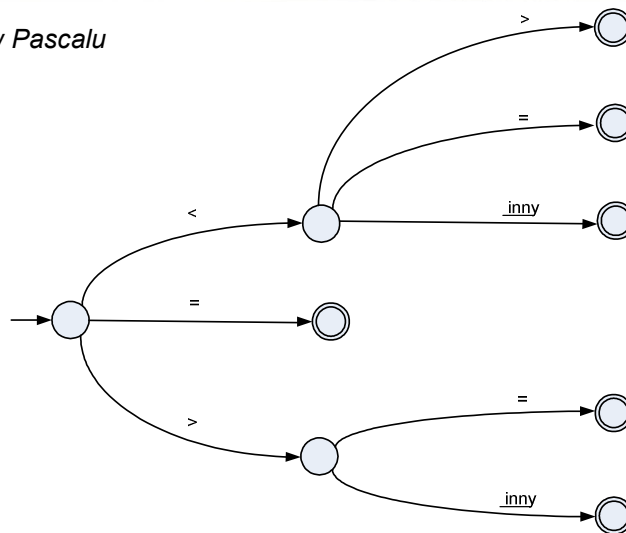
$\text{cyfra}^+ (\cdot \text{cyfra}^+)? (e|E) (+|-)? \text{cyfra}^+)?$



Diagramy przejść (2)

Operatory relacyjne w Pascalu

$< | <= | <> | = | >= | >$

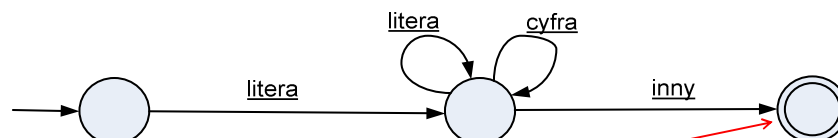




Diagramy przejść (3)

Identyfikatory

litera (litera | cyfra)*



„oddaj” ostatni przeczytany symbol na wejście;
sprawdź, czy leksem to słowo kluczowe;
jeśli tak – zwróć odpowiednie słowo kluczowe;
jeśli nie – sprawdź, czy identyfikator jest
już w tablicy symboli;
jeśli jest – zwróć adres jego pozycji;
jeśli nie ma – utwórz nową pozycję
i wpisz identyfikator do tablicy symboli.