



Syntax Analysis, Parsing

Peter Fritzson, Christoph Kessler,
IDA, Linköpings universitet, 2010.

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

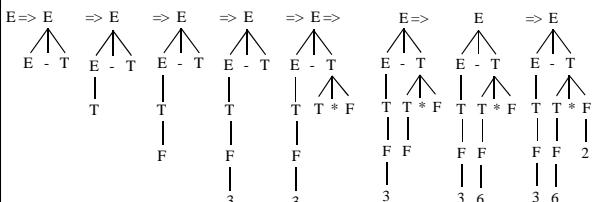
TDDD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

5.2

Top-Down Parsing

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

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



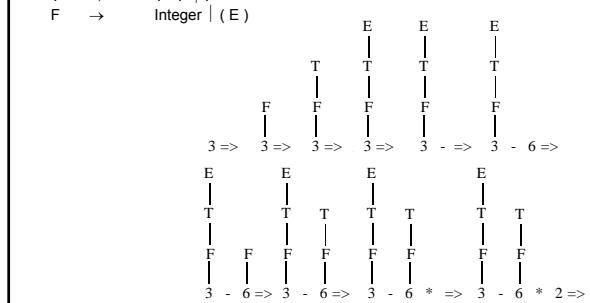
TDDD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

5.3

Bottom-up Parsing

■ Example: **Bottom-up** parsing with input: $3 - 6 * 2$ (same CFG as in previous example)

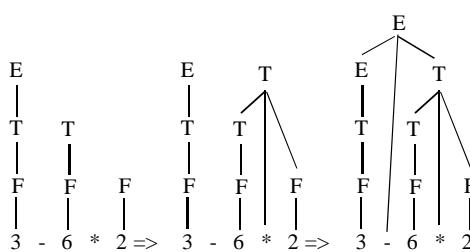
$$\begin{array}{l} E \rightarrow E - T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow \text{Integer} \mid (E) \end{array}$$



TDDD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

5.4

Bottom-up Parsing cont.



TDDD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

5.5

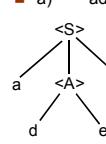
Top-Down Analysis

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

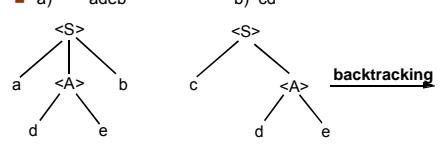
■ Example: **Top-down analysis with backtracking**

$$\begin{array}{ll} <S> \rightarrow a <A> b & 1 \text{ token lookahead works well} \\ & | c <A> \\ <A> \rightarrow d e & 1 \text{ token lookahead works well} \\ & | d \end{array} \quad \begin{array}{l} \text{test right side until something} \\ \text{fits} \end{array}$$

■ a) adeb



b) cd



TDDD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

5.6

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;
}
```

TDDD16/B44, P Fritzon, 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;
}
```

TDDD16/B44, P Fritzon, 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.

TDDD16/B44, P Fritzon, 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



TDDD16/B44, P Fritzon, 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
```

TDDD16/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

5.11

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

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

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

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

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

TDDD16/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

5.12

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:

$$\begin{array}{l} \text{arglist} \rightarrow () \\ | (\text{args}) \\ \text{or} \\ A \rightarrow a b \\ | a c \end{array}$$

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

TODD16/B44, P.Fritzson, 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: β^*
- $[\beta]$ same as the regular expression: $\beta \mid \epsilon$
- $()$ 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 β may not be preceded by A) in EBNF is rewritten:

- $A \rightarrow \beta \{\alpha\}$

TODD16/B44, P.Fritzson, C.Kessler, IDA, LIU, 2010.

5.14



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

- $A \rightarrow A \alpha \mid \beta$ (where β may not be preceded by A) is rewritten to

$$\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}$$

Generally:

- $A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \dots \mid A \alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$ (where β_1, β_2, \dots may not be preceded by A) is rewritten to:

$$\begin{array}{l} A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \\ A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon \end{array}$$

TODD16/B44, P.Fritzson, C.Kessler, IDA, LIU, 2010.

5.15



2. Left Factoring Using () or []

Original Grammar:

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

Original Grammar:

$$A \rightarrow ab \mid ac$$

Solution using EBNF:

$$A \rightarrow a(b \mid c)$$

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

Solution using rewriting:

$$<\text{stmt}> \rightarrow \text{if } <\text{expr}> \text{ then } <\text{stmt}> <\text{rest-if}>$$

$$<\text{rest-if}> \rightarrow \text{else } <\text{stmt}> \mid \epsilon$$

TODD16/B44, P.Fritzson, 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

TODD16/B44, P.Fritzson, C.Kessler, IDA, LIU, 2010.

5.17

Small Rewriting Grammar Exercise



TODD16/B44, P.Fritzson, 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> <constdef>
| <constdef>
<constdef> → id = number ;
<vardecl> → VAR <vardeflist>
| ε
<vardeflist> → <vardeflist> <idlist> : <type> ;
| <idlist> : <type> ;
<idlist> → <idlist> , id
| id
<type> → integer
| real

```

TDDD16/B44, P Fritzson, 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 ) ;

```

TDDD16/B44, P Fritzson, 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();    if (token == CONST) {
} /* declarations */        scan();
                           If (token == id)
                           constdef();
                           else
                           error("Missing id after CONST");
                           while (token == id) constdef();
                           }
} /* constdecl */

```

TDDD16/B44, P Fritzson, C. Kessler, IDA, LIU, 2010.

5.21

Pascal Declarations Parser cont 1

```

void constdef()
/* <constdef> → id = number ; */
{
  scan(); /* consume ID, get next token */
  if (token == '=') {
    scan();
    if (token == ID)
      constdef();
    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();
      scan();
    }
  }
} /* vardecl */

```

TDDD16/B44, P Fritzson, 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 ',' ");
  }
  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 */

```

TDDD16/B44, P Fritzson, C. Kessler, IDA, LIU, 2010.

5.23

TDDD16 Compilers and interpreters
TDBB44 Compiler Construction



LL Parsing Issues Beyond Recursive Descent

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

Peter Fritzson, Christoph Kessler,
IDA, Linköpings universitet, 2010.

LL(k)

■ Given:

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

■ G is in $\text{LL}(k)$ if:

for any two leftmost derivations

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

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.

TODD16/B44, P. Fritzson, 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
```

TODD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

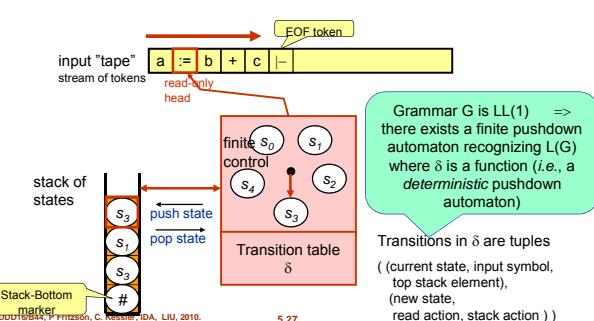
5.26



Automaton Model for Parsing Context-Free Languages

Finite pushdown automaton (FPA)

- a finite automaton with a stack of states



TODD16/B44, P. Fritzson, 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 $S' \rightarrow S \mid -$ ($\mid -$ 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 $[A \rightarrow \alpha \cdot \beta \gamma]$: read input symbol, i.e., $[A \rightarrow \alpha \cdot \beta \gamma] \rightarrow [A \rightarrow \alpha \beta \cdot \gamma]$
 - at $[A \rightarrow \alpha \cdot B \gamma]$: push $[A \rightarrow \alpha B \cdot \gamma]$, determine new production $B \rightarrow \beta$ and start from $[B \rightarrow \beta \cdot \gamma]$ Prediction!
 - at $[B \rightarrow \beta \cdot \gamma]$: pop state $[A \rightarrow \alpha B \cdot \gamma]$ to restore context (if #, error)

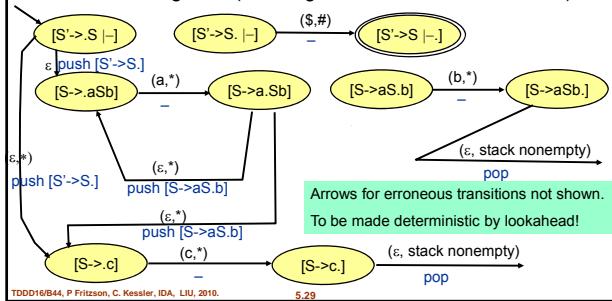
TODD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

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):



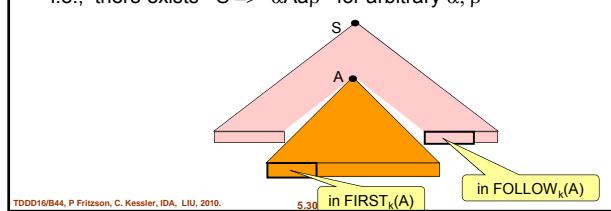
TODD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

5.29



FIRST and FOLLOW

- For a sentential form α in $(N \cup S)^*$, $\text{FIRST}(\alpha)$ denotes the set of all terminals with can be *first* in a string derived from α .
- For a nonterminal A in N , $\text{FOLLOW}(A)$ 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 α, β



TODD16/B44, P. Fritzson, C. Kessler, IDA, LIU, 2010.

5.30



Small FIRST and FOLLOW Exercise

TODD16/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

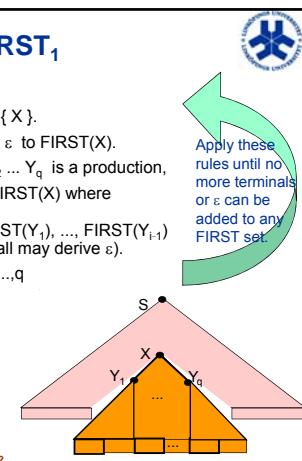
5.31

Computing FIRST = FIRST₁

For all grammar symbols X:

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

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)
TOD FIRST(S) = {a, b, 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- ϵ symbols of FIRST(X_1).
- If ϵ in FIRST(X_1), add also all non- ϵ symbols of FIRST(X_2), otherwise done.
- If ϵ also in FIRST(X_2), add also all non- ϵ symbols of FIRST(X_3), otherwise done.
- ...
- If ϵ also in FIRST(X_n), add ϵ 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}

TODD16/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

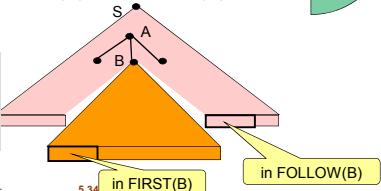
5.33

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

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' : { $S' \rightarrow S | -; S \rightarrow aSb; S \rightarrow c$ }

Diagram showing states and transitions for the FPA. States include $[S' \rightarrow S | -]$, $[S' \rightarrow S.S]$, $[S \rightarrow S | -]$, $[S \rightarrow aSb]$, $[S \rightarrow a.Sb]$, $[S \rightarrow aS.b]$, $[S \rightarrow aSb.]$, $[S \rightarrow c]$, and $[S \rightarrow c.]$. Transitions are labeled with input symbols (a, *, |-, #) and actions (push stack operations). A pink box highlights disambiguation: $\text{FIRST}_1(aSb) = \{a\}$ and $\text{FIRST}_1(c) = \{c\}$.

TODD16/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

5.35

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

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

TODD16/B44, P Fritzon, C. Kessler, IDA, LIU, 2010.

5.36

General Approach: Predictive Parsing



At any production $A \rightarrow \alpha$

- If ϵ 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., ϵ 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)

TODD16/B44, P.Fritzson, 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 FIRST and FOLLOW sets
- FIRST and FOLLOW are also useful for syntax error recovery

TODD16/B44, P.Fritzson, 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 state
- in a recursive descent parser.

TODD16/B44, P.Fritzson, C.Kessler, IDA, LIU, 2010.

5.39