

Concepts Introduced in Chapter 4

- Grammars
 - Context-Free Grammars
 - Derivations and Parse Trees
 - Ambiguity, Precedence, and Associativity
- Top Down Parsing
 - Recursive Descent, LL
- Bottom Up Parsing
 - SLR, LR, LALR
- Yacc
- Error Handling

Grammars

$$G = (N, T, P, S)$$

1. N is a finite set of nonterminal symbols
2. T is a finite set of terminal symbols
3. P is a finite subset of
$$(N \cup T)^* N (N \cup T)^* \times (N \cup T)^*$$
An element $(\alpha, \beta) \in P$ is written as
$$\alpha \rightarrow \beta$$
and is called a production.
4. S is a distinguished symbol in N and is called the start symbol.

Example of a Grammar

$\text{block} \rightarrow \underline{\text{begin}} \text{ opt_stmts } \underline{\text{end}}$
 $\text{opt_stmts} \rightarrow \text{stmt_list} \mid \epsilon$
 $\text{stmt_list} \rightarrow \text{stmt_list}; \text{stmt} \mid \text{stmt}$

Advantages of Using Grammars

- Provides a precise, syntactic specification of a programming language.
- For some classes of grammars, tools exist that can automatically construct an efficient parser.
- These tools can also detect syntactic ambiguities and other problems automatically.
- A compiler based on a grammatical description of a language is more easily maintained and updated.

Role of a Parser in a Compiler

- Detects and reports any syntax errors.
- Produces a parse tree from which intermediate code can be generated.

Conventions Used for Specifying Grammars in the Text

- terminals
 - lower case letters early in the alphabet (a, b, c)
 - punctuation and operator symbols [(,), ', +, -]
 - digits
 - boldface words (**if**, **then**)
- nonterminals
 - uppercase letters early in the alphabet (A, B, C)
 - S is the start symbol
 - lower case words

Conventions Used for Specifying Grammars in the Text (cont.)

- grammar symbols (nonterminals or terminals)
 - upper case letters late in the alphabet (X, Y, Z)
- strings of terminals
 - lower case letters late in the alphabet (u, v, ..., z)
- sentential form (string of grammar symbols)
 - lower case Greek letters (α , β , γ)

Chomsky Hierarchy

A grammar is said to be

1. regular if productions in P are all right-linear or are all left-linear

- a. right-linear

$$A \rightarrow wB \text{ or } A \rightarrow w$$

- b. left-linear

$$A \rightarrow Bw \text{ or } A \rightarrow w$$

where $A, B \in N$ and $w \in T^*$

Recognized by a finite automata (FA).

Chomsky Hierarchy (cont)

2. context-free if each production in P is of the form

$$A \rightarrow \alpha \quad \text{where } A \in N \text{ and } \alpha \in (N \cup T)^*$$

Recognized by a pushdown automata (PDA).

3. context-sensitive if each production in P is of the form

$$\alpha \rightarrow \beta \quad \text{where } |\alpha| \leq |\beta|$$

Recognized by a linear bounded automata (LBA).

4. unrestricted if each production in P is of the form

$$\alpha \rightarrow \beta \quad \text{where } \alpha \neq \epsilon$$

Recognized by a Turing machine.

Derivation

Derivation - a sequence of replacements from the start symbol in a grammar by applying productions

Example: $E \rightarrow E + E$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow - E$$

$$E \rightarrow \text{id}$$

Derive - (id) from the grammar

$$E \Rightarrow - E \Rightarrow - (E) \Rightarrow - (\text{id})$$

thus E derives - (id)

$$\text{or } E \Rightarrow - (\text{id})$$

Derivation (cont.)

leftmost derivation - each step replaces the leftmost nonterminal

Derive $\text{id} + \text{id} * \text{id}$ using leftmost derivation

$$E \Rightarrow E + E \Rightarrow \text{id} + E \Rightarrow \text{id} + E * E \Rightarrow \text{id} + \text{id} * E \Rightarrow \text{id} + \text{id} * \text{id}$$

$L(G)$ - language generated by the grammar G

sentence of G - if $S \Rightarrow^+ w$, where w is a string of terminals in $L(G)$

sentential form - if $S \Rightarrow^* \alpha$, where α may contain nonterminals

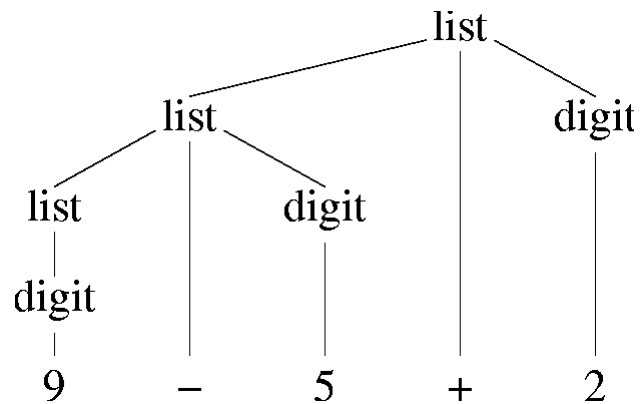
Parse Tree

A parse tree pictorially shows how the start symbol of a grammar derives a specific string in the language.

Given a context-free grammar, a parse tree has the properties:

1. The root is labeled by the start symbol.
2. Each leaf is labeled by a token or ϵ .
3. Each interior node is labeled by a nonterminal.
4. If A is a nonterminal labeling some interior node and $X_1, X_2, X_3, \dots, X_n$ are the labels of the children of that node from left to right, then $A \rightarrow X_1, X_2, X_3, \dots, X_n$ is a production of the grammar.

Example of a Parse Tree



$\text{list} \rightarrow \text{list} + \text{digit} \mid \text{list} - \text{digit} \mid \text{digit}$

Parse Tree (cont.)

Yield - the leaves of the parse tree read from left to right or the string derived from the nonterminal at the root of the parse tree.

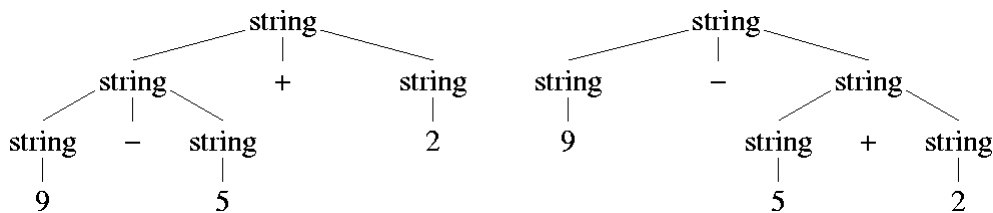
An ambiguous grammar is one that can generate two or more parse trees that yield the same string.

Example of an Ambiguous Grammar

$\text{string} \rightarrow \text{string} + \text{string}$

$\text{string} \rightarrow \text{string} - \text{string}$

$\text{string} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$



a. $\text{string} \rightarrow \text{string} + \text{string} \rightarrow \text{string} - \text{string} + \text{string}$
 $\rightarrow 9 - \text{string} + \text{string} \rightarrow 9 - 5 + \text{string} \rightarrow 9 - 5 + 2$

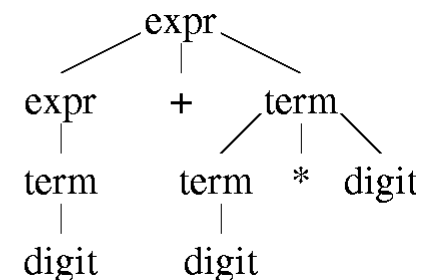
b. $\text{string} \rightarrow \text{string} - \text{string} \rightarrow 9 - \text{string}$
 $\rightarrow 9 - \text{string} + \text{string} \rightarrow 9 - 5 + \text{string} \rightarrow 9 - 5 + 2$

Precedence

By convention

$9 + 5 * 2$ $*$ has higher precedence than $+$ because it takes its operands before $+$

$\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term}$
 $\text{term} \rightarrow \text{term} * \text{digit} \mid \text{digit}$



Precedence (cont.)

Different operators have the same precedence when they are defined as alternative productions of the same nonterminal.

$\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{expr} - \text{term} \mid \text{term}$
 $\text{term} \rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor}$
 $\text{factor} \rightarrow \text{digit} \mid (\text{expr})$

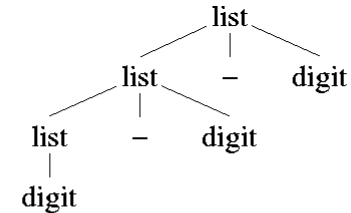
Associativity

By convention

$9 - 5 - 2$ left (operand with $-$ on both sides, the operation on the left is performed first)

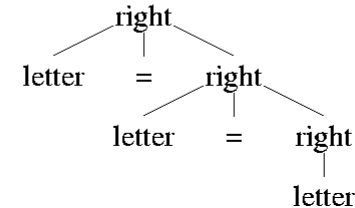
$a = b = c$ right (operand with $=$ on both sides, the operation on the right is performed first)

$\text{list} \rightarrow \text{list} - \text{digit}$
 $\text{list} \rightarrow \text{digit}$



grows to the left

$\text{right} \rightarrow \text{letter} = \text{right}$
 $\text{right} \rightarrow \text{letter}$



grows to the right

Eliminating Ambiguity

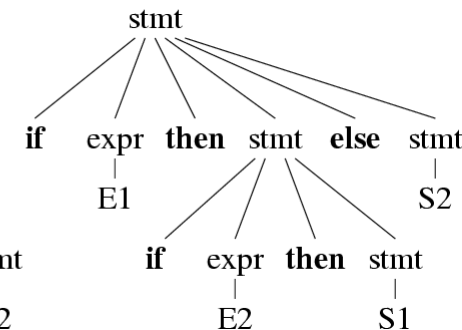
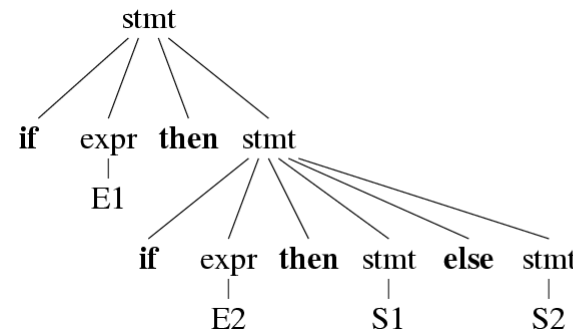
- Sometimes ambiguity can be eliminated by rewriting a grammar.

$\text{stmt} \rightarrow \text{if expr then stmt}$
 $\quad \mid \text{if expr then stmt else stmt}$
 $\quad \mid \text{other}$

- How do we parse:

if E1 then if E2 then S1 else S2

Two Parse Trees for “if E1 then if E2 then S1 else S2”



Eliminating Ambiguity (cont.)

stmt \rightarrow matched_stmt
 | unmatched_stmt

matched_stmt \rightarrow **if** expr **then** matched_stmt **else** matched_stmt
 | other

unmatched_stmt \rightarrow **if** expr **then** stmt
 | **if** expr **then** matched_stmt **else** unmatched_stmt

Parsing

universal

top-down

 recursive descent

 LL

bottom-up

 operator precedence

 LR

 SLR

 canonical LR

 LALR

Top-Down vs Bottom-Up Parsing

- top-down
 - Have to eliminate left recursion in the grammar.
 - Have to left factor the grammar.
 - Resulting grammars are harder to read and understand.
- bottom-up
 - Difficult to implement by hand, so a tool is needed.

Top-Down Parsing

Starts at the root and proceeds towards the leaves.

Recursive-Descent Parsing - a recursive procedure is associated with each nonterminal in the grammar.

Example

type \rightarrow simple | \uparrow id | array [simple] of type

simple \rightarrow integer | char | num dotdot num

Example of Recursive Descent Parsing

```
void type() {
    if ( lookahead == INTEGER || lookahead == CHAR ||
        lookahead == NUM)
        simple();
    else if (lookahead == '^') {
        match('^');
        match(ID);
    }
    else if (lookahead == ARRAY) {
        match(ARRAY);
        match('[');
        simple();
        match(']');
        match(OF);
        type();
    }
    else
        error();
}
```

Example of Recursive Descent Parsing (cont.)

```
void simple() {
    if (lookahead == INTEGER)
        match(INTEGER);
    else if (lookahead == CHAR)
        match(CHAR);
    else if (lookahead == NUM) {
        match(NUM);
        match(DOTDOT);
        match(NUM);
    }
    else
        error();
}

void match(token t)
{
    if (lookahead == t)
        lookahead = nexttoken();
    else
        error();
}
```

Top-Down Parsing (cont.)

- Predictive parsing needs to know what first symbols can be generated by the right side of a production.
- $\text{FIRST}(\alpha)$ - the set of tokens that appear as the first symbols of one or more strings generated from α . If α is ϵ or can generate ϵ , then ϵ is also in $\text{FIRST}(\alpha)$.

- Given a production

$$A \rightarrow \alpha \mid \beta$$

predictive parsing requires $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ to be disjoint.

Eliminating Left Recursion

- Recursive descent parsing loops forever on left recursion.
- Immediate Left Recursion

$$\text{Replace } A \rightarrow A\alpha \mid \beta \text{ with } \begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \epsilon \end{aligned}$$

Example:

	<u>A</u>	<u>α</u>	<u>β</u>
$E \rightarrow E + T \mid T$	E	+T	T
$T \rightarrow T * F \mid F$	T	*F	F
$F \rightarrow (E) \mid \text{id}$			

becomes

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \end{aligned}$$

Eliminating Left Recursion (cont.)

- What if a grammar is not immediately left recursive?

$$A \Rightarrow A\alpha$$

- For instance:

$$A \rightarrow B\alpha_1 \mid \alpha_4$$

$$B \rightarrow C\alpha_2$$

$$C \rightarrow A\alpha_3$$

- For example:

$$A \Rightarrow B\alpha_1 \Rightarrow C\alpha_2\alpha_1 \Rightarrow A\alpha_3\alpha_2\alpha_1$$

Eliminating Left Recursion (cont.)

In general, to eliminate left recursion given A_1, A_2, \dots, A_n

for $i = 1$ to n do

for $j = 1$ to $i-1$ do

replace each $A_i \rightarrow A_j \gamma$ with $A_i \rightarrow \delta_1 \gamma \mid \dots \mid \delta_k \gamma$

where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are the current A_j productions

end for

eliminate immediate left recursion in the A_i productions

eliminate ϵ transitions in the A_i productions

end for

This fails only if cycles ($A \Rightarrow A$) or $A \rightarrow \epsilon$ for some A .

Example of Eliminating Left Recursion

- $X \rightarrow YZ \mid a$
- $Y \rightarrow ZX \mid Xb$
- $Z \rightarrow XY \mid ZZ \mid a$

$$A_1 = X \quad A_2 = Y \quad A_3 = Z$$

$i = 1$ (eliminate immediate left recursion)
nothing to do

Example of Eliminating Left Recursion (cont.)

$i = 2, j = 1$

$$Y \rightarrow Xb \Rightarrow Y \rightarrow ZX \mid YZb \mid ab$$

now eliminate immediate left recursion

$$Y \rightarrow ZXY' \mid abY'$$

$$Y' \rightarrow ZbY' \mid \epsilon$$

now eliminate ϵ transitions

$$Y \rightarrow ZXY' \mid abY' \mid ZX \mid ab$$

$$Y' \rightarrow ZbY' \mid Zb$$

$i = 3, j = 1$

$$Z \rightarrow XY \Rightarrow Z \rightarrow YZY \mid aY \mid ZZ \mid a$$

Example of Eliminating Left Recursion (cont.)

$i = 3, j = 2$

$$Z \rightarrow YZY \Rightarrow Z \rightarrow ZXY'ZY \mid abY'ZY \mid ZXZY \\ \mid abZY \mid aY \mid ZZ \mid a$$

now eliminate immediate left recursion

$$Z \rightarrow abY'ZY Z' \mid abZY Z' \mid aYZ' \mid aZ'$$

$$Z' \rightarrow XY'ZY Z' \mid XZY Z' \mid ZZ' \mid \epsilon$$

eliminate ϵ transitions

$$Z \rightarrow abY'ZY Z' \mid abY'ZY \mid abZY Z' \mid abZY \mid aY \\ \mid aYZ' \mid aZ' \mid a$$

$$Z' \rightarrow XY'ZY Z' \mid XY'ZY \mid XZY Z' \mid XZY \mid ZZ' \\ \mid Z$$

Left-Factoring

$$A \rightarrow \alpha\beta \mid \alpha\gamma \Rightarrow A \rightarrow \alpha A' \\ A' \rightarrow \beta \mid \gamma$$

Example:

Left factor

$$\text{stmt} \rightarrow \text{if cond then stmt else stmt} \\ \mid \text{if cond then stmt}$$

becomes

$$\text{stmt} \rightarrow \text{if cond then stmt E} \\ E \rightarrow \text{else stmt} \mid \epsilon$$

Grammars must be left factored for predictive parsing so we will know which production to choose.

Nonrecursive Predictive Parsing

- Instead of recursive descent, predictive parsing can be table-driven and use an explicit stack. It uses
 1. a stack of grammar symbols (\$ on bottom)
 2. a string of input tokens (\$ on end)
 3. a parsing table [NT, T] of productions

Algorithm for Nonrecursive Predictive Parsing

1. If top == input == \$ then accept
2. If top == input then
 - pop top off the stack
 - advance to next input symbol
 - goto 1
3. If top is nonterminal
 - fetch $M[\text{top}, \text{input}]$
 - If a production
 - replace top with rhs of production
 - Else
 - parse fails
 - goto 1
4. Parse fails

First

$\text{FIRST}(\alpha)$ = the set of terminals that begin strings derived from α . If α is ϵ or generates ϵ , then ϵ is also in $\text{FIRST}(\alpha)$.

1. If X is a terminal then $\text{FIRST}(X) = \{X\}$
2. If $X \rightarrow a\alpha$, add a to $\text{FIRST}(X)$
3. If $X \rightarrow \epsilon$, add ϵ to $\text{FIRST}(X)$
4. If $X \rightarrow Y_1, Y_2, \dots, Y_k$ and $Y_1, Y_2, \dots, Y_{i-1} \Rightarrow \epsilon$ where $i \leq k$
Add every non ϵ in $\text{FIRST}(Y_i)$ to $\text{FIRST}(X)$
If $Y_1, Y_2, \dots, Y_k \Rightarrow \epsilon$, add ϵ to $\text{FIRST}(X)$

FOLLOW

$\text{FOLLOW}(A)$ = the set of terminals that can immediately follow A in a sentential form.

1. If S is the start symbol, add $\$$ to $\text{FOLLOW}(S)$
2. If $A \rightarrow \alpha B \beta$, add $\text{FIRST}(\beta) - \{\epsilon\}$ to $\text{FOLLOW}(B)$
3. If $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ and $\beta \Rightarrow \epsilon$, add $\text{FOLLOW}(A)$ to $\text{FOLLOW}(B)$

Example of Calculating FIRST and FOLLOW

Production	FIRST	FOLLOW
$E \rightarrow TE'$	$\{ (, \text{id} \}$	$\{), \$ \}$
$E' \rightarrow +TE' \mid \epsilon$	$\{ +, \epsilon \}$	$\{), \$ \}$
$T \rightarrow FT'$	$\{ (, \text{id} \}$	$\{ +,), \$ \}$
$T' \rightarrow *FT' \mid \epsilon$	$\{ *, \epsilon \}$	$\{ +,), \$ \}$
$F \rightarrow (E) \mid \text{id}$	$\{ (, \text{id} \}$	$\{ *, +,), \$ \}$

Another Example of Calculating FIRST and FOLLOW

Production	FIRST	FOLLOW
$X \rightarrow Ya$	$\{ \}$	$\{ \}$
$Y \rightarrow ZW$	$\{ \}$	$\{ \}$
$W \rightarrow c \mid \epsilon$	$\{ \}$	$\{ \}$
$Z \rightarrow a \mid bZ$	$\{ \}$	$\{ \}$

Constructing Predictive Parsing Tables

For each $A \rightarrow \alpha$ do

1. Add $A \rightarrow \alpha$ to $M[A, a]$ for each a in $\text{FIRST}(\alpha)$
2. If ϵ is in $\text{FIRST}(\alpha)$
 - a. Add $A \rightarrow \alpha$ to $M[A, b]$ for each b in $\text{FOLLOW}(A)$
 - b. If $\$$ is in $\text{FOLLOW}(A)$ add $A \rightarrow \alpha$ to $M[A, \$]$
3. Make each undefined entry of M an error.

LL(1)

- First "L" - scans input from left to right
 Second "L" - produces a leftmost derivation
 1 - uses one input symbol of lookahead at each step to make a parsing decision

A grammar whose predictive parsing table has no multiply-defined entries is LL(1).

No ambiguous or left-recursive grammar can be LL(1).

When Is a Grammar LL(1)?

A grammar is LL(1) iff for each set of productions where $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$, the following conditions hold.

1. $\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \emptyset$ where $1 \leq i \leq n$ and $1 \leq j \leq n$ and $i \neq j$

2. If $\alpha_i \Rightarrow^* \epsilon$ then

- a. $\alpha_1, \dots, \alpha_{i-1}, \alpha_{i+1}, \dots, \alpha_n$ does not $\Rightarrow^* \epsilon$
- b. $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \emptyset$ where $j \neq i$ and $1 \leq j \leq n$

Checking If a Grammar is LL(1)

Production	FIRST	FOLLOW
$S \rightarrow iEtSS' \mid a$	{ i, a }	{ e, \$ }
$S' \rightarrow eS \mid \epsilon$	{ e, ϵ }	{ e, \$ }
$E \rightarrow b$	{ b }	{ t }

Nonterminal	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow eS$			
			$S' \rightarrow \epsilon$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

So this grammar is not LL(1).

Shift-Reduce Parsing

- Shift-reduce parsing is bottom-up.
 - Attempts to construct a parse tree for an input string beginning at the leaves and working up towards the root.
- A "handle" is a substring that matches the rhs of a production.
- A "shift" moves the next input symbol on a stack.
- A "reduce" replaces the rhs of a production that is found on the stack with the nonterminal on the left of that production.
- A "viable prefix" is the set of prefixes of right sentential forms that can appear on the stack of a shift-reduce parser.
- Shift reduce parsing includes
 - operator-precedence parsing
 - LR parsing

Model of an LR Parser

- See Figure 4.35.
- Each s_i is a state.
- Each X_i is a grammar symbol (when implemented these items do not appear in the stack).
- Each a_i is an input symbol.
- All LR parsers can use the same algorithm (code).
- The action and goto tables are different for each LR parser.

Model of an LR Parser (cont.)

- A shift pushes a state on the stack and processes an input symbol.
- A reduce pops states off the stack and pushes one state back on the stack.

	terminals	nonterminals
states	action table	goto table

LR(k) Parsing

- "L" - scans input from left to right
- "R" - constructs a rightmost derivation in reverse
- "k" - uses k symbols of lookahead at each step to make a parsing decision

Uses a stack of alternating states and grammar symbols. The grammar symbols are optional. Uses a string of input symbols (\$ on end).

LR (k) Parsing (cont.)

If config == $(s_0 X_1 s_1 X_2 s_2 \dots X_m s_m, a_i a_{i+1} \dots a_n \$)$

1. if action $[s_m, a_i] == \text{shift } s$ then
new config is $(s_0 X_1 s_1 X_2 s_2 \dots X_m s_m a_i s, a_{i+1} \dots a_n \$)$
2. if action $[s_m, a_i] == \text{reduce } A \rightarrow \beta$ and
goto $[s_{m-r}, A] == s$ (where r is the length of β) then
new config is $(s_0 X_1 s_1 X_2 s_2 \dots X_{m-r} s_{m-r} A s, a_i a_{i+1} \dots a_n \$)$
3. if action $[s_m, a_i] == \text{ACCEPT}$ then stop
4. if action $[s_m, a_i] == \text{ERROR}$ then attempt recovery

Can resolve some shift-reduce conflicts with lookahead.

ex: LR(1)

Can resolve others in favor of a shift.

ex: $S \rightarrow iCtS \mid iCtSeS$

Advantages of LR Parsing

- LR parsers can recognize almost all programming language constructs expressed in context -free grammars.
- Efficient and requires no backtracking.
- Is a superset of the grammars that can be handled with predictive parsers.
- Can detect a syntactic error as soon as possible on a left-to-right scan of the input.

LR Parsing Example

1. $E \rightarrow E + T$
2. $E \rightarrow T$
3. $T \rightarrow T * F$
4. $T \rightarrow F$
5. $F \rightarrow (E)$
6. $F \rightarrow id$

See Fig 4.37.

It produces rightmost derivation in reverse:

$E \rightarrow E + T \rightarrow E + F \rightarrow E + id \rightarrow T + id \rightarrow T * F + id$
 $\rightarrow T * id + id \rightarrow F * id + id \rightarrow id * id + id$

Calculating the Sets of LR(0) Items

LR(0) item - production with a dot at some position in the rhs indicating how much has been parsed

Example:

$A \rightarrow BC$ has 3 possible LR(0) items

$A \rightarrow \cdot BC$

$A \rightarrow B \cdot C$

$A \rightarrow BC \cdot$

$A \rightarrow \epsilon$ has 1 possible item

$A \rightarrow \cdot$

3 operations required to construct the sets of LR(0) items:
(1) closure, (2) goto, and (3) augment

Example of Computing the Closure of a Set of LR(0) Items

<u>Grammar</u>	<u>Closure</u> (I_0) for $I_0 = \{E' \rightarrow \cdot E\}$
$E' \rightarrow E$	$E' \rightarrow \cdot E$
$E \rightarrow E + T \mid T$	$E \rightarrow \cdot E + T$
$T \rightarrow T * F \mid F$	$E \rightarrow \cdot T$
$F \rightarrow (E) \mid id$	$T \rightarrow \cdot T * F$
	$T \rightarrow \cdot F$
	$F \rightarrow \cdot (E)$
	$F \rightarrow \cdot id$

Calculating Goto of a Set of LR(0) Items

Calculate goto (I, X) where I is a set of items and X is a grammar symbol.

Take the closure (the set of items of the form $A \rightarrow \alpha X \cdot \beta$) where $A \rightarrow \alpha \cdot X \beta$ is in I .

<u>Grammar</u>	<u>Goto</u> ($I_1, +$) for $I_1 = \{E' \rightarrow E \cdot, E \rightarrow E \cdot + T\}$
$E' \rightarrow E$	$E \rightarrow E + \cdot T$
$E \rightarrow E + T \mid T$	$T \rightarrow \cdot T * F$
$T \rightarrow T * F \mid F$	$T \rightarrow \cdot F$
$F \rightarrow (E) \mid id$	$F \rightarrow \cdot (E)$
	$F \rightarrow \cdot id$
	<u>Goto</u> ($I_2, *$) for $I_2 = \{E \rightarrow T \cdot, T \rightarrow T \cdot * F\}$
	$T \rightarrow T * \cdot F$
	$F \rightarrow \cdot (E)$
	$F \rightarrow \cdot id$

Augmenting the Grammar

Given grammar G with start symbol S , then an augmented grammar G' is G with a new start symbol S' and new production $S' \rightarrow S$.

Analogy of Calculating the Set of LR(0) Items with Converting an NFA to a DFA

Constructing the set of items is similar to converting an NFA to a DFA. Each state in the NFA is an individual item. The closure (I) for a set of items is similar to the ϵ -closure of a set of NFA states. Each set of items is now a DFA state and goto (I, X) gives the transition from I on symbol X .

Constructing SLR Parsing Tables

Let $C = \{I_0, I_1, \dots, I_n\}$ be the parser states.

1. If $[A \rightarrow \alpha \cdot a \beta]$ is in I_i and $\text{goto}(I_i, a) = I_j$ then set action $[i, a]$ to 'shift j '.
2. If $[A \rightarrow \alpha \cdot]$ is in I_i , then set action $[i, a]$ to 'reduce $A \rightarrow \alpha$ ' for all a in the $\text{FOLLOW}(A)$. A may not be S' .
3. If $[S' \rightarrow S \cdot]$ is in I_i , then set action $[i, \$]$ to 'accept'.
4. If $\text{goto}(I_i, A) = I_j$, then set $\text{goto}[i, A]$ to j .
5. Set all other table entries to 'error'.
6. The initial state is the one holding $[S' \rightarrow \cdot S]$.

LR(1)

The unambiguous grammar

$$S \rightarrow L = R \mid R$$
$$L \rightarrow *R \mid \text{id}$$
$$R \rightarrow L$$

is not SLR.

See Fig 4.39.

action[2, =] can be a "shift 6" or "reduce $R \rightarrow L$ "
 $\text{FOLLOW}(R)$ contains "=" but no form begins with "R="

LR (1) (cont.)

Solution - split states by adding LR(1) lookahead

form of an item

$$[A \rightarrow \alpha \cdot \beta, a]$$

where $A \rightarrow \alpha \beta$ is a production and

'a' is a terminal or endmarker $\$$

Closure(I) is now slightly different

repeat

for each item $[A \rightarrow \alpha \cdot B \beta, a]$ in I ,

each production $B \rightarrow \gamma$ in the grammar,

and each terminal b in $\text{FIRST}(\beta a)$ do

add $[B \rightarrow \cdot \gamma, b]$ to I (if not there)

until no more items can be added to I

Start the construction of the set of LR(1) items by computing the closure of $\{[S' \rightarrow \cdot S, \$]\}$.

LR(1) Example

(0) 1. $S' \rightarrow S$

(1) 2. $S \rightarrow CC$

(2) 3. $C \rightarrow cC$

(3) 4. $C \rightarrow d$

$I_0:$ $[S' \rightarrow \cdot S, \$]$ $\text{goto}(S) = I_1$

$[S \rightarrow \cdot CC, \$]$ $\text{goto}(C) = I_2$

$[C \rightarrow \cdot cC, c/d]$ $\text{goto}(c) = I_3$

$[C \rightarrow \cdot d, c/d]$ $\text{goto}(d) = I_4$

$I_1:$ $[S' \rightarrow S \cdot, \$]$

$I_2:$ $[S \rightarrow C \cdot C, \$]$ $\text{goto}(C) = I_5$

$[C \rightarrow \cdot cC, \$]$ $\text{goto}(c) = I_6$

$[C \rightarrow \cdot d, \$]$ $\text{goto}(d) = I_7$

LR(1) Example (cont.)

I_3 : $[C \rightarrow c \cdot C, c/d]$ goto (C) = I_8
 $[C \rightarrow \cdot cC, c/d]$ goto (c) = I_3
 $[C \rightarrow \cdot d, c/d]$ goto (d) = I_4
 I_4 : $[C \rightarrow d \cdot, c/d]$
 I_5 : $[S \rightarrow CC \cdot, \$]$
 I_6 : $[C \rightarrow c \cdot C, \$]$ goto (C) = I_9
 $[C \rightarrow \cdot cC, \$]$ goto (c) = I_6
 $[C \rightarrow \cdot d, \$]$ goto (d) = I_7
 I_7 : $[C \rightarrow d \cdot, \$]$
 I_8 : $[C \rightarrow cC \cdot, c/d]$
 I_9 : $[C \rightarrow cC \cdot, \$]$

Constructing the LR(1) Parsing Table

Let $C = \{I_0, I_1, \dots, I_n\}$

1. If $[A \rightarrow \alpha \cdot a \beta, b]$ in I_i and $\text{goto}(I_i, a) = I_j$ then set $\text{action}[i, a]$ to “shift j”.
2. If $[A \rightarrow \alpha \cdot, a]$ is in I_i , then set $\text{action}[i, a]$ to ‘reduce $A \rightarrow \alpha$ ’. A may not be S' .
3. If $[S' \rightarrow S \cdot, \$]$ is in I_i , then set $\text{action}[i, \$]$ to “accept.”
4. If $\text{goto}(I_i, A) = I_j$, then set $\text{goto}[i, A]$ to j.
5. Set all other table entries to error.
6. The initial state is the one holding $[S' \rightarrow \cdot S, \$]$

Constructing LALR Parsing Tables

- Combine LR(1) sets with the same sets of the first parts (ignore lookahead).
- Table is the same size as SLR.
- Will not introduce shift-reduce conflicts since shifts depend only on the core and don't use lookahead.
- May introduce reduce-reduce conflicts but seldom do for grammars describing programming languages.
- Last example collapses to table shown in Fig 4.43.
- Algorithms exist that skip constructing all the LR(1) sets of items.

Compaction of LR Parsing Tables

- A typical programming language may have 50 to 100 terminals and over 100 productions. This can result in several hundred states and a very large action table.
- One technique to save space is to recognize that many rows of the action table are identical. Can create a pointer for each state with the same actions so that it points to the same location.
- Could save further space by creating a list for the actions of each state, where the list consists of terminal-symbol/action pairs. This would eliminate the blank or error entries in the action table. While this technique would save a lot of space, the parser would be much slower.

Using Ambiguous Grammars

1. $E \rightarrow E + E$ $E \rightarrow E + T \mid T$
2. $E \rightarrow E * E$ instead of $T \rightarrow T * F \mid F$
3. $E \rightarrow (E)$ $F \rightarrow (E) \mid \text{id}$
4. $E \rightarrow \text{id}$

See Figure 4.48.

Advantages:

Grammar is easier to read.

Parser is more efficient.

Using Ambiguous Grammars (cont.)

Can use precedence and associativity to solve the problem.

See Fig 4.49.

shift / reduce conflict in state $\text{action}[7,+]=(s4,r1)$

$s4 = \text{shift } 4$ or $E \rightarrow E \cdot + E$

$r1 = \text{reduce } 1$ or $E \rightarrow E + E \cdot$

$\text{id} + \text{id} + \text{id}$

↑ cursor here

$\text{action}[7,*]=(s5,r1)$

$\text{action}[8,+]=(s4,r2)$ $\text{action}[8,*]=(s5,r2)$

Another Ambiguous Grammar

0. $S' \rightarrow S$
1. $S \rightarrow \text{iSeS}$
2. $S \rightarrow \text{iS}$
3. $S \rightarrow \text{a}$

See Figure 4.50.

$\text{action}[4,e]=(s5,r2)$

Ambiguities from Special-Case Productions

$E \rightarrow E \text{ sub } E \text{ sup } E$

$E \rightarrow E \text{ sub } E$

$E \rightarrow E \text{ sup } E$

$E \rightarrow \{ E \}$

$E \rightarrow c$

Ambiguities from Special-Case Productions (cont)

1. $E \rightarrow E \text{ sub } E \text{ sup } E$ $\text{FIRST}(E) = \{ \{, c \}$
2. $E \rightarrow E \text{ sub } E$ $\text{FOLLOW}(E) = \{ \text{sub}, \text{sup}, \}', \$ \}$
3. $E \rightarrow E \text{ sup } E$
4. $E \rightarrow \{ E \}$ sub, sup have equal precedence
5. $E \rightarrow c$ and are right associative

action[7,sub]=(s4,r2)	action[7,sup]=(s10,r2)
action[8,sub]=(s4,r3)	action[8,sup]=(s5,r3)
action[11,sub]=(s4,r1,r3)	action[11,sup]=(s5,r1,r3)
action[11,}]= (r1,r3)	action[11,\$]= (r1,r3)

YACC

Yacc source program

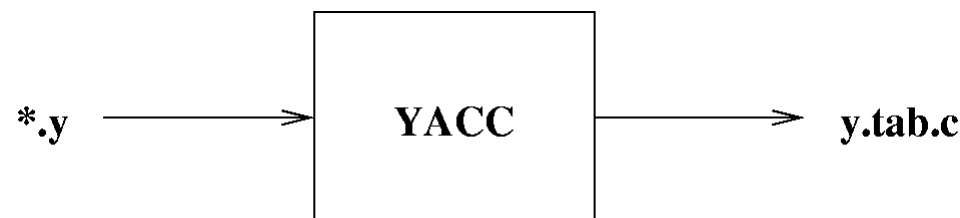
declarations

%%

translation rules

%%

supporting C-routines



YACC Declarations

- In declarations:
 - Can put ordinary C declarations in

```
%{  
    ...  
}%
```
 - Can declare tokens using
 - %token
 - %left
 - %right
 - Precedence is established by the order the operators are listed (low to high).

YACC Translation Rules

- Form

```
A : Body ;
```

where A is a nonterminal and Body is a list of nonterminals and terminals.
- Semantic actions can be enclosed before or after each grammar symbol in the body.
- Yacc chooses to shift in a shift/reduce conflict.
- Yacc chooses the first production in a reduce/reduce conflict.

Yacc Translation Rules (cont.)

- When there is more than one rule with the same left hand side, a ' | ' can be used.

A : B C D ;

A : E F ;

A : G ;

=>

A : B C D

| E F

| G

;

Example of a Yacc Specification

```
%token IF ELSE NAME      /* defines multicharacter tokens */
%right '='                /* low precedence, a=b=c shifts */
%left '+' '-'             /* mid precedence, a-b-c reduces */
%left '*' '/'             /* high precedence, a/b/c reduces */
%%
stmt  : expr ';'
      | IF '(' expr ')' stmt
      | IF '(' expr ')' stmt ELSE stmt
      ; /* prefers shift to reduce in shift/reduce conflict */
expr  : NAME '=' expr      /* assignment */
      | expr '+' expr
      | expr '-' expr
      | expr '*' expr
      | expr '/' expr
      | '-' expr %prec '*' /* can override precedence */
      | NAME
      ;
%% /* definitions of yylex, etc. can follow */
```

Yacc Actions

- Actions are C code segments enclosed in { } and may be placed before or after any grammar symbol in the right hand side of a rule.
- To return a value associated with a rule, the action can set \$\$.
- To access a value associated with a grammar symbol on the right hand side, use \$i, where i is the position of that grammar symbol.
- The default action for a rule is

```
{ $$ = $1; }
```

Syntax Error Handling

- Errors can occur at many levels
 - lexical - unknown operator
 - syntactic - unbalanced parentheses
 - semantic - variable never declared
 - logical - dereference a null pointer
- Goals of error handling in a parser
 - detect and report the presence of errors
 - recover from each error to be able to detect subsequent errors
 - should not slow down the compilation of correct programs

Syntax Error Handling (cont.)

- Viable–prefix property - detect an error as soon as the parser sees a prefix of the input that is not a prefix of any string in the language.

Error-Recovery Strategies

- Panic-mode - skip until one of a synchronizing set of tokens is found (e.g. ';', "end"). Is very simple to implement but may miss detection of some errors (when more than one error in a single statement).
- Phrase-level - replace prefix of remaining input by a string that allows the parser to continue. Hard for the compiler writer to anticipate all error situations.
- Error productions - augment the grammar of the source language to include productions for common errors. When production is used, an appropriate error diagnostic would be issued. Feasible to only handle a limited number of errors.

Error-Recovery Strategies (cont)

- Global correction - choose minimal sequence of changes to allow a least-cost correction. Often considered too costly to actually be implemented in a parser. Also the closest correct program may not be what the programmer intended.

Error-Recovery in Predictive Parsing

- It is easier to recover from an error in a nonrecursive predictive parser than using recursive descent.
- Panic-mode recovery
 - Assume the nonterminal A is on the stack when we encounter an error. As a starting point can place all symbols in FOLLOW(A) into the synchronizing set for the nonterminal A. May also wish to add symbols that begin higher-level constructs to the synchronizing set of lower-level constructs. If a terminal is on top of the stack, then can pop the terminal and issue a message stating that the terminal was discarded.

Error-Recovery in Predictive Parsing (cont.)

Phrase-level recovery

- Can be implemented by filling in the blank entries in the predictive parsing table with pointers to error routines. The compiler writer would attempt each situation appropriately (issue error message and update input symbols and pop from the stack).

Error-Recovery in LR Parsing

- Canonical LR Parser - will never make a single reduction before recognizing an error.
- SLR & LALR Parsers - may make extra reductions but will never shift an erroneous input symbol on the stack.
- Panic-mode recovery - scan down stack until a state with a goto on a particular nonterminal representing a major program construct (e.g. expression, statement, block, etc.) is found. Input symbols are discarded until one is found that is in the FOLLOW of the nonterminal. The parser then pushes on the state in goto. Thus, it attempts to isolate the construct containing the error.

Error-Recovery in LR Parsing (cont)

- Phrase-level recovery - Implement an error recovery routine for each error entry in the table.
- Error productions - Used in YACC. Pops symbols until topmost state has an error production, then shifts error onto stack. Then discards input symbols until it finds one that allows parsing to continue. The semantic routine with an error production can just produce a diagnostic message.