

TDDD55 Compilers and interpreters

TDDE66 Compiler Construction



# Syntax Analysis, Parsing

# Parser

- A parser for a CFG (*Context-Free Grammar*) is a program which determines whether a string  $w$  is part of the language  $L(G)$ .

## □ Function

- Produces a parse tree if  $w \in L(G)$ .
- Calls semantic routines.
- Manages syntax errors, generates error messages.

## □ Input:

- String (finite sequence of tokens)
- Input is read from left to right.

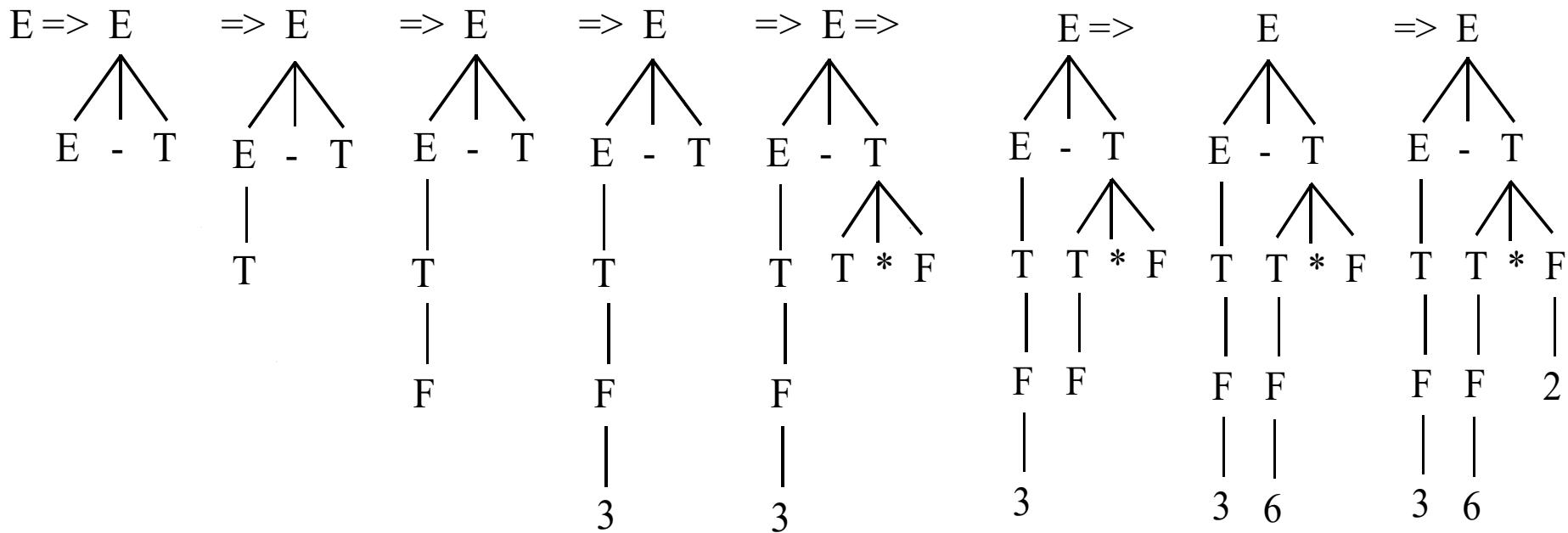
## □ Output:

- Parse tree / error messages

# Top-Down Parsing

- Example: **Top-down** parsing with input:  $3 - 6 * 2$

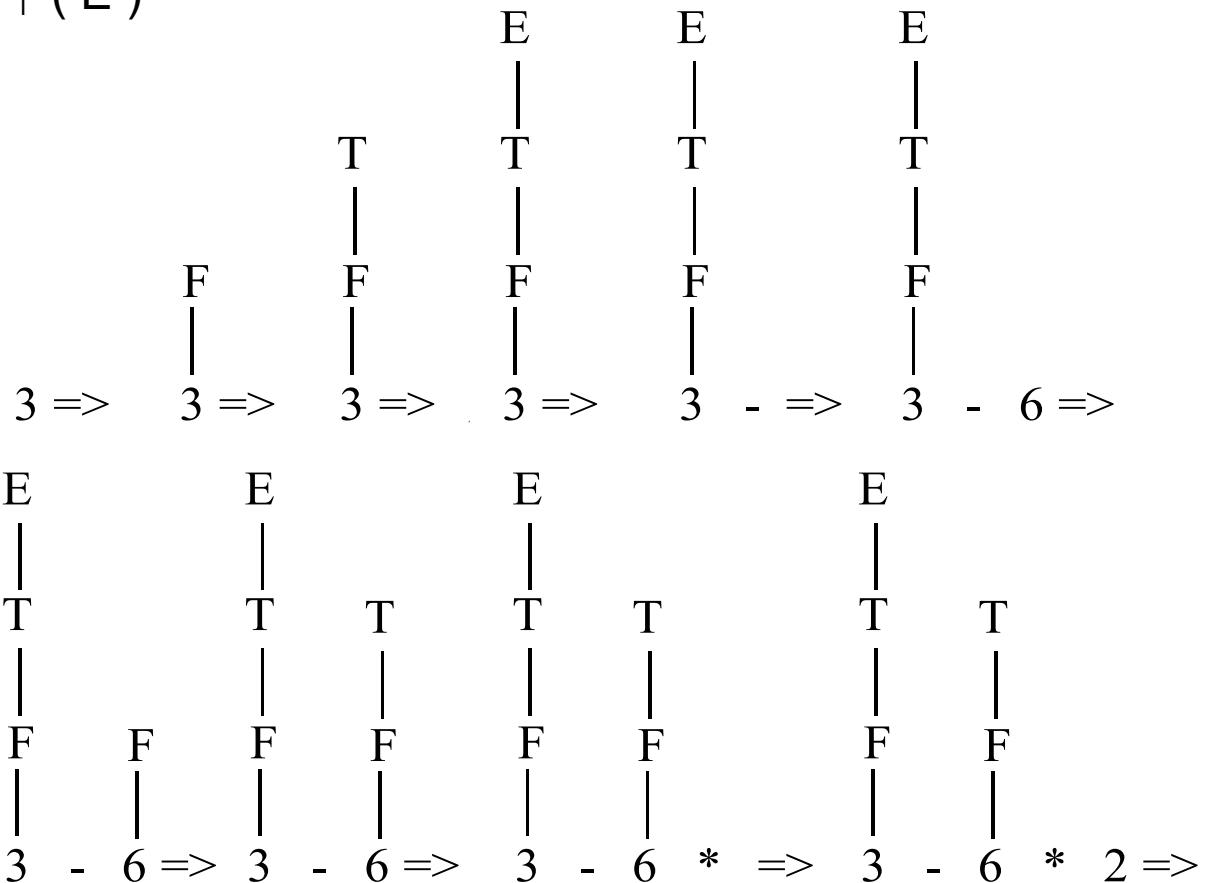
$$\begin{array}{l}
 E \rightarrow E - T \quad | \quad T \\
 T \rightarrow T * F \quad | \quad F \\
 F \rightarrow \text{Integer} \quad | \quad ( E )
 \end{array}$$



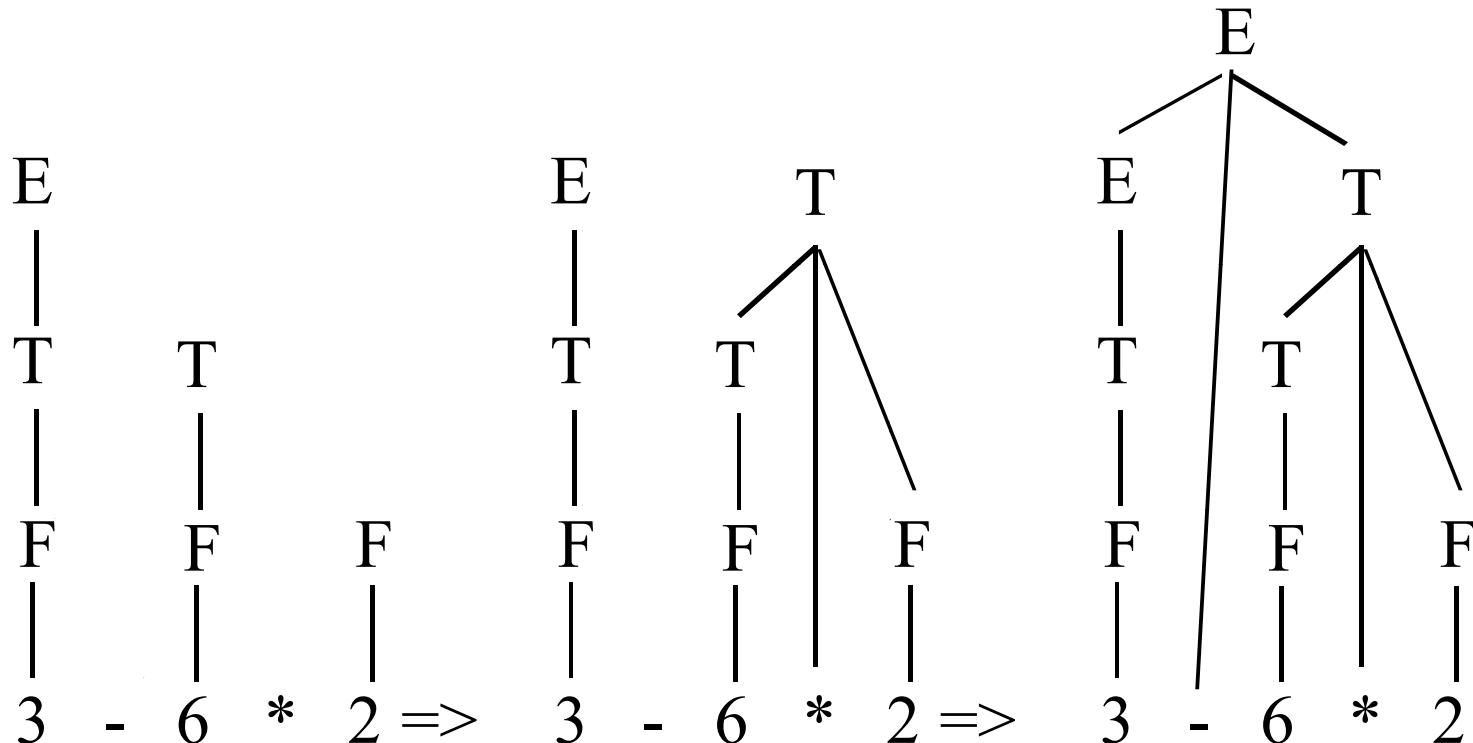
# Bottom-up Parsing

- ❑ Example: **Bottom-up** parsing with input: 3 - 6 \* 2

$E \rightarrow E - T \mid T$       (same CFG as in previous example)  
 $T \rightarrow T * F \mid F$   
 $F \rightarrow \text{Integer} \mid ( E )$



# Bottom-up Parsing cont.



# Top-Down Analysis

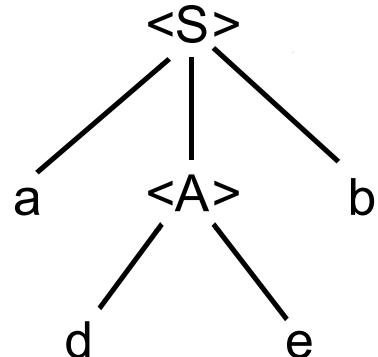
- How do we know in which order the string is to be derived?
  - Use one or more tokens lookahead.

- Example: **Top-down analysis with backtracking**

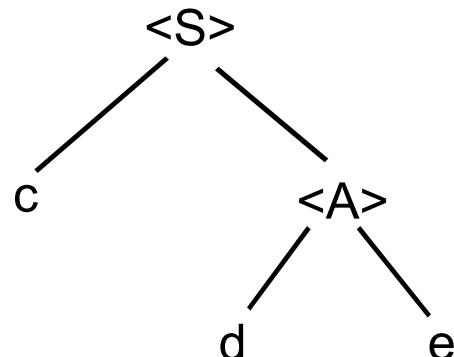
$$\begin{array}{lcl}
 <S> \rightarrow & a <A> b \\
 & | \\
 & c <A> \\
 <A> \rightarrow & d e \\
 & | \\
 & d
 \end{array}$$

*1 token lookahead works well*  
*1 token lookahead works well*  
*test right side until something fits*

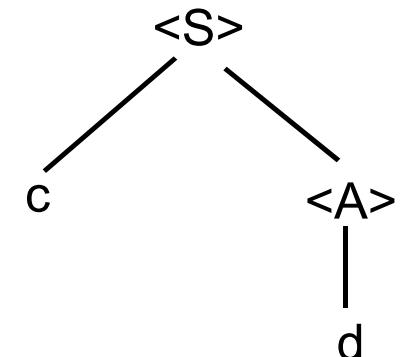
- a) adeb



- b) cd



backtracking



# Top-down Analys with Backtracking, cont.

- Top-down analysis with backtracking is implemented by writing a **procedure or a function for each nonterminal** whose task is to find one of its right sides:

```
extern char scan(); // set by scan: static char *inpptr=0;
bool A() { /* A -> d e | d */
    char* savep;
    savep = inpptr;
    if (*inpptr == 'd') {
        scan(); /* Get next token, move inpptr a step */
        if (*inpptr == 'e') {
            scan();
            return true; /* 'de' found */
        }
    }
    inpptr = savep;
    /* 'de' not found, backtrack and try 'd'*/
    if (*inpptr == 'd') {
        scan(); return true; /* 'd' found, OK */
    }
    return false;
}
```

$\begin{array}{l} \langle S \rangle \rightarrow a \langle A \rangle b \\ \qquad\qquad\qquad | \qquad c \langle A \rangle \\ \langle A \rangle \rightarrow d e \\ \qquad\qquad\qquad | \qquad d \end{array}$

# Top-down Analys with Backtracking, cont.

```
bool S() { /* S -> a A b | c A */
    if (*inpptr == 'a') {
        scan();
        if (A()) {
            if (*inpptr == 'b') {
                scan();
                return true;
            } else return false;
        } else return false;
    }
    else if (*inpptr == 'c') {
        scan();
        if (A()) return true; else return false;
    }
    else return false;
}
```

$\begin{array}{l} \langle S \rangle \rightarrow a \langle A \rangle b \\ \quad \quad \quad | \\ \quad \quad \quad c \langle A \rangle \\ \langle A \rangle \rightarrow d e \\ \quad \quad \quad | \\ \quad \quad \quad d \end{array}$

# Construction of a top-down parser

- Write a procedure for each nonterminal.
- Call scan directly *after each token is consumed.*
  - Reason: The look-ahead token should be available
- Start by calling the procedure for the start symbol.

At each step check the leftmost non-treated vocabulary symbol.

- If it is a *terminal symbol*
  - Match it with the current token and read the next token.
- If it is a *nonterminal symbol*
  - Call the routine for this nonterminal.
- In case of error call the error management routine.

# Example: An LL(1) grammar which describes binary numbers

$S \rightarrow \text{BinaryDigit } \text{BinaryNumber}$

$\text{BinaryNumber} \rightarrow \text{BinaryDigit } \text{BinaryNumber}$

|  $\epsilon$

$\text{BinaryDigit} \rightarrow 0 \mid 1$

# Sketch of a Top-Down Parser (recursive descent)

```
void TopDown(input, output)
{
/* main program */
scan();
S();
if not eof then error(...);
}
```

## Grammar:

$S \rightarrow \text{BinaryDigit } \text{BinaryNumber}$   
 $\text{BinaryNumber} \rightarrow \text{BinaryDigit}$   
                         $\text{BinaryNumber}$   
                         $\mid \epsilon$   
 $\text{BinaryDigit} \rightarrow 0 \mid 1$

```
void BinaryDigit()
{
    if (token==0 || token==1) scan();
    else error(...);
} /* BinaryDigit */

void BinaryNumber()
{
    if (token==0 || token==1)
    {
        BinaryDigit();
        BinaryNumber();
    } /* OK for the case with  $\epsilon$  */
} /* BinaryNumber */

void S()
{
    BinaryDigit();
    BinaryNumber();
} /* S */
```

# A Top-Down Parser that does not Work, Infinite Recursion:

```
void TopDown(input, output)
{
/* main program */
scan();
S();
if not eof then error(...);
}
```

## Grammar:

```
S → BinaryDigit BinaryNumber
BinaryNumber → BinaryNumber
          BinaryDigit
          | ε
BinaryDigit → 0 | 1
```

```
void BinaryDigit()
{
    if (token==0 || token==1) scan();
    else error(...);
} /* BinaryDigit */

void BinaryNumber()
{
    if (token==0 || token==1)
    {
        BinaryNumber(); /* Infinite Recursion */
        BinaryDigit();
    } /* OK for the case with ε */
} /* BinaryNumber */

void S()
{
    BinaryDigit();
    BinaryNumber();
} /* S */
```

# Non-LL(1) Structures in a Grammar:

- Left recursion, example:

$$\begin{array}{l} E \rightarrow E - T \\ | \\ T \end{array}$$

- Productions for a nonterminal with the same prefix in two or more right-hand sides, example:

$$\begin{array}{l} \text{args} \rightarrow () \\ | \\ (\text{ args }) \end{array}$$

*or*

$$\begin{array}{l} A \rightarrow ab \\ | \\ ac \end{array}$$

- The problem can be solved in most cases by rewriting the grammar to an LL(1) grammar

# Convert a grammar for top-down parsing?

## 1. Eliminate left recursion

### a) Transform the grammar to iterative form

EBNF (*Extended BNF*) Notation:

- $\{\beta\}$  zero or more, same as:  $\beta^*$
- $[\beta]$  at least once, same as:  $\beta \mid \epsilon$
- ( ) left factoring, e.g.  $A \rightarrow ab \mid ac$  in EBNF is rewritten:  
 $A \rightarrow a(b \mid c)$

Transform the grammar to be iterative using EBNF

- $A \rightarrow A \alpha \mid \beta$  (where  $\beta$  may not be preceded by  $A$ )  
in EBNF is rewritten:
- $A \rightarrow \beta \{ \alpha \}$

# 1b) Transform the Grammar to Right Recursive Form Using a Rewrite Rule:

- $A \rightarrow A\alpha | \beta$  (where  $\beta$  may not be preceded by  $A$ )

is rewritten to

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' | \varepsilon$$

Generally:

- $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | \beta_1 | \beta_2 | \dots | \beta_n$   
(where  $\beta_1, \beta_2, \dots$  may not be preceded by  $A$ )

is rewritten to:

$$A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_n A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A' | \varepsilon$$

## 2. Left Factoring Using ( ) or [ ]

*Original Grammar:*

$$\begin{aligned} <\text{stmt}> \rightarrow & \text{ if } <\text{expr}> \text{ then } <\text{stmt}> \\ | & \text{ if } <\text{expr}> \text{ then } <\text{stmt}> \text{ else } <\text{stmt}> \end{aligned}$$

*Solution using EBNF:*

$$\begin{aligned} <\text{stmt}> \rightarrow & \text{ if } <\text{expr}> \text{ then } <\text{stmt}> \\ & [ \text{ else } <\text{stmt}> ] \end{aligned}$$

*Solution using rewriting:*

$$\begin{aligned} <\text{stmt}> \rightarrow & \text{ if } <\text{expr}> \text{ then } <\text{stmt}> <\text{rest-if}> \\ <\text{rest_if}> \rightarrow & \text{ else } <\text{stmt}> \mid \epsilon \end{aligned}$$

*Original Grammar*

□ A → ab | ac

*Solution using EBNF:*

□ A → a (b | c)

# Summary LL(1) and Recursive Descent

## Summary of the LL(1) grammar:

- Many CFGs are not LL(1)
- Some can be rewritten to LL(1)
  - The underlying structure is lost (because of rewriting).

## Two main methods for writing a top-down parser

- Table-driven, LL(1)
- Recursive descent

LL(1)	Recursive Descent
Table-driven	Hand-written
+ Fast	- Much coding, + fast
+ Good error management and restart	+ Easy to include semantic actions; good error mgmt

# Small Rewriting Grammar Exercise

See 00-LectureExercises

# Example: A recursive Descent Parser for Pascal Declarations, Orig. Grammar

$\langle \text{declarations} \rangle \rightarrow \langle \text{constdecl} \rangle \langle \text{vardecl} \rangle$

$\langle \text{constdecl} \rangle \rightarrow \text{CONST} \langle \text{consdeflist} \rangle$

|  $\epsilon$

$\langle \text{consdeflist} \rangle \rightarrow \langle \text{consdeflist} \rangle \langle \text{constdef} \rangle$

|  $\langle \text{constdef} \rangle$

$\langle \text{constdef} \rangle \rightarrow \text{id} = \text{number} ;$

$\langle \text{vardecl} \rangle \rightarrow \text{VAR} \langle \text{vardeflist} \rangle$

|  $\epsilon$

$\langle \text{vardeflist} \rangle \rightarrow \langle \text{vardeflist} \rangle \langle \text{idlist} \rangle : \langle \text{type} \rangle ;$

|  $\langle \text{idlist} \rangle : \langle \text{type} \rangle ;$

$\langle \text{idlist} \rangle \rightarrow \langle \text{idlist} \rangle , \text{id}$

|  $\text{id}$

$\langle \text{type} \rangle \rightarrow \text{integer}$

|  $\text{real}$

# Rewrite in EBNF so that a Recursive Descent Parser can be Written

$\langle \text{declarations} \rangle \rightarrow \langle \text{constdecl} \rangle \langle \text{vardecl} \rangle$

$\langle \text{constdecl} \rangle \rightarrow \text{CONST} \langle \text{consdef} \rangle \{ \langle \text{consdef} \rangle \}$

|  $\epsilon$

$\langle \text{constdef} \rangle \rightarrow \text{id} = \text{number} ;$

$\langle \text{vardecl} \rangle \rightarrow \text{VAR} \langle \text{vardef} \rangle \{ \langle \text{vardef} \rangle \}$

|  $\epsilon$

$\langle \text{vardef} \rangle \rightarrow \text{id} \{ , \text{id} \} : ( \text{integer} | \text{real} ) ;$

# A Recursive Descent Parser for the New Pascal Declarations Grammar in EBNF

- We have one character lookahead.
- scan should be called when we have consumed a character.

```
void declarations()
/*<declarations>→<constdecl><vardecl> */
{
    constdecl();
    vardecl();
} /* declarations */
```

```
void constdecl()
/* <constdecl> → CONST <consdef>
{ <consdef> }
| ε */
{
    if (token == CONST) {
        scan();
        If (token == id)
            constdef();
        else
            error("Missing id after CONST");
    }
    while (token == id) constdef();
}
} /* constdecl */
```

# Pascal Declarations Parser cont 1

```
void constdef()
/* <constdef> → id = number ; */
{
    scan(); /* consume ID, get next token */
    if (token == '=')
        scan();
    else
        error("Missing '=' after id");
    if (token == NUMBER) then
        scan();
    else
        error("Missing number");

    if (token == ';')
        scan(); /* consume ';;', get next token */
    else
        error("Missing ';' after const decl");
} /* constdef */
```

```
void vardecl()
/* <vardecl> → VAR <vardef> { <vardef> }
|ε */
{
    If (token == VAR) {
        scan();

        if (token == ID)
            vardef();
        else
            error("Missing id after VAR");

    while (token == ID) {
        vardef();
    }
}
} /* vardecl */
```

# Pascal Declarations Parser cont 2

```
void vardef() /* <vardef> → id { , id } : ( integer | real ) ; */  
{  
    scan();  
    while (token == ',') {  
        scan();  
        if (token == ID)  
            scan();  
        else error("id expected after ',' ");  
    } /* while */  
  
    if (token == ':') {  
        scan();  
        if ((token == INTEGER) || (token == REAL))  
            scan();  
        else error("Incorrect type of variable");  
        if (token ==';')  
            scan();  
        else error("Missing ';' in variable decl.");  
  
    } else error("Missing ':' in var. decl.");  
} /* vardef */
```

```
{ /* main */  
    scan(); /* lookahead token */  
    declarations();  
    if (token!=eof_token) then error(...);  
} /* main */
```

# LL Parsing Issues

## Beyond Recursive Descent

**LL( $k$ )**

**LL items**

**Finite pushdown automaton**

**FIRST and FOLLOW**

**Table-driven Predictive Parser**

# LL( $k$ )

- Given:
    - Context-free grammar  $G = (N, \Sigma, P, S)$
    - Integer  $k > 0$

- G is (in) **LL( $k$ )** if:  
for any two leftmost derivations

- $S \Rightarrow^*_{\text{Im}} uY\alpha \Rightarrow u\beta\alpha \Rightarrow^* ux$  and
  - $S \Rightarrow^*_{\text{Im}} uY\alpha \Rightarrow u\gamma\alpha \Rightarrow^* uy$

with  $x[1:k] = y[1:k]$

it holds  $\beta = \gamma$ .

the  $k$  first tokens of  $x$  and  $y$   
are equal

- ❑ That is, for fixed left context  $u$ , the choice for the "right" production to apply to  $Y$  is uniquely determined by the next  $k$  input tokens.

# Example

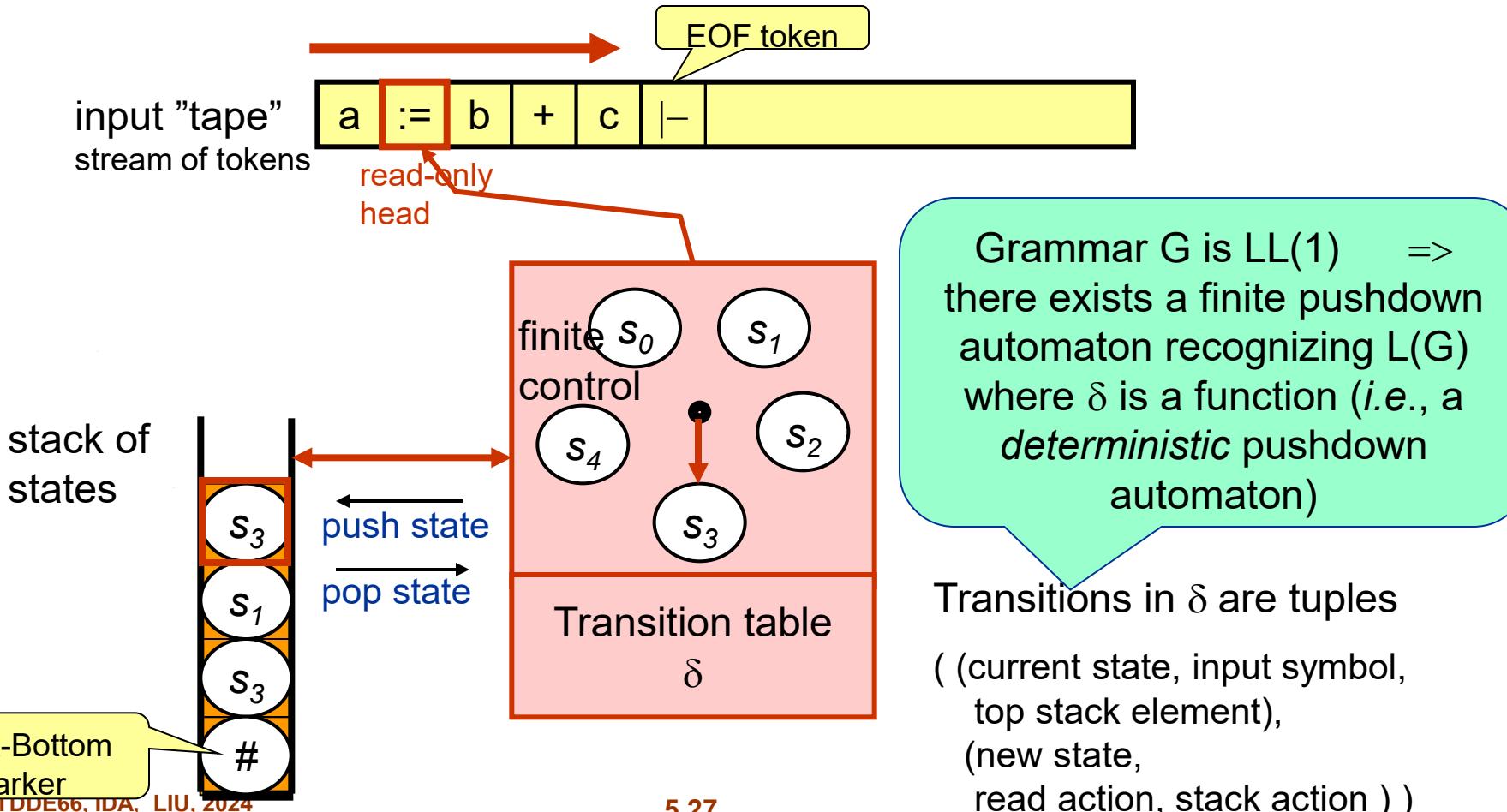
- The following grammar is LL(1)  
(terminals are bold-face):

```
S -> if ident then S else S fi
      | while ident do S od
      | begin S end
      | ident := ident
```

# Automaton Model for Parsing Context-Free Languages

## Finite pushdown automaton (FPA)

- ❑ a finite automaton with a stack of states



# Context-Free Items

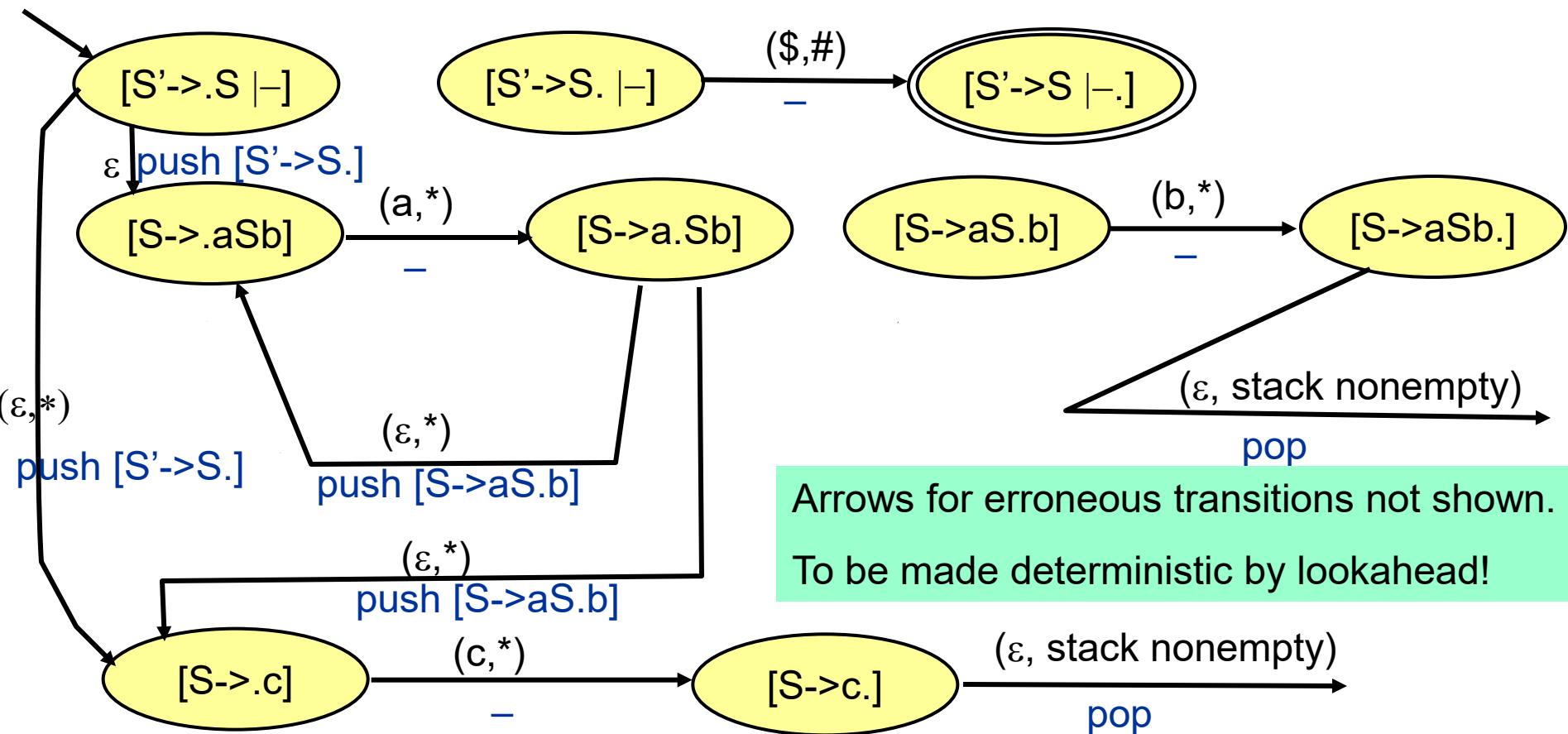
Given CFG G, construct states of the finite pushdown automaton:

- Add new start symbol  $S'$  with  $S' \rightarrow S \mid -$  ( $\mid -$  means End-of-Input)
- For each production  $A \rightarrow \alpha_1 \dots \alpha_k$  e.g.  $A \rightarrow aBc$   
create  $k+1$  **context-free items** (= states)
  - e.g.,  $[A \rightarrow .aBc]$ ,  $[A \rightarrow a.Bc]$ ,  $[A \rightarrow aB.c]$ ,  $[A \rightarrow aBc.]$
- Construct a **predictive parser** as finite pushdown automaton:
  - start in state  $[S' \rightarrow .S \mid -]$  with empty stack (#)
  - halt and accept in state  $[S' \rightarrow S \mid -]$  with empty stack (#)
  - at  $[A \rightarrow \alpha.b\gamma]$ : read input symbol, i.e.,  $[A \rightarrow \alpha.b\gamma] \rightarrow [A \rightarrow \alpha b.\gamma]$
  - at  $[A \rightarrow \alpha.B\gamma]$ : push  $[A \rightarrow \alpha B.\gamma]$ ,  
determine new production  $B \rightarrow \beta$   
and start from  $[B \rightarrow .\beta]$
  - at  $[B \rightarrow \beta.]$ : pop state  $[A \rightarrow \alpha B.\gamma]$  to restore context (if #, error)


 Prediction!

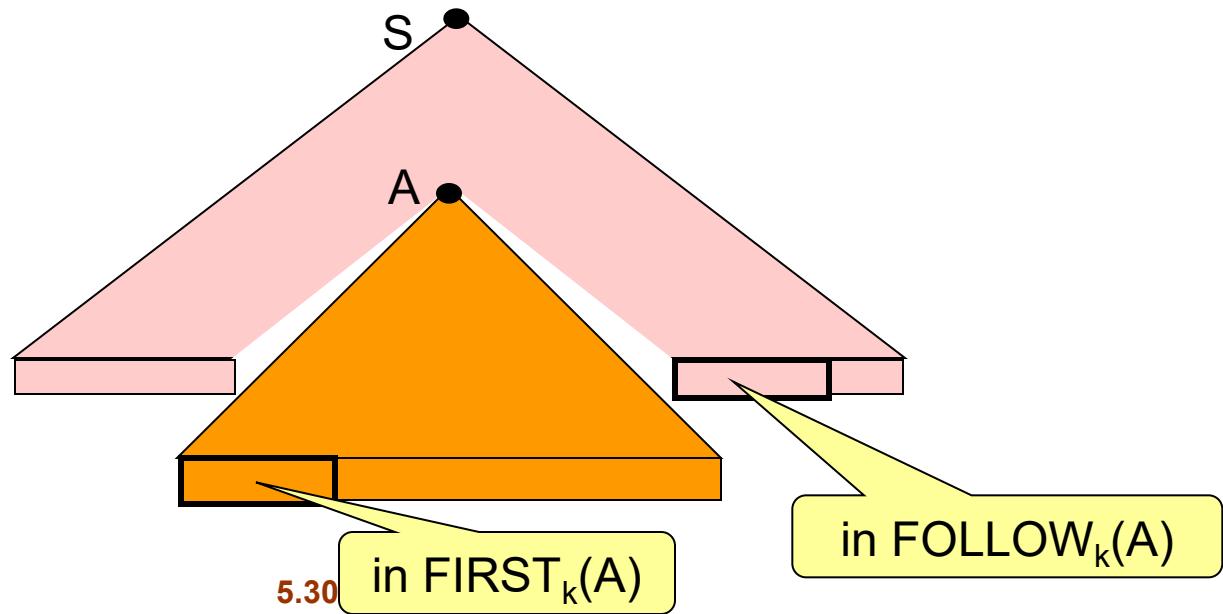
# Example

- Grammar with productions  $\{ S \rightarrow aSb \mid c \}$
- Add new start symbol  $S'$ :  $\{ S' \rightarrow S; \quad S \rightarrow aSb; \quad S \rightarrow c \}$
- Transition diagram (showing **stack actions** below arrows):



# FIRST and FOLLOW

- For a sentential form  $\alpha$  in  $(N \cup S)^+$ ,  
**FIRST**( $\alpha$ ) denotes the set of all terminals  
 which can be ***first*** in a string derived from  $\alpha$ .
- For a nonterminal  $A$  in  $N$ ,  
**FOLLOW**( $A$ ) denotes the set of all terminals (e.g.  $a$ )  
 that could appear immediately ***after***  $A$  in a sentential form  
 i.e., there exists  $S \Rightarrow^* \alpha A a \beta$  for arbitrary  $\alpha, \beta$



# Small FIRST and FOLLOW Exercise

See 00-LectureExercises

# Computing FIRST = FIRST<sub>1</sub>

For all grammar symbols X:

- If X is a terminal, then FIRST(X) = { X }.
- If  $X \rightarrow \varepsilon$  is a production, then add  $\varepsilon$  to FIRST(X).
- If X is a nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_q$  is a production,
  - then place all those  $a$  of  $\Sigma$  in FIRST(X) where for some  $i$ ,  $a$  is in FIRST( $Y_i$ ) and  $\varepsilon$  is in all of FIRST( $Y_1$ ), ..., FIRST( $Y_{i-1}$ ) (that is,  $Y_1$ , ...,  $Y_{i-1}$  all may derive  $\varepsilon$ ).
  - If  $\varepsilon$  is in FIRST( $Y_j$ ) for all  $j=1,2,\dots,q$  then add  $\varepsilon$  to FIRST(X).

Apply these rules until no more terminals or  $\varepsilon$  can be added to any FIRST set.

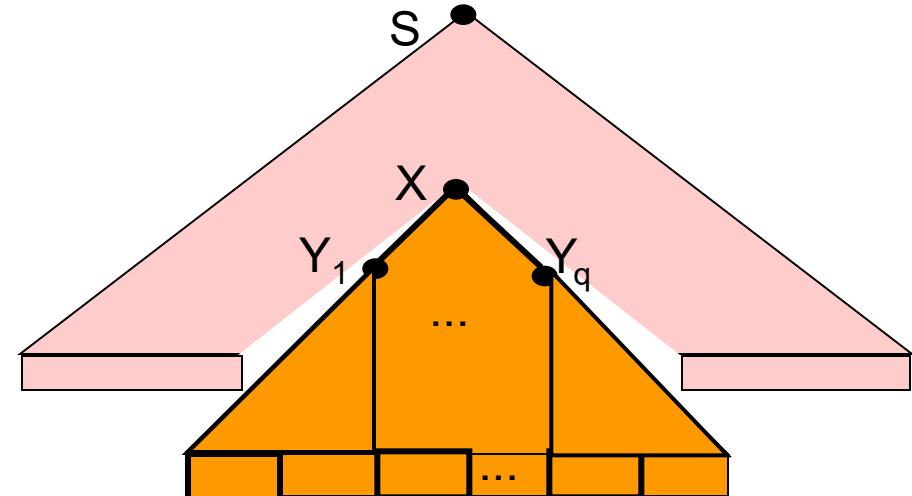
For the example grammar

$S' \rightarrow S; \quad S \rightarrow aSb; \quad S \rightarrow c$

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

$\text{FIRST}(S') = \text{FIRST}(S)$

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



# Computing FIRST (cont.)

For any string  $X_1 X_2 \dots X_n$  of grammar symbols:

- Add to  $\text{FIRST}(X_1 X_2 \dots X_n)$  all non- $\epsilon$  symbols of  $\text{FIRST}(X_1)$ .
- If  $\epsilon$  in  $\text{FIRST}(X_1)$ , add also all non- $\epsilon$  symbols of  $\text{FIRST}(X_2)$ , otherwise done.
- If  $\epsilon$  also in  $\text{FIRST}(X_2)$ , add also all non- $\epsilon$  symbols of  $\text{FIRST}(X_3)$ , otherwise done.
- ...
- If  $\epsilon$  also in  $\text{FIRST}(X_n)$ , add  $\epsilon$  to  $\text{FIRST}(X_1 X_2 \dots X_n)$

For the example grammar

$S' \rightarrow S; \quad S \rightarrow aSb; \quad S \rightarrow c$

$\text{FIRST}(abc) = \{a\}$

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

# Computing FOLLOW

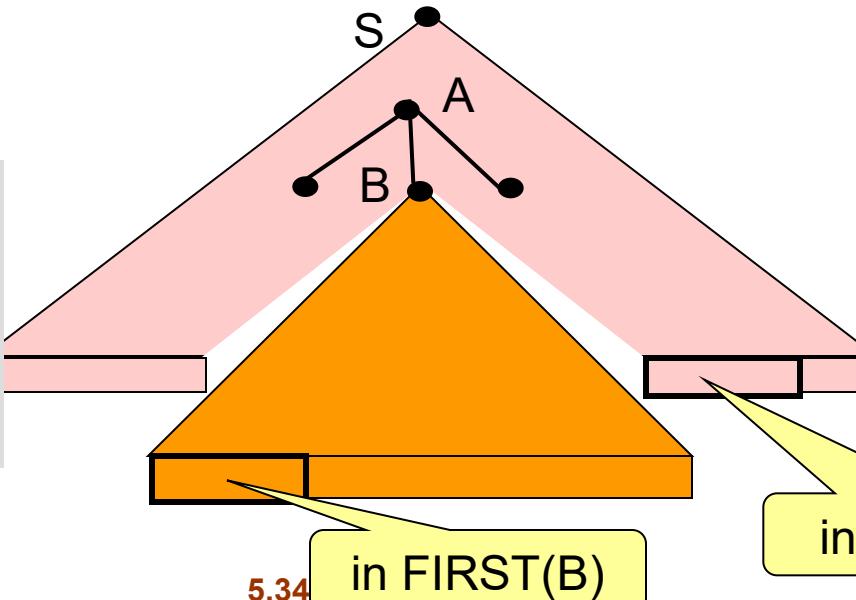
Compute  $\text{FOLLOW}(B)$  for each nonterminal  $B$ :

- Add  $|-$  to  $\text{FOLLOW}(S)$
- If there is a production  $A \rightarrow^* \alpha B \beta$  for arbitrary  $\alpha, \beta$   
then add all of  $\text{FIRST}(\beta)$  except  $\epsilon$  to  $\text{FOLLOW}(B)$
- If there is a production  $A \rightarrow \alpha B$ ,  
or a production  $A \rightarrow \alpha B \beta$  where  $\epsilon$  in  $\text{FIRST}(\beta)$ , i.e.  $\beta \Rightarrow^* \epsilon$ ,  
then add all of  $\text{FOLLOW}(A)$  to  $\text{FOLLOW}(B)$ .

Apply these rules until no more terminals or  $\epsilon$  can be added to any FOLLOW set.

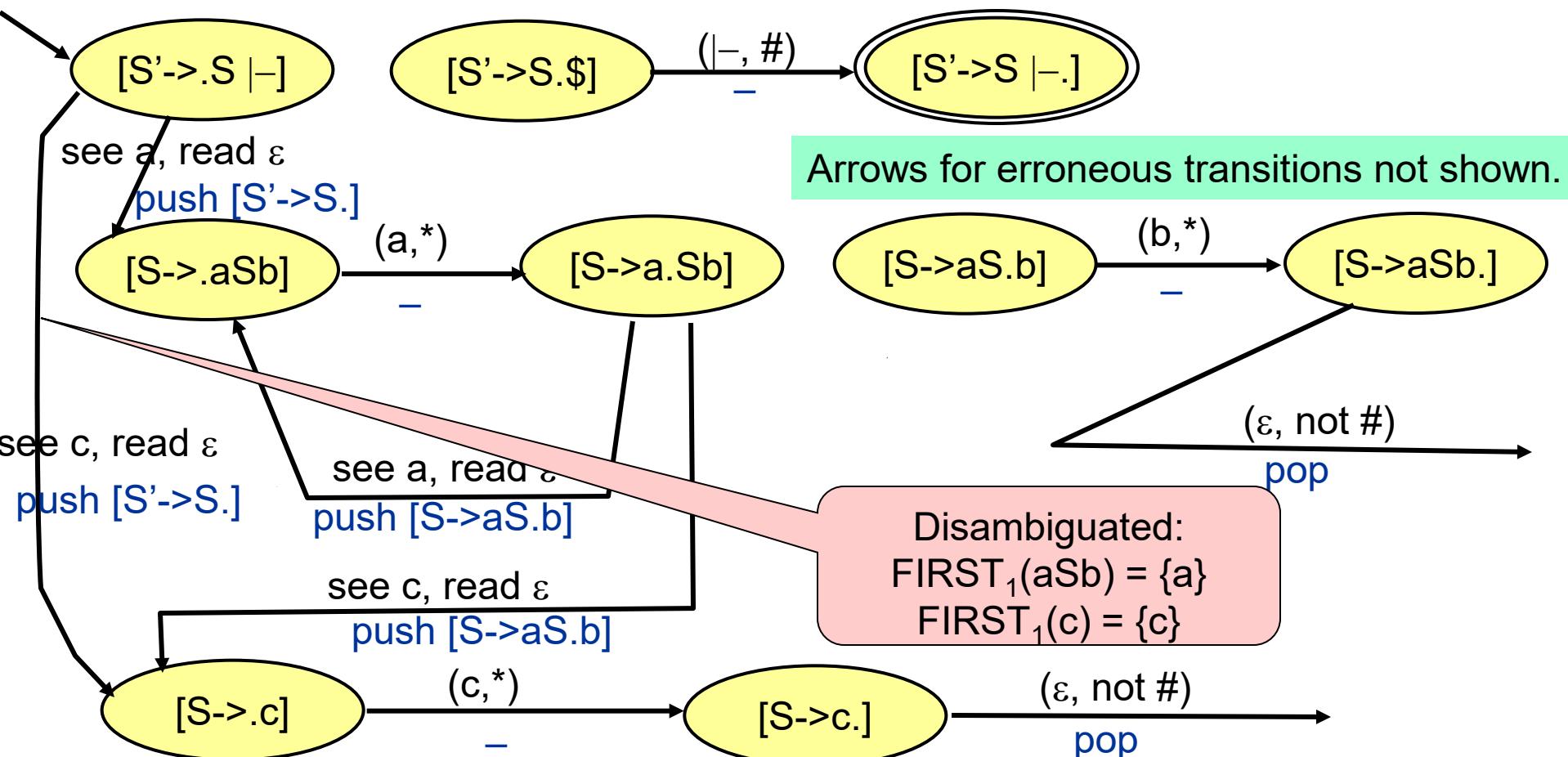
For the example grammar  
 $S \rightarrow aSb; S \rightarrow c$

$\text{FOLLOW}(S) = \{|-, b\}$



# Example Cont.: Finite Pushdown Automaton (FPA) Made Deterministic

- Grammar with productions {  $S \rightarrow aSb \mid c$  }
- Added new start symbol  $S'$ : {  $S' \rightarrow S \mid -; \quad S \rightarrow aSb; \quad S \rightarrow c$  }



# Example (cont.): Transition table ( $k=1$ )

state	final ?	lookahead a	lookahead b	lookahead c	lookahead  -
[S'->.S  -]	no	push [S'->S.\$]; [S->.aSb]	[Error]	push [S'->S.\$]; [S->.c]	[Error]
[S'->S.  -]	no	[Error]	[Error]	[Error]	read  -; [S'->S  -.]
[S'->S  -.]	yes				
[S->.aSb]	no	read a; [S->a.Sb]	[Error]	[Error]	[Error]
[S->a.Sb]	no	push [S->aS.b]; [S->.aSb]	[Error]	push [S->aS.b]; [S->.c]	[Error]
[S->aS.b]	no	[Error]	read b; [S->aSb.]	[Error]	[Error]
[S->aSb.]	no	[Error]	pop state	[Error]	pop state
[S->.c]	no	[Error]	[Error]	read c; [S->c.]	[Error]
[S->c.]	no	[Error]	pop state	[Error]	pop state

# General Approach: Predictive Parsing

At any production  $A \rightarrow \alpha$

- If  $\varepsilon$  is not in  $\text{FIRST}(\alpha)$ :
  - Parser expands by production  $A \rightarrow \alpha$  if current lookahead input symbol is in  $\text{FIRST}(\alpha)$ .
- otherwise (i.e.,  $\varepsilon$  in  $\text{FIRST}(\alpha)$ ):
  - Expand by production  $A \rightarrow \alpha$  if current lookahead symbol is in  $\text{FOLLOW}(A)$  or if it is  $|-$  and  $|-$  is in  $\text{FOLLOW}(A)$ .

Use these rules to fill the transition table.

(pseudocode: see [ASU86] p. 190, [ALSU06] p. 224)

# Summary: Parsing LL( $k$ ) Languages

## ❑ Predictive LL parser

- iterative, based on finite pushdown automaton
- transition-table-driven
- can be generated automatically

## ❑ Recursive-descent parser

- recursive
- manually coded
- easier to fix intermediate code generation, error handling

## ❑ Both require lookahead (or backtracking)

to predict the next production to apply

- Removes nondeterminism
- Necessary checks derived from FIRST and FOLLOW sets
- FIRST and FOLLOW are also useful for syntax error recovery

# Homework

- ❑ Now, read again the part on recursive descent parsers and find the equivalent of
  - Context-free items (pushdown automaton (PDA) states)
  - The stack of states
  - Pushing a state to stack
  - Popping a state from stack
  - Start state, final statein a recursive descent parser.

# Thank you! Questions?

Next lecture: LR Parsing, Part 1