West Chester University Dr. Richard Burns Fall 2023

CSC 416/565: DESIGN AND CONSTRUCTION OF COMPILERS

LAST CLASS

Predictive Parsing Algorithm: Recursive Descent

TODAY

Left Recursion

 Motivation: Recursive Descent parsing algorithm does not work with grammars that are <u>left</u> recursive.

Recursive Descent parsing algorithm does not work with grammars that are left recursive.

Unable to use this grammar w/ recursive descent:

```
S \rightarrow E\$
E \rightarrow E + T \mid E - T \mid T
T \rightarrow T * F \mid T / F \mid F
F \rightarrow id \mid num \mid (E)
```

Why? Infinite loop

```
boolean E_1() { return E() && advance(PLUS) && T(); }
boolean E() { ...; return E_1() || ...}
```

Definition: A <u>left-recursive grammar</u> has some non-terminal S, such that $S \rightarrow +S\alpha$, where α is any combination of zero or more terminals or non-terminals.

- A production is left recursive if its LHS symbol is also the first symbol of its RHS.
- → means *one or more* derivation steps...
- "...eventually derives to"

It is still possible to make derivations with a grammar that is left recursive.

 Often easier to understand left-recursive grammars than those modified to not be left-recursive.

Grammar:

$$S \rightarrow S\alpha \mid \beta$$

- What is the language defined by this left-recursive grammar?
- Rewrite <u>left-recursive grammar</u> so that it uses <u>right recursion</u> instead.

A general rule to eliminate left recursion

• A left-recursive grammar has the form:

$$\begin{array}{ccc} X & \to & X \gamma \\ X & \to & \alpha \end{array}$$

- where α does not start with X
- All strings derived from grammar will start with α , and continue with zero or more γ .
- Can rewrite using right-recursion as:

$$X \rightarrow \alpha X'$$
 $X' \rightarrow \gamma X' \mid \epsilon$

Example: expression-term grammar

Left-recursive grammar:

Right-recursive grammar:

$$E \rightarrow E + T$$

$$E \rightarrow T$$

Can work for *multiple* productions

• Left recursive:

$$N \rightarrow N \alpha_1 \mid N \alpha_2 \dots \mid N \alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

Transformation, introducing new nonterminal N'

$$N \rightarrow \beta_1 N' \mid \beta_2 N' \mid \dots \mid \beta_n N'$$

 $N' \rightarrow \alpha_1 N' \mid \alpha_2 N' \mid \dots \mid \alpha_m N' \mid \epsilon$

More examples in [Thain]

• Removing left-recursion that spans multiple productions is trickier.

- left-recursive because: $S \rightarrow^+ S\beta\alpha$
- ... more complicated general algorithms exist

Predictive Parsing

- New Idea: in <u>predictive parsing</u>, parser can *predict* which production to use, simply by looking at the next few tokens:
- No backtracking! (as we had in <u>Recursive Descent</u> in the last class)
- Works only when the next few tokens provide *enough information* about which production to use.
 - no more choices to make
 - bad/incorrect choices not possible

Predictive Parsing

Predictive parsers accept LL(k) class of grammars:

Will parse the input:

- left-to-right ("first L")
- leftmost derivation ("second L")
- using k "tokens of lookahead"
 - k=1: if we can look at next token and always know which production to use

Grammar: $E \rightarrow T \mid T + E$ $T \rightarrow int \mid int * T \mid (E)$

Some issues:

- With 1 symbol of lookahead, we may know that our next token is <int>, but not enough information to choose which T production rule to apply.
- Sometimes referred to as <u>prediction conflicts</u> caused by *common prefixes*
- *Idea*: Left factor the grammar when two productions for the same nonterminal start with the same symbols.

Left Factoring

- "Take the allowable 'endings' and make a new nonterminal to stand for them."
- "Re-write the productions to defer the decision about which production to use in predictive parsing until enough of the input has been seem that we can make the right choice."

Grammar: $E \rightarrow T \mid T + E$ $T \rightarrow int \mid int * T \mid (E)$

Apply left factoring.

Another Example

Grammar:

```
S \rightarrow if E then S else S

S \rightarrow if E then S
```

Apply left factoring.

Note: Can't simply just increase the lookahead. Rewriting the grammar is necessary.

LL(1) Parsing Tables

- From a left-factored grammar, going to create a parsing table.
- Parser will use this table to create parse tree, by performing one leftto-right pass on the input, with a leftmost derivation, using 1 symbol of lookahead.
- Components:
 - left column: current leftmost nonterminal in parse tree
 - top row: next input token
 - each table cell entry: information about which production to use to expand the current nonterminal

1	$S \rightarrow A C $$
2	$C \rightarrow c$
3	λ
4	$A \rightarrow a B C d$
5	B Q
6	$B \rightarrow b B$
7	λ
8	$Q \rightarrow q$
9	λ

	Lookahead						
Nonterminal	а	b	С	d	q	\$	
S	1	1	1		1	1	•
С			2	3		3	
Α	4	5	5		5	5	
В		6	7	7	7	7	
Q			9		8	9	
	01						

Figure 5.10: LL(1) table. The blank entries should trigger error actions in the parser.

Grammar is augmented with the $S \rightarrow E\$$ production.

$$S \rightarrow E \$$$

$$T \rightarrow F T'$$

$$E \rightarrow T E'$$

$$T' \rightarrow *F T'$$

$$F \rightarrow \text{id}$$

$$T' \rightarrow *F T'$$

$$F \rightarrow \text{num}$$

$$E' \rightarrow +T E'$$

$$E' \rightarrow -T E'$$

$$T' \rightarrow F T'$$

$$T' \rightarrow F T'$$

$$F \rightarrow (E)$$

GRAMMAR 3.15.

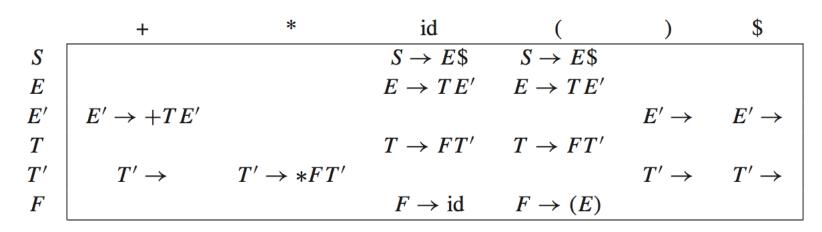


TABLE 3.17. Predictive parsing table for Grammar 3.15. We omit the columns for num, /, and -, as they are similar to others in the table.

Alg. for processing input stream using parsing table

- Maintain frontier (fringe) of parse tree as a stack.
 - stack will contain non-terminals that have yet to be expanded and terminals yet to be matched against the input
 - top of stack: leftmost terminal or nonterminal
- Initialize stack to start symbol.
- Iterate until:
 - reject on error state
 - accept on end of input and empty stack

Input String:

id+id\$

$$S \rightarrow E \$$$
 $T \rightarrow F T'$
 $E \rightarrow T E'$
 $T' \rightarrow *F T'$
 $F \rightarrow \text{id}$
 $T' \rightarrow *F T'$
 $F \rightarrow \text{num}$
 $E' \rightarrow +T E'$
 $E' \rightarrow -T E'$
 $T' \rightarrow F T'$
 $T' \rightarrow F T'$
 $F \rightarrow (E)$
 $F \rightarrow (E)$

GRAMMAR 3.15.

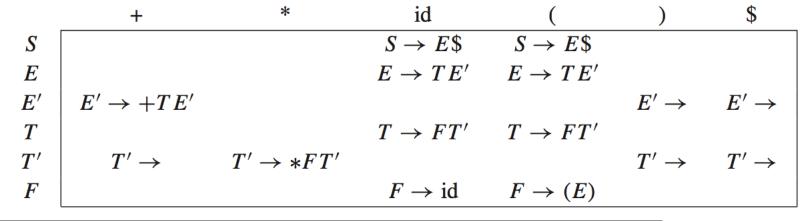


TABLE 3.17. Predictive parsing table for Grammar 3.15. We omit the columns for num, /, and -, as they are similar to others in the table.

- Grammars parsable with LL(1) parsing tables are called LL(1) grammars.
- Grammars parsable with LL(2) parsing tables are called LL(2) grammars, and so on...
- for LL(k), columns are every sequence of k terminals.
 - tables grow very large
 - LL(1) only table used in practice

- Duplicate entries in a parsing table ⇔ grammar is ambiguous.
- No ambiguous grammar is LL(K) for any k.

• How to construct a parsing table for a LL(1) grammar.

After the Break

BEFORE THE BREAK

 finished with example that used a LL(1) parsing table to perform derivation on input token stream

NOW

- How can we construct LL(1) parsing table?
- *Intuition:* during the parsing process, we have:
 - the leftmost nonterminal on the fringe that we're ready to expand
 - 2. the next token

Parsing Table Construction

Motivation: FIRST and FOLLOW Sets

- Example table entry in LL(1) Parsing Table: T[X, t] = Y
- Under what conditions should the parser use the production rule: $X \rightarrow Y$

Scenario #1

• Y can derive t in the first position.

$$Y \rightarrow^* tZ$$

Using this move would be a good idea ...
 eventually the parser could match the t.

$$t \in FIRST(Y)$$

- The terminal/symbol t is one of the things that Y can produce in the very first position.
- There may be others that Y could also produce.

Scenario #2

• Y reduces to ε , t follows Y in a derivation.

$$Y \rightarrow * t$$

t ∉ FIRST(Y)

 $Y \rightarrow * ε$

• This production rule is still useful if parser can get rid of X by deriving Y (and by deriving ε).

$$t \in FOLLOW(X)$$

 Not talking about X deriving t, just that t appears in a derivation after X.

$$S \rightarrow^* W X t Z$$

Formal Definition: FIRST

FIRST (X) - the set of all terminals, potentially including ϵ , that can begin strings derived from X

$$FIRST(X) = \{t | X \to^* tY\} \cup \{\epsilon | X \to^* \epsilon\}$$

Alg. for computing FIRST sets

- 1. $FIRST(t) = \{t\}$
- first set of a terminal includes itself
- for each terminal symbol t
- 2. $\varepsilon \in FIRST(X)$
- if $X \rightarrow \epsilon$
- if $X \rightarrow Y_1$, ... Y_k and $\varepsilon \in FIRST(Y_i)$ for $1 \le i \le k$
 - // Y₁, ... Y_k are all "nullable" and each can derive the empty string
- 3. $FIRST(Z) \subseteq FIRST(X)$
- if $X \rightarrow Y_1 \dots Y_k Z$ and $\varepsilon \in FIRST(Y_i)$ for $1 \le i \le k$
 - $Y_1 \dots Y_k$ can all derive empty string, therefore $X \rightarrow^* Z$

Compute FIRST sets for the left-factored grammar from before the break:

```
Grammar: E \rightarrow TX
X \rightarrow + E \mid \epsilon
T \rightarrow int Y \mid (E)
Y \rightarrow * T \mid \epsilon
```

Formal Definition: FOLLOW

- FOLLOW(X) set of *terminals* that can immediately follow X $FOLLOW(X) = \{t | S \rightarrow^* WXtZ\}$
- t ∈ FOLLOW(X) if there is any derivation containing Xt
 - tricky: this can also occur if the derivation contains XYZt where both Y and Z derive ε

FOLLOW Intuition #1

Intuition #1:

- if $X \rightarrow YZ$, then FIRST(Z) $\{\epsilon\} \subseteq FOLLOW(Y)$
 - (ε never appears in FOLLOW sets)
 - if two symbols are adjacent somewhere in the grammar, the "FIRST of the second symbol is in the FOLLOW of the first symbol"

Grammar:

 $X \rightarrow YZ$ $Z \rightarrow ab$

 $a \in FIRST(Z)$ $a \in FOLLOW(Y)$ $X \rightarrow Yab$

FOLLOW Intuition #2

- Intuition #2:
- if $X \rightarrow YZ$, then $FOLLOW(X) \subseteq FOLLOW(Z)$
 - anything that occurs at the end of a production, its FOLLOW set will include the FOLLOW of the symbol on the LHS.

 $a \in FOLLOW(X)$

 $a \in FOLLOW(Z)$

FOLLOW Intuition #3

- Intuition #3:
- Generalizing: What is the end of the production?
- What if $Z \rightarrow^* \epsilon$?
- if $Z \rightarrow *\varepsilon$, then FOLLOW(X) \subseteq FOLLOW(Y)
 - Need to recursively check rightmost symbols for potential ε reduction

 $a \in FOLLOW(X)$

 $a \in FOLLOW(Y)$

Alg. for computing FOLLOW sets

Rules from before:

- 1. $$ \in FOLLOW(S_1)$
- After we reduce start symbol, expect to see EOF.
- 2. $FIRST(Y) \{\epsilon\} \subseteq FOLLOW(X)$
- Not interested in ε: FOLLOW sets are sets of terminals.
- For each production: $S \rightarrow aXY$
 - Optional a terminal(s)
 - Y = terminal or nonterminal
- 3. $FOLLOW(S) \subseteq FOLLOW(X)$
- For each production: $S \rightarrow aXY$, where $\varepsilon \in FIRST(Y)$
 - Y "falls out" of the derivation
- 4. $FOLLOW(S) \subseteq FOLLOW(Y)$
 - If Y does not "fall out"

Compute FOLLOW sets for the left-factored grammar from before the break:

Grammar: $E \rightarrow TX$ $X \rightarrow + E \mid \epsilon$ $T \rightarrow int Y \mid (E)$ $Y \rightarrow * T \mid \epsilon$

... now going to <u>construct</u> a parsing table for our left-factored LL(1) grammar!

Parsing Table Construction Alg.

For each production in G: $X \rightarrow Y$

- For each terminal in $t \in FIRST(Y)$, add T[X, t] = Y
- If $\varepsilon \in FIRST(Y)$, for each $t \in FOLLOW(X)$, add T[X, t] = Y
- If $\varepsilon \in FIRST(Y)$ and $\varphi \in FOLLOW(X)$, add $T[X, \varphi] = Y$

```
Grammar:

E \rightarrow TX

X \rightarrow + E \mid \epsilon

T \rightarrow int Y \mid (E)

Y \rightarrow * T \mid \epsilon
```

Another Example

Grammar:

 $S \rightarrow Sa \mid b$

- Demonstrating that left-recursive grammars are not LL(1)!
- Recall that the language for this grammar is strings of a b followed by zero or more a's
- Try building a parse table.

Conclusion

• One way to check if grammar is LL(1) is to build parsing table and see if there are no duplicate entries.

- Classes of grammars that are *not* LL(1):
 - Not left factored
 - Left recursive
 - Ambiguous
 - Plus some others (i.e. those that require more than 1 token of lookahead)

Next Class

Bottom-Up Parsing

