

# Teoretyczne podstawy informatyki

## Wykład 11: Opis wzorców: gramatyki

Opis wzorców polegający na wykorzystaniu modelu definicji rekurencyjnych, nazywamy gramatyką bezkontekstową (ang. context-free grammar). Jednym z ważnych zastosowań gramatyk są specyfikacje języków programowania. Gramatyki stanowią zwięzłą notację opisu ich składni. Istnieje możliwość mechanicznego przekonwertowania gramatyki typowego języka programowania na analizator składniowy (ang. parser), który stanowi jeden z kluczowych elementów kompilatora takiego języka. Analizator składniowy pozwala na zidentyfikowanie struktury programu źródłowego, często w postaci drzewa wyrażień dla każdej instrukcji programu.

# Gramatyki bezkontekstowe

Wyrażenia arytmetyczne można w naturalny sposób zdefiniować rekurencyjnie.

Weźmy pod uwagę wyrażenia arytmetyczne zawierające:

- (a) Cztery operatory dwuargumentowe  $+$ ,  $-$ ,  $*$ ,  $/$
- (b) Nawiasy służące do grupowania podwyrażeń
- (c) Operandy które są liczbami

Tradycyjna definicja takich wyrażeń stanowi indukcje:

## Podstawa:

Liczba jest wyrażeniem.

## Indukcja:

Jeżeli  $E$  oznacza dowolne wyrażenie, to wyrażeniami są także wszystkie z poniższych elementów:

- (1)  $(E)$ . Oznacza to że wyrażenie można umieścić w nawiasach w wyniku czego otrzymuje się nowe wyrażenie.
- (2)  $E + E$ . Oznacza to że dwa wyrażenie połączone znakiem plus stanowią wyrażenie.
- (3)  $E - E$ .
- (4)  $E * E$ .
- (5)  $E / E$ .

Powyższa indukcja indukuje język czyli zbiór ciągów znaków. Podstawa określa że każda liczba należy do tego języka. Reguła (1) określa, że jeżeli  $s$  jest ciągiem znaków należącym do omawianego języka, to także ciąg znaków objęty nawiasami należy do tego języka. Taki ciąg  $s$  jest poprzedzony znakiem nawiasu otwierającego, zaś po nim występuje znak nawiasu zamykającego. Reguły (2),(3),(4),(5) określają, że jeżeli  $s$  i  $t$  są dwoma ciągami znaków należącymi do języka, to należą do niego również ciągi znaków  $s+t$ ,  $s-t$ ,  $s*t$ ,  $s/t$ .

Gramatyki pozwalają na zapisywanie takich reguł w sposób zwięzły i precyzyjny. Tak by wyglądał zapis definicji wyrażeń arytmetycznych

$\langle \text{Wyrażenie} \rangle \rightarrow \text{liczba}$

$\langle \text{Wyrażenie} \rangle \rightarrow ( \langle \text{Wyrażenie} \rangle )$

$\langle \text{Wyrażenie} \rangle \rightarrow \langle \text{Wyrażenie} \rangle + \langle \text{Wyrażenie} \rangle$

$\langle \text{Wyrażenie} \rangle \rightarrow \langle \text{Wyrażenie} \rangle - \langle \text{Wyrażenie} \rangle$

$\langle \text{Wyrażenie} \rangle \rightarrow \langle \text{Wyrażenie} \rangle * \langle \text{Wyrażenie} \rangle$

$\langle \text{Wyrażenie} \rangle \rightarrow \langle \text{Wyrażenie} \rangle / \langle \text{Wyrażenie} \rangle$

<Wyrażenie>

Ten symbol nosi nazwę kategorii syntaktycznej (ang. syntactic category). Zastępuje on dowolny ciąg znaków należący do języka wyrażen arytmetycznych.

→

Ten symbol należy traktować jako zwrot „może się składać z”

liczba

To pewien abstrakcyjny symbol zastępczy dla dowolnego ciągu znaków.

## Terminologia

Istnieją trzy rodzaje symboli wykorzystywanych w gramatykach

metasymbol - sam w sobie nie ma żadnego znaczenia, oddziela zdefiniowaną kategorię syntaktyczną od opisu sposobu w jaki ciąg znaków może tworzyć daną kategorię.

kategoria syntaktyczna - reprezentuje zbiór zdefiniowanych ciągów znaków.

symbol terminalny - może być znakami, np. (, +, -) lub symbolem abstrakcyjnym

Gramatyka składa się z jednej lub większej liczby produkcji (ang. productions). Każda produkcja składa się z trzech części:

- (1) Części nagłówkowej (ang. head), która jest kategoria syntaktyczna umieszczona po lewej stronie strzałki
- (2) Metasymbolu (np. strzałki)
- (3) Części zasadniczej (ang. body)

<Cyfra> → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

<Liczba> → <Cyfra>

<Liczba> → <Liczba> <Cyfra>

<Wyrażenie> → <Liczba>

<Wyrażenie> → ( <Wyrażenie> )

<Wyrażenie> → <Wyrażenie> + <Wyrażenie>

<Wyrażenie> → <Wyrażenie> - <Wyrażenie>

<Wyrażenie> → <Wyrażenie> \* <Wyrażenie>

<Wyrażenie> → <Wyrażenie> / <Wyrażenie>

Gramatyka wyrażeń  
w której liczby  
zdefiniowano przy  
pomocy konstrukcji  
gramatycznych.

Gramatyka ciągów znaków o zbilansowanej liczbie nawiasów.

**<Zbilansowany> →  $\epsilon$**

**<Zbilansowany> → (<Zbilansowany>) <Zbilansowany>**

Gramatyka ciągów znaków o zbilansowanej liczbie nawiasów zdefiniowana na podstawie gramatyki wyrażeń arytmetycznych.

**<ZbilansowaneW> →  $\epsilon$**

**<ZbilansowaneW> → (<ZbilansowaneW>)**

**<ZbilansowaneW> → <ZbilansowaneW> <ZbilansowaneW>**

Jeżeli kategorie syntaktyczna <Wyrażenie> zastąpi się bardziej odpowiednią nazwą <ZbilansowaneW>, to otrzymuje się nową gramatykę, opisującą ciągi znaków o zbilansowanej liczbie nawiasów. Produkcje są naturalne. Gramatyki przedstawione powyżej definiują ten sam zbiór ciągów znaków.

## Produkcje definiujące część instrukcji języka C

**<Instrukcja> → while (warunek) <Instrukcja>**  
**<Instrukcja> → if (warunek) <Instrukcja>**  
**<Instrukcja> → if (warunek) <Instrukcja> else <Instrukcja>**  
**<Instrukcja> → {<ListaInstr>;}**  
**<Instrukcja> → prostaInstr;**  
**<ListaInstr> → ε**  
**<ListaInstr> → <ListaInstr> <Instrukcja>**

Można opisywać gramatycznie strukturę przebiegu sterowania występującą w językach takich jak C. Załóżmy istnienie abstrakcyjnych symboli terminalnych warunek oraz instrProsta. Pierwszy z nich oznacza wyrażenie warunkowe i można go zastąpić kategorią syntaktyczną <Warunek>. Symbol terminalny instrProsta określa instrukcję nie zawierającą zagnieżdżonych struktur sterujących, takich jak instrukcja przypisania, wywołania funkcji, odczytu, zapisu i skoku. Można zastąpić symbol terminalny kategorią syntaktyczną oraz rozszerzającymi ją produkcjami. Jako kategorii syntaktycznej instrukcji języka C będziemy używać kategorii <Instrukcja>.

# Języki gramatyk

Gramatyka to definicja indukcyjna zawierająca zbiory znaków. W przypadku gramatyk jest rzeczą normalną definiowanie kilku kategorii syntaktycznych za pomocą jednej gramatyki. Dla każdej kategorii syntaktycznej  $\langle S \rangle$  danej gramatyki, definiuje się język  $L(\langle S \rangle)$  w sposób opisany poniżej.

## Podstawa:

W przypadku każdej kategorii syntaktycznej  $\langle S \rangle$  danej gramatyki, język  $L(\langle S \rangle)$  jest zbiorem pustym.

## Indukcja:

Założmy, że gramatyka posiada produkcję  $\langle S \rangle \rightarrow X_1 X_2 \dots X_m$  danej gramatyki, gdzie każdy element  $X_i$ , dla  $i=1, 2, \dots, n$ , jest albo kategorią syntaktyczną, albo symbolem terminalnym. Dla każdego  $i=1, 2, \dots, n$  wybieramy ciąg znaków  $s_i$  dla  $X_i$  w sposób następujący.

(1) Jeżeli  $X_i$  jest symbolem terminalnym, to  $X_i$  można użyć jedynie jako ciągu znaków  $s_i$ .

(2) Jeżeli  $X_i$  jest kategorią syntaktyczną, to wybieramy jako  $s_i$  dowolny ciąg znaków, o którym wiadomo już że należy do języka  $L(X_i)$ .

Wówczas złożenie  $s_1 s_2 \dots s_n$  stanowi ciąg znaków należący do języka  $L(\langle S \rangle)$ .



Jednym z metodycznych sposobów zaimplementowania takiej definicji jest wykonanie sekwencyjnego przebiegu przez produkcję gramatyki. W każdym przebiegu następuje uaktualnienie języka każdej kategorii syntaktycznej przy użyciu reguły indukcyjnej na wszystkie możliwe sposoby, tzn. dla każdego  $X_i$  będącego kategorią syntaktyczną wybieramy ciągi znaków ze zbioru  $L(\langle X_i \rangle)$  na wszystkie możliwe sposoby.

- (1)  $\langle \text{Instrukcja} \rangle \rightarrow \text{while (warunek) } \langle \text{Instrukcja} \rangle$
- (2)  $\langle \text{Instrukcja} \rangle \rightarrow \{ \langle \text{ListaInstr} \rangle \};$
- (3)  $\langle \text{Instrukcja} \rangle \rightarrow \text{prostaInstr};$
- (4)  $\langle \text{ListaInstr} \rangle \rightarrow \langle \text{ListaInstr} \rangle \langle \text{Instrukcja} \rangle$
- (5)  $\langle \text{ListaInstr} \rangle \rightarrow \varepsilon$

**Uproszczona  
gramatyka instrukcji**

- $$\langle I \rangle \rightarrow w c \langle I \rangle$$
- $$\langle I \rangle \rightarrow \{ \langle L \rangle \}$$
- $$\langle I \rangle \rightarrow s;$$
- 
- $$\langle L \rangle \rightarrow \langle L \rangle \langle I \rangle$$
- $$\langle L \rangle \rightarrow \varepsilon$$

**Uproszczona  
notacja**

## Uproszczona gramatyka instrukcji

$\langle I \rangle \rightarrow w c \langle I \rangle$   
 $\langle I \rangle \rightarrow \{ \langle L \rangle \}$   
 $\langle I \rangle \rightarrow s ;$   
  
 $\langle L \rangle \rightarrow \langle L \rangle \langle I \rangle$   
 $\langle L \rangle \rightarrow \varepsilon$

Język definiowany przez gramatykę może być nieskończony, czyli nie ma możliwości wypisania wszystkich należących do niego znaków.

## Nowe ciągi znaków dodawane w pierwszych trzech przebiegach

	I	L
Przebieg 1.	s ;	$\varepsilon$
Przebieg 2.	wcs ; }	s ;
Przebieg 3.	wcwcs ; ws { {s ;}	wcs ; { s ; s ; s ; wcs ; s ; {

## Drzewa rozbioru

Można pokazać że ciąg  $s$  należy do języka  $L(\langle S \rangle)$  pewnej kategorii syntaktycznej  $\langle S \rangle$  w wyniku powtarzalnego stosowania produkcji. Rozpoczyna się od pewnych ciągów znaków wynikających z produkcji bazowych, czyli tych, które w swojej części zasadniczej nie posiadają żadnych kategorii syntaktycznych. Następnie należy „zastosować” produkcję wobec ciągów znaków które już otrzymano dla różnych kategorii syntaktycznych. Każde takie zastosowanie polega na podstawieniu ciągów znaków za występujące w części zasadniczej produkcji różne kategorie syntaktyczne, stąd skonstruowanie ciągu znaków należącego do kategorii syntaktycznej stanowiącej część nagłówkową. W końcu ciąg  $s$  konstruuje się poprzez zastosowanie produkcji posiadającej w części nagłówkowej kategorię  $\langle S \rangle$ . Możemy ilustrować przynależność  $s$  do  $L(\langle S \rangle)$  w formie drzewa, zwanego drzewem rozbioru lub drzewem analizy składniowej (ang. parse tree). Wierzchołki drzewa rozbioru etykietuje się albo symbolami terminalnymi, albo kategoriami syntaktycznymi, albo symbolem  $\varepsilon$ . Liście są etykietowane jedynie symbolami terminalnymi lub symbolem  $\varepsilon$ , zaś wierzchołki wewnętrzne są etykietowane jedynie kategoriami syntaktycznymi. Każdy wierzchołek wewnętrzny  $v$  reprezentuje zastosowanie produkcji. Tzn. kategoria syntaktyczna etykietująca  $v$  wierzchołek stanowi część nagłówkową produkcji. Etykiety potomków wierzchołka  $v$ , od strony lewej do prawej, tworzą część zasadniczą tej produkcji.

<Cyfra> → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

<Liczba> → <Cyfra>

<Liczba> → <Liczba> <Cyfra>

<Wyrażenie> → <Liczba>

<Wyrażenie> → ( <Wyrażenie> )

<Wyrażenie> → <Wyrażenie> + <Wyrażenie>

<Wyrażenie> → <Wyrażenie> - <Wyrażenie>

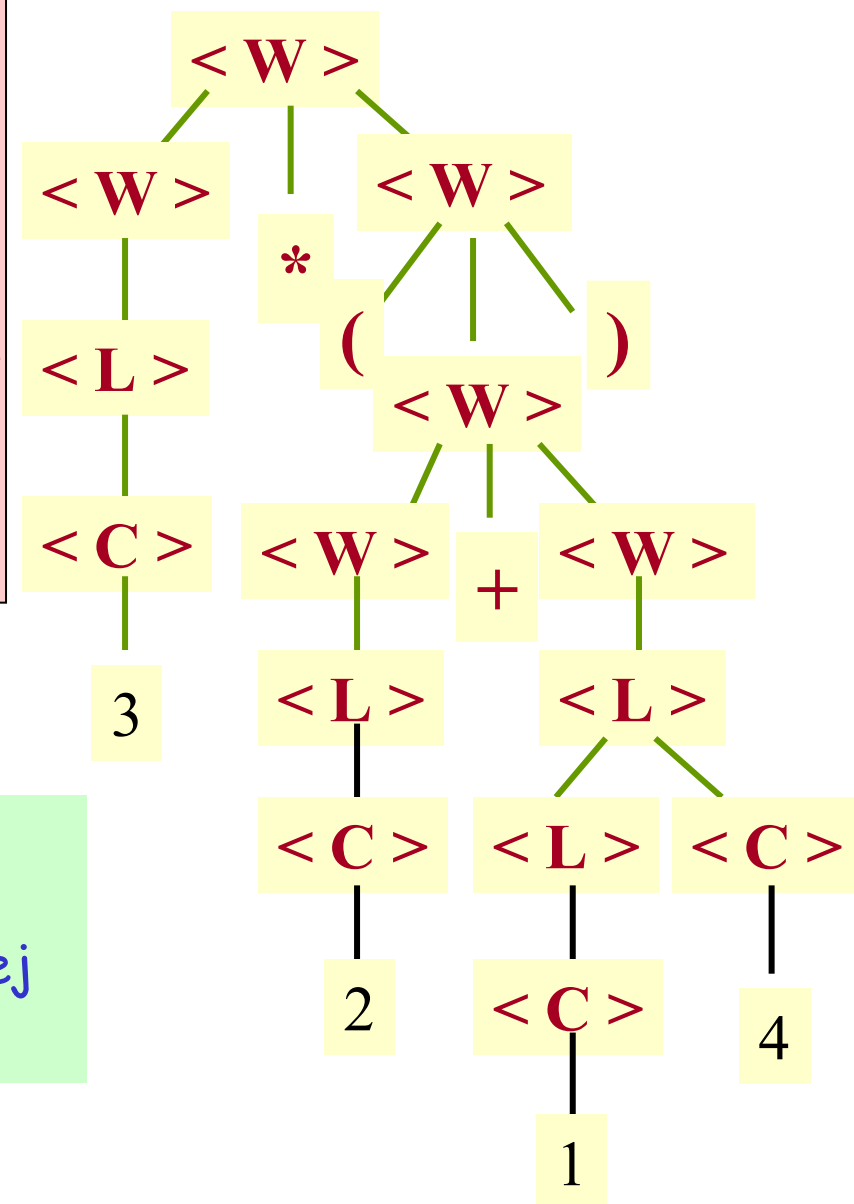
<Wyrażenie> → <Wyrażenie> \* <Wyrażenie>

<Wyrażenie> → <Wyrażenie> / <Wyrażenie>

Drzewo rozbioru dla ciągu znaków

$3 * (2 + 14)$

przy użyciu gramatyki zdefiniowanej  
powyżej.



# Konstruowanie drzew rozbioru

Każde drzewo rozbioru reprezentuje ciąg symboli terminalnych  $s$ , który nosi nazwę wyniku (ang. yield) drzewa. Ciąg  $s$  składa się z etykiet liści drzewa ułożonych w kolejności od strony lewej do prawej. Jeżeli drzewo posiada jeden wierzchołek, wierzchołek ten jest etykietowany symbolem terminalnym lub symbolem  $\varepsilon$ , ponieważ jest liściem. Jeżeli drzewo posiada więcej niż jeden wierzchołek, to korzeń zostaje zaetykietowany kategorią syntaktyczną, gdyż korzeń drzewa posiadającego dwa lub więcej wierzchołków jest zawsze wierzchołkiem wewnętrznym. Ta kategoria syntaktyczna zawsze zawiera wśród swoich ciągów znaków także wynik drzewa.

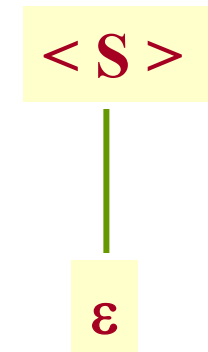
# Definicja indukcyjna drzew rozbioru dla danej gramatyki

## Podstawa:

Dla każdego symbolu terminalnego, np.  $x$ , danej gramatyki istnieje drzewo posiadające tylko jeden wierzchołek zaetykietowany jako  $x$ . Wynikiem takiego drzewa jest  $x$ .

## Indukcja:

Założmy, że istnieje produkcja  $\langle S \rangle \rightarrow X_1 X_2 \dots X_m$ , gdzie każdy z symboli  $X$  oznacza albo symbol terminalny, albo kategoria syntaktyczna. Jeżeli  $n=0$ , czyli produkcja ma postać  $\langle S \rangle \rightarrow \varepsilon$ , to istnieje drzewo którego wynikiem jest  $\varepsilon$ , a korzeniem  $\langle S \rangle$ . Z uwagi na te produkcje ciąg  $\varepsilon$  należy do języka  $L(\langle S \rangle)$ .

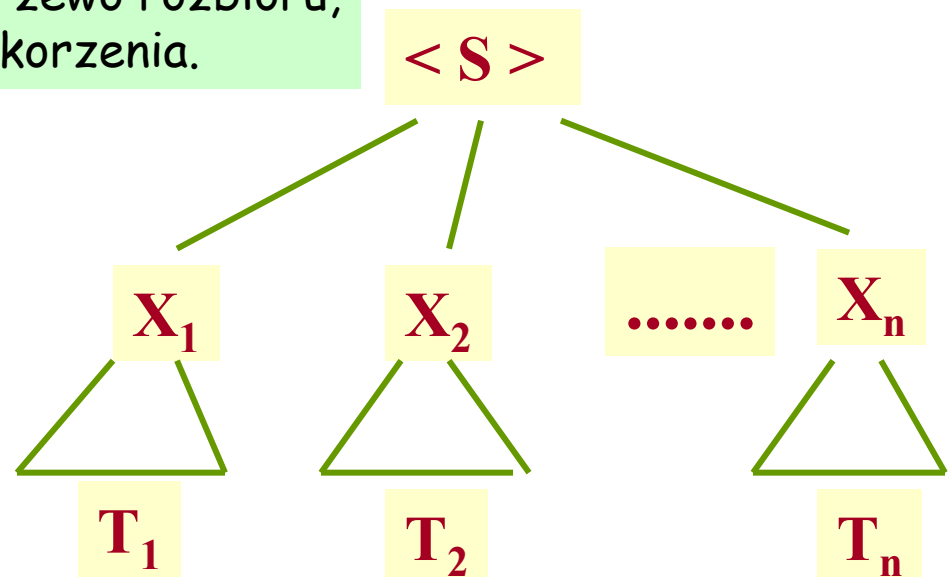


Teraz załóżmy, że  $\langle S \rangle \rightarrow X_1 X_2 \dots X_n$  a  $n \geq 1$ . Dla każdego  $X_i$  możemy wybrać drzewo  $T_i$  w sposób następujący:

(1) Jeśli  $X_i$  jest symbolem terminalnym, musimy wybrać dla każdego wystąpienia tego symbolu terminalnego 1-wierzchołkowe drzewo zaetykietowane jako  $X_i$ . Jeżeli jeden lub więcej symboli  $X$  jest tym samym symbolem terminalnym, to musimy wybrać różne jednowierzchołkowe drzewa o tej samej etykiecie

(2) Jeżeli  $X_i$  jest kategorią syntaktyczną, możemy wybrać dowolne już skonstruowane drzewo rozbioru, takie, że posiada on  $X_i$  jako etykietę korzenia.

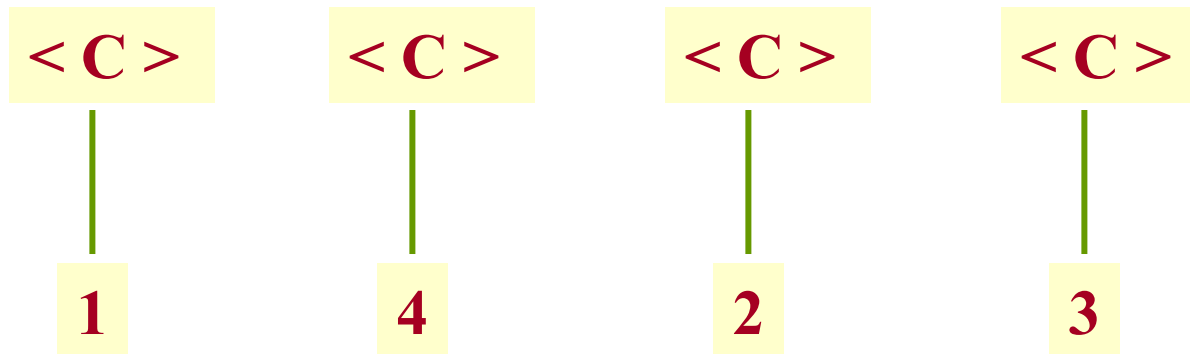
Następnie konstruujemy drzewo, tworząc korzeń zaetykietowany kategorią syntaktyczną  $\langle S \rangle$  umieszczoną w części nagłówkowej produkcji oraz przypisujemy mu jako potomków korzenie drzew wybranych dla  $X_1, X_2, \dots, X_n$ .



# Etapy konstruowania drzewa rozbioru 3\* (2+14)

(a) Konstruujemy jednowierzchołkowe drzewo dla każdego symbolu terminalnego w drzewie. Używamy produkcji

**<Cyfra> → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9**

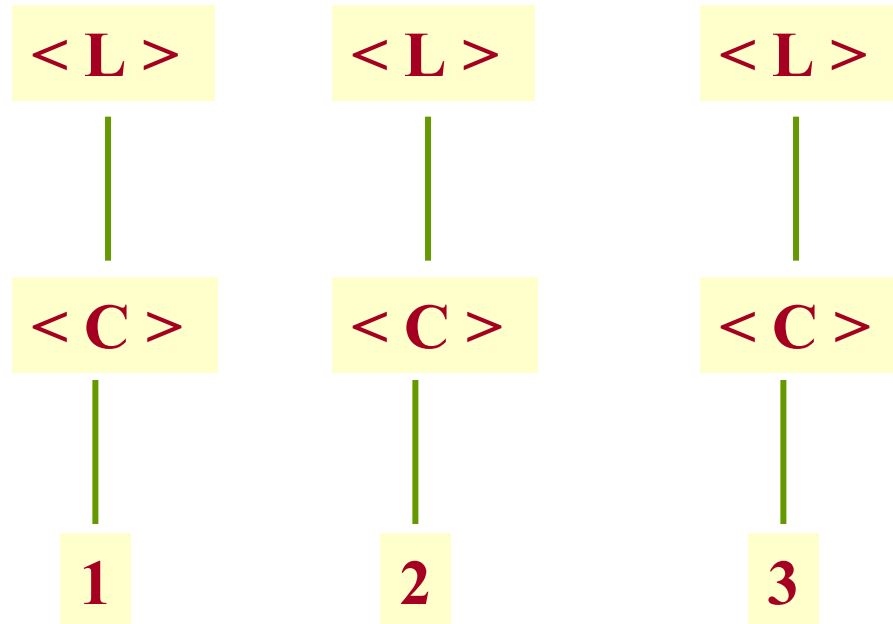




# Etapy konstruowania drzewa rozbioru 3\* (2+14)

(b) Określamy że cyfry są liczbami. Używamy produkcji

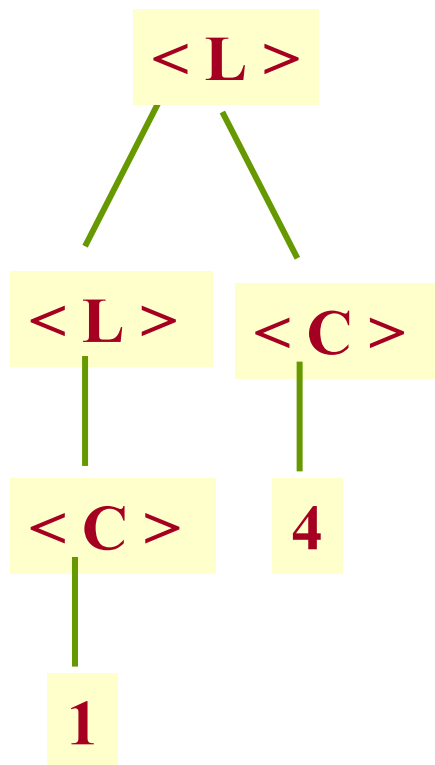
**<Liczba> → <Cyfra>**



# Etapy konstruowania drzewa rozbioru 3\* (2+14)

(c) Określamy że cyfr i liczba jest liczba. Wynikiem drzewa jest 14.  
Używamy produkcji

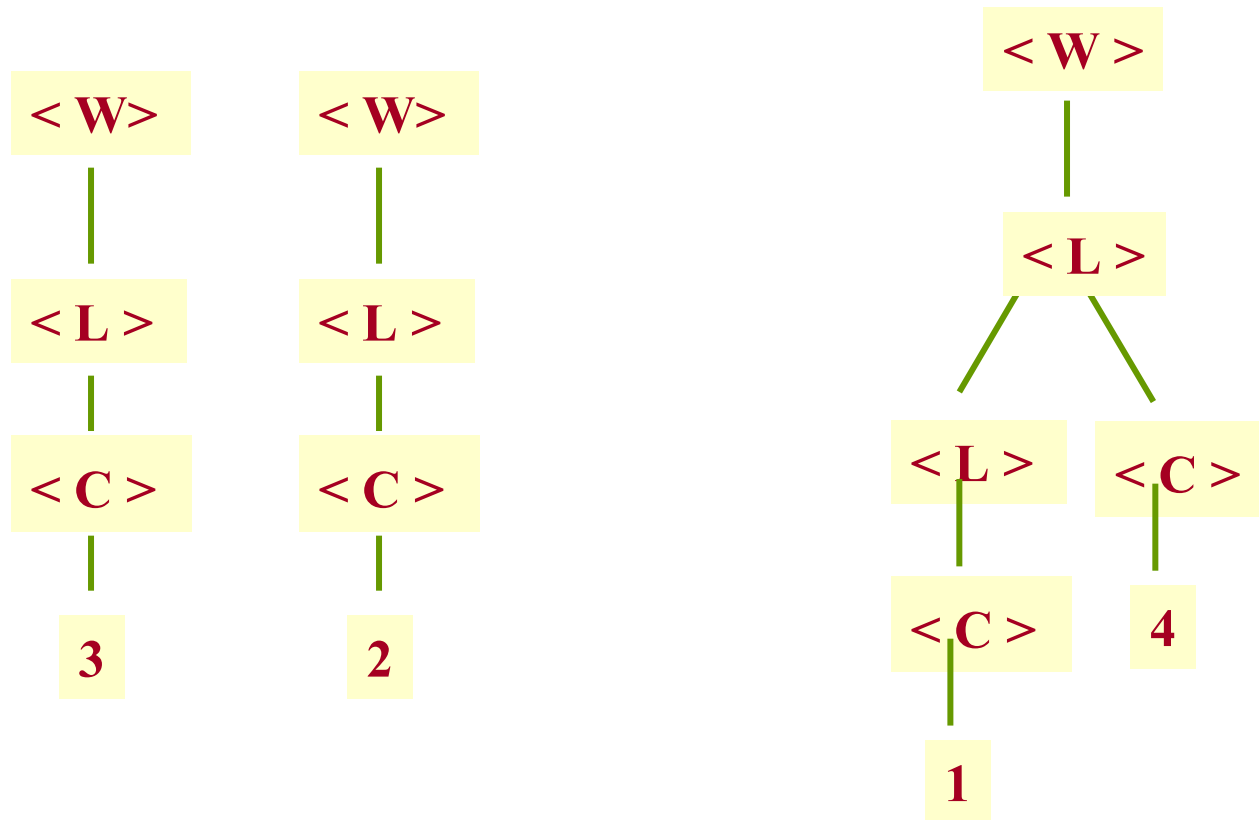
$\langle \text{Liczba} \rangle \rightarrow \langle \text{Liczba} \rangle \langle \text{Cyfra} \rangle$



# Etapy konstruowania drzewa rozbioru $3 * (2 + 14)$

(d) Tworzymy drzewa rozbioru dla wyrażeń 3, 2, 14. Używamy produkcji

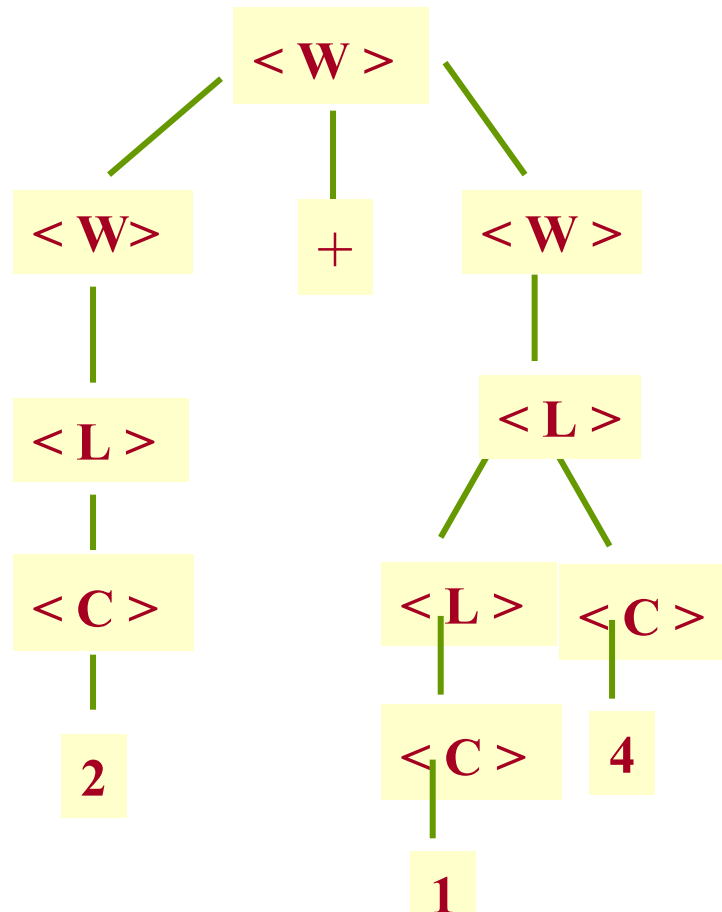
**<Wyrażenie> → <Liczba>**



# Etapy konstruowania drzewa rozbioru $3 * (2+14)$

(e) Tworzymy drzewo dla sumy  $2 + 14$ . Używamy produkcji

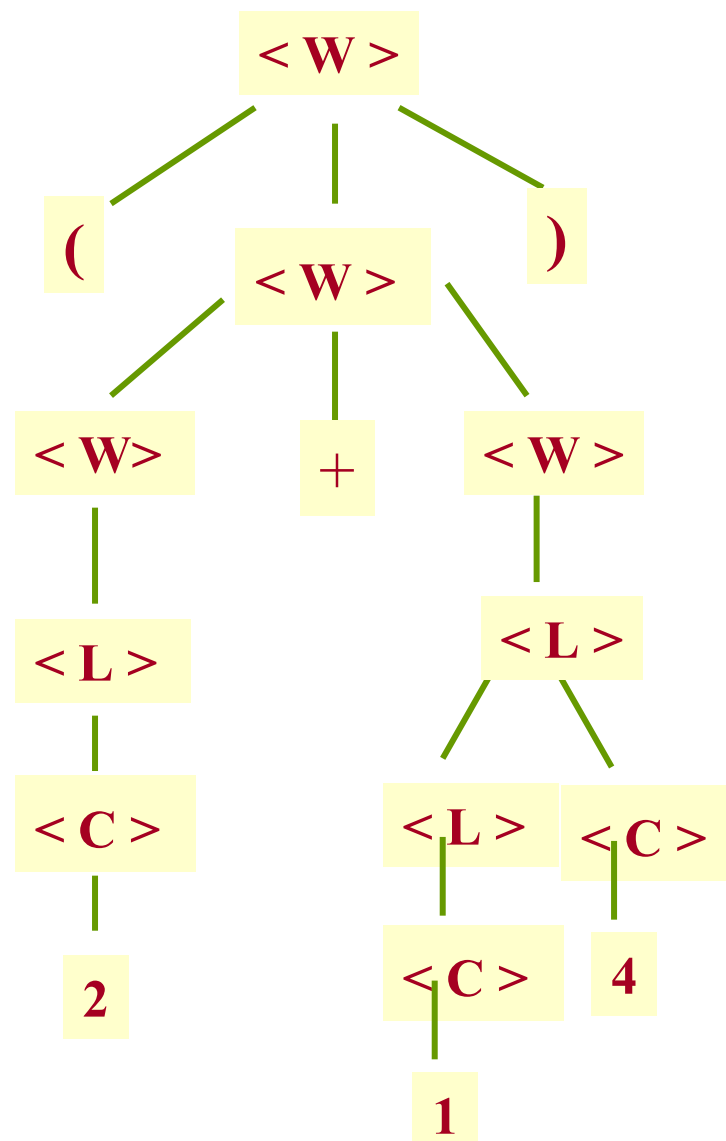
$\langle \text{Wyrażenie} \rangle \rightarrow \langle \text{Wyrażenie} \rangle + \langle \text{Wyrażenie} \rangle$



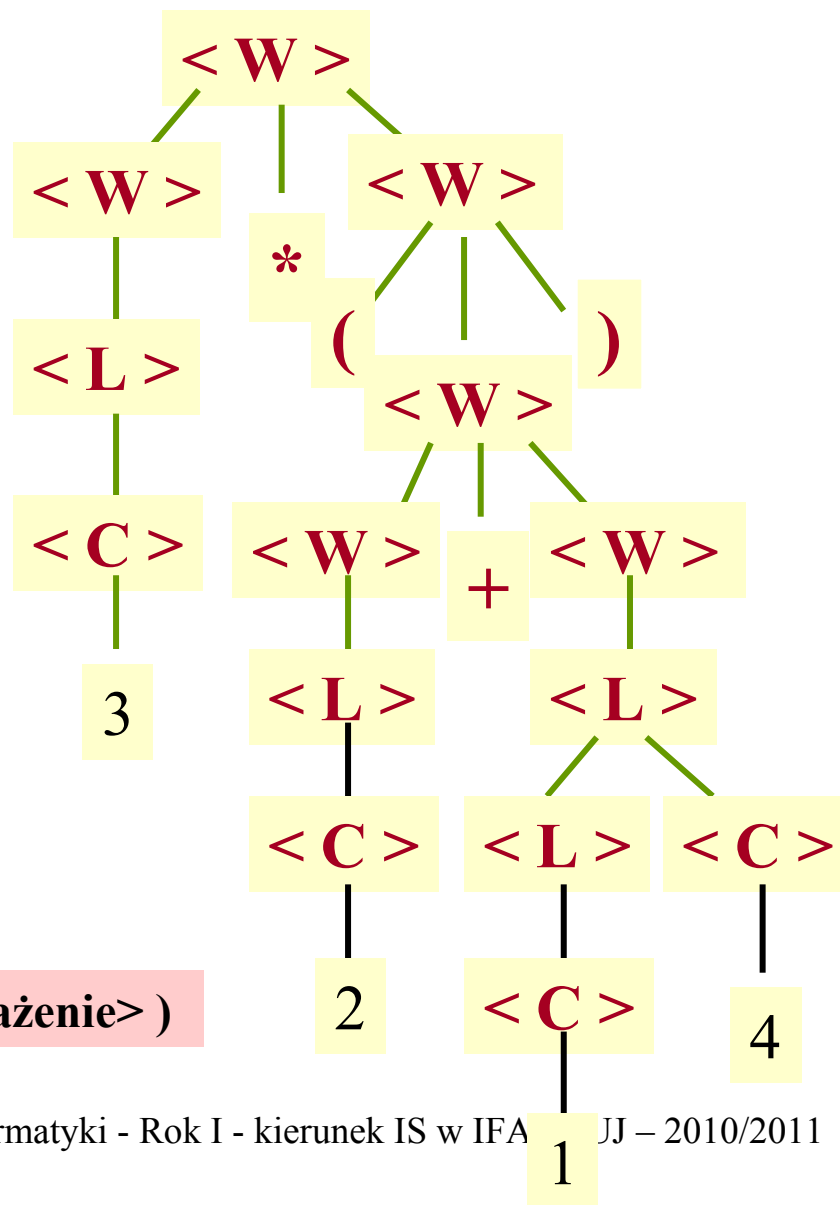
# Etapy konstruowania drzewa rozbioru $3 * (2 + 14)$

(f) Tworzymy drzewo dla sumy  $(2 + 14)$ . Używamy produkcji

$\langle \text{Wyrażenie} \rangle \rightarrow ( \langle \text{Wyrażenie} \rangle )$



# Etapy konstruowania drzewa rozbioru 3\* (2+14)



(g) Tworzymy drzewo  $3 * (2 + 14)$ .  
Używamy produkcji

$\langle \text{Wyrażenie} \rangle \rightarrow \langle \text{Wyrażenie} \rangle * (\langle \text{Wyrażenie} \rangle )$

# Uzasadnienie poprawności konstrukcji drzew rozbioru

Można udowodnić za pomocą dwóch prostych indukcji, że wyniki drzew rozbioru o korzeniu  $\langle S \rangle$  odpowiadają dokładnie ciągom znaków języka  $L(\langle S \rangle)$  dla dowolnej kategorii syntaktycznej  $\langle S \rangle$ .

To znaczy:

(1) Jeżeli  $T$  jest drzewem rozbioru o korzeniu zaetykietowanym przez  $\langle S \rangle$  i daje ono wynik  $s$ , to ciąg  $s$  należy do języka  $L(\langle S \rangle)$ .

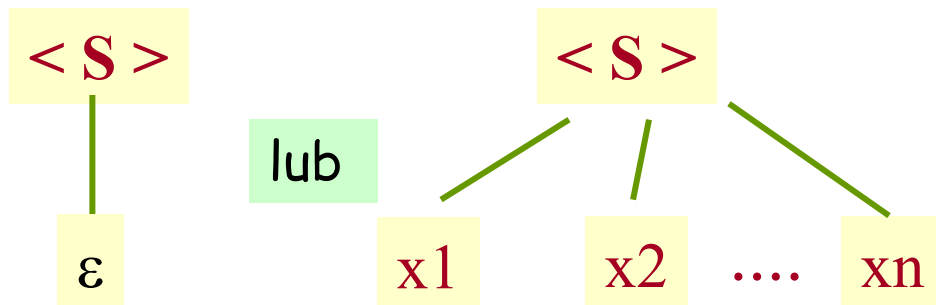
(2) Jeżeli ciąg  $s$  należy do języka  $L(\langle S \rangle)$ , to istnieje drzewo rozbioru dające wynik  $s$  oraz posiadające korzeń zaetykietowany przez  $\langle S \rangle$ .

Drzewa rozbioru są tworzone z mniejszych drzew rozbioru w ten sam sposób, w jaki dłuższe ciągi znaków składa się z krótszych, przy użyciu podstawień za kategorie syntaktyczne w częściach zasadniczych produkcji.

Dowód rozpoczynamy od części (1), której dowodzimy za pomocą indukcji zupełnej względem wysokości drzewa.

## Podstawa:

Zakładamy że wysokość drzewa rozbioru wynosi 1. Wówczas drzewo ma postać



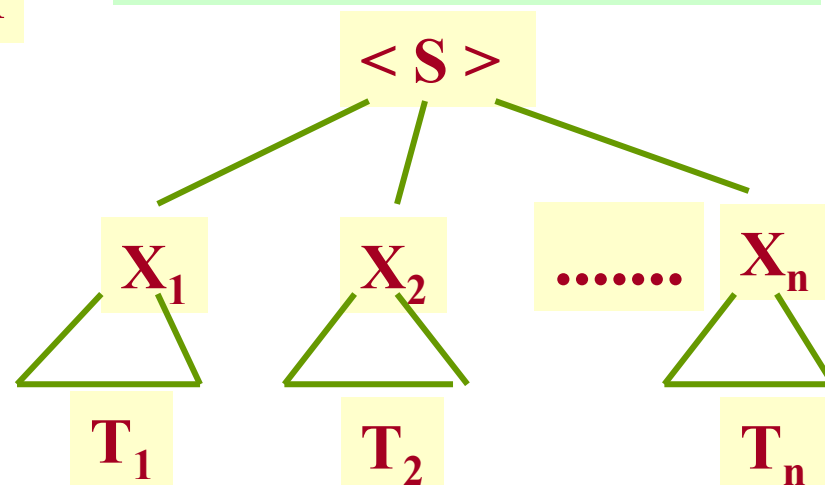
gdzie każde  $x$  jest symbolem terminalnym. Zatem  $x_1x_2\dots x_n$  stanowi ciąg znaków należący do języka  $L(\langle S \rangle)$

## Indukcja:

Zakładamy, że (1) jest prawdziwe dla wszystkich drzew o wysokości  $k$  lub mniejszej.

Bierzemy pod uwagę drzewo o wysokości  $k+1$ . Każde poddrzewo  $T$  ma wysokość co najwyżej  $k$ . Jeżeli  $X_i$  jest kategorią syntaktyczną, to wynik drzewa  $T_i$ , na przykład  $s_i$ , należy do języka  $L(X_i)$ . Jeżeli  $X_i$  jest symbolem terminalnym, ciąg  $s_i$  definiujemy jako  $X_i$ . Wówczas wynikiem drzewa jest  $s_1s_2\dots s_n$ .

Na podstawie reguły indukcyjnej dla języka kategorii  $\langle S \rangle$ ,  $s_1s_2\dots s_n$  należy do  $L(\langle S \rangle)$ .





Teraz należy udowodnić część (2) mówiącą, że każdy ciąg  $s$  kategorii syntaktycznej  $\langle S \rangle$  posiada drzewo rozbioru o korzeniu  $\langle S \rangle$  oraz wyniku  $s$ . Na początku należy zauważyć, że dla każdego symbolu terminalnego  $x$  istnieje drzewo rozbioru, w którym  $x$  jest zarówno korzeniem jak i wynikiem.

### Podstawa:

Zakładamy, że  $s$  wymaga jednego zastosowania kroku indukcyjnego w celu wykazania, że  $s$  należy do  $L(\langle S \rangle)$ . Wówczas musi istnieć produkcja postaci  $\langle S \rangle \rightarrow x_1x_2\dots x_n$ , gdzie wszystkie symbole  $x$  są symbolami terminalnymi a  $s=x_1x_2\dots x_n$ . Wiadomo, że istnieje jednowierzchołkowe drzewo rozbioru z etykietą  $x_i$ , dla  $i=1,2,\dots,n$ . A zatem istnieje drzewo rozbioru o wyniku  $s$  i korzeniu zaetykietowanym przez  $\langle S \rangle$ . (Patrz drzewa na str. 24.)

### Indukcja:

Zakładamy, że dowolny ciąg znaków  $t$  należący do języka dowolnej kategorii syntaktycznej  $\langle T \rangle$  w wyniku zastosowania kroku indukcyjnego  $k$  lub mniejsza liczbę razy posiada drzewo rozbioru z  $t$  jako wynikiem oraz  $\langle T \rangle$  jako korzeniem. Bierzemy pod uwagę ciąg  $s$  należący do języka kategorii syntaktycznej  $\langle S \rangle$  w wyniku zastosowania kroku indukcyjnego  $k+1$  lub mniejsza liczbę razy.

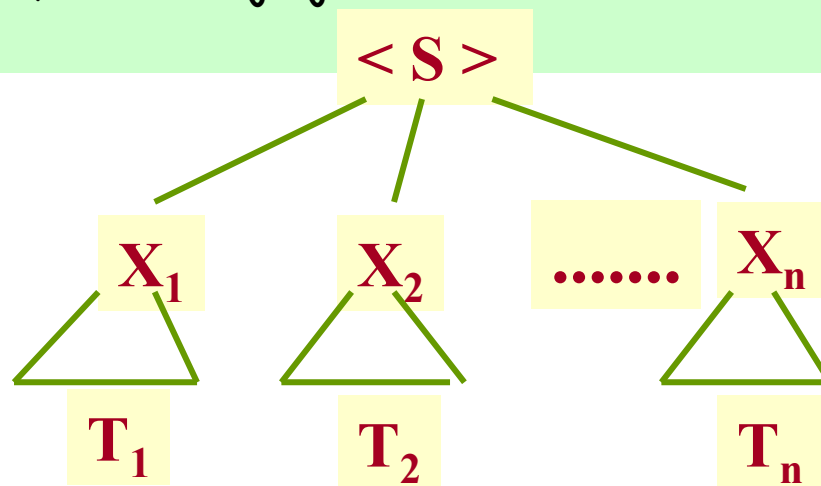
## Indukcja (cont):

Wówczas istnieje produkcja  $\langle S \rangle \rightarrow X_1 X_2 \dots X_n$  oraz zachodzi  $s = s_1 s_2 \dots s_n$ , gdzie każdy podciąg  $s_i$  jest:

- (1) Równy  $X_i$ , jeżeli  $X_i$  jest symbolem terminalnym.
- (2) Pewnym ciągiem, o którym wiadomo, że należy do języka  $L(X_i)$  w wyniku zastosowania reguły indukcyjnej co najwyżej  $k$  razy, jeżeli  $X_i$  jest kategorią syntaktyczną.

A zatem dla każdego  $i$  można określić drzewo  $T_i$  posiadające wynik  $s_i$  oraz korzeń zaetykietowany przez  $X_i$ . Jeżeli  $X_i$  jest kategorią syntaktyczną, wykorzystujemy hipotezę indukcyjną w celu stwierdzenia, że  $T_i$  istnieje, zaś jeżeli  $X_i$  jest symbolem terminalnym, nie ma potrzeby wykorzystywania hipotezy indukcyjnej w celu stwierdzenia, że istnieje jednowierzchołkowe drzewo zaetykietowane przez  $X_i$ .

A zatem drzewo skonstruowane jak to po lewej stronie, posiada wynik  $s$  oraz korzeń zaetykietowany przez  $\langle S \rangle$ .



## Drzewa rozbioru i drzewa wyrażeń

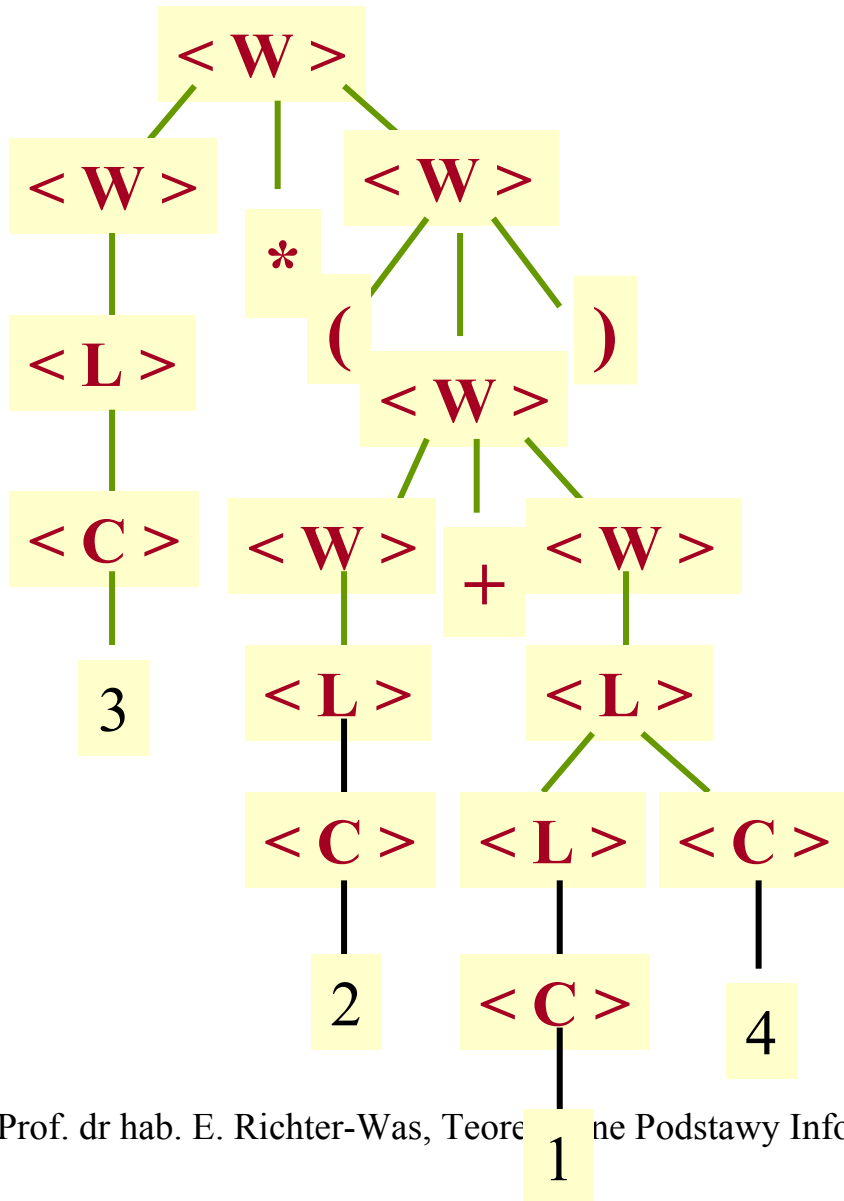
Mając sformułowaną gramatykę wyrażeń możemy drzewa rozbioru przekonwertować na drzewa wyrażeń, dokonując trzech transformacji:

(1) Wierzchołki związane z poszczególnymi operandami niepodzielnymi są łączone w jeden wierzchołek zaetykietowany danym operandem

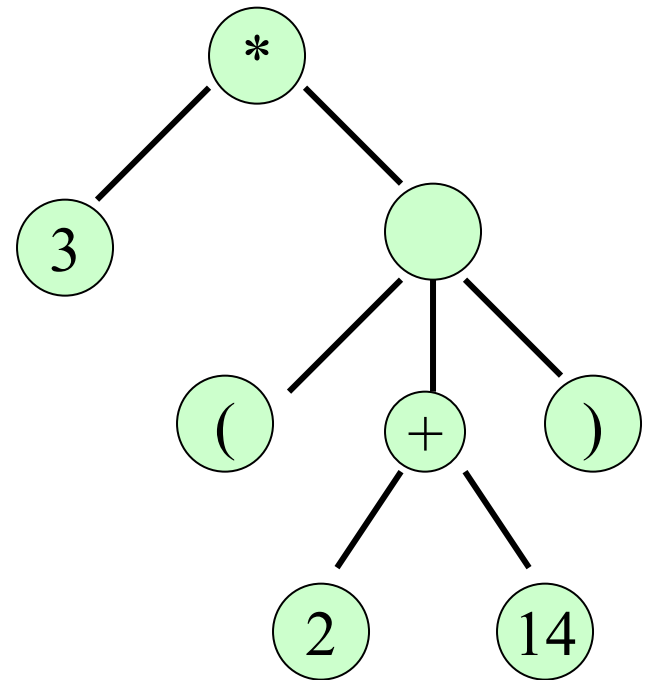
(2) Operatory zostają przesunięte z liści do ich wierzchołków nadrzędnych. To znaczy symbol operatora, taki jak +, staje się etykietą wierzchołka umieszczonego nad nim, który wcześniej był zaetykietowany kategorią syntaktyczną „wyrażenia”.

(3) Wierzchołki wewnętrzne, których etykietami wciąż są „wyrażenia” zostają usunięte.

## drzewo rozbioru



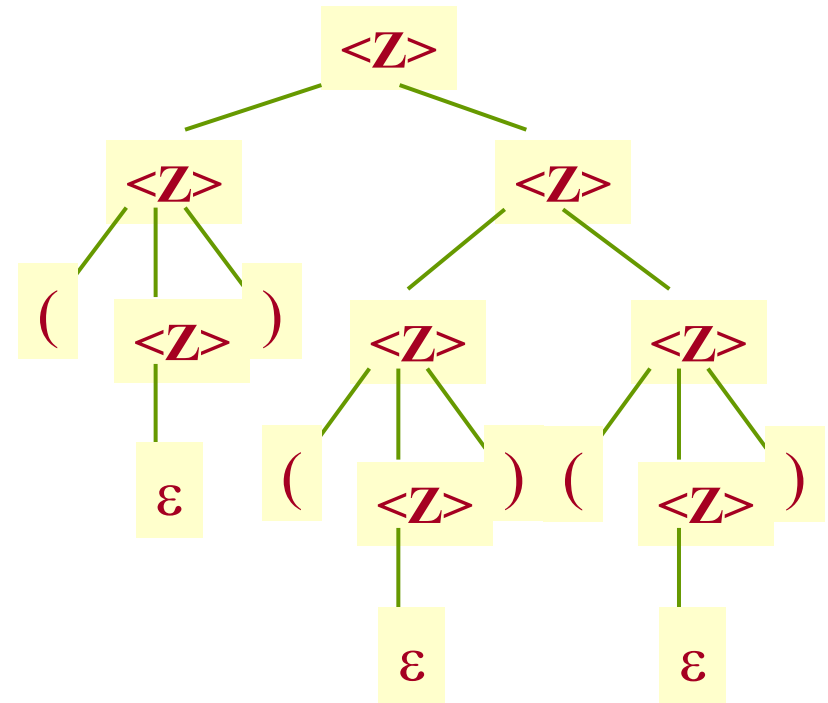
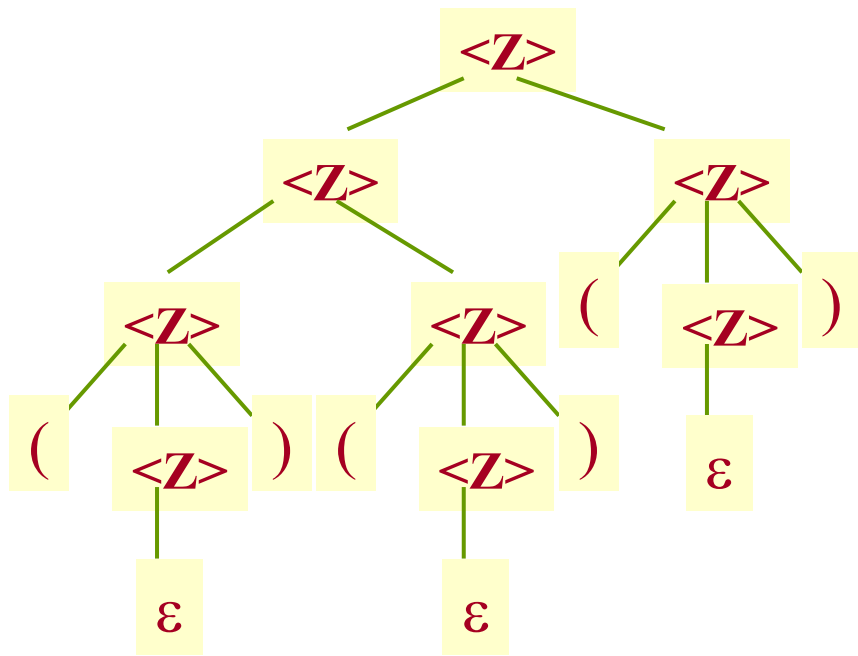
## drzewo wyrażeń



# Niejednoznaczność i projektowanie gramatyk

Rozpatrzmy gramatykę zbilansowanych nawiasów.  
Chcemy utworzyć drzewo rozbioru dla ciągu znaków  
( ) ( ). Można utworzyć dwa takie drzewa.

$\langle Z \rangle \rightarrow \varepsilon$   
 $\langle Z \rangle \rightarrow (\langle Z \rangle)$   
 $\langle Z \rangle \rightarrow \langle Z \rangle \langle Z \rangle$

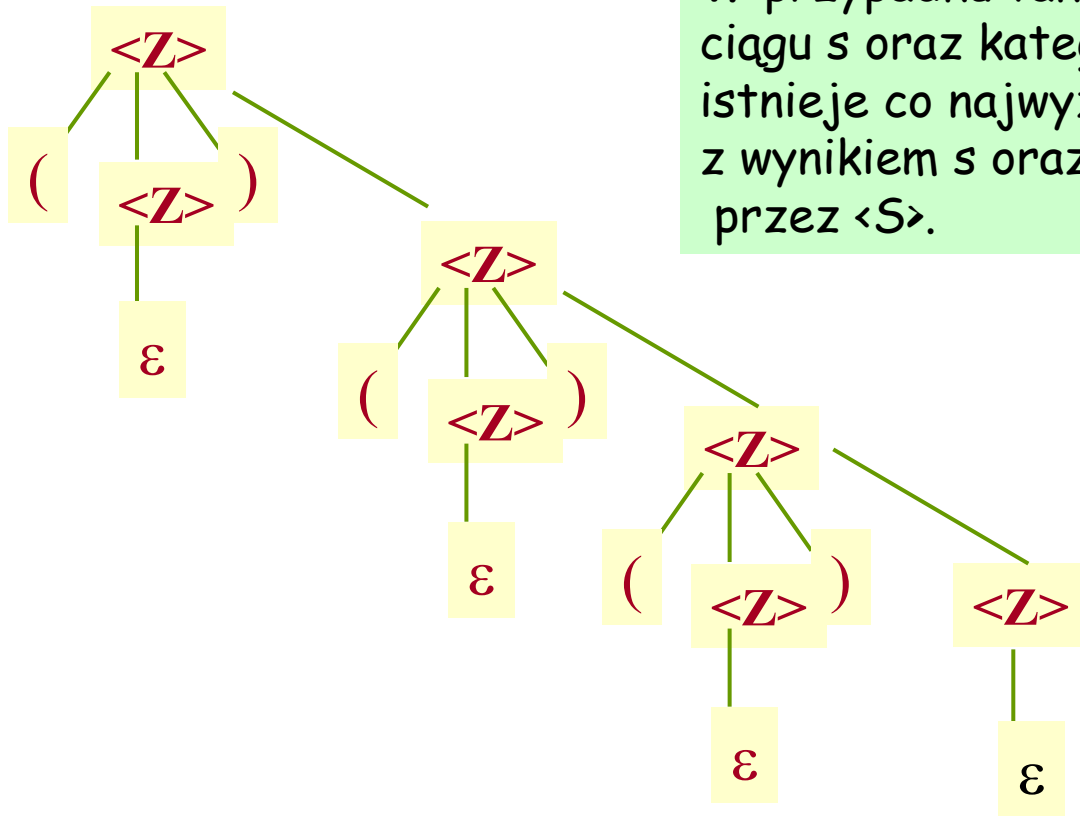


Gramatyka w której istnieją dwa lub więcej drzewa rozbioru o tym samym wyniku oraz tej samej kategorii syntaktycznej etykietującej korzeń jest nazywana gramatyka niejednoznaczna. (ang. *ambiguous*). Wystarczy żeby istniał choć jeden taki ciąg który jest niejednoznaczny.

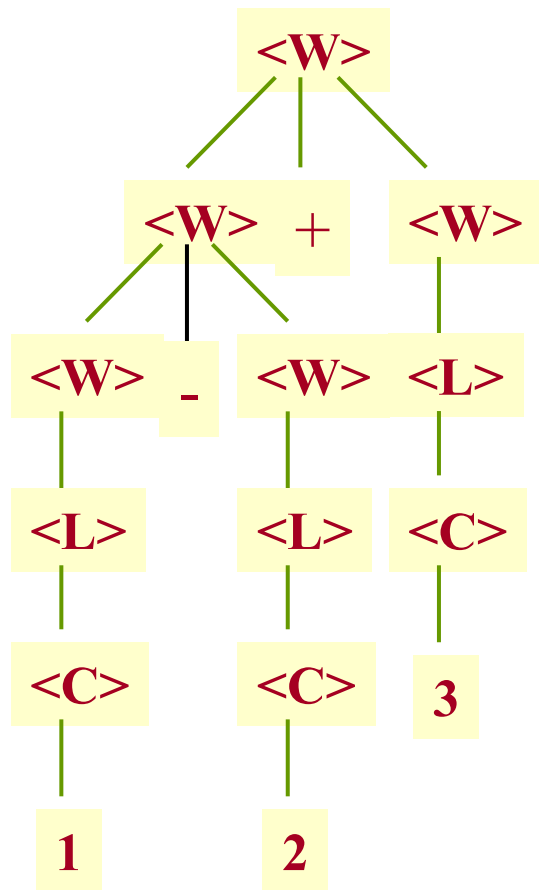
Rozpatrzmy inną gramatykę zbilansowanych nawiasów. Chcemy utworzyć drzewo rozbioru dla ciągu znaków  $()()()$ . Można utworzyć tylko jedno takie drzewo.

$\langle Z \rangle \rightarrow \varepsilon$   
 $\langle Z \rangle \rightarrow (\langle Z \rangle) \langle Z \rangle$

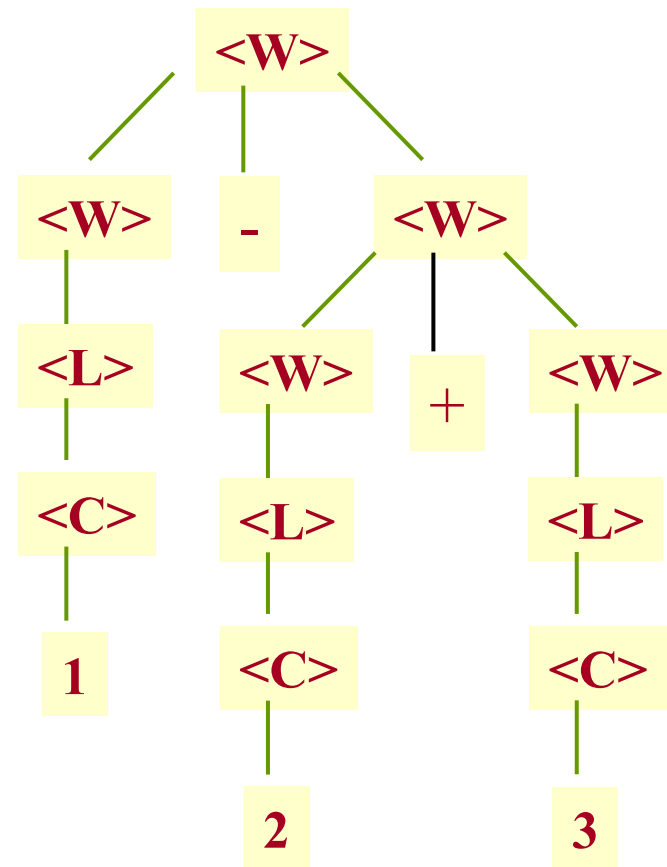
Gramatyka, która nie jest niejednoznaczna nosi nazwę jednoznacznej (ang. *unambiguous*). W przypadku takiej gramatyki dla każdego ciągu  $s$  oraz kategorii syntaktycznej  $\langle S \rangle$  istnieje co najwyżej jedno drzewo rozbioru z wynikiem  $s$  oraz korzeniem zaetykietowanym przez  $\langle S \rangle$ .



Niejednoznaczność gramatyk wyrażen może być poważnym problemem.  
 Niektóre drzewa rozbioru mogą dawać złe wartości dla wyrażen.  
 Dwa drzewa rozbioru dla wyrażenia: 1-2+3

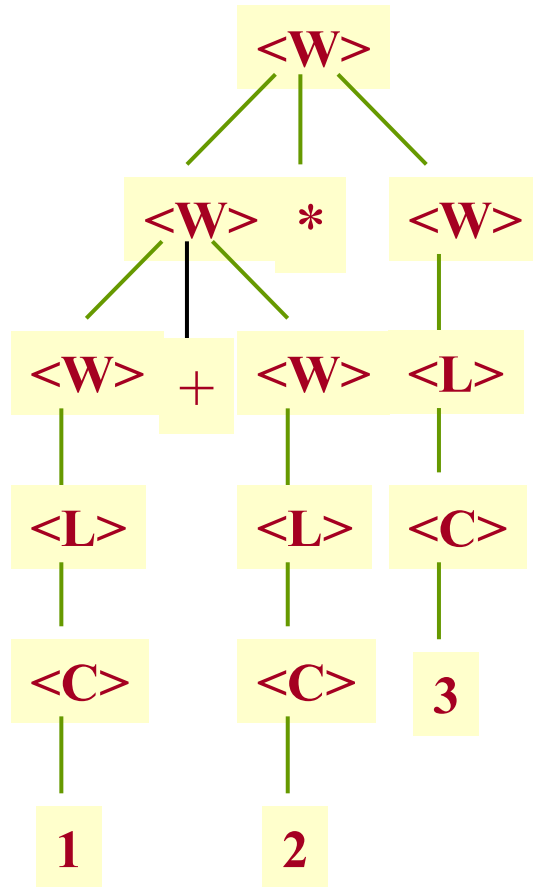


**Poprawne drzewo rozbioru**  
 $1-2+3 = 2$

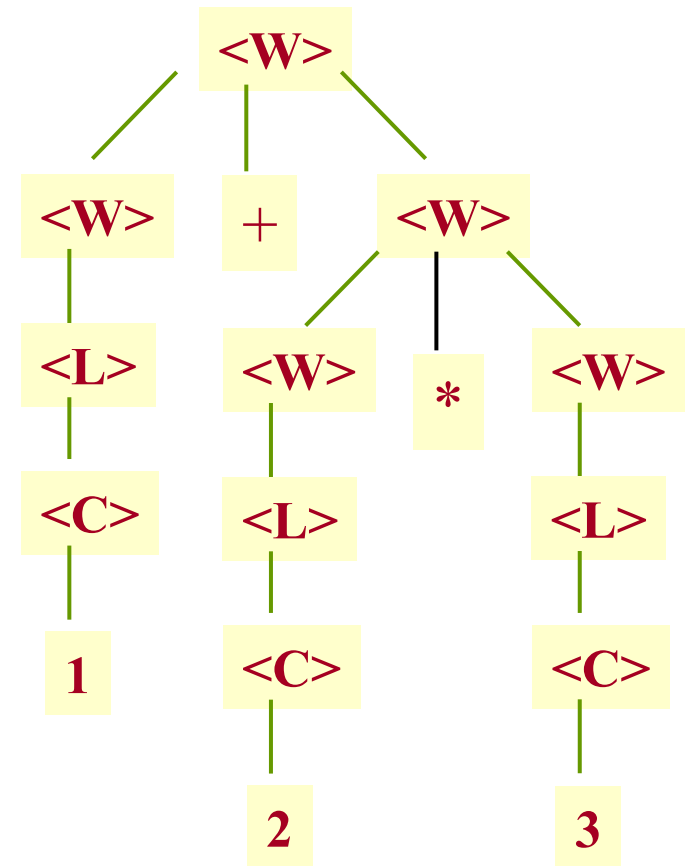


**Niepoprawne drzewo rozbioru**  
 $1-2+3=-4$

Niejednoznaczność gramatyk wyrażeń może być poważnym problemem.  
 Niektóre drzewa rozbioru mogą dawać złe wartości dla wyrażeń.  
 Dwa drzewa rozbioru dla wyrażenia: 1+2\*3



Niepoprawne drzewo rozbioru  
 $1 + 2 * 3 = 9$



Poprawne drzewo rozbioru  
 $1 + 2 * 3 = 7$



# Istota niejednoznaczności

Analizator składniowy który konstruuje drzewa rozbioru dla programów stanowi podstawowy element kompilatora. Jeżeli gramatyka opisująca język programowania jest niejednoznaczna, oraz jeżeli jej niejednoznaczności są wyprowadzone jednostronnie, to w przypadku przynajmniej części programów istnieje więcej niż jedno drzewo rozbioru. Jeżeli gramatyka programu jest niejednoznaczna, kompilator nie może podjąć prawidłowej decyzji odnośnie do drzewa rozbioru dla pewnych programów, a w związku z tym nie może zdecydować, jakie działania powinien wykonać program w języku maszynowym. Kompilatory muszą korzystać ze specyfikacji które są jednoznaczne.

# Jednoznaczne gramatyki wyrażeń

Konstrukcja jednoznacznej gramatyki polega na zdefiniowaniu trzech kategorii syntaktycznych o następującym znaczeniu:

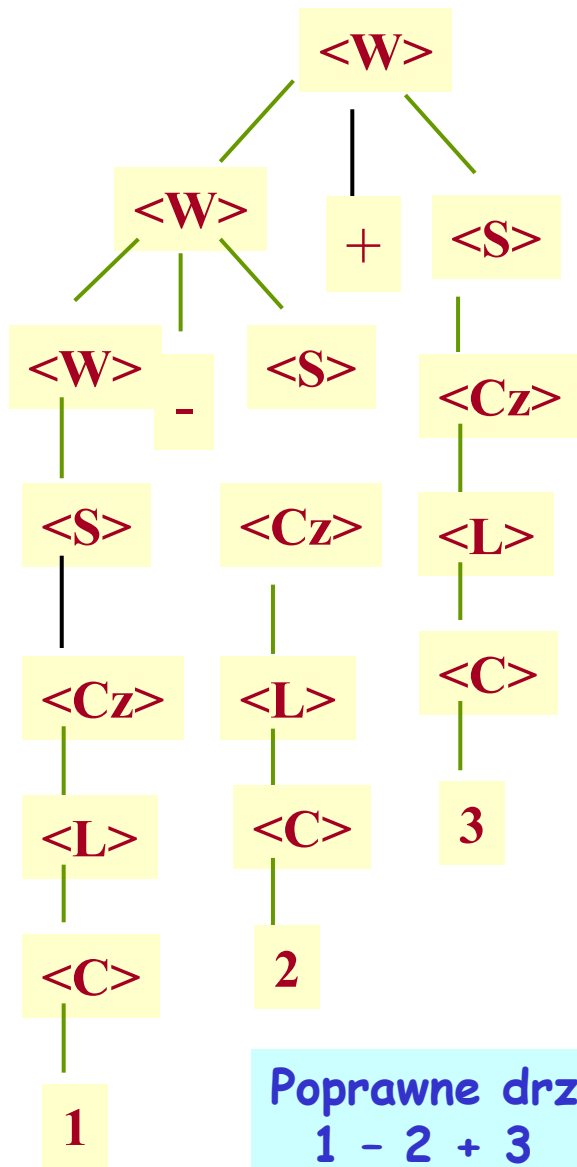
**<Czynnik>** - generuje wyrażenia, które nie mogą zostać rozdzielone, to znaczy czynnik jest albo pojedynczym operandem, albo dowolnym wyrażeniem umieszczonym w nawiasie.

**<Składnik>** - generuje iloczyn lub iloraz czynników. Pojedynczy czynnik jest składnikiem, a więc stanowi ciąg czynników rozdzielonych operatorami \* lub /. Przykładami składników są 12 lub  $12/3*45$ .

**<Wyrażenie>** - generuje różnicę lub sumę jednego lub większej liczby składników. Pojedynczy składnik jest wyrażeniem, a więc stanowi sekwencję składników rozdzielonych operatorami + lub -. Przykładami wyrażeń są 12,  $12/3*45$  lub  $12+3*45-6$ .

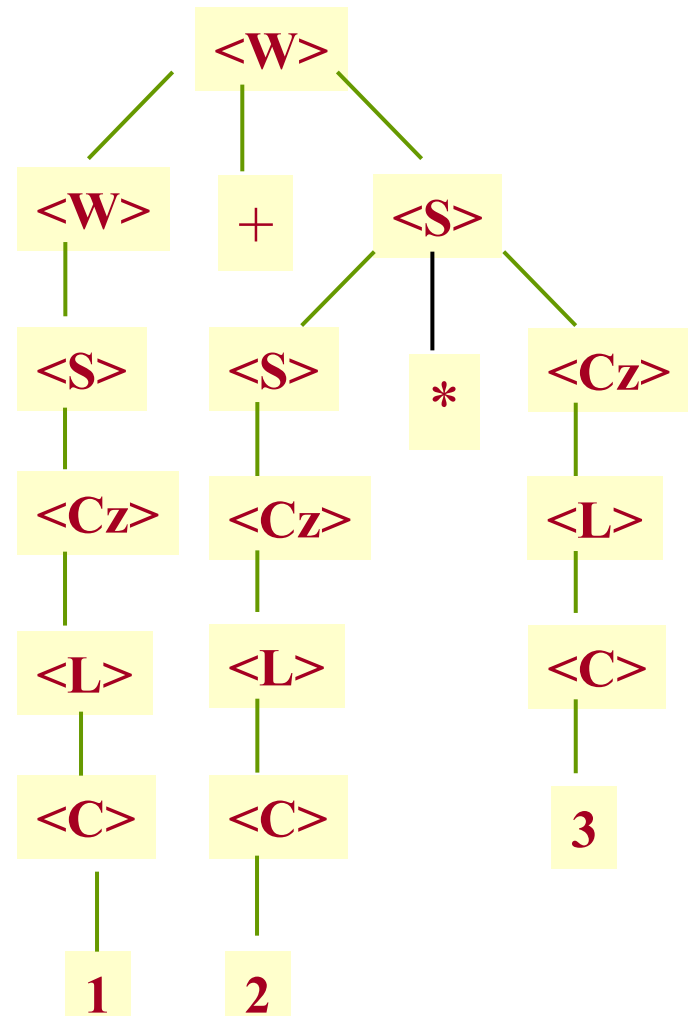
- (1)  $\langle W \rangle \rightarrow \langle W \rangle + \langle S \rangle \mid \langle W \rangle - \langle S \rangle \mid \langle S \rangle$
- (2)  $\langle S \rangle \rightarrow \langle S \rangle * \langle Cz \rangle \mid \langle S \rangle / \langle Cz \rangle \mid \langle Cz \rangle$
- (3)  $\langle Cz \rangle \rightarrow ( \langle W \rangle ) \mid \langle L \rangle$
- (4)  $\langle L \rangle \rightarrow \langle L \rangle \langle C \rangle \mid \langle C \rangle$
- (5)  $\langle C \rangle \rightarrow 0 \mid 1 \mid \dots \mid 9$

**gramatyka jednoznaczna  
wyrażeń arytmetycznych**



**Poprawne drzewo rozbioru**  
1 - 2 + 3

**Poprawne drzewo rozbioru**  
1 + 2 \* 3



Rozróżnienie między wyrażeniami, składnikami i czynnikami wymusza poprawne grupowanie wyrażzeń na różnych poziomach pierwszeństwa działań.

Gramatyki przypominają wyrażenia regularne tym, że obie notacje opisują języki, ale nie definiują bezpośrednio algorytmu określania, czy dany ciąg znaków należy do definiowanego języka. W przypadku wyrażzeń regularnych pokazaliśmy jak można konwertować wyrażenia regularne najpierw na automat niedeterministyczny, a potem na automat deterministyczny. Ten drugi można bezpośrednio implementować jako program.

Analogiczny proces można opisać w przypadku gramatyk. Konwersja gramatyki na automat deterministyczny jest niemożliwa. Istnieje możliwość przekonwertowania gramatyki na program, który podobnie jak automat odczytuje dane wejściowe i określa czy dany ciąg wejściowy należy do języka gramatyki. Najważniejsza z takich technik nosi nazwę rozbiór lewostronny (ang. LR parsing).

# Analiza składniowa i konstrukcja drzew rozbioru

Technika zwana *schodzeniem rekurencyjnym (ang. recursive descent)*, w przypadku której gramatyka jest zastępowana kolekcją wzajemnie rekurencyjnych funkcji, z których każda odpowiada jednej kategorii syntaktycznej gramatyki. Celem działania funkcji  $S$ , która odpowiada kategorii syntaktycznej  $\langle S \rangle$ , jest odczytanie ciągu znaków wejściowych, które tworzą ciąg należący do języka  $L(\langle S \rangle)$  oraz zwrócenie wskaźnika do korzenia drzewa rozbioru tego ciągu.

Część zasadniczą produkcji można traktować jako sekwencję warunków - symboli terminalnych i kategorii syntaktycznych - które muszą zostać spełnione, aby móc określić ciąg znaków występujących w części nagłówkowej produkcji.

Przykładowo można rozpatrzyć gramatykę jednoznaczna wyrażen o zbilansowanej liczbie nawiasów.

(1)  $\langle Z \rangle \rightarrow \varepsilon$

(2)  $\langle Z \rangle \rightarrow (\langle Z \rangle) \langle Z \rangle$

Produkcja (2) określa, że jednym ze sposobów znalezienia ciągu o zbilansowanej liczbie nawiasów jest:

- (1) znalezienie znaku (
- (2) znalezienie ciągu z bilansowanej liczbie nawiasów
- (3) znalezienie ciągu )
- (4) znalezienie kolejnego ciągu o zbilansowanej liczbie nawiasów

Warunek dotyczący symbolu terminalnego jest spełniony wówczas, gdy okaże się, że dany symbol terminalny jest kolejnym symbolem wejściowym, ale warunek nie może zostać spełniony, jeżeli symbolem wejściowym jest pewien inny znak. Aby określić, czy kategoria syntaktyczna określona w części zasadniczej produkcji jest spełniona, należy wywołać funkcję dla tej kategorii syntaktycznej.

Chcemy określić czy sekwencja symboli terminalnych  $X_1X_2\dots X_n$  jest ciągiem należącym do kategorii syntaktycznej  $\langle S \rangle$ , a jeżeli tak, to znaleźć jego drzewo rozbioru. Do pliku wejściowego wstawiane są symbole  $X_1X_2\dots X_n$  ENDM gdzie ENDM to tzw. znacznik końcowy, (ang. end marker). Kursor wejściowy (ang. input cursor) oznacza symbol terminalny który ma zostać poddany przetworzeniu, czyli bieżący symbol terminalny. Jeżeli dane wejściowe stanowią ciąg znaków, to kursor może być wskaźnikiem znaku. Analizę składniową programu rozpoczyna się od wywołania funkcji  $S$  dla początkowej kategorii syntaktycznej  $\langle S \rangle$ , przy kursorze wejściowym wskazującym na początek ciągu wejściowego.

## Inicjalizacja programu wykrywającego $\langle S \rangle$ na wyjściu.

---



=> Gdy rozpatrywana jest część zasadnicza produkcji i dochodzi się w produkcji do symbolu terminalnego **a**, należy wyszukać odpowiadający mu symbol terminalny **a** na pozycji wskazywanej przez kursor wejściowy. Jeżeli zostanie znaleziony symbol **a**, kursor przechodzi do kolejnego symbolu terminalnego danych wejściowych. Jeżeli bieżący symbol wejściowy różni się od **a**, to dopasowanie kończy się niepowodzeniem i nie można znaleźć drzewa rozbioru dla ciągu wejściowego.

=> Gdy rozpatrywana jest część zasadnicza produkcji i dochodzi się w produkcji do kategorii syntaktycznej **<T>**, należy wywołać funkcję **T** dla **<T>**. Jeżeli wywołanie kończy się niepowodzeniem, cały proces analizy składniowej kończy się niepowodzeniem, a ciąg wejściowy można określić jako nie należący do analizowanego języka. Jeżeli wywołanie funkcji **T** kończy się powodzeniem, to następuje przesunięcie kursora wejściowego o pewną ilość znaków (odpowiadających funkcji **T**) oraz funkcja **T** zwraca drzewo, które jest drzewem rozbioru przetwarzanej części ciągu wejściowego.

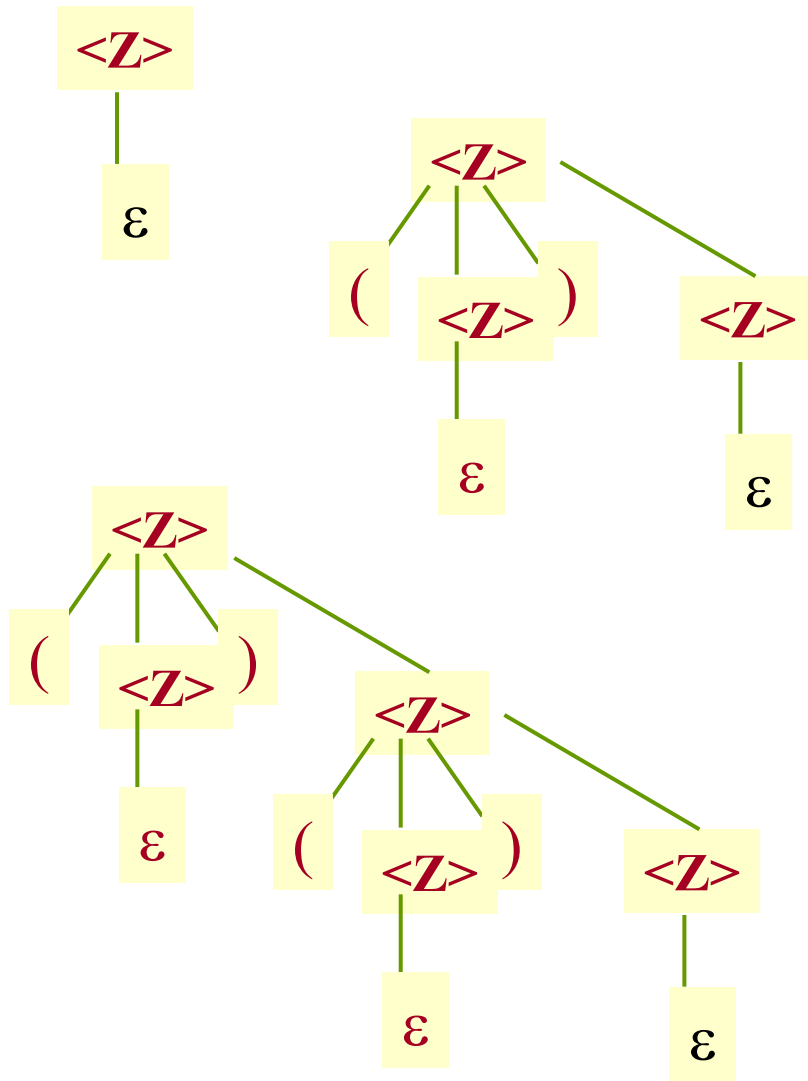
Kiedy opisywane działania zakończą się sukcesem w przypadku każdego symbolu określonego w części zasadniczej produkcji, należy zestawić drzewo rozbioru dla fragmentu danych wejściowych reprezentowanych przez tę produkcję.



Funkcja konstruuująca drzewa rozbioru dla ciągów zawierających zbilansowaną liczbę nawiasów.

```
TREE Z()
{
TREE fusrtZ, secondZ;
if ( *nextTerminal == '(' /*zgodnie z produkcja 2*/{
    nextTerminal++;
    firstZ = Z();
    if(secondZ == FAILED)
        return FAILED;
    else
        return makeNode4('Z',makeNode0('('),firstZ,makeNode0(')'),secondZ);
    }
    else /* pierwsze wołanie Z zakończone niepowodzeniem */
        return FAILED;
}
else /* zgodnie z produkcja 1 */
    return makeNode1('Z', makeNode0('e'));
}
```

Drzewa skonstruowane w wyniku rekurencyjnego wywoływania funkcji.



Jeżeli kategoria syntaktyczna  $\langle S \rangle$  reprezentuje język, którego ciągi znaków należy rozpoznawać i analizować, proces analizy rozpoczyna się od ustawienia kursora wejściowego na pierwszym wejściowym symbolu terminalnym. Wywołanie funkcji  $S$  powoduje skonstruowanie drzewa rozbioru (o ile takie istnieje) dla danych wejściowych i kończy się niepowodzeniem, jeżeli ciąg wejściowy nie należy do języka  $L(\langle S \rangle)$ .

## Konstrukcja analizatorów składniowych schodzenia rekurencyjnego

Kluczowe jest zapewnienie, aby dla każdej kategorii syntaktycznej  $\langle S \rangle$ , jeżeli istnieje więcej niż jedna produkcja posiadająca jako część nagłóvkową  $\langle S \rangle$ , można było tylko przez zbadanie bieżącego symbolu terminalnego (często określanego mianem symbolu antycypowanego (ang. lookahead) określić jedną produkcję  $\langle S \rangle$ , którą należy wypróbować.

Nie jest możliwe stwierdzenie czy dla danej gramatyki istnieje algorytm, który zawsze będzie podejmował poprawne decyzje. Jeżeli dysponuje się strategią którą uważa się za poprawną, to dla każdej kategorii syntaktycznej  $\langle S \rangle$  można zaprojektować funkcje  $S$ , której działanie będzie polegało na:

(1) Zbadaniu symbolu antycypowanego i zdecydowaniu która produkcję wypróbować. Zakładamy, że zasadnicza część wybranej produkcji ma postać  $X_1X_2\dots X_n$ .

(2) Dla  $i=1,2,\dots,n$  wykonaniu następujących działań w przypadku elementu  $X_i$ :

(a) jeżeli  $X_i$  jest symbolem terminalnym, należy sprawdzić, czy jest symbolem antycypowanym. Jeżeli tak, należy przesunąć kursor wejściowy. Jeżeli nie, dane wywołanie funkcji  $S$  kończy się niepowodzeniem.

(b) jeżeli  $X_i$  jest kategorią syntaktyczną, na przykład  $\langle T \rangle$ , należy wywołać funkcje  $T$  odpowiadającą danej kategorii syntaktycznej. Jeżeli  $T$  kończy się niepowodzeniem, to niepowodzeniem kończy się również wywołanie funkcji  $S$ . Jeżeli  $T$  kończy się sukcesem należy zachować zwrócone drzewo.

# Algorytm analizy składniowej oparty na tabeli

Rekurencyjne wywołania funkcji implementuje się zwykle za pomocą stosu zapisów aktywacji. Funkcje analizatora składniowego schodzenia rekurencyjnego wykonują bardzo specyficzne działania. Istnieje możliwość zastąpienia ich pojedynczą funkcją, która bada tabele i samodzielnie manipuluje stosem.

Funkcja  $S$  kategorii syntaktycznej  $\langle S \rangle$  w pierwszej kolejności decyduje o tym, jakiej produkcji należy użyć, a następnie przechodzi sekwencje kolejnych etapów - jeden etap dla każdego symbolu należącego do części zasadniczej wybranej produkcji. Zatem można utworzyć stos symboli gramatycznych który z grubsza będzie przypominał stos zapisów aktywacji. Na stos zostają odkładane zarówno symbole terminalne jak i kategorie syntaktyczne. Kiedy kategoria syntaktyczna  $\langle S \rangle$  znajduje się na szczycie stosu najpierw określa się odpowiednią produkcję. Następnie należy zastąpić  $\langle S \rangle$  częścią zasadniczą produkcji (lewy koniec na szczycie stosu). Kiedy na szczycie stosu mamy symbol terminalny i odpowiada on bieżącemu symbolowi wejściowemu należy zdjąć element ze stosu i przesunąć kursor wejściowy.

# Tabele analizy składniowej

Alternatywą dla pisania zbioru funkcji rekurencyjnych jest skonstruowanie tabeli analizy składniowej (ang. parsing table), której wiersze odpowiadają kategoriom syntaktycznym, zaś kolumny odpowiadają możliwym symbolom antycypowanym. Wartość umieszczona w polu określonym przez wiersz kategorii syntaktycznej  $\langle S \rangle$  oraz kolumnę symbolu antycypowanego  $X$  jest numerem produkcji, której częścią nagłówkową jest  $\langle S \rangle$ , i która musi zostać wykorzystana w celu rozszerzenia  $\langle S \rangle$  w przypadku, gdy symbolem antycypowanym jest  $X$ .

## Przykład:

### Gramatyka

(1)  $\langle Z \rangle \rightarrow \varepsilon$

(2)  $\langle Z \rangle \rightarrow (\langle Z \rangle) \langle Z \rangle$

### Tabela analizy składniowej

	(	)	ENDM
$\langle Z \rangle$	2	1	1

## Gramatyka

- (1)  $\langle I \rangle \rightarrow w c \langle I \rangle$
- (2)  $\langle I \rangle \rightarrow \{ \langle D \rangle$
- (3)  $\langle I \rangle \rightarrow s ;$
- (4)  $\langle D \rangle \rightarrow \langle I \rangle \langle D \rangle$
- (5)  $\langle D \rangle \rightarrow \}$

## Tabela analizy składniowej

	w	c	{	}	s	;	ENDM
$\langle I \rangle$	1		2		3		
$\langle D \rangle$	4		4	5	4		

Postać gramatyki przedstawionej powyżej umożliwia jej analizę składniową za pomocą schodzenia rekurencyjnego lub za pomocą analizy składniowej opartej na tabeli.

$\langle D \rangle$  -kategoria syntaktyczna „dokończenie”.

## Etapy działania analizatora składniowego dla ciągu: {w c s ; s ; } ENDM

STOS	SYMBOL ANTYCYP.	RESZTA
<I>	{	wcs;s;}ENDM
{<D>	{	wcs;s;}ENDM
<D>	w	cs;s;}ENDM
<I><D>	w	cs;s;}ENDM
wc<I><D>	w	cs;s;}ENDM
c<I><D>	c	s;s;}ENDM
<I><D>	s	;s;} ENDM
s;<D>	s	;s;} ENDM
;<D>	;	s;} ENDM
<D>	s	;} ENDM
<I><D>	s	;} ENDM
s;<D>	s	;} ENDM
;<D>	;	} ENDM
<D>	}	ENDM
}	}	ENDM
ε	ENDM	ε

# Konstruowanie drzewa rozbioru

Opisany algorytm określa czy dany ciąg znaków należy do danej kategorii syntaktycznej, ale nie tworzy drzewa rozbioru. Istnieje jednak możliwość wprowadzenia prostej modyfikacji algorytmu, pozwalającej również na utworzenie drzewa rozbioru, kiedy okaże się że ciąg wejściowy należy do kategorii syntaktycznej za pomocą której zainicjalizowano stos.

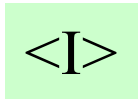
Analizator składniowy schodzenia rekurencyjnego, tworzy drzewo rozbioru wg. konwencji wstępującej (ang. bottom-up), tzn. rozpoczynając od liści i łącząc je w coraz większe poddrzewa w miarę kolejnych powrotów z wywołań funkcji.

W przypadku analizatora składniowego opartego na tabeli odpowiedniejszym sposobem jest budowanie drzewa według konwencji zstępującej (ang. top-down). Oznacza to rozpoczęcie konstrukcji od korzenia i w miarę wybierania kolejnych produkcji, za pomocą których mają być rozszerzane kategorie syntaktyczne na szczycie stosu, jednocześnie tworzy się potomków pewnego wierzchołka należącego do konstruowanego drzewa. Potomkowie ci odpowiadają symbolom należącym do części zasadniczej wybranej produkcji.

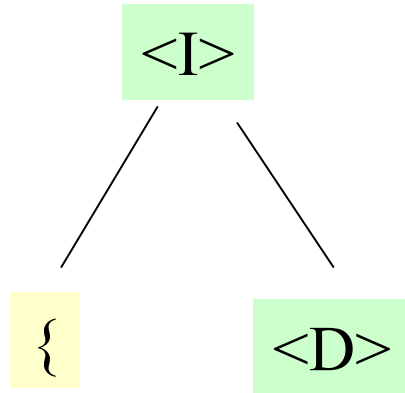


# Etapy konstruowania drzewa rozbioru dla ciągu: {w c s ; s ; } ENDM

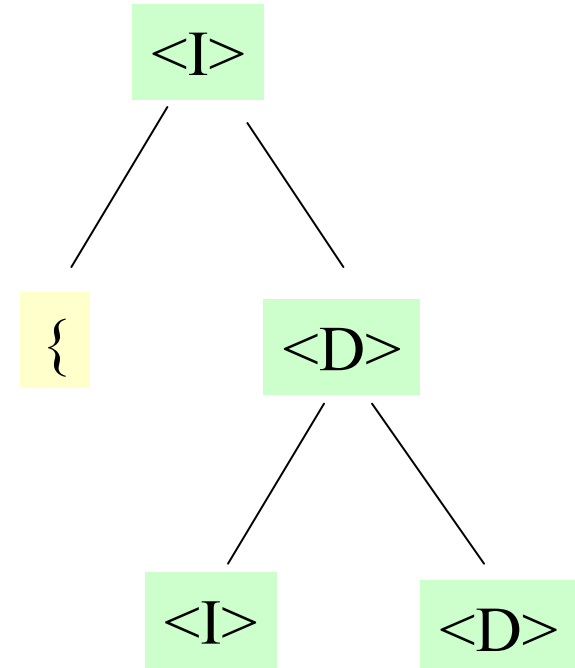
(a)



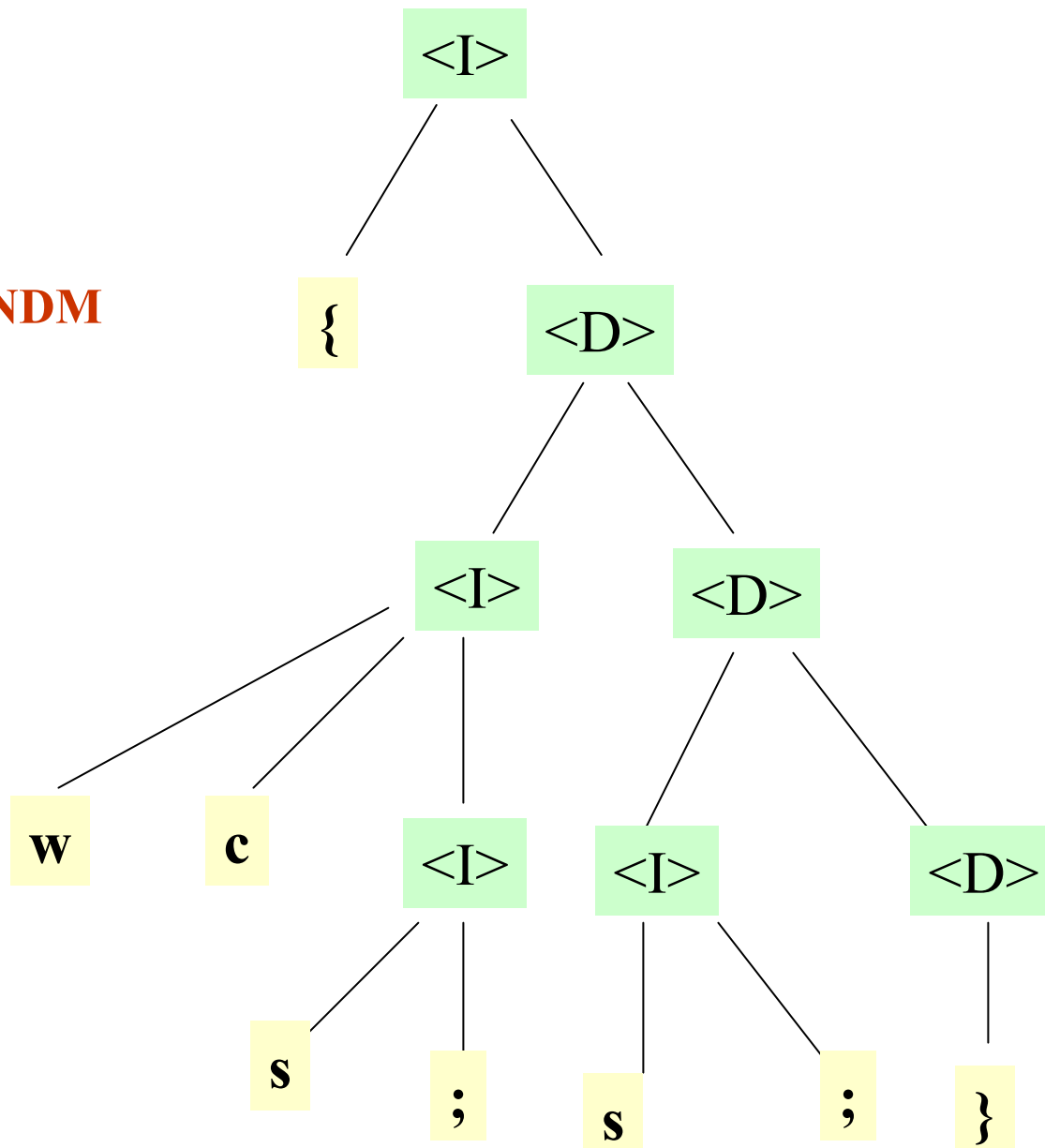
(b)



(c)



**Pełne drzewo rozbioru  
dla analizy składniowej  
dla ciągu: {w c s ; s ; } ENDM**



# Konwertowanie gramatyk na możliwe do analizy składniowej.

Wiele gramatyk wymaga wprowadzenia modyfikacji aby można było przeprowadzając analizę składniową za pomocą metody schodzenia rekurencyjnego lub opartej na tabeli.

a) eliminujemy rekurencje lewostronna, tzn.

$$\langle L \rangle \rightarrow \langle L \rangle \langle I \rangle \mid \varepsilon$$

zamieniamy na

$$\langle L \rangle \rightarrow \langle I \rangle \langle L \rangle \mid \varepsilon$$

b) stosujemy lewostronny rozkład na czynniki, tzn. tworzymy nowe kategorie syntaktyczne  $\langle X \rangle \rightarrow \langle C \rangle \langle D \rangle$  jeżeli dwie produkcje dla  $\langle X \rangle$  zaczynają się wspólnym symbolem  $C$ .

# Gramatyki a wyrażenia regularne

Zarówno gramatyki jak i wyrażenia regularne są notacjami służącymi do opisywania języków.

- => Dotychczas pokazaliśmy że notacja wyrażeń regularnych jest równoważna z dwiema innymi notacjami - automatami deterministycznymi oraz niedeterministycznymi.
- => Gramatyki dają większą możliwość opisu od notacji wyrażeń regularnych. Każdy język możliwy do opisu przez wyrażenia regularne można też opisać przy pomocy gramatyk. Istnieją natomiast języki które można opisać za pomocą gramatyk, ale nie można za pomocą wyrażeń regularnych.

# Symulowanie wyrażeń regularnych za pomocą gramatyk

## Twierdzenie:

Dla każdego wyrażenia regularnego  $R$  istnieje gramatyka, taka, że dla jednej z należących do niej kategorii syntaktycznych  $\langle S \rangle$  zachodzi związek

$$L(\langle S \rangle) = L(R).$$

## Podstawa:

Przypadek podstawowy to  $n=0$ , gdzie wyrażenie regularne  $R$  posiada zero wystąpień operatorów. Wówczas  $R$  jest albo pojedynczym symbolem, np.  $x$ , albo jest  $\varepsilon$  lub  $\emptyset$ . Tworzymy nową kategorię syntaktyczną  $\langle S \rangle$ . Gdy  $R=x$ , tworzymy również produkcje  $\langle S \rangle \rightarrow x$ . Zatem  $L(\langle S \rangle) = \{x\}$ , zaś  $L(R)$  jest tym samym językiem zawierającym jeden ciąg znaków. Jeżeli  $R$  jest równe  $\varepsilon$ , w podobny sposób tworzymy produkcję  $\langle S \rangle \rightarrow \varepsilon$  dla  $\langle S \rangle$ , a jeśli  $R = \emptyset$ , nie tworzymy dla  $\langle S \rangle$  w ogóle żadnej produkcji. Wówczas  $L(\langle S \rangle)$  to  $\{\varepsilon\}$ , kiedy  $R$  jest  $\varepsilon$ , oraz  $L(\langle S \rangle)$  jest  $\emptyset$ , kiedy  $R$  jest  $\emptyset$ .

## Indukcja:

Założmy, że hipoteza indukcyjna jest spełniona w przypadku wyrażeń regularnych o  $n$  lub mniejszej liczbie wystąpień operatorów. Niech  $R$  będzie wyrażeniem regularnym o  $n+1$  wystąpieniach operatorów. Istnieją trzy przypadki, w zależności od tego, czy ostatnim operatorem użytym do skonstruowania wyrażenia regularnego  $R$  jest operator sumy, złożenia czy domknięcia.

Zakładamy, że mamy gramatykę  $G_1$  z kategorią syntaktyczną  $\langle S_1 \rangle$  oraz gramatykę  $G_2$  z kategorią syntaktyczną  $\langle S_2 \rangle$ , takie, że  $L(\langle S_1 \rangle) = L(R_1)$  oraz  $L(\langle S_2 \rangle) = L(R_2)$ .

(1)  $R = R_1 \mid R_2$ . Tworzymy nową kategorię syntaktyczną  $\langle S \rangle$  oraz do produkcji dodajemy  $\langle S \rangle \rightarrow \langle S_1 \rangle \mid \langle S_2 \rangle$ . Wówczas  $L(\langle S \rangle) = L(R_1) \cup L(R_2) = L(R)$ .

(2)  $R = R_1 R_2$ . Tworzymy nową kategorię syntaktyczną  $\langle S \rangle$  oraz do produkcji dodajemy  $\langle S \rangle \rightarrow \langle S_1 \rangle \langle S_2 \rangle$ . Wówczas  $L(\langle S \rangle) = L(R_1) \cup L(R_2) = L(R)$ .

(3)  $R = R_1^*$ . Tworzymy nową kategorię syntaktyczną  $\langle S \rangle$  oraz do produkcji dodajemy  $\langle S \rangle \rightarrow \langle S_1 \rangle \langle S \rangle \mid \varepsilon$ . Wówczas  $L(\langle S \rangle) = L(\langle S_1 \rangle)^*$ , ponieważ  $\langle S \rangle$  generuje ciągi znaków zawierające zero lub więcej kategorii  $\langle S_1 \rangle$ .

# Konstrukcja gramatyki dla wyrażenia regularnego: $a \mid bc^*$

1. Tworzymy kategorie syntaktyczne dla trzech symboli, które pojawiają się w tym wyrażeniu:

$\langle A \rangle \rightarrow a$        $\langle B \rangle \rightarrow b$        $\langle C \rangle \rightarrow c$

2. Tworzymy gramatykę dla  $c^*$ :  $\langle D \rangle \rightarrow \langle C \rangle \langle D \rangle \mid \varepsilon$   
Wówczas  $L(\langle D \rangle) = L(\langle C \rangle)^* = c^*$

3. Tworzymy gramatykę dla  $bc^*$ :  $\langle E \rangle \rightarrow \langle B \rangle \langle D \rangle$

4. Tworzymy gramatykę dla całego wyrażenia regularnego  $a \mid bc^*$ :  
 $\langle F \rangle \rightarrow \langle A \rangle \mid \langle E \rangle$

**końcowa postać gramatyki**

$\langle F \rangle \rightarrow \langle A \rangle \mid \langle E \rangle$   
 $\langle E \rangle \rightarrow \langle B \rangle \langle D \rangle$   
 $\langle D \rangle \rightarrow \langle C \rangle \langle D \rangle \mid \varepsilon$   
 $\langle A \rangle \rightarrow a$   
 $\langle B \rangle \rightarrow b$   
 $\langle C \rangle \rightarrow c$

# Język posiadający gramatykę ale nie posiadający wyrażenia regularnego

Język  $E$  będzie zbiorem znaków składających się z jednego lub większej liczby symboli 0, po których występuje ta sama liczba symboli 1, to znaczy:

$$E = \{ 01, 0011, 000111, \dots \}$$

W celu opisanego ciągów znaków języka  $E$  można użyć przydatnej notacji opartej na wykładnikach. Niech  $s^n$ , gdzie  $s$  jest ciągiem znaków, zaś  $n$  liczba całkowita, oznacza  $ss\dots s$  ( $n$  razy), to znaczy  $s$  złożone ze sobą  $n$  razy.

Wówczas:

$$E = \{0^11^1, 0^21^2, 0^31^3, \dots\} \quad \text{lub} \quad E = \{0^n1^n \mid n \geq 1\}$$

Język  $E$  można zapisać za pomocą gramatyki:

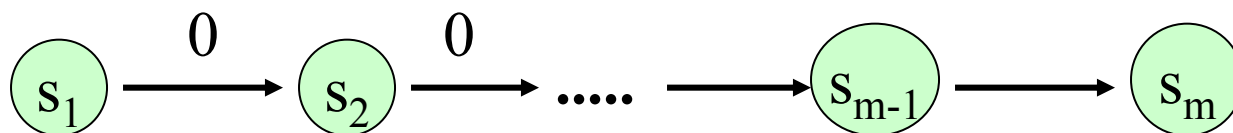
$$\begin{aligned} \langle S \rangle &\rightarrow 0 \langle S \rangle 1 \\ \langle S \rangle &\rightarrow 0 1 \end{aligned}$$

Ponieważ nie istnieją żadne inne ciągi znaków możliwe do utworzenia na podstawie tych dwóch produkcji,  $E = L(\langle S \rangle)$ .



## Dowód niedefiniowalności języka E za pomocą dowolnego wyrażenia regularnego.

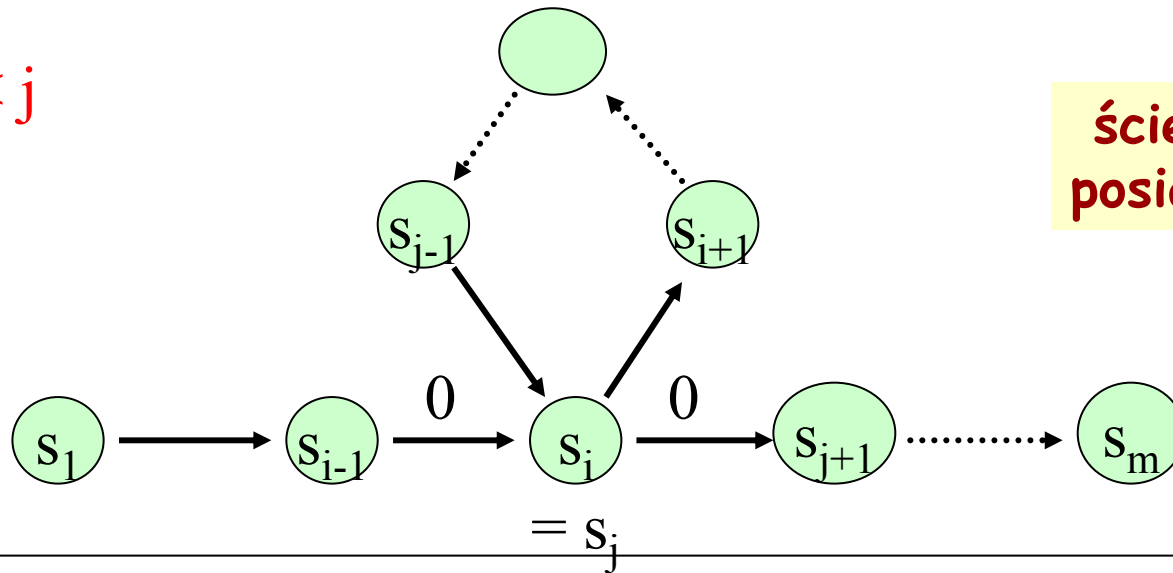
Dowód niedefiniowalności języka E poprzez wyrażenie regularne przeprowadza się wykazując brak możliwości zbudowania deterministycznego automatu skończonego (a więc brak wyrażenia regularnego który można by zamienić na równoważny automat deterministyczny).



**wprowadzanie  
symboli 0 do  
automatu A**

Założmy, że język E jest językiem pewnego deterministycznego automatu skończonego A. Wówczas A posiada pewną liczbę stanów, np. m. Jeżeli automat A otrzymuje na wejściu ciąg 000..... Stan początkowy nieznanego automatu nosi nazwę  $s_0$ . Automat musi posiadać przejście z  $s_0$  do  $s_1$  itd.itd. Odczytanie stanu i powoduje przejście do  $s_i$ . Automat jednak nie może pamiętać ile symboli 0 odczytano, niektóre stany mogą być zdegenerowane (pętle).

$i < j$



ścieżka musi posiadać pętle

Automat nie może „pamiętać” ile symboli 0 odczytano. Jeżeli znajduje się w stanie  $s_m$ , odczytanych mogło zostać dokładnie  $m$  symboli 0. Wówczas musi być prawda, że jeżeli rozpocznie się w stanie  $m$  i wprowadzi do automatu  $A$  dokładnie  $m$  symboli 1, to automat  $A$  dochodzi do stanu akceptującego.

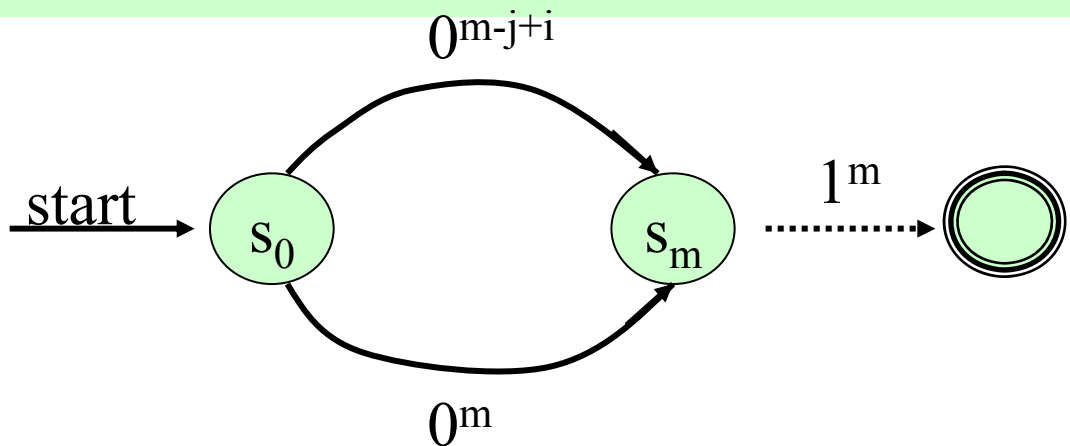
Założmy, że do automatu wprowadzono  $m-j+i$  symboli 0.  $i$  symboli 0 powoduje przejście automatu  $A$  ze stanu  $s_0$  do stanu  $s_i$ , który jest tym samym stanem, co  $s_j$ . Widać, że  $m-j$  symboli 0 powoduje przejście automatu  $A$  ze stanu  $s_j$  do stanu  $s_m$ . Zatem  $m-j+i$  symboli 0 powoduje przejście automatu  $A$  ze stanu  $s_0$  do stanu  $s_m$ .

Stąd  $m-j+i$  symboli 0, po których następuje  $m$  symboli 1 powoduje przejście automatu  $A$  ze stanu  $s_0$  do stanu akceptującego. Innymi słowy, ciąg znaków  $0^{m-j+i} 1^m$  należy do języka automatu  $A$ . Jednak z uwagi na fakt że  $j > i$ , ciąg ten posiada więcej symboli 1 niż symboli 0 i nie należy do języka  $E$ . Można wyciągnąć wniosek, że język automatu  $A$  nie odpowiada dokładnie językowi  $E$ , co jest sprzeczne z przyjętym założeniem.

Wystartowaliśmy z założenia że język  $E$  posiada deterministyczny automat skończony i doszliśmy do sprzeczności. A więc założenie było fałszywe. Stąd wnioskujemy że język  $E$  nie posiada też wyrażenia regularnego.

Język  $E = \{0^n 1^n \mid n \geq 1\}$  to jeden z przykładów nieskończonej liczby języków, które można określić za pomocą gramatyki, ale nie można za pomocą wyrażenia regularnego.

**Automat  $A$  nie ma możliwości stwierdzenia, czy odczytano  $m$  symboli 0 lub  $m-j+i$  symboli 0**



# Tylko definicje.....

Gramatyka jest prawostronnie liniowa, jeżeli każda produkcja ma postać:  $A \rightarrow w B$  lub  $A \rightarrow w$ .

Gramatyka jest lewostronnie liniowa, jeżeli każda produkcja ma postać:  $A \rightarrow Bw$  lub  $A \rightarrow w$ .

Gramatyka która jest lewostronnie liniowa, lub prawostronnie liniowa to gramatyka regularna.

Gramatyka nieograniczona to taka, która dopuszcza produkcje o postaci  $\alpha \rightarrow \beta$ , gdzie  $\alpha, \beta$  są dowolnymi łańcuchami symboli tej gramatyki, przy czym  $\alpha \neq \beta$ .

Gramatyka kontekstowa to taka gramatyka nieograniczona, dla której  $\beta$  jest co najmniej tak długie jak  $\alpha$ .

## Posumowanie

- ⇒ Gramatyka bezkontekstowa wykorzystuje model funkcji rekurencyjnych.
- ⇒ Jednym z ważnych zastosowań gramatyk są specyfikacje języków programowania. Gramatyki stanowią zwięzłą notację opisu ich składni.
- ⇒ Drzewa analizy składniowej (drzewa rozbioru) stanowią formę reprezentacji, która przedstawia strukturę ciągu znaków zgodną z daną gramatyką.
- ⇒ Niejednoznaczność - to problem, który pojawia się w sytuacji gdy ciąg znaków posiada dwa lub więcej odrębnych drzew analizy składniowej, przez co nie posiada unikatowej struktury zgodnie z daną gramatyką

## Posumowanie

⇒ Metoda zamiany gramatyki na analizator składniowy to algorytm pozwalający stwierdzić, czy dany ciąg znaków należy do pewnego języka.

⇒ Gramatyki posiadają większe możliwości w zakresie opisu języków niż wyrażenia regularne. Gramatyki oferują co najmniej tak samo duże możliwości opisu języków, jak wyrażenia regularne przez przedstawienie sposobu symulowania wyrażeń regularnych za pomocą gramatyk. Istnieją jednakże języki które można wyrazić za pomocą gramatyk, ale nie można za pomocą wyrażeń regularnych.