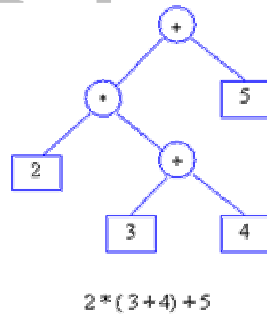


Chapter 4

(b) parsing

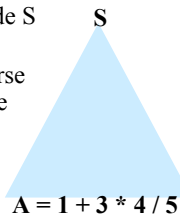


Parsing

- A grammar describes the strings of tokens that are syntactically legal in a PL
- A *recogniser* simply accepts or rejects strings.
- A generator produces sentences in the language described by the grammar
- A *parser* constructs a derivation or parse tree for a sentence (if possible)
- Two common types of parsers:
 - bottom-up or data driven
 - top-down or hypothesis driven
- A *recursive descent parser* is a way to implement a top-down parser that is particularly simple.

Top down vs. bottom up parsing

- The parsing problem is to connect the root node S with the tree leaves, the input
- **Top-down parsers:** starts constructing the parse tree at the top (root) of the parse tree and move down towards the leaves. Easy to implement by hand, but work with restricted grammars. examples:
 - Predictive parsers (e.g., LL(k))
- **Bottom-up parsers:** build the nodes on the bottom of the parse tree first. Suitable for automatic parser generation, handle a larger class of grammars. examples:
 - shift-reduce parser (or LR(k) parsers)
- Both are general techniques that can be made to work for all languages (but not all grammars!).



Top down vs. bottom up parsing

- Both are general techniques that can be made to work for all languages (but not all grammars!).
- Recall that a given language can be described by several grammars.
- Both of these grammars describe the same language

$E \rightarrow E + Num$	$E \rightarrow Num + E$
$E \rightarrow Num$	$E \rightarrow Num$
- The first one, with its left recursion, causes problems for top down parsers.
- For a given parsing technique, we may have to transform the grammar to work with it.

Parsing complexity

- How hard is the parsing task?
 - Parsing an arbitrary Context Free Grammar is $O(n^3)$, e.g., it can take time proportional the cube of the number of symbols in the input. This is bad!
 - If we constrain the grammar somewhat, we can always parse in linear time. This is good!
 - Linear-time parsing
 - LL parsers
 - Recognize LL grammar
 - Use a top-down strategy
 - LR parsers
 - Recognize LR grammar
- Use a bottom-up strategy

- **LL(n) : Left to right, Leftmost derivation, look ahead at most n symbols.**
- **LR(n) : Left to right, Right derivation, look ahead at most n symbols.**

Top Down Parsing Methods

- Simplest method is a full-backup *recursive descent* parser.
- Write recursive recognizers (subroutines) for each grammar rule
 - If rules succeeds perform some action (I.e., build a tree node, emit code, etc.)
 - If rule fails, return failure. Caller may try another choice or fail
 - On failure it “backs up”
- Problems
 - When going forward, the parser consumes tokens from the input, so what happens if we have to back up?
 - Backup is, in general, inefficient
 - Grammar rules which are left-recursive lead to non-termination

Recursive Decent Parsing Example

Example: For the grammar:

```
<term> -> <factor> {(*|/)<factor>}
```

We could use the following recursive descent parsing subprogram (this one is written in C)

```
void term() {  
    factor(); /* parse first factor*/  
    while (next_token == ast_code ||  
           next_token == slash_code) {  
        lexical(); /* get next token */  
        factor(); /* parse next factor */  
    }  
}
```

Informal recursive descent parsing

Problems

- Some grammars cause problems for top down parsers.
- Top down parsers do not work with left-recursive grammars.
 - E.g., one with a rule like: $E \rightarrow E + T$
 - We can transform a left-recursive grammar into one which is not.
- A top down grammar can limit backtracking if it only has one rule per non-terminal
 - The technique of factoring can be used to eliminate multiple rules for a non-terminal.

Left-recursive grammars

- A grammar is left recursive if it has rules like
 $X \rightarrow X \beta$
Or if it has indirect left recursion, as in
 $X \rightarrow X \beta$
- Why is this a problem?
- Consider
 $E \rightarrow E + \text{Num}$
 $E \rightarrow \text{Num}$
- We can manually or automatically rewrite a grammar to remove left-recursion, making it suitable for a top-down parser.

Elimination of Left Recursion

- Consider the left-recursive grammar
 $S \rightarrow S \alpha \mid \beta$
- S generates all strings starting with a β and followed by a number of α
- Can rewrite using right-recursion
 $S \rightarrow \beta S'$
 $S' \rightarrow \alpha S' \mid \epsilon$

More Elimination of Left-Recursion

- In general
 $S \rightarrow S \alpha_1 \mid \dots \mid S \alpha_n \mid \beta_1 \mid \dots \mid \beta_m$
- All strings derived from S start with one of β_1, \dots, β_m and continue with several instances of $\alpha_1, \dots, \alpha_n$
- Rewrite as
 $S \rightarrow \beta_1 S' \mid \dots \mid \beta_m S'$
 $S' \rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \epsilon$

General Left Recursion

- The grammar
$$S \rightarrow A \alpha \mid \delta$$
$$A \rightarrow S \beta$$
is also left-recursive because
$$S \rightarrow^+ S \beta \alpha$$
where \rightarrow^+ means “can be rewritten in one or more steps”
- This indirect left-recursion can also be automatically eliminated

Summary of Recursive Descent

- Simple and general parsing strategy
 - Left-recursion must be eliminated first
 - ... but that can be done automatically
- Unpopular because of backtracking
 - Thought to be too inefficient
- In practice, backtracking is eliminated by restricting the grammar, allowing us to successfully *predict* which rule to use.

Predictive Parser

- A **predictive parser** uses information from the **first terminal symbol** of each expression to decide which production to use.
- A predictive parser is also known as an **LL(*k*)** parser because it does a **Left-to-right parse**, a **Leftmost-derivation**, and ***k*-symbol lookahead**.
- A grammar in which it is possible to decide which production to use examining only the first token (as in the previous example) are called **LL(1)**
- LL(1) grammars are widely used in practice.
 - The syntax of a PL can be adjusted to enable it to be described with an LL(1) grammar.

Predictive Parser

Example: consider the grammar

```
S → if E then S else S
S → begin S L
S → print E
L → end
L → ; S L
E → num = num
```

An *S* expression starts either with an IF, BEGIN, or PRINT token, and an *L* expression start with an END or a SEMICOLON token, and an *E* expression has only one production.

LL(k) and LR(k) parsers

- Two important classes of parsers are called LL(k) parsers and LR(k) parsers.
- The name LL(k) means:
 - L - Left-to-right scanning of the input
 - L - Constructing leftmost derivation
 - k – max number of input symbols needed to select a parser action
- The name LR(k) means:
 - L - Left-to-right scanning of the input
 - R - Constructing rightmost derivation in reverse
 - k – max number of input symbols needed to select a parser action
- So, a LL(1) parser never needs to “look ahead” more than one input token to know what parser production to apply.

Predictive Parsing and Left Factoring

- Consider the grammar
$$E \rightarrow T + E \mid T$$
$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$
- Hard to predict because
 - For T, two productions start with *int*
 - For E, it is not clear how to predict which rule to use
- A grammar must be left-factored before use for predictive parsing
- Left-factoring involves rewriting the rules so that, if a non-terminal has more than one rule, each begins with a terminal.

Left-Factoring Example

- Consider the grammar
$$E \rightarrow T + E \mid T$$
$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$
- Factor out common prefixes of productions
$$E \rightarrow T X$$
$$X \rightarrow + E \mid \epsilon$$
$$T \rightarrow (E) \mid \text{int} Y$$
$$Y \rightarrow * T \mid \epsilon$$

Left Factoring

- Consider a rule of the form
$$A \rightarrow a B_1 \mid a B_2 \mid a B_3 \mid \dots a B_n$$
- A top down parser generated from this grammar is not efficient as it requires backtracking.
- To avoid this problem we left factor the grammar.
 - collect all productions with the same left hand side and begin with the same symbols on the right hand side
 - combine the common strings into a single production and then append a new non-terminal symbol to the end of this new production
 - create new productions using this new non-terminal for each of the suffixes to the common production.
- After left factoring the above grammar is transformed into:
$$A \rightarrow a A_1$$
$$A_1 \rightarrow B_1 \mid B_2 \mid B_3 \dots B_n$$

Using Parsing Tables

- LL(1) means that for each non-terminal and token there is only one production
- Can be specified via 2D tables
 - One dimension for current non-terminal to expand
 - One dimension for next token
 - A table entry contains one production
- Method similar to recursive descent, except
 - For each non-terminal S
 - We look at the next token a
 - And chose the production shown at [S,a]
- We use a stack to keep track of pending non-terminals
- We reject when we encounter an error state
- We accept when we encounter end-of-input

LL(1) Parsing Table Example

- Left-factored grammar
 - $E \rightarrow T X$ $X \rightarrow + E \mid \epsilon$
 - $T \rightarrow (E) \mid \text{int } Y$ $Y \rightarrow * T \mid \epsilon$
- The LL(1) parsing table:

	int	*	+	()	\$
E	$T X$			$T X$		
X			$+ E$		ϵ	ϵ
T	$\text{int } Y$			(E)		
Y		$* T$	ϵ		ϵ	ϵ

LL(1) Parsing Table Example

- Consider the [E, int] entry
 - “When current non-terminal is E and next input is *int*, use production $E \rightarrow T X$
 - This production can generate an *int* in the first place
- Consider the [Y, +] entry
 - “When current non-terminal is Y and current token is +, get rid of Y”
 - Y can be followed by + only in a derivation in which $Y \rightarrow \epsilon$
- Blank entries indicate error situations
 - Consider the [E,*] entry
 - “There is no way to derive a string starting with * from non-terminal E”

LL(1) Parsing Algorithm

```

initialize stack = <S $> and next
repeat
  case stack of
    <X, rest> : if  $T[X, \text{next}] = Y_1 \dots Y_n$ 
                then stack  $\leftarrow \langle Y_1 \dots Y_n, \text{rest} \rangle$ ;
    rest>;
                else error ();
    <t, rest> : if  $t == \text{next} ++$ 
                then stack  $\leftarrow \langle \text{rest} \rangle$ ;
                else error ();
until stack == < >
    
```

LL(1) Parsing Example

Stack	Input	Action
E \$	int * int \$	T X
T X \$	int * int \$	int Y
int Y X \$	int * int \$	terminal
Y X \$	* int \$	* T
* T X \$	* int \$	terminal
T X \$	int \$	int Y
int Y X \$	int \$	terminal
Y X \$	\$	ϵ
X \$	\$	ϵ
\$	\$	ACCEPT

Constructing Parsing Tables

- LL(1) languages are those defined by a parsing table for the LL(1) algorithm
- No table entry can be multiply defined
- We want to generate parsing tables from CFG
- If $A \rightarrow \alpha$, where in the line of A we place α ?
- In the column of t where t can start a string derived from α
 - $\alpha \rightarrow^* t \beta$
 - We say that $t \in \text{First}(\alpha)$
- In the column of t if α is ϵ and t can follow an A
 - $S \rightarrow^* \beta A t \delta$
 - We say $t \in \text{Follow}(A)$

Computing First Sets

Definition: $\text{First}(X) = \{t \mid X \rightarrow^* t\alpha\} \cup \{\epsilon \mid X \rightarrow^* \epsilon\}$

Algorithm sketch (see book for details):

1. for all terminals t do $\text{First}(t) \leftarrow \{t\}$
2. for each production $X \rightarrow \epsilon$ do $\text{First}(X) \leftarrow \{\epsilon\}$
3. if $X \rightarrow A_1 \dots A_n \alpha$ and $\epsilon \in \text{First}(A_i)$, $1 \leq i \leq n$ do
 - add $\text{First}(\alpha)$ to $\text{First}(X)$
4. for each $X \rightarrow A_1 \dots A_n$ s.t. $\epsilon \in \text{First}(A_i)$, $1 \leq i \leq n$ do
 - add ϵ to $\text{First}(X)$
5. repeat steps 4 & 5 until no First set can be grown

First Sets. Example

- Recall the grammar

$$\begin{array}{ll} E \rightarrow T X & X \rightarrow + E \mid \epsilon \\ T \rightarrow (E) \mid \text{int } Y & Y \rightarrow * T \mid \epsilon \end{array}$$

- First sets

$$\begin{array}{ll} \text{First}(()) = \{ (\} & \text{First}(T) = \{ \text{int}, (\} \\ \text{First}()) = \{) \} & \text{First}(E) = \{ \text{int}, (\} \\ \text{First}(\text{int}) = \{ \text{int} \} & \text{First}(X) = \{ +, \epsilon \} \\ \text{First}(+) = \{ + \} & \text{First}(Y) = \{ *, \epsilon \} \\ \text{First}(*) = \{ * \} & \end{array}$$

Computing Follow Sets

- Definition:

$$\text{Follow}(X) = \{ t \mid S \rightarrow^* \beta X t \delta \}$$

- Intuition

- If S is the start symbol then $\$ \in \text{Follow}(S)$
- If $X \rightarrow A B$ then $\text{First}(B) \subseteq \text{Follow}(A)$ and $\text{Follow}(X) \subseteq \text{Follow}(B)$
- Also if $B \rightarrow^* \epsilon$ then $\text{Follow}(X) \subseteq \text{Follow}(A)$

Computing Follow Sets

Algorithm sketch:

1. $\text{Follow}(S) \leftarrow \{ \$ \}$
2. For each production $A \rightarrow \alpha X \beta$
 - add $\text{First}(\beta) - \{ \epsilon \}$ to $\text{Follow}(X)$
3. For each $A \rightarrow \alpha X \beta$ where $\epsilon \in \text{First}(\beta)$
 - add $\text{Follow}(A)$ to $\text{Follow}(X)$
 - repeat step(s) ___ until no Follow set grows

Follow Sets. Example

- Recall the grammar

$$E \rightarrow T X \qquad X \rightarrow + E \mid \epsilon$$

$$T \rightarrow (E) \mid \text{int } Y \qquad Y \rightarrow * T \mid \epsilon$$

- Follow sets

$$\text{Follow}(+) = \{ \text{int}, (\}$$

$$\text{Follow}(*) = \{ \text{int}, (\}$$

$$\text{Follow}(() = \{ \text{int}, (\}$$

$$\text{Follow}(E) = \{ \}, \$ \}$$

$$\text{Follow}(X) = \{ \$,) \}$$

$$\text{Follow}(T) = \{ +,) , \$ \}$$

$$\text{Follow}()) = \{ +,) , \$ \}$$

$$\text{Follow}(Y) = \{ +,) , \$ \}$$

$$\text{Follow}(\text{int}) = \{ *, +,) , \$ \}$$

Constructing LL(1) Parsing Tables

- Construct a parsing table T for CFG G
- For each production $A \rightarrow \alpha$ in G do:
 - For each terminal $t \in \text{First}(\alpha)$ do
 - $T[A, t] = \alpha$
 - If $\epsilon \in \text{First}(\alpha)$, for each $t \in \text{Follow}(A)$ do
 - $T[A, t] = \alpha$
 - If $\epsilon \in \text{First}(\alpha)$ and $\$ \in \text{Follow}(A)$ do
 - $T[A, \$] = \alpha$

Notes on LL(1) Parsing Tables

- If any entry is multiply defined then G is not LL(1)
 - If G is ambiguous
 - If G is left recursive
 - If G is not left-factored
- Most **programming language** grammars are not LL(1)
- There are tools that build LL(1) tables

Bottom-up Parsing

- YACC uses bottom up parsing. There are two important operations that bottom-up parsers use. They are namely shift and reduce.
 - (In abstract terms, we do a simulation of a Push Down Automata as a finite state automata.)
- Input: given string to be parsed and the set of productions.
- Goal: Trace a rightmost derivation in reverse by starting with the input string and working backwards to the start symbol.

Algorithm

1. Start with an empty stack and a full input buffer. (The string to be parsed is in the input buffer.)
 2. Repeat until the input buffer is empty and the stack contains the start symbol.
 - a. **Shift** zero or more input symbols onto the stack from input buffer until a handle (beta) is found on top of the stack. If no handle is found report syntax error and exit.
 - b. **Reduce** handle to the nonterminal A. (There is a production $A \rightarrow \text{beta}$)
 3. **Accept** input string and return some representation of the derivation sequence found (e.g., **parse tree**)
- The four key operations in bottom-up parsing are **shift, reduce, accept** and **error**.
 - Bottom-up parsing is also referred to as shift-reduce parsing.
 - Important thing to note is to know when to shift and when to reduce and to which reduce.

Example of Bottom-up Parsing

STACK	INPUT BUFFER	ACTION	
\$	num1+num2*num3\$	shift	
\$num1	+num2*num3\$	reduc	E -> E+T
\$F	+num2*num3\$	reduc	T
\$T	+num2*num3\$	reduc	E-T
\$E	+num2*num3\$	shift	T -> T*F
\$E+	num2*num3\$	shift	F
\$E+num2	*num3\$	reduc	T/F
\$E+F	*num3\$	reduc	F -> (E)
\$E+T	*num3\$	shift	id
E+T*	num3\$	shift	-E
E+T*num3	\$	reduc	num
E+T*F	\$	reduc	
E+T	\$	reduc	
E	\$	accept	