



AKADEMIA GÓRNICZO-HUTNICZA
IM. STANISŁAWA STASZICA W KRAKOWIE

Analiza leksykalna – 1

Teoria kompilacji

Dr inż. Janusz Majewski
Katedra Informatyki

Zadanie analizy leksykalnej



Przykład:

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

Wy: $\underline{id}_1 := (\underline{id}_2 + \underline{id}_3) * \underline{num}_4$

Tablica symboli:

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

Zadanie analizy leksykalnej

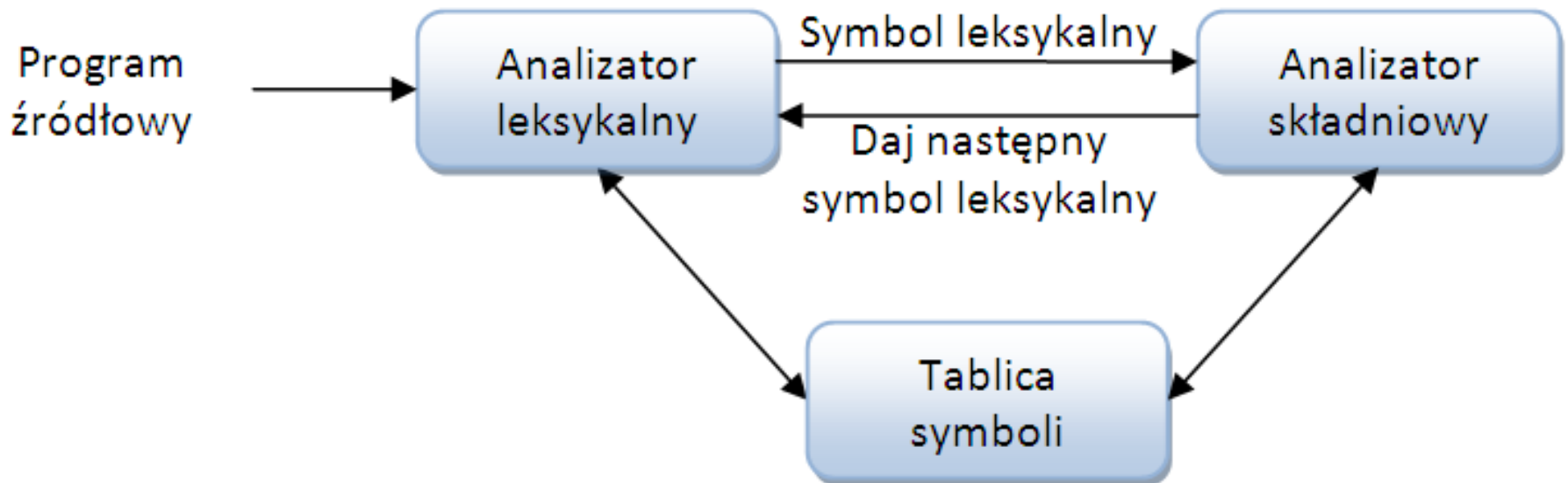
Przykład:

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

Wy: $\underline{id}_1 \underline{:=} (\underline{id}_2 + \underline{id}_3) * \underline{num}_4$

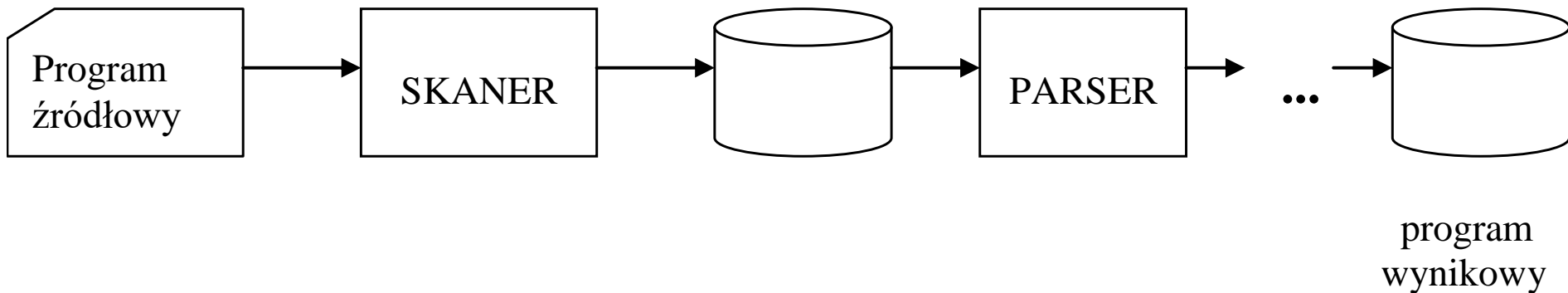
Wyjście skanera:	<=	Wejście skanera:
(<u>id</u> , 1)		COST
(<u>equ</u> ,)		:=
(<u>left-par</u> ,)		(
(<u>id</u> , 2)		PRICE
(<u>plus</u> ,)		+
(<u>id</u> , 3)		TAX
(<u>right-par</u> ,))
(<u>mult</u> ,)		*
(<u>num</u> , 4)		0.98

Współpraca z parserem



Architektura kompilatora (1)

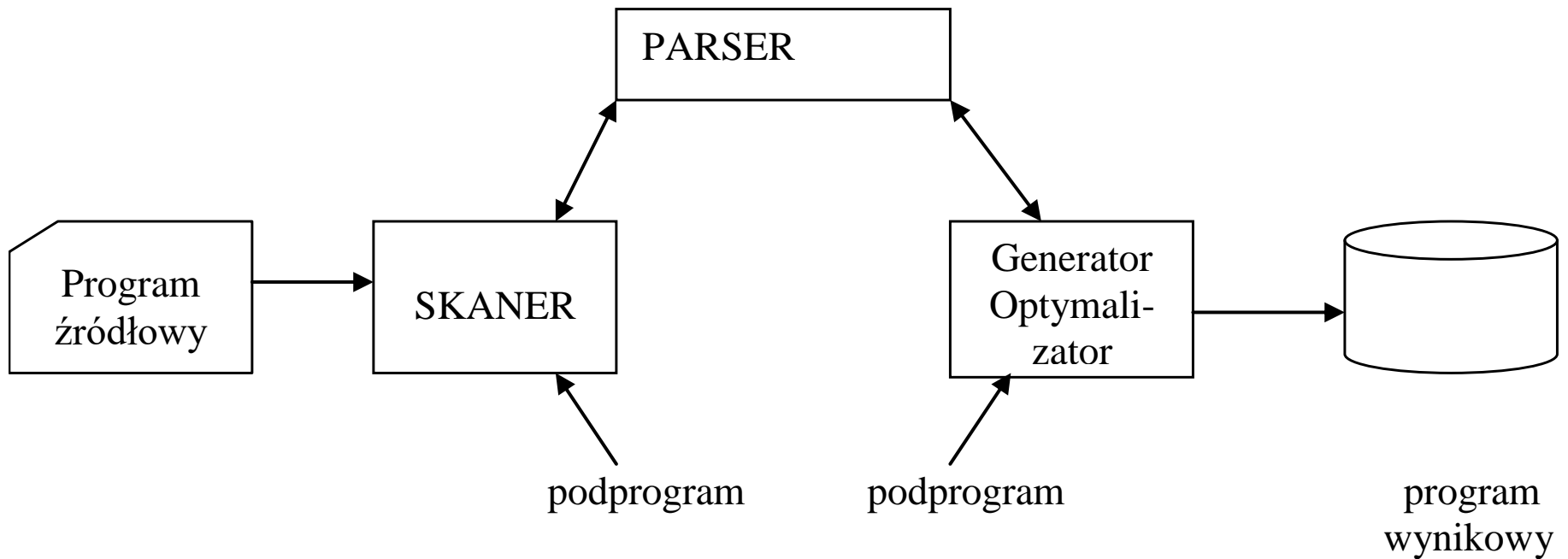
a) Realizacja szeregową



- translator wieloprzebiegowy
- translator nakładkowy
- mniejsza zajętość PaO
- dłuższy czas kompilacji

Architektura kompilatora (2)

b) Realizacja równoległa



- translator jednoprzebiegowy
- większa zajętość PaO
- krótszy czas kompilacji



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>ws</u>	□	<i>biały_znak</i> ⁺
<u>id</u>	pi2	<i>litera</i> (<i>litera</i> <i>cyfra</i>)*
<u>relop</u>	=	< > <= >= = <>
<u>num</u>	6.2832	<i>cyfra</i> ⁺ (. <i>cyfra</i> ⁺)? ((E e) (+ -)? <i>cyfra</i> ⁺)?

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



Podział gramatyki G opisującej dany język programowania

Gramatyka G języka źródłowego (najczęściej bezkontekstowa):

$$G = \langle V, \Sigma, P, S \rangle$$

gdzie na ogół:

$$\Sigma = \{A, B, \dots, Z, a, b, \dots, z, 0, 1, \dots, 9, +, -, *, /, \dots, (,), \dots\}$$

Przebudowujemy gramatykę G na rodzinę gramatyk:

$G_S = \langle V_S, \Sigma_S, P_S, S_S \rangle$ - gramatyka syntaktyczna

$$\Sigma_S = \{ \Sigma_{S_i} \mid \Sigma_{S_i} \text{ jest_tokenem} \} \quad \text{- tokeny są terminalami w gramatyce}$$

syntaktycznej

$G_i = \langle V_i, \Sigma_i, P_i, S_i \rangle$ - gramatyki leksykalne, jedna dla każdego tokenu

$$S_i = \Sigma_{S_i}$$

$$\Sigma_{S_i} \in V_i, \Sigma_i \subset \Sigma$$

i -ty token Σ_{S_i} jest symbolem początkowym (nieterminalem) i -tej gramatyki leksykalnej.

Podział gramatyki G opisującej dany język programowania

Podział gramatyki języka źródłowego $G = \langle V, \Sigma, P, S \rangle$

- Gramatyka syntaktyczna: $G_S = \langle V_S, \Sigma_S, P_S, S_S \rangle$
- Rodzina gramatyk leksykalnych $G_i = \langle V_i, \Sigma_i, P_i, S_i \rangle$ po jednej dla każdego tokenu

$$V = V_S \cup \bigcup_{i=1}^n V_i$$

$$\Sigma = \bigcup_{i=1}^n \Sigma_i$$

$$P = P_S \cup \bigcup_{i=1}^n P_i$$

$$S = S_S$$

Gramatyka języka źródłowego G jest więc efektem „podstawienia” gramatyk leksykalnych G_i do gramatyki syntaktycznej G_S .

Podział gramatyki G opisującej dany język programowania

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

Wzorce są opisami sposobu wyodrębniania tokenów.

Wzorce pełnią rolę produkcji w gramatykach leksykalnych.

Tokeny są symbolami nieterminalnymi w gramatykach leksykalnych.

Tokeny są terminalami w gramatyce syntaktycznej.



Trudności w budowaniu analizatora leksykalnego

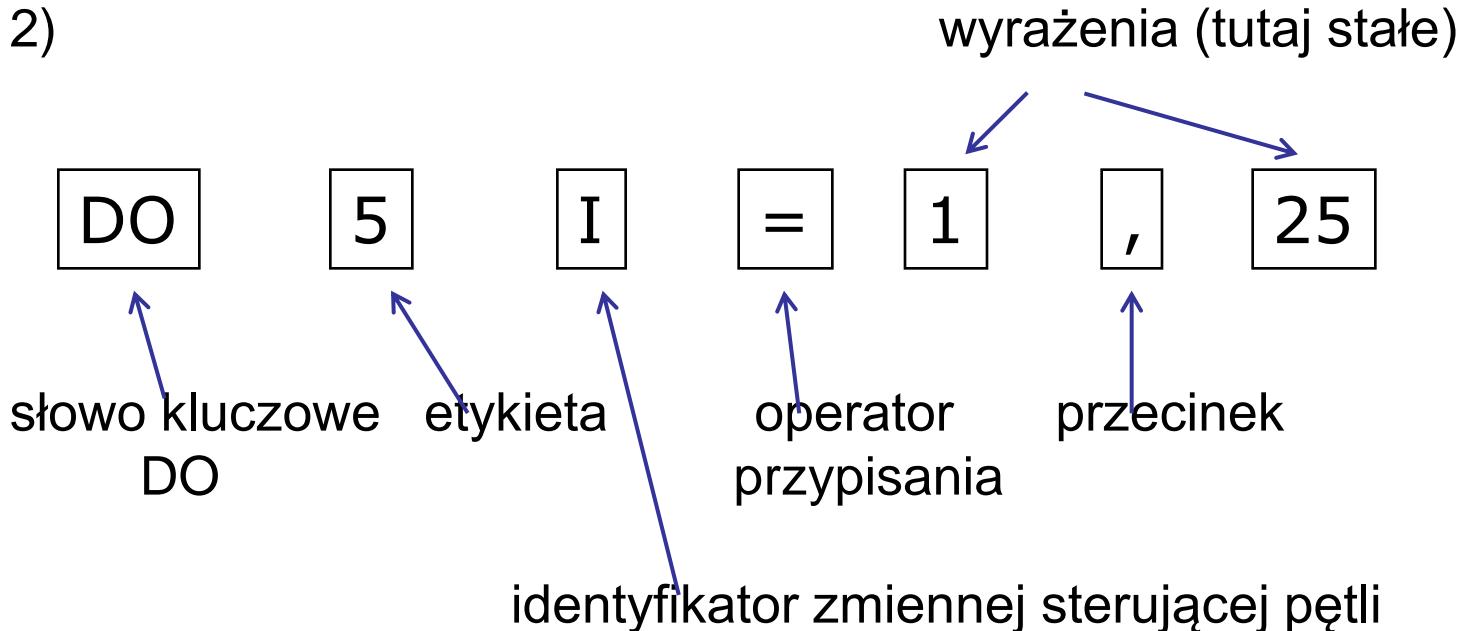
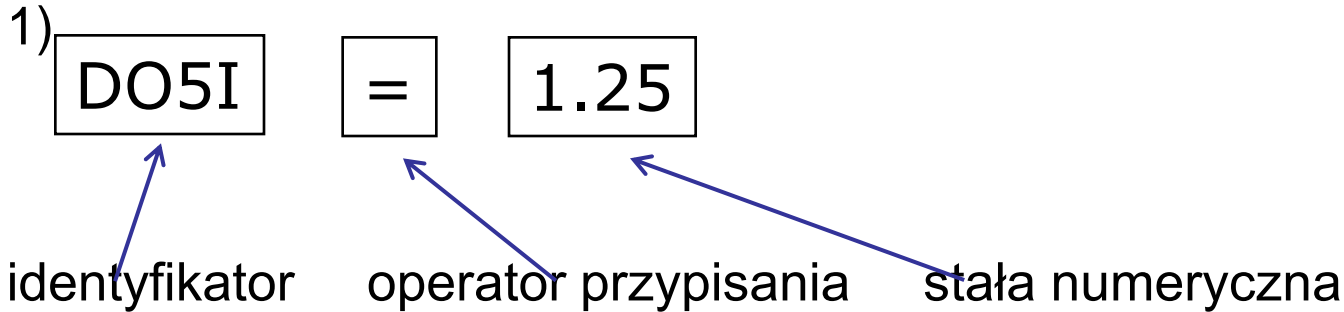
Przykład:

W niektórych językach programowania słowa kluczowe nie są zastrzeżone (FORTRAN, PL/I). W języku PL/I poprawny jest zapis:

```
IF THEN THEN THEN = ELSE ; ELSE  
ELSE = THEN;
```


Trudności w budowaniu analizatora leksykalnego


Przykład c.d.:



Trudności w budowaniu analizatora leksykalnego

Przykład c.d.:

DO 5 l = 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:

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

cyfry \rightarrow cyfra cyfra *

część-ułamkowa \rightarrow . cyfry $\mid \varepsilon$

wykładnik $\rightarrow (E \mid e) (+ \mid - \mid \varepsilon)$ cyfry $\mid \varepsilon$

num \rightarrow cyfry część-ułamkowa 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

Rozpoznawanie tokenów

Przykładowa gramatyka syntaktyczna:

$stmt \rightarrow \underline{if} \ expr \ \underline{then} \ stmt \ / \ \underline{if} \ expr \ \underline{then} \ stmt \ \underline{else} \ stmt \ / \ \varepsilon$

$expr \rightarrow term \ \underline{relop} \ term \ / \ term$

$term \rightarrow \underline{id} \ / \ \underline{num}$

Definicje regularne:

$\underline{delim} \rightarrow \square \ / \ \backslash t \ / \ \backslash n$

$\underline{ws} \rightarrow \underline{delim}^+$

$\underline{letter} \rightarrow a \ / \ b \ / \ \dots \ / \ z \ / \ A \ / \ B \ / \ \dots \ / \ Z$

$\underline{digit} \rightarrow 0 \ / \ 1 \ / \ \dots \ / \ 9$

$\underline{if} \rightarrow if$

$\underline{then} \rightarrow then$

$\underline{else} \rightarrow else$

$\underline{relop} \rightarrow < \ / \ <= \ / \ <> \ / \ > \ / \ >= \ / \ =$

$\underline{id} \rightarrow \underline{letter} \ (\ \underline{letter} \ / \ \underline{digit} \)^*$

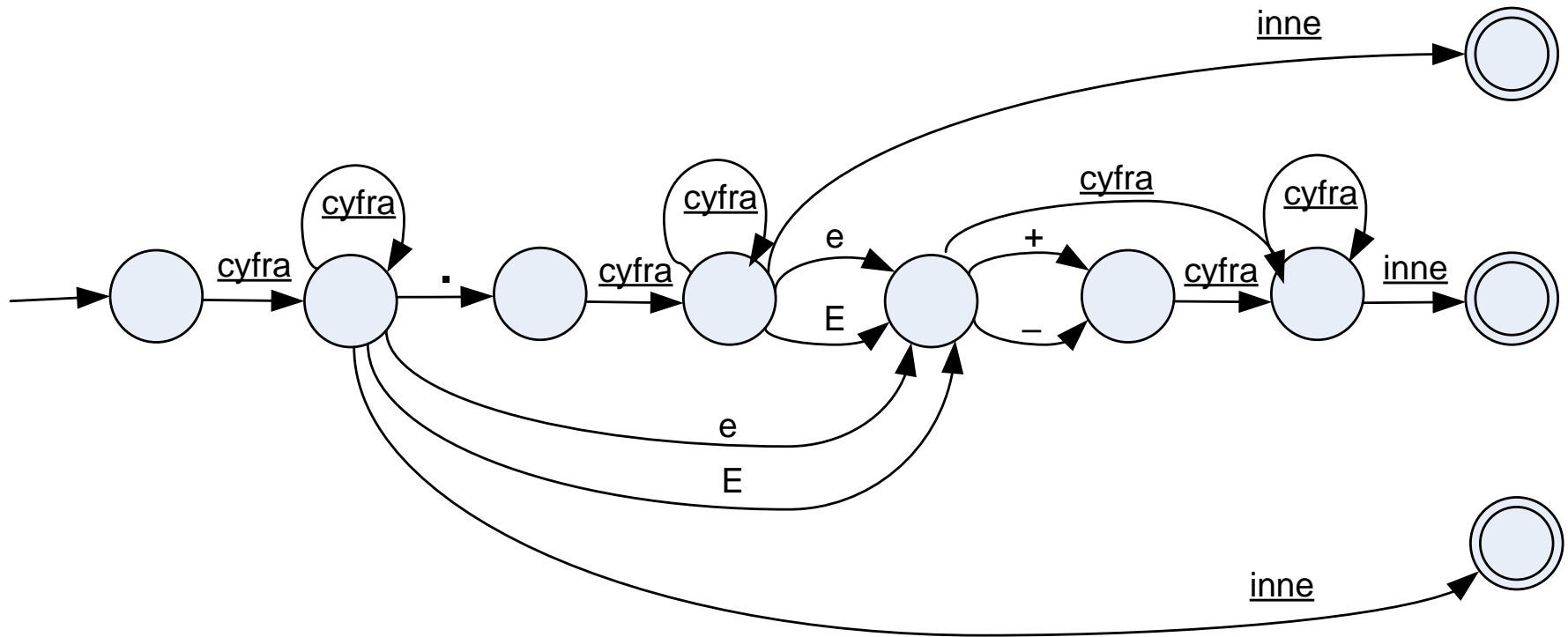
$\underline{num} \rightarrow \underline{digit}^+ \ (\ . \ \underline{digit}^+ \)? \ (\ E \ (\ + \ / \ - \)? \ \underline{digit}^+ \)?$

Skaner ma rozpoznawać...

	token	atrybut tokenu
<u>ws</u>	-	-
<u>if</u>	<u>if</u>	-
<u>then</u>	<u>then</u>	-
<u>else</u>	<u>else</u>	-
<u>id</u>	<u>id</u>	wskaźnik do tablicy symboli
<u>num</u>	<u>num</u>	wskaźnik do tablicy symboli
<	<u>relop</u>	LT
<=	<u>relop</u>	LE
<>	<u>relop</u>	NE
=	<u>relop</u>	EQ
>	<u>relop</u>	GT
>=	<u>relop</u>	GE

Diagramy przejść (1)

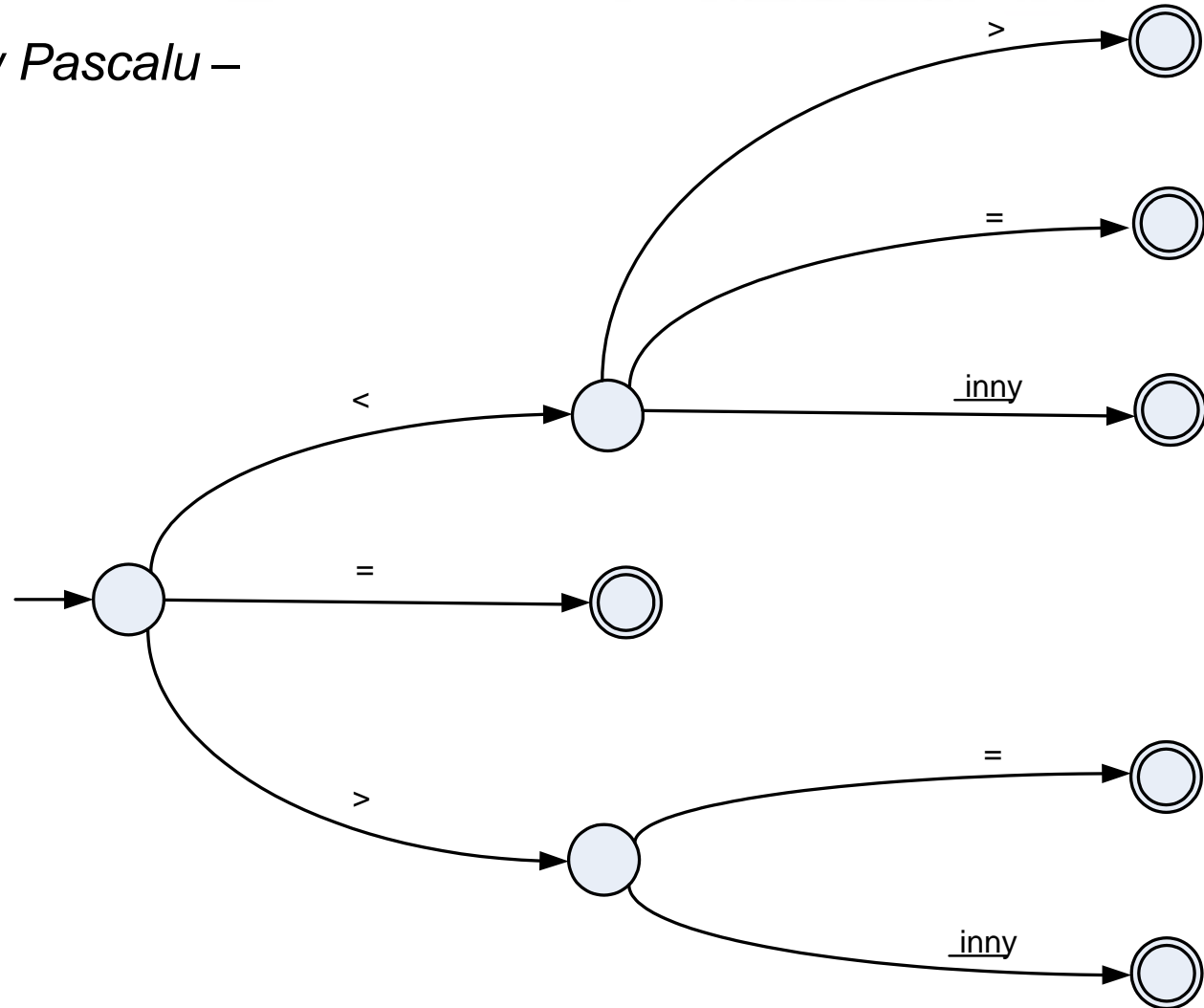
*Liczba bez znaku w Pascalu – token **num***
cyfra⁺ (cyfra⁺)? ((e|E) (+|-)? cyfra⁺)?



Diagramy przejść (2)

Operatory relacyjne w Pascalu –
token relop

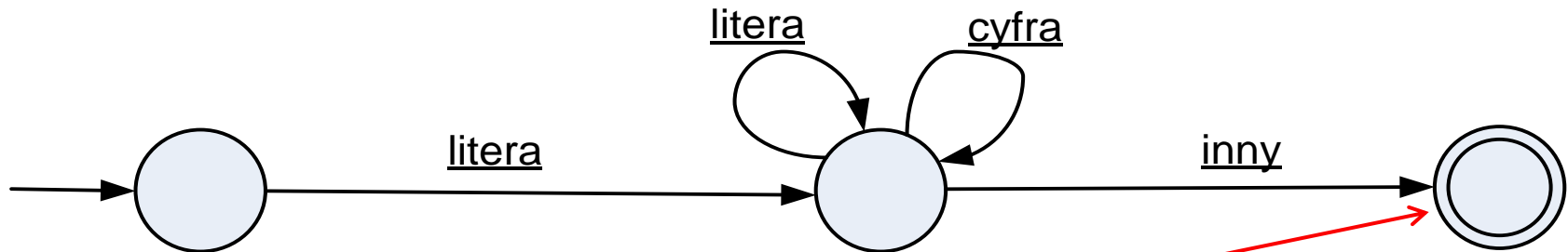
< | <= | <> | = | >= | >



Diagramy przejść (3)

Identyfikator – token id

litera (litera | cyfra)*

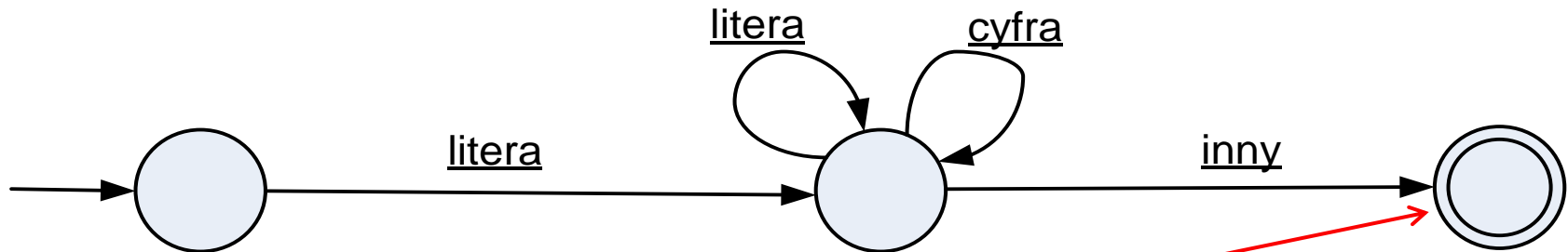


„oddaj” ostatni przeczytany symbol na wejście;
sprawdź, czy leksem to słowo kluczowe;
jeśli tak – zwróć odpowiednie słowo kluczowe;
jeśli nie – sprawdź, czy identyfikator jest
już w tablicy symboli;
jeśli jest – zwróć adres jego pozycji;
jeśli nie ma – utwórz nową pozycję, wpisz
identyfikator do tablicy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

Diagramy przejść (4)

Identyfikatory

litera (litera | cyfra)*



„oddaj” ostatni przeczytany symbol na wejście;
sprawdź, czy leksem to słowo kluczowe...

Nie zawsze stosuje się metodę umieszczania słów kluczowych w tablicy i sprawdzania każdego identyfikatora, czy nie jest on słowem kluczowym. Można budować dla słów kluczowych odrębne wzorce i konstruować automaty rozpoznające. Wówczas jednak liczba stanów analizatora gwałtownie rośnie.

