

Table-driven LL(1) parsing algorithm

References: Michael Scott's Programming Languages, Tucker-Noonan's Programming Languages

Website: <https://www.cs.princeton.edu/courses/archive/spring20/cos320/LL1/>

Grammar in “augmented form”:

S is an artificial start state, whereas E is the actual start state. We treat \$ as a terminal symbol denoting the end of input. Must remove left recursion from BNF and do left factoring. Example below: E stands for Expr, T for Term, and F for Factor. Removing left recursion from $E \rightarrow E + T$ gives us $E \rightarrow T E'$, $E' \rightarrow + T E'$, and $E' \rightarrow \epsilon$.

$S \rightarrow E \$$
$E \rightarrow T E'$
$E' \rightarrow + T E'$
$E' \rightarrow \epsilon$
$T \rightarrow F T'$
$T' \rightarrow * F T'$
$T' \rightarrow \epsilon$
$F \rightarrow (E)$
$F \rightarrow id$

$$FIRST(\alpha) \equiv \{a : \alpha \Rightarrow^* a\beta\} \cup (\text{if } \alpha \Rightarrow^* \epsilon \text{ then } \{\epsilon\} \text{ else } \emptyset)$$

$$FOLLOW(A) \equiv \{a : S \Rightarrow^+ A a\beta\} \cup (\text{if } S \Rightarrow^* A \text{ then } \{\epsilon\} \text{ else } \emptyset)$$

A: nonterminal, a: terminal, ϵ : empty string, other Greek letters: sentential form (i.e., a string of terminal and/or non-terminal symbols).

FIRST(α) = Set of terminal symbols (tokens) occurring on the left of any sentential form derived from α .

Special rules: Add ϵ to $FIRST(\alpha)$ if α is nullable (i.e., α derives ϵ). For a terminal symbol a , $FIRST(a) = \{a\}$.

FOLLOW(A) = Set of terminal symbols appearing after A in any sentential form derived from S.

Special rule: Add ϵ to $FOLLOW(A)$ when S derives a sentential form ending with A.

Non-terminal	Nullable?	FIRST	FOLLOW
S	✗	(, id	
E	✗	(, id), \$
E'	✓	+ , ϵ), \$
T	✗	(, id	+ ,), \$
T'	✓	* , ϵ	+ ,), \$
F	✗	(, id	+ , * ,), \$

LL(1) parse table construction:

Consider Table[A, a] and each of A's production rules $A \rightarrow \alpha$.

- If $FIRST(\alpha)$ contains a, enter $A \rightarrow \alpha$ into Table[A, a]
- If $FIRST(\alpha)$ contains ϵ , enter $A \rightarrow \epsilon$ into Table[A, b] for each terminal b in $FOLLOW(A)$.

	\$	+	*	()	id
S				$S \rightarrow E \$$		$S \rightarrow E \$$
E				$E \rightarrow T E'$		$E \rightarrow T E'$
E'	$E' \rightarrow \epsilon$	$E' \rightarrow + T E'$			$E' \rightarrow \epsilon$	
T				$T \rightarrow F T'$		$T \rightarrow F T'$
T'	$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$	$T' \rightarrow * F T'$		$T' \rightarrow \epsilon$	
F				$F \rightarrow (E)$		$F \rightarrow id$

Multiple productions in a cell means predictive parsing like LL(1) is impossible for the given grammar.

Parsing Algorithm:

First, push S into the stack. Repeat the following steps, where X is the current top of the stack and a is the current input token:

- If $X = a = \$ \rightarrow$ Successful parsing!
- Else if X is a terminal symbol:
 - $X = a$: pop X off the stack and get the next input token.
 - $X \neq a$: Error!
- Else if X is a nonterminal symbol:
 - Table[X, a] has nothing \rightarrow Syntax Error!
 - Table[X, a] contains $X \rightarrow \alpha$: pop X and push α with the leftmost symbol of α on top.