



## Syntax Analysis, Parsing

## Parser

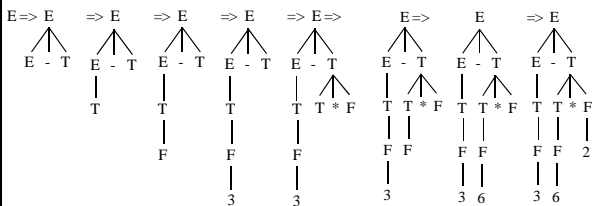


- A parser for a CFG (*Context-Free Grammar*) is a program which determines whether a string  $w$  is part of the language  $L(G)$ .
- **Function**
  - Produces a parse tree if  $w \in L(G)$ .
  - Calls semantic routines.
  - Manages syntax errors, generates error messages.
- **Input:**
  - String (finite sequence of tokens)
  - Input is read from left to right.
- **Output:**
  - Parse tree / error messages

## Top-Down Parsing



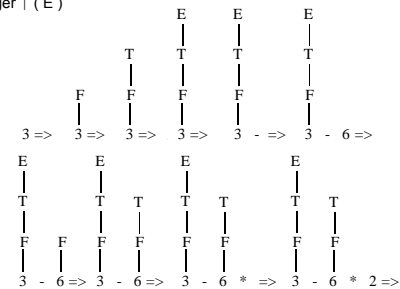
- Example: **Top-down** parsing with input:  $3 - 6 * 2$
- $E \rightarrow E - T \mid T$
- $T \rightarrow T * F \mid F$
- $F \rightarrow \text{Integer} \mid ( E )$



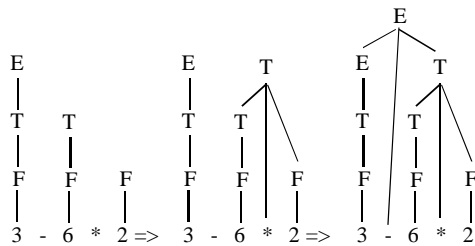
## Bottom-up Parsing



- 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.

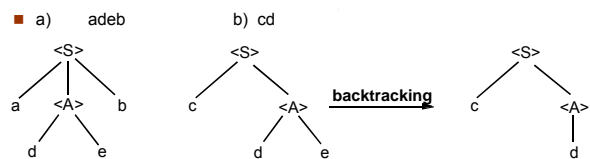


## Top-Down Analysis



- 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**
- $\langle S \rangle \rightarrow a \langle A \rangle b$  *1 token lookahead works well*
- $\quad \quad \quad c \langle A \rangle$  *1 token lookahead works well*
- $\langle A \rangle \rightarrow d e$  *test right side until something fits*
- $\quad \quad \quad d$



## Top-down Analys with Backtracking, cont.

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

```
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;
}
```

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010.

5.7

## 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;
}
```

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010.

5.8

## 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.

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010.

5.9

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

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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010.

5.10

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

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

### Grammar:

```
S → BinaryDigit BinaryNumber
BinaryNumber → BinaryDigit
                BinaryNumber
                | ε
BinaryDigit → 0 | 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 ε */
} /* B' */
```

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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010.

5.11

## 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 here */
        BinaryDigit();
    } /* OK for the case with ε */
} /* B' */
```

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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010.

5.12

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

- Left recursion, example:
 
$$E \rightarrow E - T$$

$$| T$$
- Productions for a nonterminal with the same prefix in two or more right-hand sides, example:
 
$$\text{arglist} \rightarrow ( )$$

$$| ( \text{args} )$$

or

$$A \rightarrow a b$$

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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.13

## Convert a grammar for top-down parsing?

### 1. Eliminate left recursion

#### a) Transform the grammar to iterative form

EBNF (*Extended BNF*) Notation:

- $\{\beta\}$  same as the regular expression:  $\beta^*$
- $[\beta]$  same as the regular expression:  $\beta | \epsilon$
- $()$  left factoring, e.g.  $A \rightarrow ab | ac$  in EBNF is rewritten:  $A \rightarrow a (b | c)$

Transform the grammar to be iterative using EBNF

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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.14

### 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
 
$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' | \epsilon$$

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

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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.15

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

Original Grammar:

$\langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle$   
 $| \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

Solution using EBNF:

$\langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle$   
 $[ \text{else } \langle \text{stmt} \rangle ]$

Solution using rewriting:

$\langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \langle \text{rest-if} \rangle$   
 $\langle \text{rest-if} \rangle \rightarrow \text{else } \langle \text{stmt} \rangle | \epsilon$

Original Grammar

- $A \rightarrow ab | ac$

Solution using EBNF:

- $A \rightarrow a (b | c)$

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.16

## 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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.17

## Small Rewriting Grammar Exercise

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.18

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



```

<declarations> → <constdecl> <vardecl>
<constdecl> → CONST <consdeflist>
                | ε
<consdeflist> → <consdeflist> <consdef>
                | <consdef>
<consdef> → id = number ;
<vardecl> → VAR <vardeflist>
                | ε
<vardeflist> → <vardeflist> <idlist> : <type> ;
                | <idlist> : <type> ;
<idlist> → <idlist> , id
                | id
<type> → integer
                | real
    
```

TDD55/B44, P Fritzon, C. Kessler, IDA, LIU, 2010. 5.19

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



```

<declarations> → <constdecl> <vardecl>
<constdecl> → CONST <consdef> { <consdef> }
                | ε
<consdef> → id = number ;
<vardecl> → VAR <vardef> { <vardef> }
                | ε
<vardef> → id { , id } : ( integer | real ) ;
    
```

TDD55/B44, P Fritzon, C. Kessler, IDA, LIU, 2010. 5.20

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

TDD55/B44, P Fritzon, C. Kessler, IDA, LIU, 2010. 5.21

## Pascal Declarations Parser cont 1



```

void constdef()              void vardecl()
/* <constdef> → id = number ; */ /* <vardecl> → VAR <vardef> { <vardef> } */
{                               { [ε */
    scan(); /* consume ID, get next token */ {
    if (token == '=')           if (token == VAR) {
        scan();                 scan();
    else                         if (token == ID)
        error("Missing '=' after id");         vardef();
    if (token == NUMBER) then    else
        scan();                 error("Missing id after VAR");
    else                           while (token == ID) {
        error("Missing number");         vardef();
    }                               }
    if (token == ';')           scan(); /* consume ';', get next token */
        scan(); /* consume ';', get next token */
    else                           error("Missing ';' after const decl");
    } /* constdef */           } /* vardecl */
    
```

TDD55/B44, P Fritzon, C. Kessler, IDA, LIU, 2010. 5.22

## 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 */
    
```

TDD55/B44, P Fritzon, C. Kessler, IDA, LIU, 2010. 5.23

TDD55 Compilers and interpreters  
TDD55 Compiler Construction



## LL Parsing Issues Beyond Recursive Descent

**LL(k)**  
**LL items**  
**Finite pushdown automaton**  
**FIRST and FOLLOW**  
**Table-driven Predictive Parser**

Peter Fritzon, Christoph Kessler,  
IDA, Linköping universitet, 2011.

## 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
    - $S \Rightarrow_{lm}^* uY\alpha \Rightarrow^* u\beta\alpha \Rightarrow^* ux$  and
    - $S \Rightarrow_{lm}^* uY\alpha \Rightarrow^* u\gamma\alpha \Rightarrow^* uy$
 with  $x[1:k] = y[1:k]$ 
    - the  $k$  first tokens of  $x$  and  $y$  are equal
 it holds  $\beta = \gamma$ .
- 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.

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.25

## Example

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

```

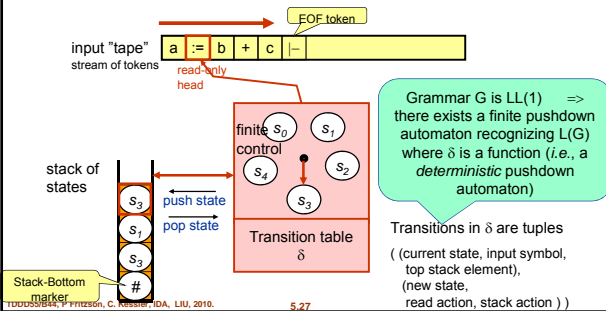
S -> if ident then S else S fi
    | while ident do S od
    | begin S end
    | ident := ident
    
```

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.26

## Automaton Model for Parsing Context-Free Languages

### Finite pushdown automaton (FPA)

- a finite automaton with a stack of states



TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.27

## Context-Free Items

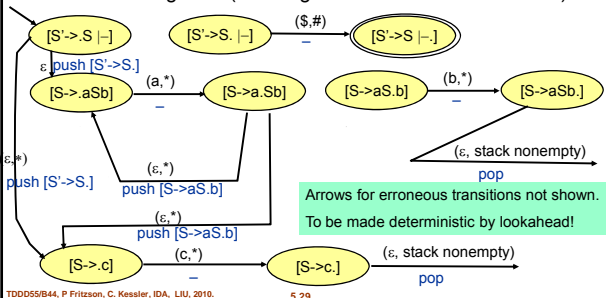
Given CFG  $G$ , construct states of the finite pushdown automaton:

- Add new start symbol  $S'$  with  $S' \rightarrow S \mid \_$  ( $\_$  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.,  $[A \rightarrow \cdot aBc]$ ,  $[A \rightarrow a \cdot Bc]$ ,  $[A \rightarrow aB \cdot c]$ ,  $[A \rightarrow aBc \cdot]$
- Construct a **predictive parser** as finite pushdown automaton:
  - start in state  $[S' \rightarrow \cdot S \mid \_]$  with empty stack ( $\#$ )
  - halt and accept in state  $[S' \rightarrow S \mid \_]$  with empty stack ( $\#$ )
  - at  $[A \rightarrow \alpha \cdot b\gamma]$ : read input symbol, i.e.,  $[A \rightarrow \alpha \cdot b\gamma] \rightarrow [A \rightarrow \alpha b \cdot \gamma]$
  - at  $[A \rightarrow \alpha \cdot B\gamma]$ : push  $[A \rightarrow \alpha B \cdot \gamma]$ , determine new production  $B \rightarrow \beta$  and start from  $[B \rightarrow \cdot \beta]$  **Prediction!**
  - at  $[B \rightarrow \beta \cdot]$ : pop state  $[A \rightarrow \alpha B \cdot \gamma]$  to restore context (if  $\#$ , error)

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.28

## Example

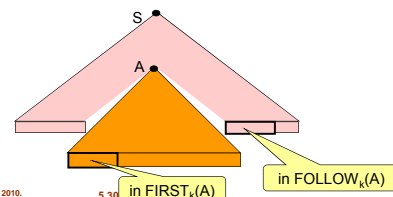
- Grammar with productions  $\{ S \rightarrow aSb \mid c \}$
- Add new start symbol  $S'$ :  $\{ S' \rightarrow S; S \rightarrow aSb; S \rightarrow c \}$
- Transition diagram (showing **stack actions** below arrows):



TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.29

## FIRST and FOLLOW

- For a sentential form  $\alpha$  in  $(N \cup \Sigma)^+$ ,  $FIRST(\alpha)$  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  $S \Rightarrow^* \alpha A \beta$  for arbitrary  $\alpha, \beta$



TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.30

## Small FIRST and FOLLOW Exercise

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.31

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

For all grammar symbols X:

- If X is a terminal, then  $FIRST(X) = \{ X \}$ .
- If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  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  $\epsilon$  is in all of  $FIRST(Y_1), \dots, FIRST(Y_{i-1})$  (that is,  $Y_1, \dots, Y_{i-1}$  all may derive  $\epsilon$ ).
  - If  $\epsilon$  is in  $FIRST(Y_j)$  for all  $j=1,2,\dots,q$  then add  $\epsilon$  to  $FIRST(X)$ .

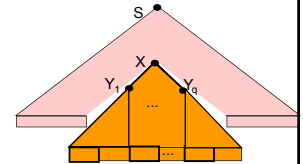
Apply these rules until no more terminals or  $\epsilon$  can be added to any FIRST set.

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

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

$FIRST(S') = FIRST(S)$

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



## Computing FIRST (cont.)

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

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

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

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

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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.33

## Computing FOLLOW

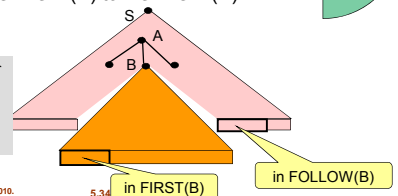
Compute  $FOLLOW(B)$  for each nonterminal B:

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

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

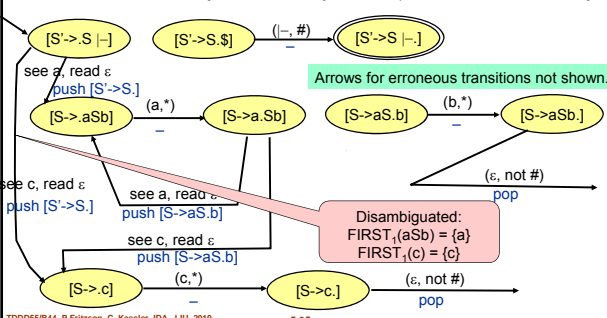
For the example grammar  
 $S \rightarrow aSb; S \rightarrow c$

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



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

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



TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.35

## 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 -> aSb]	[Error]	[Error]	[Error]
[S -> aSb]	no	push [S -> aSb]; [S -> aSb]	[Error]	push [S -> aSb]; [S -> c]	[Error]
[S -> aSb]	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

TDD055/B44, P Fritzzon, C. Kessler, IDA, LIU, 2010. 5.36

## General Approach: Predictive Parsing



At any production  $A \rightarrow \alpha$

- If  $\epsilon$  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.,  $\epsilon$  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)

TDD055/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

5.37

## 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  $\text{FIRST}$  and  $\text{FOLLOW}$  sets
  - $\text{FIRST}$  and  $\text{FOLLOW}$  are also useful for syntax error recovery

TDD055/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

5.38

## 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.

TDD055/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

5.39