

# Syntax Analysis:

## Context-free Grammars, Pushdown Automata and Parsing

### Part - 3

Y.N. Srikant

Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- What is syntax analysis? (covered in lecture 1)
- Specification of programming languages: context-free grammars (covered in lecture 1)
- Parsing context-free languages: push-down automata (covered in lectures 1 and 2)
- Top-down parsing: LL(1) and recursive-descent parsing
- Bottom-up parsing: LR-parsing

# Testable Conditions for LL(1)

- We call strong LL(1) as LL(1) from now on and we will not consider lookaheads longer than 1
- The classical condition for LL(1) property uses *FIRST* and *FOLLOW* sets
- If  $\alpha$  is any string of grammar symbols ( $\alpha \in (N \cup T)^*$ ), then
$$FIRST(\alpha) = \{a \mid a \in T, \text{ and } \alpha \Rightarrow^* ax, x \in T^*\}$$
$$FIRST(\epsilon) = \{\epsilon\}$$
- If  $A$  is any nonterminal, then
$$FOLLOW(A) = \{a \mid S \Rightarrow^* \alpha A a \beta, \alpha, \beta \in (N \cup T)^*,$$
$$a \in T \cup \{\$\}\}$$
- $FIRST(\alpha)$  is determined by  $\alpha$  alone, but  $FOLLOW(A)$  is determined by the “context” of  $A$ , i.e., the derivations in which  $A$  occurs

# FIRST and FOLLOW Computation Example

- Consider the following grammar  
 $S' \rightarrow S\$, S \rightarrow aAS \mid c, A \rightarrow ba \mid SB, B \rightarrow bA \mid S$
- $FIRST(S') = FIRST(S) = \{a, c\}$  because  
 $S' \Rightarrow S\$ \Rightarrow \underline{c}\$,$  and  $S' \Rightarrow S\$ \Rightarrow \underline{a}AS\$ \Rightarrow \underline{aba}S\$ \Rightarrow \underline{abac}\$$
- $FIRST(A) = \{a, b, c\}$  because  
 $A \Rightarrow \underline{ba},$  and  $A \Rightarrow SB,$  and therefore all symbols in  $FIRST(S)$  are in  $FIRST(A)$
- $FOLLOW(S) = \{a, b, c, \$\}$  because  
 $S' \Rightarrow \underline{S}\$,$   
 $S' \Rightarrow^* \underline{a}AS\$ \Rightarrow \underline{a}SBS\$ \Rightarrow \underline{aSb}AS\$,$   
 $S' \Rightarrow^* \underline{a}SBS\$ \Rightarrow \underline{aS}SS\$ \Rightarrow \underline{aS}aASS\$,$   
 $S' \Rightarrow^* \underline{aS}SS\$ \Rightarrow \underline{aS}cS\$$
- $FOLLOW(A) = \{a, c\}$  because  
 $S' \Rightarrow^* \underline{a}AS\$ \Rightarrow \underline{aA}aAS\$,$   
 $S' \Rightarrow^* \underline{a}AS\$ \Rightarrow \underline{aA}c$

# Computation of *FIRST*: Terminals and Nonterminals

```
{  
  for each ( $a \in T$ )  $FIRST(a) = \{a\}$ ;  $FIRST(\epsilon) = \{\epsilon\}$ ;  
  for each ( $A \in N$ )  $FIRST(A) = \emptyset$ ;  
  while (FIRST sets are still changing) {  
    for each production  $p$  {  
      Let  $p$  be the production  $A \rightarrow X_1 X_2 \dots X_n$ ;  
       $FIRST(A) = FIRST(A) \cup (FIRST(X_1) - \{\epsilon\})$ ;  
       $i = 1$ ;  
      while ( $\epsilon \in FIRST(X_i) \ \&\& \ i \leq n - 1$ ) {  
         $FIRST(A) = FIRST(A) \cup (FIRST(X_{i+1}) - \{\epsilon\})$ ;  $i++$ ;  
      }  
      if ( $i == n$ )  $\&\& (\epsilon \in FIRST(X_n))$   
         $FIRST(A) = FIRST(A) \cup \{\epsilon\}$   
    }  
  }  
}
```

# Computation of $FIRST(\beta)$ : $\beta$ , a string of Grammar Symbols

```
{ /* It is assumed that FIRST sets of terminals and nonterminals
   are already available */
  FIRST( $\beta$ ) =  $\emptyset$ ;
  while (FIRST sets are still changing) {
    Let  $\beta$  be the string  $X_1 X_2 \dots X_n$ ;
    FIRST( $\beta$ ) = FIRST( $\beta$ )  $\cup$  (FIRST( $X_1$ ) -  $\{\epsilon\}$ );
     $i = 1$ ;
    while ( $\epsilon \in FIRST(X_i)$  &&  $i \leq n - 1$ ) {
      FIRST( $\beta$ ) = FIRST( $\beta$ )  $\cup$  (FIRST( $X_{i+1}$ ) -  $\{\epsilon\}$ );  $i++$ ;
    }
    if ( $i == n$ ) && ( $\epsilon \in FIRST(X_n)$ )
      FIRST( $\beta$ ) = FIRST( $\beta$ )  $\cup$   $\{\epsilon\}$ 
  }
}
```

# FIRST Computation: Algorithm Trace - 1

- Consider the following grammar

$$S' \rightarrow S\$, \quad S \rightarrow aAS \mid \epsilon, \quad A \rightarrow ba \mid SB, \quad B \rightarrow cA \mid S$$

- Initially,  $\text{FIRST}(S) = \text{FIRST}(A) = \text{FIRST}(B) = \emptyset$
- Iteration 1
  - $\text{FIRST}(S) = \{a, \epsilon\}$  from the productions  $S \rightarrow aAS \mid \epsilon$
  - $\text{FIRST}(A) = \{b\} \cup \text{FIRST}(S) - \{\epsilon\} \cup \text{FIRST}(B) - \{\epsilon\} = \{b, a\}$   
from the productions  $A \rightarrow ba \mid SB$   
(since  $\epsilon \in \text{FIRST}(S)$ ,  $\text{FIRST}(B)$  is also included;  
since  $\text{FIRST}(B) = \emptyset$ ,  $\epsilon$  is not included)
  - $\text{FIRST}(B) = \{c\} \cup \text{FIRST}(S) - \{\epsilon\} \cup \{\epsilon\} = \{c, a, \epsilon\}$   
from the productions  $B \rightarrow cA \mid S$   
( $\epsilon$  is included because  $\epsilon \in \text{FIRST}(S)$ )

# FIRST Computation: Algorithm Trace - 2

- The grammar is

$$S' \rightarrow S\$, \quad S \rightarrow aAS \mid \epsilon, \quad A \rightarrow ba \mid SB, \quad B \rightarrow cA \mid S$$

- From the first iteration,

$$\text{FIRST}(S) = \{a, \epsilon\}, \quad \text{FIRST}(A) = \{b, a\}, \quad \text{FIRST}(B) = \{c, a, \epsilon\}$$

- Iteration 2

(values stabilize and do not change in iteration 3)

- $\text{FIRST}(S) = \{a, \epsilon\}$  (no change from iteration 1)
- $\text{FIRST}(A) = \{b\} \cup \text{FIRST}(S) - \{\epsilon\} \cup \text{FIRST}(B) - \{\epsilon\} \cup \{\epsilon\}$   
 $= \{b, a, c, \epsilon\}$  (changed!)
- $\text{FIRST}(B) = \{c, a, \epsilon\}$  (no change from iteration 1)

# Computation of FOLLOW

```
{ for each ( $X \in N \cup T$ ) FOLLOW( $X$ ) =  $\emptyset$ ;  
  FOLLOW( $S$ ) =  $\{\$$ \}; /*  $S$  is the start symbol of the grammar */  
  repeat {  
    for each production  $A \rightarrow X_1 X_2 \dots X_n$  { /*  $X_i \neq \epsilon$  */  
      FOLLOW( $X_n$ ) = FOLLOW( $X_n$ )  $\cup$  FOLLOW( $A$ );  
      REST = FOLLOW( $A$ );  
      for  $i = n$  downto 2 {  
        if ( $\epsilon \in \text{FIRST}(X_i)$ ) { FOLLOW( $X_{i-1}$ ) =  
          FOLLOW( $X_{i-1}$ )  $\cup$  (FIRST( $X_i$ ) -  $\{\epsilon\}$ )  $\cup$  REST;  
          REST = FOLLOW( $X_{i-1}$ );  
        } else { FOLLOW( $X_{i-1}$ ) = FOLLOW( $X_{i-1}$ )  $\cup$  FIRST( $X_i$ );  
          REST = FOLLOW( $X_{i-1}$ ); }  
      }  
    }  
  } until no FOLLOW set has changed  
}
```

# FOLLOW Computation: Algorithm Trace

- Consider the following grammar  
 $S' \rightarrow S\$, S \rightarrow aAS \mid \epsilon, A \rightarrow ba \mid SB, B \rightarrow cA \mid S$
- Initially,  $follow(S) = \{\$\}$ ;  $follow(A) = follow(B) = \emptyset$   
 $first(S) = \{a, \epsilon\}$ ;  $first(A) = \{a, b, c, \epsilon\}$ ;  $first(B) = \{a, c, \epsilon\}$ ;
- Iteration 1 /\* In the following,  $x \cup = y$  means  $x = x \cup y$  \*/
  - $S \rightarrow aAS$ :  $follow(S) \cup = \{\$\}$ ;  $rest = follow(S) = \{\$\}$   
 $follow(A) \cup = (first(S) - \{\epsilon\}) \cup rest = \{a, \$\}$
  - $A \rightarrow SB$ :  $follow(B) \cup = follow(A) = \{a, \$\}$   
 $rest = follow(A) = \{a, \$\}$   
 $follow(S) \cup = (first(B) - \{\epsilon\}) \cup rest = \{a, c, \$\}$
  - $B \rightarrow cA$ :  $follow(A) \cup = follow(B) = \{a, \$\}$
  - $B \rightarrow S$ :  $follow(S) \cup = follow(B) = \{a, c, \$\}$
  - At the end of iteration 1  
 $follow(S) = \{a, c, \$\}$ ;  $follow(A) = follow(B) = \{a, \$\}$

# FOLLOW Computation: Algorithm Trace (contd.)

- $first(S) = \{a, \epsilon\}$ ;  $first(A) = \{a, b, c, \epsilon\}$ ;  $first(B) = \{a, c, \epsilon\}$ ;
- At the end of iteration 1  
 $follow(S) = \{a, c, \$\}$ ;  $follow(A) = follow(B) = \{a, \$\}$
- Iteration 2
- $S \rightarrow aAS$ :  $follow(S) \cup = \{a, c, \$\}$ ;  
 $rest = follow(S) = \{a, c, \$\}$   
 $follow(A) \cup = (first(S) - \{\epsilon\}) \cup rest = \{a, c, \$\}$  (changed!)
- $A \rightarrow SB$ :  $follow(B) \cup = follow(A) = \{a, c, \$\}$  (changed!)  
 $rest = follow(A) = \{a, c, \$\}$   
 $follow(S) \cup = (first(B) - \{\epsilon\}) \cup rest = \{a, c, \$\}$  (no change)
- At the end of iteration 2  
 $follow(S) = follow(A) = follow(B) = \{a, c, \$\}$ ;
- The *follow* sets do not change any further

# LL(1) Conditions

- Let  $G$  be a context-free grammar
- $G$  is LL(1) iff for every pair of productions  $A \rightarrow \alpha$  and  $A \rightarrow \beta$ , the following condition holds
  - $dirsymp(\alpha) \cap dirsymp(\beta) = \emptyset$ , where  
 $dirsymp(\gamma) =$  if  $(\epsilon \in first(\gamma))$  then  
 $((first(\gamma) - \{\epsilon\}) \cup follow(A))$  else  $first(\gamma)$   
( $\gamma$  stands for  $\alpha$  or  $\beta$ )
  - $dirsymp$  stands for “direction symbol set”
- An equivalent formulation (as in ALSU’s book) is as below
  - $first(\alpha.follow(A)) \cap first(\beta.follow(A)) = \emptyset$
- Construction of the LL(1) parsing table

for each production  $A \rightarrow \alpha$

for each symbol  $s \in dirsymp(\alpha)$

/\*  $s$  may be either a terminal symbol or \$ \*/

add  $A \rightarrow \alpha$  to  $LLPT[A, s]$

Make each undefined entry of  $LLPT$  as *error*

# LL(1) Table Construction using *FIRST* and *FOLLOW*

for each production  $A \rightarrow \alpha$   
  for each terminal symbol  $a \in first(\alpha)$   
    add  $A \rightarrow \alpha$  to  $LLPT[A, a]$   
  if  $\epsilon \in first(\alpha)$  {  
    for each terminal symbol  $b \in follow(A)$   
      add  $A \rightarrow \alpha$  to  $LLPT[A, b]$   
    if  $\$ \in follow(A)$   
      add  $A \rightarrow \alpha$  to  $LLPT[A, \$]$   
  }  
Make each undefined entry of  $LLPT$  as *error*

- After the construction of the LL(1) table is complete (following any of the two methods), if any slot in the LL(1) table has two or more productions, then the grammar is NOT LL(1)

# Simple Example of LL(1) Grammar

- P1:  $S \rightarrow \text{if } (a) S \text{ else } S \mid \text{while } (a) S \mid \text{begin } SL \text{ end}$   
P2:  $SL \rightarrow S S'$   
P3:  $S' \rightarrow ; SL \mid \epsilon$
- {if, while, begin, end, a, (, ), ;} are all terminal symbols
- Clearly, all alternatives of P1 start with distinct symbols and hence create no problem
- P2 has no choices
- Regarding P3,  $\text{dirsymp}(;SL) = \{;\}$ , and  $\text{dirsymp}(\epsilon) = \{\text{end}\}$ , and the two have no common symbols
- Hence the grammar is LL(1)

# LL(1) Table Construction Example 1

LL(1) Parsing Table for the original grammar

	if	id	else	a	\$
S'	S' → S\$			S' → S\$	
S	S → if id S S → if id S else S			S → a	

Original Grammar

Grammar is not LL(1)

S' → S\$  
S → if id S |  
if id S else S |  
a

tokens: if, id, else, a

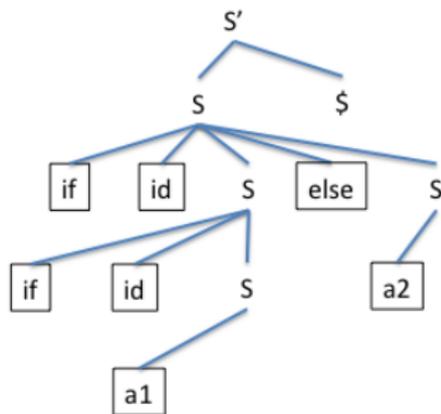
$\text{dirsyimb}(S\$) = \{\text{if}, \text{a}\}$ ;  $\text{dirsyimb}(a) = \{\text{a}\}$   
 $\text{dirsyimb}(\text{if id S}) = \{\text{if}\}$   
 $\text{dirsyimb}(\text{if id S else S}) = \{\text{if}\}$

$$\text{dirsyimb}(\text{if id S}) \cap \text{dirsyimb}(a) = \emptyset$$

$$\text{dirsyimb}(\text{if id S else S}) \cap \text{dirsyimb}(a) = \emptyset$$

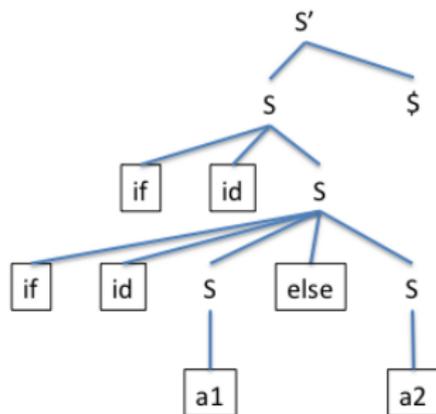
$$\text{dirsyimb}(\text{if id S}) \cap \text{dirsyimb}(\text{if id S else S}) \neq \emptyset$$

# LL(1) Table Problem Example 1



string: if id (if id a1) else a2

parentheses are not part of the string



string: if id (if id a1 else a2)

parentheses are not part of the string

# LL(1) Table Construction Example 2

Original Grammar

$S' \rightarrow S\$$   
 $S \rightarrow \text{if id } S \mid$   
 $\quad \text{if id } S \text{ else } S \mid$   
 $\quad a$

LL(1) Parsing Table for modified grammar

	if	else	a	\$
$S'$	$S' \rightarrow S\$$		$S' \rightarrow S\$$	
$S$	$S \rightarrow \text{if id } S \ S1$		$S \rightarrow a$	
$S1$		$S1 \rightarrow \epsilon$ $S1 \rightarrow \text{else } S$		$S1 \rightarrow \epsilon$

$\text{dirsymb}(S\$) = \{\text{if}, a\}$ ;  $\text{dirsymb}(a) = \{a\}$   
 $\text{dirsymb}(\text{if id } S \ S1) = \{\text{if}\}$   
 $\text{dirsymb}(\text{else } S) = \{\text{else}\}$   
 $\text{dirsymb}(\epsilon) = \{\text{else}, \$\}$

Grammar is not LL(1)

Left-Factored Grammar

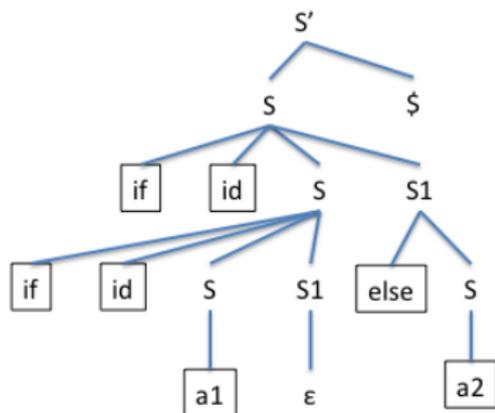
$S' \rightarrow S\$$   
 $S \rightarrow \text{if id } S \ S1 \mid a$   
 $S1 \rightarrow \epsilon \mid \text{else } S$

tokens: if, id, else, a

$$\text{dirsymb}(\text{if id } S \ S1) \cap \text{dirsymb}(a) = \emptyset$$

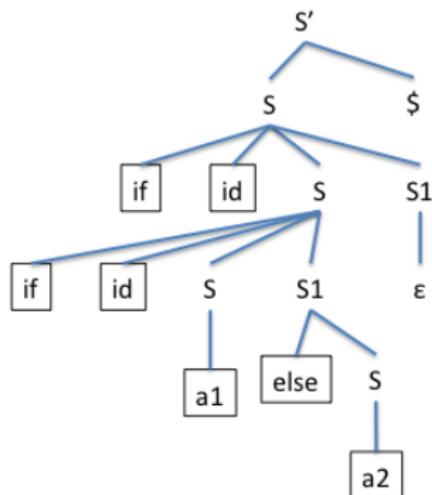
$$\text{dirsymb}(\epsilon) \cap \text{dirsymb}(\text{else } S) \neq \emptyset$$

# LL(1) Table Problem Example 2



string: if id (if id a1) else a2

parentheses are not part of the string



string: if id (if id a1 else a2)

parentheses are not part of the string

# LL(1) Table Construction Example 3

$S' \rightarrow S\$$

$S \rightarrow aAS \mid c$

$A \rightarrow ba \mid SB$

$B \rightarrow bA \mid S$

Grammar is LL(1)

LL(1) Parsing Table

	a	b	c	\$
S'	$S' \rightarrow S\$$		$S' \rightarrow S\$$	
S	$S \rightarrow aAS$		$S \rightarrow c$	
A	$A \rightarrow SB$	$A \rightarrow ba$	$A \rightarrow SB$	
B	$B \rightarrow S$	$B \rightarrow bA$	$B \rightarrow S$	

$\text{first}(S) = \{a, c\}$

$\text{first}(A) = \{a, b, c\}$

$\text{first}(B) = \{a, b, c\}$

$\text{dirsymp}(aAS) \cap \text{dirsymp}(c) = \emptyset$

$\text{dirsymp}(ba) \cap \text{dirsymp}(SB) = \emptyset$

$\text{dirsymp}(bA) \cap \text{dirsymp}(S) = \emptyset$

$\text{follow}(S) = \{a, b, c, \$\}$

$\text{follow}(A) = \{a, c\}$

$\text{follow}(B) = \{a, c\}$

$\text{dirsymp}(S\$) = \{a, c\}$

$\text{dirsymp}(aAS) = \{a\}$

$\text{dirsymp}(c) = \{c\}$

$\text{dirsymp}(ba) = \{b\}$

$\text{dirsymp}(SB) = \{a, c\}$

$\text{dirsymp}(bA) = \{b\}$

$\text{dirsymp}(S) = \{a, c\}$

# LL(1) Table Construction Example 4

Left-Recursive Grammar  
for Statement List

$$\begin{aligned} S' &\rightarrow SL \$ \\ SL &\rightarrow SL S \mid S \\ S &\rightarrow a \end{aligned}$$

$$\begin{aligned} \text{dirsymp}(SL \$) &= \{a\} \\ \text{dirsymp}(a) &= \{a\} \\ \text{dirsymp}(SL S) &= \{a\} \\ \text{dirsymp}(S) &= \{a\} \end{aligned}$$

$$\text{dirsymp}(SL S) \cap \text{dirsymp}(S) \neq \emptyset$$

$$\begin{aligned} \text{dirsymp}(SL \$) &= \{a\} \\ \text{dirsymp}(a) &= \{a\} \\ \text{dirsymp}(S A) &= \{a\} \\ \text{dirsymp}(\epsilon) &= \{\$ \} \end{aligned}$$

LL(1) Parsing Table for  
Left-Recursive Grammar

	a
S'	$S' \rightarrow SL \$$
SL	$SL \rightarrow SL S$ $SL \rightarrow S$
S	$S \rightarrow a$

Grammar is not LL(1)

Right-Recursive Grammar  
for Statement List

$$\begin{aligned} S' &\rightarrow SL \$ \\ SL &\rightarrow S A \\ A &\rightarrow S A \mid \epsilon \\ S &\rightarrow a \end{aligned}$$

LL(1) Parsing Table for  
Right-Recursive Grammar

	a	\$
S'	$S' \rightarrow SL \$$	
SL	$SL \rightarrow S A$	
A	$A \rightarrow S A$	$A \rightarrow \epsilon$
S	$S \rightarrow a$	

$$\text{dirsymp}(S A) \cap \text{dirsymp}(\epsilon) = \emptyset$$

# Elimination of Useless Symbols

Now we study the *grammar transformations*, elimination of useless symbols, elimination of left recursion and left factoring

- Given a grammar  $G = (N, T, P, S)$ , a non-terminal  $X$  is *useful* if  $S \Rightarrow^* \alpha X \beta \Rightarrow^* w$ , where,  $w \in T^*$   
Otherwise,  $X$  is useless
- Two conditions have to be met to ensure that  $X$  is useful
  - 1  $X \Rightarrow^* w$ ,  $w \in T^*$  ( $X$  derives some terminal string)
  - 2  $S \Rightarrow^* \alpha X \beta$  ( $X$  occurs in some string derivable from  $S$ )
- Example:  $S \rightarrow AB \mid CA$ ,  $B \rightarrow BC \mid AB$ ,  $A \rightarrow a$ ,  
 $C \rightarrow aB \mid b$ ,  $D \rightarrow d$ 
  - 1  $A \rightarrow a$ ,  $C \rightarrow b$ ,  $D \rightarrow d$ ,  $S \rightarrow CA$
  - 2  $S \rightarrow CA$ ,  $A \rightarrow a$ ,  $C \rightarrow b$

# Testing for $X \Rightarrow^* w$

$G' = (N', T', P', S')$  is the new grammar

$N\_OLD = \phi$ ;

$N\_NEW = \{X \mid X \rightarrow w, w \in T^*\}$

while  $N\_OLD \neq N\_NEW$  do {

$N\_OLD = N\_NEW$ ;

$N\_NEW = N\_OLD \cup \{X \mid X \rightarrow \alpha, \alpha \in (T \cup N\_OLD)^*\}$

}

$N' = N\_NEW$ ;  $T' = T$ ;  $S' = S$ ;

$P' = \{p \mid \text{all symbols of } p \text{ are in } N' \cup T'\}$

# Testing for $S \Rightarrow^* \alpha X \beta$

$G' = (N', T', P', S')$  is the new grammar

$N' = \{S\}$ ;

Repeat {

for each production  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$  with  $A \in N'$  do

add all nonterminals of  $\alpha_1, \alpha_2, \dots, \alpha_n$  to  $N'$  and

all terminals of  $\alpha_1, \alpha_2, \dots, \alpha_n$  to  $T'$

} until there is no change in  $N'$  and  $T'$

$P' = \{p \mid \text{all symbols of } p \text{ are in } N' \cup T'\}$ ;  $S' = S$