

Compiler Construction

Lent Term 2020

Lecture 15

Parsing Part II : Table-based Methods

1. Top-down parsing (LL grammars)
 1. Predictive parsing
 2. Bottom-up parsing (LR grammars)
 1. SLR(1) (Simple LR)
 2. LR(1)

**We will NOT cover
LR(1) this year!**

Timothy G. Griffin

tgg22@cam.ac.uk

**Computer Laboratory
University of Cambridge**

Leftmost, rightmost derivations

$$w \in T^*, \alpha, \beta \in (N \cup T)^*$$

Given : $wA\beta$ and a production $A \rightarrow \gamma$

a leftmost derivation step is written as

$$wA\beta \Rightarrow_{lm} w\gamma\beta$$

Given : αAw and a production $A \rightarrow \gamma$

a rightmost derivation step is written as

$$\alpha Aw \Rightarrow_{rm} \alpha\gamma w$$

LL(k) and LR(k)

- **LL(k)** : (L)eft-to-right parse, (L)eft-most derivation, k-symbol lookahead. Based on looking at the next k tokens, an LL(k) parser must *predict* the next production. We have been looking at LL(1).
- **LR(k)** : (L)eft-to-right parse, (R)ight-most derivation, k-symbol lookahead. Postpone production selection until *the entire* right-hand-side has been seen (and as many as k symbols beyond). LR parsers perform a rightmost derivation backwards!

For LL(1), augment Grammar with end-of-input

$$G'_3 = (N'_3, T_3, P'_3, S)$$

$$N'_3 = \{E, E', T, T', F, S\} \quad T_3 = \{+, *, (,), \text{id}, \$\}$$

P'_3 :

$S \rightarrow E\$$ (\$ is end of input marker)

$E \rightarrow T E'$

$E' \rightarrow +T E' \mid \varepsilon$

$T \rightarrow F T'$

$T' \rightarrow *F T' \mid \varepsilon$

$F \rightarrow (E) \mid \text{id}$

Can top-down parsing be automated?

$S \Rightarrow_{lm} E\$$ (derivation of $(x + y)$)

$\Rightarrow_{lm} TE'\$$

$\Rightarrow_{lm} FT' E'\$$

$\Rightarrow_{lm} (E)T' E'\$$

$\Rightarrow_{lm} (TE')T' E'\$$

$\Rightarrow_{lm} (FT' E')T' E'\$$

$\Rightarrow_{lm} (xT' E')T' E'\$$

$\Rightarrow_{lm} (xE')T' E'\$$

$\Rightarrow_{lm} (x + TE')T' E'\$$

$\Rightarrow_{lm} (x + FT' E')T' E'\$$

$\Rightarrow_{lm} (x + yT' E')T' E'\$$

$\Rightarrow_{lm} (x + yE')T' E'\$$

$\Rightarrow_{lm} (x + y)T' E'\$$

$\Rightarrow_{lm} (x + y)E'\$$

$\Rightarrow_{lm} (x + y)\$$

Idea : If $S \Rightarrow_{lm}^+ w\alpha\$$ then w has been read from the input and α is on the stack.

This looks promising. But

input	stack	via production
$(x + y)\$$	S	$S \rightarrow E\$$
$(x + y)\$$	$E\$$	$E \rightarrow TE'$
$(x + y)\$$	$TE'\$$	$T \rightarrow FT'$
$(x + y)\$$	$FT' E'\$$	$F \rightarrow (E)$
$(x + y)\$$	$(E)T' E'\$$	match
$x + y)\$$	$E)T' E'\$$	$E \rightarrow TE'$
$x + y)\$$	$TE')T' E'\$$	$T \rightarrow FT'$
$x + y)\$$	$FT' E')T' E'\$$	$F \rightarrow id$
$x + y)\$$	$idT' E')T' E'\$$	match
$+ y)\$$	$T' E')T' E'\$$	$T' \rightarrow \varepsilon$

.... how do we automate selection of the production to use at each step?

input	stack	via production
+ y)\$	$E')T' E' \$$	$E' \rightarrow +TE'$
+ y)\$	$+TE')T' E' \$$	match
y)\$	$TE')T' E' \$$	$T \rightarrow FT'$
y)\$	$FT' E')T' E' \$$	$F \rightarrow id$
y)\$	$idT' E')T' E' \$$	match
)\$	$T' E')T' E' \$$	$T' \rightarrow \varepsilon$
)\$	$E')T' E' \$$	$E' \rightarrow \varepsilon$
)\$	$)T' E' \$$	match
\$	$T' E' \$$	$T' \rightarrow \varepsilon$
\$	$E' \$$	$E' \rightarrow \varepsilon$
\$	\$	accept!

FIRST (we will see how to compute later)

$$\text{FIRST}(A) = \left\{ a \in T / \exists \beta \in (N \cup T)^*, A \Rightarrow^+ a\beta \right\} \\ \cup \left\{ \varepsilon / A \Rightarrow^+ \varepsilon \right\}$$

$$S \rightarrow E\$ \quad \text{FIRST}(S) = \{ (, id \}$$

$$E \rightarrow T E' \quad \text{FIRST}(E) = \{ (, id \}$$

$$E' \rightarrow +T E' / \varepsilon \quad \text{FIRST}(E') = \{ +, \varepsilon \}$$

$$T \rightarrow F T' \quad \text{FIRST}(T) = \{ (, id \}$$

$$T' \rightarrow * F T' | \varepsilon \quad \text{FIRST}(T') = \{ *, \varepsilon \}$$

$$F \rightarrow (E) | id \quad \text{FIRST}(F) = \{ (, id \}$$

FOLLOW (we will see how to compute later)

$$\text{FOLLOW}(A) = \{a / \exists \alpha \beta, S \Rightarrow^+ \alpha A a \beta\}$$

$$S \rightarrow E\$$$

$$E \rightarrow T E'$$

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

$$E' \rightarrow +T E' / \varepsilon$$

$$\text{FOLLOW}(E') = \{), \$ \}$$

$$T \rightarrow F T'$$

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

$$T' \rightarrow *F T' / \varepsilon$$

$$\text{FOLLOW}(T') = \{ +,), \$ \}$$

$$F \rightarrow (E) \mid \text{id}$$

$$\text{FOLLOW}(F) = \{ +, *,), \$ \}$$

" $)$ " \in FOLLOW(E)?

$$S \Rightarrow E\$ \Rightarrow TE'\$ \Rightarrow FT' E'\$ \Rightarrow (E)T' E'\$$$

The LL(1) Parsing table M

for all $A \in N$, $a \in T$, $M[A, a] = \{\}$

for each $A \in N$

for each production $A \rightarrow \alpha$

if $a \in \text{FIRST}(\alpha)$ and $a \in T$

then $M[A, a] = M[A, a] \cup \{A \rightarrow \alpha\}$

else if $\varepsilon \in \text{FIRST}(\alpha)$

then for each $b \in \text{FOLLOW}(A)$

$M[A, b] = M[A, b] \cup \{A \rightarrow \alpha\}$

See slide 22 for
definition of
 $\text{FIRST}(\alpha)$

Table M for grammar G_3'

	id	$+$	$*$	$($	$)$	$\$$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \varepsilon$	$E' \rightarrow \varepsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \varepsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \varepsilon$	$T' \rightarrow \varepsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

The LL(1) Parsing Algorithm

$a := \text{LexNextToken}()$

$X := \text{TopOfStack}()$

while ($X \neq \$$)

 if $X = a$

 then pop; $a := \text{LexNextToken}()$

 else if $M[X, a] = \{X \rightarrow \alpha\}$

 then pop; push α (leftmost symbol on top)

$X := \text{TopOfStack}()$

Now use M to parse $(x+y) \dots$

input	stack	action
$(x + y)\$$	S	$M[S, (] = \{S \rightarrow E\$ \}$
$(x + y)\$$	$E\$$	$M[E, (] = \{E \rightarrow TE' \}$
$(x + y)\$$	$TE'\$$	$M[T, (] = \{T \rightarrow FT' \}$
$(x + y)\$$	$FT' E'\$$	$M[F, (] = \{F \rightarrow (E) \}$
$(x + y)\$$	$(E)T' E'\$$	<i>match</i>
$x + y)\$$	$E)T' E'\$$	$M[E, id] = \{E \rightarrow TE' \}$
$x + y)\$$	$TE')T' E'\$$	$M[T, id] = \{T \rightarrow FT' \}$
$x + y)\$$	$FT' E')T' E'\$$	$M[F, id] = \{F \rightarrow id \}$
$x + y)\$$	$idT' E')T' E'\$$	<i>match</i>
$+ y)\$$	$T' E')T' E'\$$	$M[T', +] = \{T' \rightarrow \varepsilon \}$

... kachunk, kachunk, kachunk ...

input	stack	action
+ y)\$	$E')T' E' \$$	$M[E', +] = \{E' \rightarrow +TE'\}$
+ y)\$	$+TE')T' E' \$$	<i>match</i>
y)\$	$TE')T' E' \$$	$M[T, id] = \{T \rightarrow FT'\}$
y)\$	$FT' E')T' E' \$$	$M[F, id] = \{F \rightarrow id\}$
y)\$	$idT' E')T' E' \$$	<i>match</i>
)\$	$T' E')T' E' \$$	$M[T',)] = \{T' \rightarrow \varepsilon\}$
)\$	$E')T' E' \$$	$M[E',)] = \{E' \rightarrow \varepsilon\}$
)\$	$)T' E' \$$	<i>match</i>
\$	$T' E' \$$	$M[T', \$] = \{T' \rightarrow \varepsilon\}$
\$	$E' \$$	$M[E', \$] = \{E' \rightarrow \varepsilon\}$
\$	\$	<i>accept</i>

Computing FIRST: Easy when there are no epsilon productions!

for all $a \in T$, $\text{FIRST}(a) := \{a\}$

for all $A \in N$, $\text{FIRST}(A) := \{\}$

while FIRST changes

if $A \rightarrow X\alpha$ is a production ($X \in N \cup T$)

then $\text{FIRST}(A) := \text{FIRST}(A) \cup \text{FIRST}(X)$

Computing FOLLOW: Easy when there are no epsilon productions!

for all $A \in N$, $\text{FOLLOW}(A) := \{\}$

$\text{FOLLOW}(S) := \{\$ \}$ (S is the start symbol)

while FOLLOW changes

if $A \rightarrow \alpha B \beta$ is a production ($B \in N, \beta \neq \varepsilon$)

then $\text{FOLLOW}(B) := \text{FOLLOW}(B) \cup \text{FIRST}(\beta)$

if $A \rightarrow \alpha B$ is a production ($B \in N$)

then $\text{FOLLOW}(B) := \text{FOLLOW}(B) \cup \text{FOLLOW}(A)$

With epsilon productions **NULLABLE** is a useful function

$\text{NULLABLE} : (N \cup T)^* \rightarrow \{true, false\}$

$\text{NULLABLE}(\varepsilon) = true$

$\text{NULLABLE}(a\beta) = false$

$\text{NULLABLE}(A\beta) = (A \Rightarrow^* \varepsilon) \wedge \text{NULLABLE}(\beta)$

Computing FIRST with epsilon productions!

for all $a \in T$, $\text{FIRST}(a) := \{a\}$

for all $A \in N$, $\text{FIRST}(A) := \{\}$

while FIRST changes

if $A \rightarrow \varepsilon$ is a production

then $\text{FIRST}(A) := \text{FIRST}(A) \cup \{\varepsilon\}$

if $A \rightarrow X_1 X_2 \cdots X_k$ is a production

then $j = 1$; $\text{done} := \text{false}$

while not done and $j \leq k$

$\text{FIRST}(A) := \text{FIRST}(A) \cup (\text{FIRST}(X_j) - \{\varepsilon\})$

if $\text{NULLABLE}(X_j)$

then $j := j + 1$

else $\text{done} := \text{true}$

if $j = k + 1$ then $\text{FIRST}(A) := \text{FIRST}(A) \cup \{\varepsilon\}$

Computing FOLLOW with epsilon productions!

for all $A \in N$, $\text{FOLLOW}(A) := \{ \}$

$\text{FOLLOW}(S) := \{ \$ \}$ (S is the start symbol)

while FOLLOW changes

if $A \rightarrow \alpha B \beta$ is a production ($B \in N, \beta \neq \varepsilon$)

then $\text{FOLLOW}(B) := \text{FOLLOW}(B) \cup (\text{FIRST}(\beta) - \{ \varepsilon \})$

if $A \rightarrow \alpha B \beta$ is a production and $\varepsilon \in \text{FIRST}(\beta)$

then $\text{FOLLOW}(B) := \text{FOLLOW}(B) \cup \text{FOLLOW}(A)$

if $A \rightarrow \alpha B$ is a production ($B \in N$)

then $\text{FOLLOW}(B) := \text{FOLLOW}(B) \cup \text{FOLLOW}(A)$

Many grammars cannot be parsed LL(1)

$S \rightarrow d \mid XY S$		FIRST	FOLLOW
$Y \rightarrow c \mid \varepsilon$	S	$\{a, c, d\}$	$\{\}$
$X \rightarrow Y \mid a$	Y	$\{c\}$	$\{a, c, d\}$
	X	$\{a, c\}$	$\{a, c, d\}$

$$M[S, d] = \{ S \rightarrow d, S \rightarrow XY S \}$$

This is ambiguity!

Grammar is not LL(1)!



Bottom-up (LR) parsing to the rescue!

$$G_2 = (N_2, T_1, P_2, E)$$

$$N_2 = \{E, T, F\}$$

$$T_1 = \{+, *, (,), \text{id}\}$$

$$E \rightarrow E + T \mid T \quad T \rightarrow T * F \mid F \quad F \rightarrow (E) \mid \text{id}$$

With LR parsing we no longer have to eliminate left recursion from the grammar!



Define $\text{FIRST}(\alpha)$ for $\alpha \in (N \cup T)^*$

$$\text{FIRST}'(\varepsilon) = \{\varepsilon\}$$

$$\text{FIRST}'(X) = \text{FIRST}(X) \quad (X \in N \cup T)$$

$$\text{FIRST}'(a\beta) = \{a\} \quad (a \in T)$$

$$\text{FIRST}'(A\beta) = \quad (A \in N)$$

if $\text{NULLABLE}(A)$

then $\text{FIRST}'(\beta) \cup (\text{FIRST}(A) - \{\varepsilon\})$

else $\text{FIRST}(A)$

Now, abuse notation and just write

$\text{FIRST}(\alpha)$ instead of $\text{FIRST}'(\alpha)$ 22