

# **Formalni jezici i jezični procesori I**

## **PARSERI**

prof. dr. sc. Sanda Martinčić - Ipšić

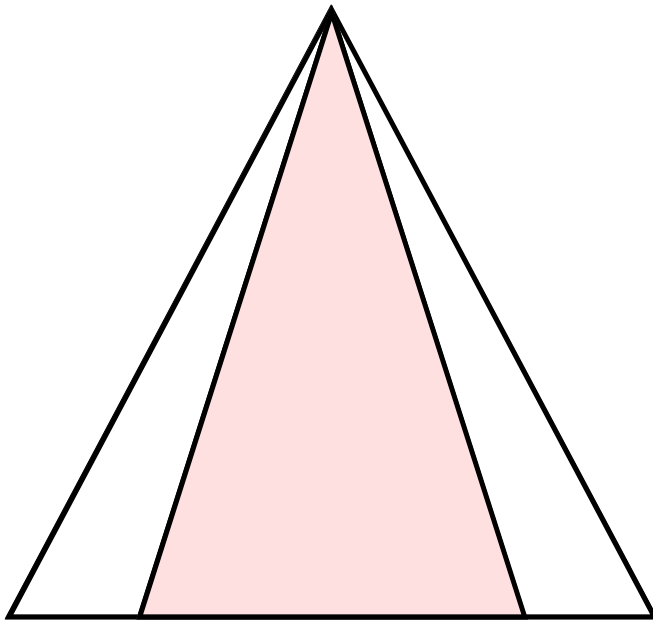
smart@uniri.hr

# Parsiranje

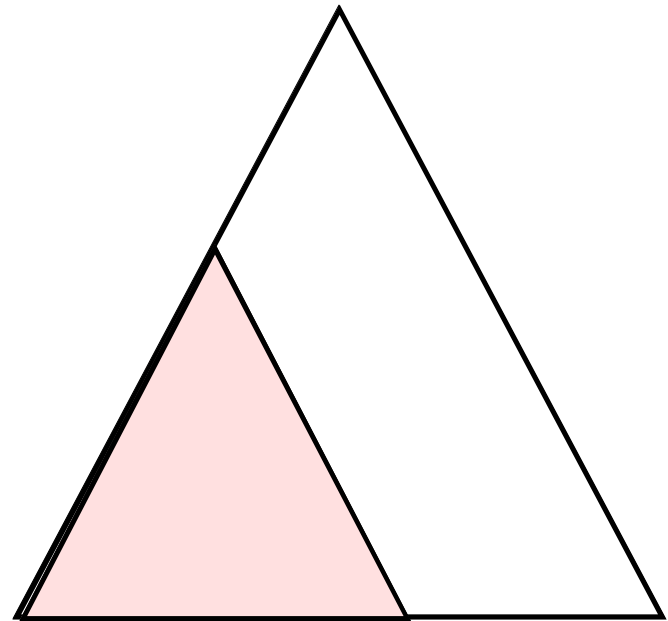
- postupak prepoznavanja niza i gradnja generativnog stabla na temelju zadanih produkcija kontekstno neovisne gramatike
- želimo odrediti je li ulazni niz (program) generirala gramatika (je li u skladu sa sintaksnim pravilima zapisanima u gramatici)

# Parsiranje Od-vrha-prema-dnu vs. Od-dna-prema-vrhu

- veličina stabla parsiranja



**Top-down**



**Bottom-up**

# Parsiranje niza

- prepoznavanje niza – određivanje
  - pripada li niz w jeziku  $L(G)$ ?
  - generira li gramatika  $G$  niz  $w$ ?
  - $w$  – C program, SQL query, XML dokument, ...
- **generativno stablo** (stablo parsiranja) koristi se za interpretaciju niza (u listovima isključivo završni znakovi)
- **parsiranje** – prepoznavanje niza + gradnja generativnog stabla (za niz  $w$  i gramatiku  $G$ )
  - od vrha prema dnu
  - od dna prema vrhu

# Generativno stablo za gramatiku $G=(V, T, P, S)$

- **čvorovi** stabla:  $V \cup T \cup \{\varepsilon\}$
- **korijen** stabla:  $S$  – početni nezavršni znak
- **unutrašnji čvorovi**:  $A \in V$  (primjena produkcija  $P$  iz  $G$ )
  - čvor  $n$  ima djecu:  $n_1, n_2, \dots, n_k$  i čvor  $n$  je označen s  $A$  a čvorovi  $n_1, n_2, \dots, n_k$  sa  $X_1, X_2, \dots, X_k$  onda je  $A \rightarrow X_1 X_2 \dots X_k$  produkcija iz  $P$
- $\varepsilon$ : isključivo označava listove stabla (to je jedino dijete svojih roditelja – unutrašnjih čvorova)
- listovi stabla:  $T \cup \{\varepsilon\}$ 
  - ako ih čitamo s lijeva na desno dobijemo niz  $w$  iz  $L(G)$
  - produkcije  $P$  se primjenjuju dok u čvorovima nisu isključivo završni znakovi niza  $w$

# Parsiranje od vrha prema dnu

- početak: korijen je S – početni nezavršni znak
- primijenimo produkcije P za gradnju ostalih čvorova
- postupak ponavljamo sve dok listovi nisu označeni isključivo završnim znakovima
- učinkovito i jednostavno za programiranje
- LL(k) parsiranje-
  - 1. L : od lijeva na desno čita se niz
  - 2. L : produkcije se primjenjuju na krajnje lijevi nezavršni znak u generiranom međunizu
  - k : odluka o primjeni produkcije donosi se na osnovu k pročitanih znakova;  $k \geq 1$

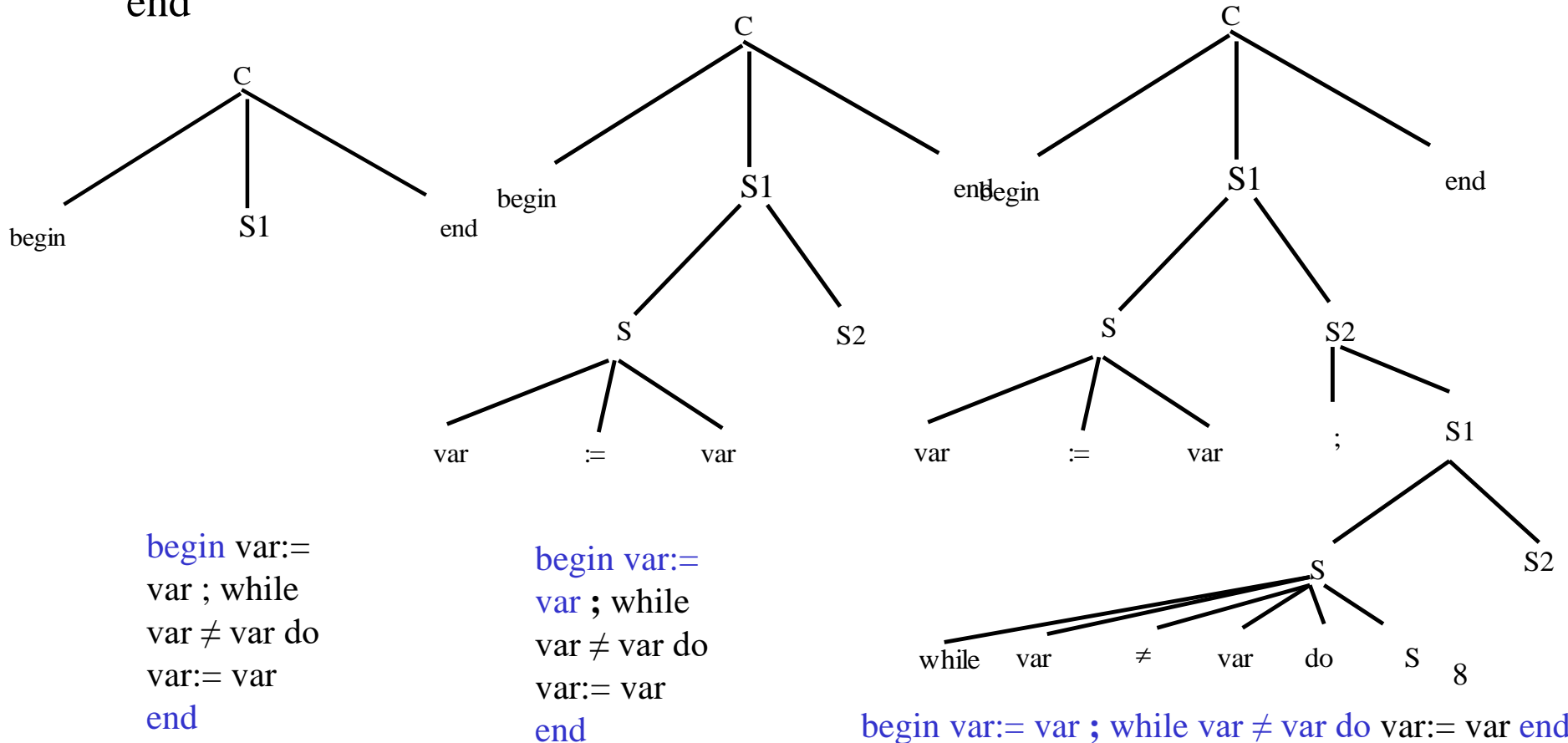
# Jednoznačná gramatika

- $G = (\{C, S, S1, S2\}, \{\text{begin, end, while, do, }, :=, \neq, \text{var}\}, P, C)$
- $P = \{ C \rightarrow \text{begin } S1 \text{ end}$   
     $S1 \rightarrow S \ S2$   
     $S2 \rightarrow ; \ S1 \mid \varepsilon$   
     $S \rightarrow \text{var} := \text{var} \mid \text{while } \text{var} \neq \text{var} \text{ do } S \mid \text{begin } S1 \text{ end } \}$
- **niz w:**  
    begin var:= var;  
        while var  $\neq$  var do  
            var:= var  
        end  
    end

# Generativno stablo

**niz:**begin var:= var;  
 while var ≠ var do  
 var:= var  
 end

$\mathbf{P} = \{ C \rightarrow \text{begin } S1 \text{ end}$   
 $S1 \rightarrow S \ S2$   
 $S2 \rightarrow ; \ S1 \mid \varepsilon$   
 $S \rightarrow \text{var} := \text{var} \mid \text{while } \text{var} \neq \text{var} \text{ do } S \mid \text{begin } S1 \text{ end} \}$

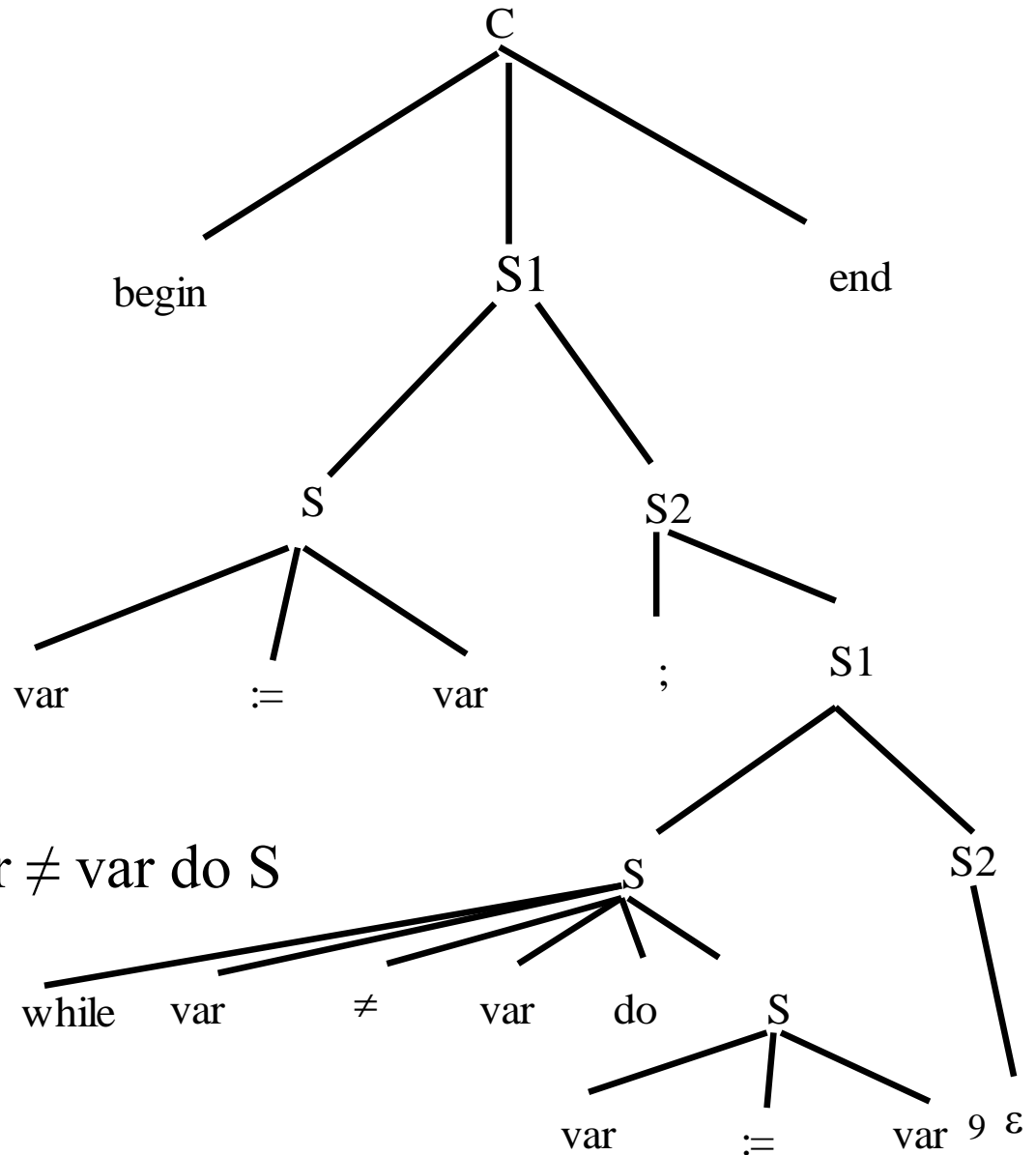




# Generativno stablo

```
begin var:= var;
  while var ≠ var
do
  var:= var
end
```

$C \rightarrow \text{begin } S1 \text{ end}$   
 $S1 \rightarrow S \ S2$   
 $S2 \rightarrow ; \ S1 \mid \varepsilon$   
 $S \rightarrow \text{var } := \text{var} \mid \text{while var } \neq \text{var do } S$   
 $\mid \text{begin } S1 \text{ end}$



# LL parser i LL gramatika

- prvi L: ulazni niz se čita s lijeva na desno (left-to-right)
- drugi L: produkcije se primjenjuju na krajnje lijevi nezavršni znak u generiranom međunizu (leftmost derivation)
- LL(1) – parsira se čitanjem 1 znaka u nizu
- LL(k) – parsira se čitanjem k znakova u nizu

# Primjer izgradnje parsera za gramatiku

$S \rightarrow NP VP$

$Det \rightarrow that \mid this \mid a$

$S \rightarrow Aux NP VP$

$Noun \rightarrow book \mid flight \mid meal \mid money$

$S \rightarrow VP$

$Verb \rightarrow book \mid include \mid prefer$

$NP \rightarrow Det Nominal$

$Aux \rightarrow does$

$Nominal \rightarrow Noun$

$Proper-Noun \rightarrow Houston \mid TWA$

$Nominal \rightarrow Noun Nominal$

$Prep \rightarrow from \mid to \mid on$

$NP \rightarrow Proper-Noun$

$VP \rightarrow Verb$

$VP \rightarrow Verb NP$

$Nominal \rightarrow Nominal PP$

*Book that flight.*

# Generativno stablo

*Book that flight*

$S \rightarrow VP$

$VP \rightarrow \text{Verb NP}$

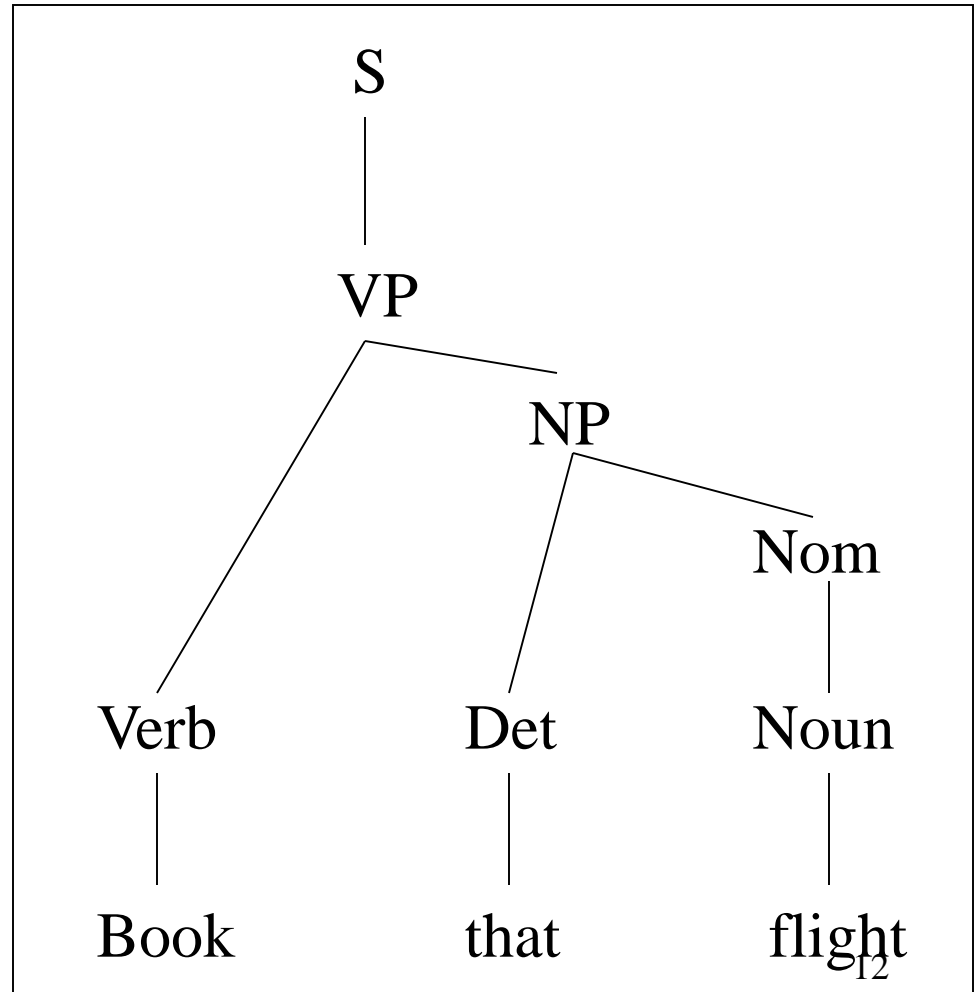
$NP \rightarrow \text{Det Nominal}$

$\text{Nominal} \rightarrow \text{Noun}$

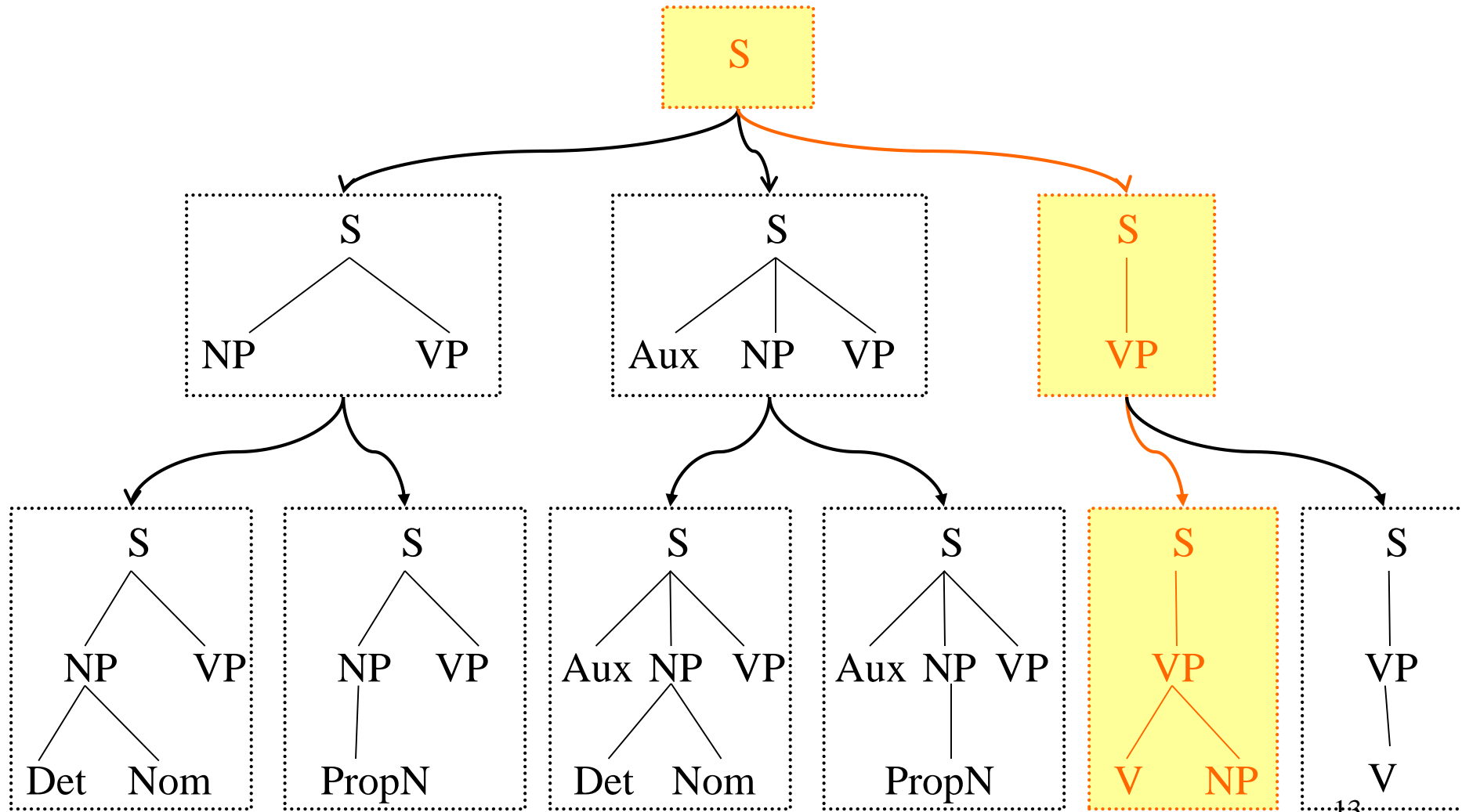
$\text{Verb} \rightarrow \text{book}$

$\text{Det} \rightarrow \text{that}$

$\text{Noun} \rightarrow \text{flight}$



# Primjer parsiranja od vrha prema dnu



# Rekurzivni spust I

- uporaba rekurzije za parsiranje od vrha prema dnu
- nezavršnim znakovima pridružuju se potprogrami koji ispituju da li desna strana produkcije odgovara pročitanoj nizu
- završni znakovi desnih strana produkcija uspoređuju se sa nizom
- $\perp$  - oznaka za kraj niza

## *Programska realizacija algoritma:*

1. u glavnom programu pročita se prvi znak niza  $w$  i pozove se potprogram za početni nezavršni znak gramatike
  - ako se pri izvođenju pročita oznaka za kraj niza  $\perp$  niz se prihvata u suprotnom se ne prihvata

# Rekurzivni spust II

```
Glavni program( )  
  {  
    ulaz = prvi znak iz niza w;  
    potprogram pridružen početnim  
      nezavršnom znaku S;  
    ako (ulaz != ⊥)  
      ispiši ("w ∉ L(G)");  
    inače  
      ispiši ("w ∈ L(G)");  
  }
```

# Rekurzivni spust III

2. za svaku produkciju svakog nezavršnog znaka gradi se poseban potprogram  $A \rightarrow a_1\beta_1 \mid a_2\beta_2 \mid \dots \mid a_n\beta_n$

```
A ( )  
{  
  slučaj (ulaz)  
  {  
    a1:  
      ispitaј slijedeće znakove prema β1;  
    a2:  
      ispitaј slijedeće znakove prema β2;  
    ....  
    an:  
      ispitaј slijedeće znakove prema βn;  
    svi ostali znakovi:  
      ispiši ("w ∉ L(G)");  
  }  
}
```



# Rekurzivni spust IV

3. za desne stane produkcija  $\beta_i : A \rightarrow a_1\beta_1 \mid a_2\beta_2 \mid \dots \mid a_n\beta_n$
- za bilo koji **završni znak**  $b$  na desnoj stani produkcije  $A \rightarrow a\alpha by$  koji nije na krajnje lijevom mjestu

$A( )$

```
{  
    ...  
    ulaz = slijedeći znak iz ulaznog niza w;  
    ako (ulaz != b)  
        ispiši ("w ∉ L(G)");  
    ... }
```

- za bilo koji **nezavršni znak**  $B$  na desnoj stani produkcije  $A \rightarrow a\alpha B\gamma$  koji nije na krajnje lijevom mjestu

$A( )$

```
{  
    ...  
    ulaz = slijedeći znak iz ulaznog niza w;  
    B( );  
    ...  
}
```

# Rekurzivni spust V

- ako je nezavršni znak B na na krajnje lijevom mjestu produkcije  $A \rightarrow Ba\alpha\gamma$

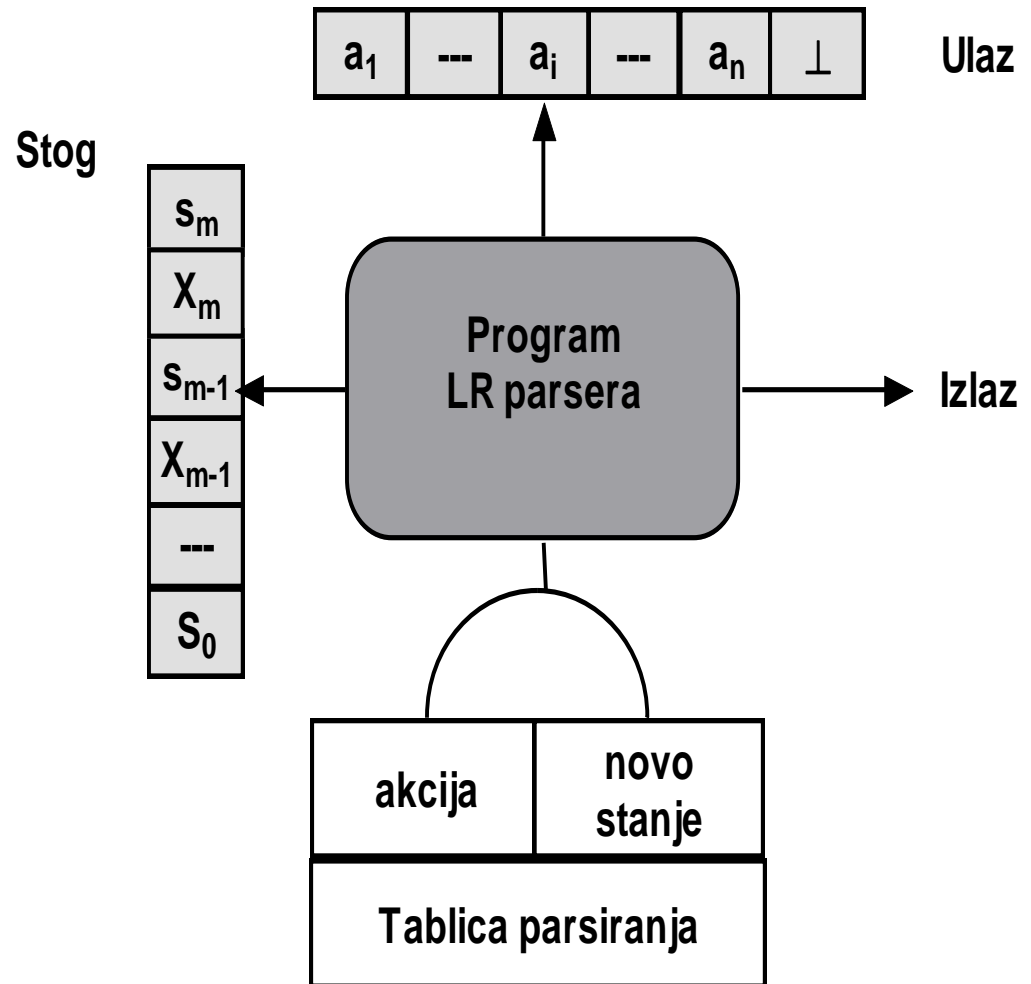
```
A ( )  
{  
    ...  
    B ( ) ;  
    ...  
}
```

# Parsiranje od dna prema vrhu

- gradnja započinje od listova, završnih znakova gramatike
- ako je međuniz niza  $w$  jednak desnoj strani produkcije onda se zamijeni lijevom
- na taj način gradi se stablo uključujući i korijen
- broj znakova se postepeno smanjuje prema vrhu zato produkcije nazivamo redukcijama
- metoda primjerena za generatore parsera
- LR(k) parsiranje –
  - L niz čitamo s lijeva na desno
  - R stablo gradimo s desna na lijevo
  - $k$  : najviše  $k$  znakova treba pročitati da se primjeni redukcija;  
 $k \geq 1$

# Model LR parsera

- LR parser gradi se posebnim generatorom
- LR parseri različitih gramatika razlikuju se samo u **tablici parsiranja**
- u radu koristi potisni stog (LIFO)



# Rad LR parsera I

- tablica parsiranja: tablica **Akcija** i tablica **Novo stanje**
- čita se **vrh stoga**  $s_m$  i **znak na ulazu**  $a_i$   $[s_m, a_i]$
- na temelju  $[s_m, a_i]$  se odredi akcija iz tablice:
  - Pomak (odredi stanje),
  - Reduciraj  $A \rightarrow \beta$  (primjeni produkciju gramatike),
  - Prihvati,
  - Odbaci
- **konfiguracija LR parsera** se dobiva **sadržajem stoga i nepročitanim dijelom ulaznoga niza** ( $X_1$  su znakovi gramatike)
  - sadržaj stoga
  - nepročitani dio ulaznoga niza

$(s_0 X_1 s_1 X_2 s_2 \dots X_m s_m$

$a_i a_{i+1} \dots a_n \perp)$

# Rad LR parsera II

- **Akcija  $[s_m, a_i]$ =Pomak s** promjeni konfiguraciju LR parsera u  $(s_0 X_1 s_1 X_2 s_2 \dots X_m s_m a_i s, a_{i+1} \dots a_n \perp)$
- **Akcija  $[s_m, a_i]$ = Reduciraj  $A \rightarrow \beta$** 
  - s vrha stoga uzme se  $2r$  znakova (ako je  $r$  broj znakova u  $\beta$ )
  - na vrhu stoga ostaje stanje  $s_{m-r}$ , i na stog se stavi  $A$  (nezavršni znak produkcije) a ulazni niz ostaje nepromijenjen
  - promjeni konfiguraciju LR parsera u

$$(s_0 X_1 s_1 X_2 s_2 \dots X_{m-r} s_{m-r} A s, a_i a_{i+1} \dots a_n \perp)$$

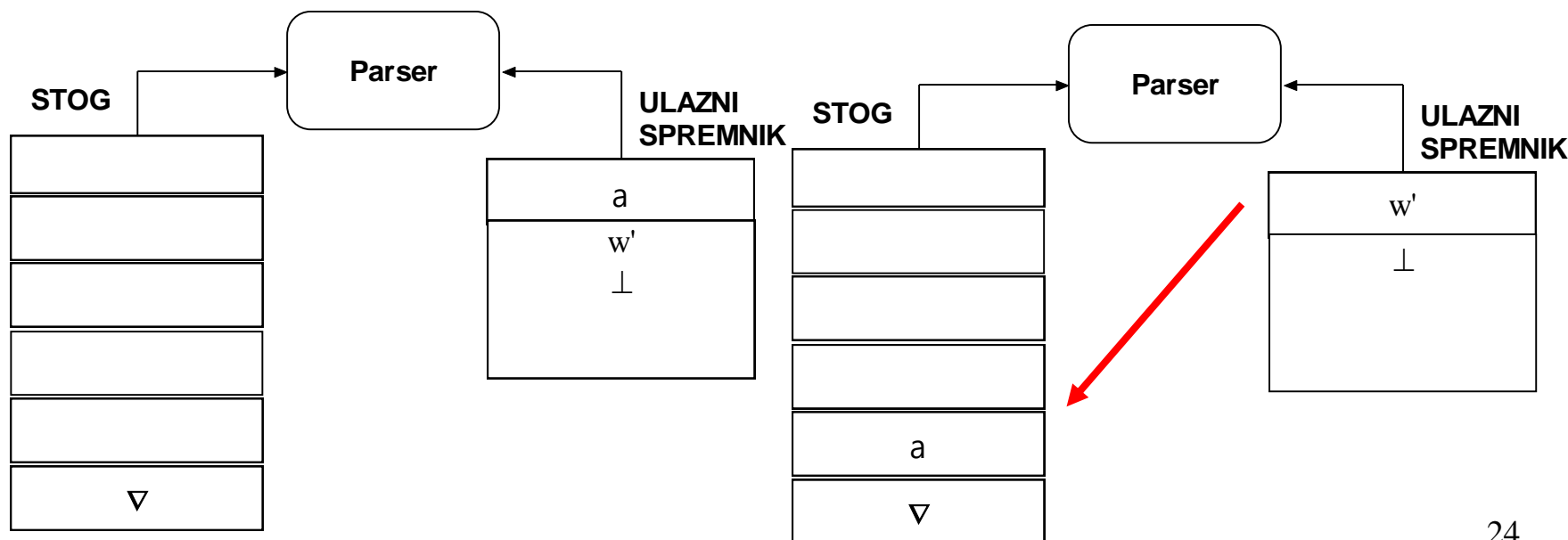
- **Akcija  $[s_m, a_i]$ = Prihvati**
  - prihvća se niz  $w$ , zaustavlja parsiranje
- **Akcija  $[s_m, a_i]$ = Odbaci**
  - odbacuje se niz  $w$ , zaustavlja parsiranje

# Rad bottom-up parsera

- koristi se potisni stog i ulazni spremnik
  - na jednom koraku se čita više znakova stoga (ne samo vrh)
  - na stogu su završni i nezavršni nizovi
  - u ulaznom spremniku je niz  $w$  i svi međunizovi koji se parsiraju
    - oznaka dna stoga:  $\Delta$
    - oznaka kraja niza:  $\perp$
- akcije:
  - *Pomakni; Reduciraj; Prihvati; Odbaci*

# Rad bottom-up parsera: *Pomakni*

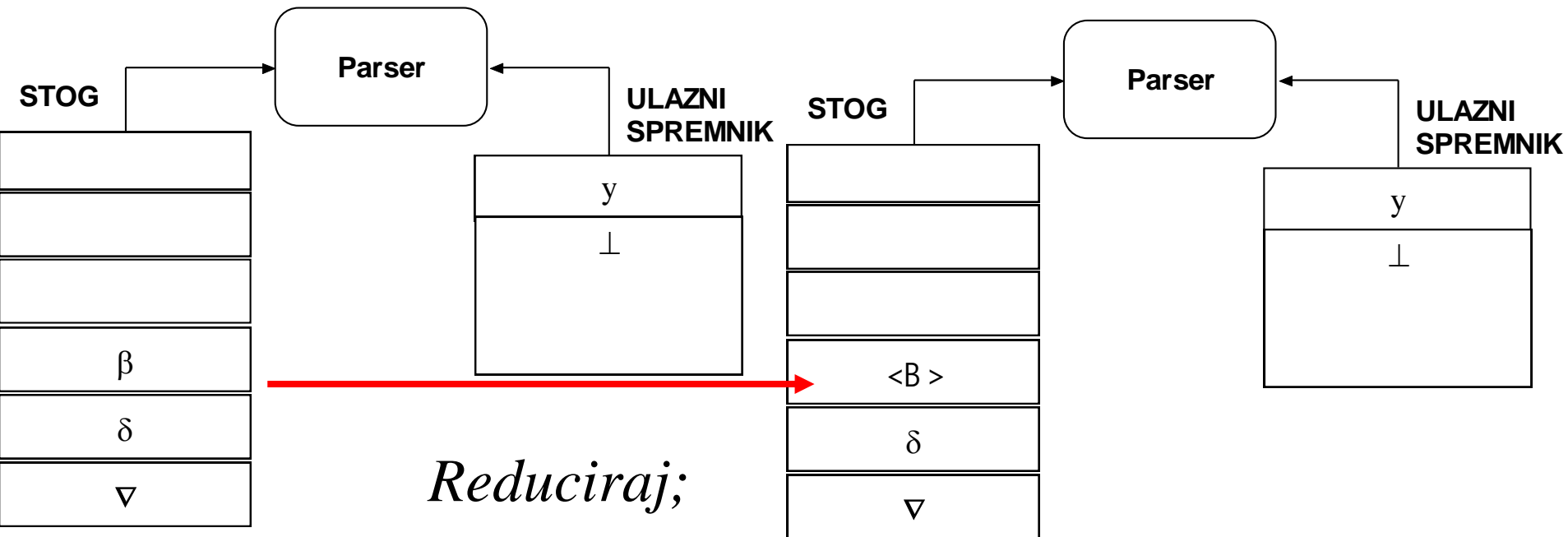
- na početku je *stog prazan*
- čita se prvi znak ulaznog niza
- *Pomakni*; prvi znak ulaznog niza se zapisuje na vrh stoga
- kazaljka na ulaznom nizu se postavlja za jedan znak desno
  - sljedeći znak





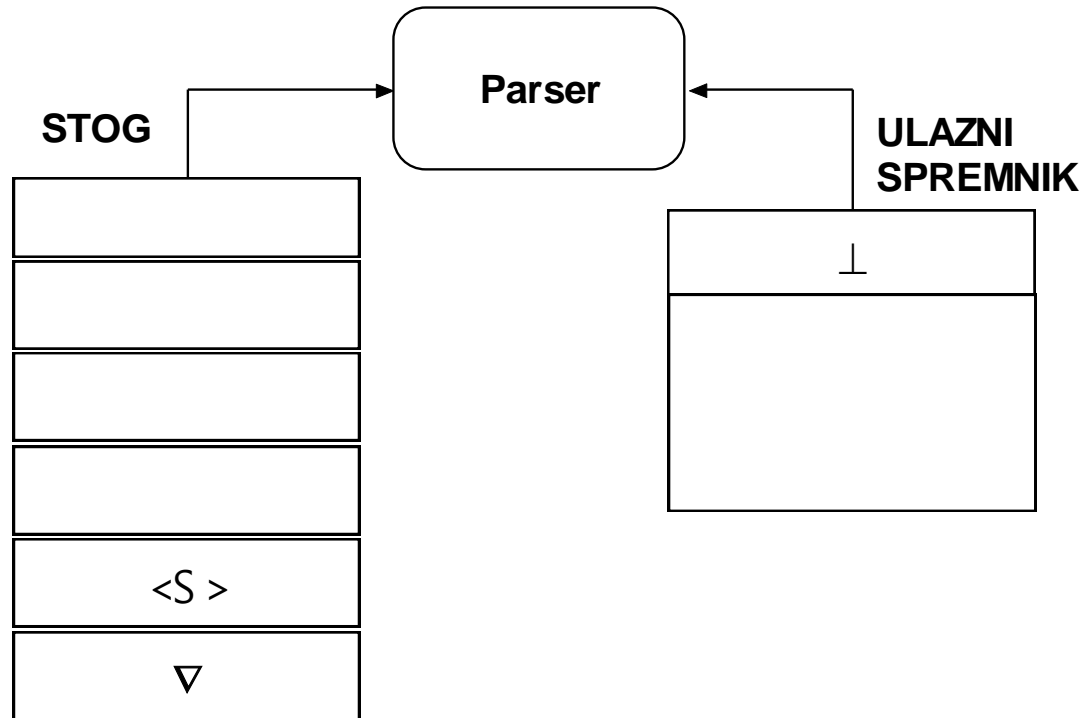
# Rad bottom-up parsera: *Reduciraj*;

- na vrhu stoga je  $\beta$  i postoji produkcija  $\langle B \rangle \rightarrow \beta$
- na vrhu stoga se zamijeni  $\beta$  s nezavršnim znakom  $\langle B \rangle$
- međuniz  $\delta\beta\gamma \leq \delta\langle B \rangle\gamma$



# Rad bottom-up parsera: *Prihvati;*

- Pročitani svi znakovi ulaznog niza: do  $\perp$
- na stogu je početni nezavršni znak  $\langle S \rangle$  i oznaka kraja  $\Delta$



# Primjer: rad parsera

gramatika: 1.  $E \rightarrow E+T$  2.  $E \rightarrow T$  3.  $T \rightarrow T * F$  4.  $T \rightarrow F$  5.  $F \rightarrow (E)$  6.  $F \rightarrow \text{var}$

ulazni niz: **var+var\*var**

## Tablica parsera:

akcije **si**: stavi stanje si na stog **rj**: reduciraj produkcijom j **prihvati**: odbaci: sve  
prazne ćelije **početno stanje 0**

	<i>Akcija</i>						<i>NovoStanje</i>		
	var	+	*	(	)	$\perp$	E	T	F
<b>0</b>	s5			s4			1	2	3
<b>1</b>		s6				prihvati			
<b>2</b>		r2	s7		r2	r2			
<b>3</b>		<b>r4</b>	<b>r4</b>		r4	r4			
<b>4</b>	s5			s4			8	2	3
<b>5</b>		r6	r6		r6	r6			
<b>6</b>	s5			s4				9	3
<b>7</b>	s5			s4					10
<b>8</b>		s6			s11				
<b>9</b>		r1	s7		r1	r1			
<b>10</b>		r3	r3		r3	r3			
<b>11</b>		r5	r5		r5	r5			

# Primjer: rad parsera II

Ulazni niz	Stog	Akcija
var+var*var⊥	$\Delta 0$	s5
+var*var⊥	$\Delta 0 \text{var} 5$	r6: redukcija $F \rightarrow \text{var}$ , skidam 2 znaka, stanje 0, F=3
+var*var⊥	$\Delta 0 F 3$	r4: redukcija $T \rightarrow F$ , skidam 2 znaka, stanje 0, T=2
+var*var⊥	$\Delta 0 T 2$	r2: redukcija $E \rightarrow T$ , skidam 2 znaka, stanje 0, E=1
+var*var⊥	$\Delta 0 E 1$	s6: stavi + i stanje 6
var*var⊥	$\Delta 0 E 1 + 6$	s5: stavi var i stanje 5
*var⊥	$\Delta 0 E 1 + 6 \text{var} 5$	r6: redukcija $F \rightarrow \text{var}$ , skidam 2 znaka, stanje 6, F=3
*var⊥	$\Delta 0 E 1 + 6 F 3$	r4: redukcija $T \rightarrow F$ , skidam 2 znaka, stanje 6, T=9
*var⊥	$\Delta 0 E 1 + 6 T 9$	s7: stavi * i stanje 7
var⊥	$\Delta 0 E 1 + 6 T 9 * 7$	s5: stavi var i stanje 5
⊥	$\Delta 0 E 1 + 6 T 9 * 7 \text{var} 5$	r6: redukcija $F \rightarrow \text{var}$ , skidam 2 znaka, stanje 7, F=10
⊥	$\Delta 0 E 1 + 6 T 9 * 7 F 10$	r3: redukcija $T \rightarrow T * F$ , skidam 6 znakova, stanje 6, T=9
⊥	$\Delta 0 E 1 + 6 T 9$	r1: redukcija $E \rightarrow E + T$ , skidam 6 znakova, stanje 0, E=1
⊥	$\Delta 0 E 1$	Prihvati

# Algoritam LR parsera

- *ulaz*: niz  $w$  i tablica parsiranja
- *izlaz*: ili se  $w$  prihvaća ili odbaci
- *postupak*: na početku rada na stogu je stanje  $s_0$ , a na ulazu  $w \perp$ , parser izvodi program parsera sve dok se ne izvede akcija prihvati ili odbaci

## Program LR parsera:

ProgramLRParsera()

```
{   postavi kazaljku na 1. znak u nizu  $w \perp$ ;  
  
    dok(1)           //ponavljaj zauvijek  
    {               //s je stanje na vrhu stoga  
                   //kazaljka pokazuje na znak  $a$  u nizu  $w$   
    }  
    (nastavak)
```

# Program LR parsera

slučaj (akcija[s,a])

```
{      pomak s' :  
        stavi a na stog;  
        stavi s' na stog;  
        pomakni kazaljku na slijedeći znak iz w;  
reduciraj  $A \rightarrow \beta$ :  
        skini sa stoga  $2 * |\beta|$  znakova;  
        // s' stanje na vrhu nakon skidanja  
        stavi A na stog;  
        stavi NovoStanje[s',A] na stog;  
        ispiši("primjenjena produkcija  $A \rightarrow \beta$ ");  
prihvati:  
        ispiši("niz w se prihvća");  
        završi;  
odbaci:  
        ispiši("niz w se ne prihvća");  
        završi;  
}}}
```

# Primjer izgradnje parsera za gramatiku

$S \rightarrow NP VP$

$Det \rightarrow that \mid this \mid a$

$S \rightarrow Aux NP VP$

$Noun \rightarrow book \mid flight \mid meal \mid money$

$S \rightarrow VP$

$Verb \rightarrow book \mid include \mid prefer$

$NP \rightarrow Det Nominal$

$Aux \rightarrow does$

$Nominal \rightarrow Noun$

$Proper-Noun \rightarrow Houston \mid TWA$

$Nominal \rightarrow Noun Nominal$

$Prep \rightarrow from \mid to \mid on$

$NP \rightarrow Proper-Noun$

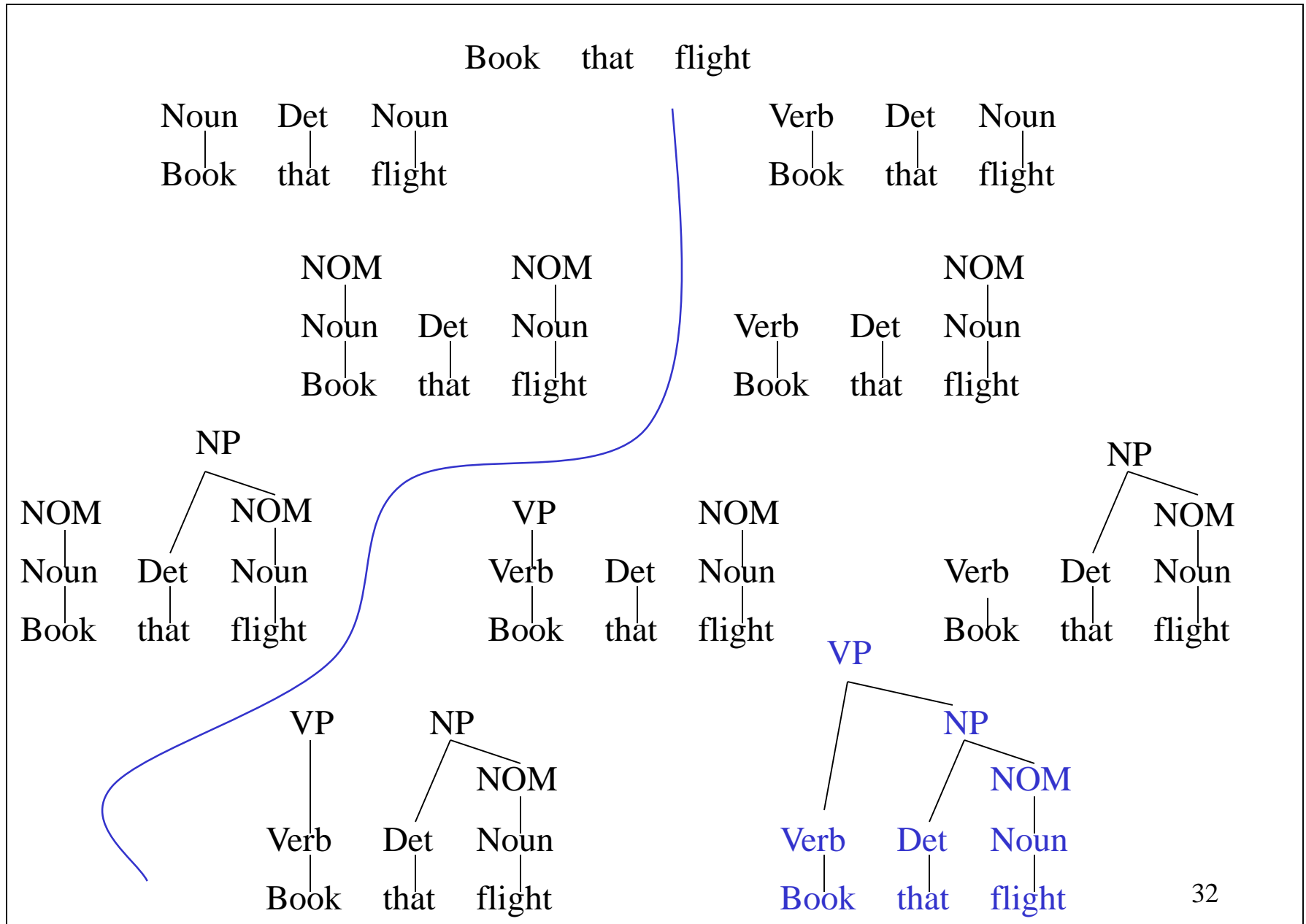
$VP \rightarrow Verb$

$VP \rightarrow Verb NP$

$Nominal \rightarrow Nominal PP$

*Book that flight.*

# Primjer parsiranja od dna prema vrhu





# Parsiranje: Formalni II

- Parsiranje od vrha prema dnu
  - prvo u dubinu / prvo u širinu
  - S-gramatika / Q-gramatika
  - LL(1)-gramatika /rekurzivni spust
- Parsiranje od dna prema vrhu
  - Pomakni-pronađi / Pomakni-reduciraj
  - Prednosti operatora
  - LR parsiranje
  - Earlyev algoritam
  - CYK (Cocke-Younger-Kasami) algoritam

# LITERATURA

- S. Srbljić: *Jezični procesori I + II*, Element, Zagreb, 2002.
- J.E. Hopcroft, J.D. Ullman: *Introduction to Automata Theory, Languages and Computation*, Addison-Wesley, USA, 1979.
- A.V. Aho, R. Sethi, J.D. Ullman: *Compilers Principles, Techniques and Tools*, 1987.
- Michael Sipser, [\*Introduction to the Theory of Computation\*](#), second edition, Course Technology, MIT, 2005.