

# LL parser

From Wikipedia, the free encyclopedia

In computer science, an **LL parser** is a top-down parser for a subset of context-free languages. It parses the input from **L**eft to **r**ight, performing **L**eftmost derivation of the sentence.

An LL parser is called an LL( $k$ ) parser if it uses  $k$  tokens of lookahead when parsing a sentence. If such a parser exists for a certain grammar and it can parse sentences of this grammar without backtracking then it is called an LL( $k$ ) grammar. LL( $k$ ) grammars can generate more languages the higher the number  $k$  of lookahead tokens.<sup>[1]</sup> A corollary of this is that not all context-free languages can be recognized by an LL( $k$ ) parser. An LL parser is called an LL(\*) parser (an LL-regular parser<sup>[2]</sup>) if it is not restricted to a finite  $k$  tokens of lookahead, but can make parsing decisions by recognizing whether the following tokens belong to a regular language (for example by means of a Deterministic Finite Automaton).

LL grammars, particularly LL(1) grammars, are of great practical interest, as parsers for these grammars are easy to construct, and many computer languages are designed to be LL(1) for this reason. LL parsers are table-based parsers, similar to LR parsers. LL grammars can also be parsed by recursive descent parsers.

## Contents

- 1 Overview
- 2 Parser
- 3 Concrete example
  - 3.1 Set up
  - 3.2 Parsing procedure
  - 3.3 Parser implementation in C++
  - 3.4 Parser implementation in Python
- 4 Remarks
- 5 Constructing an LL(1) parsing table
- 6 Constructing an LL( $k$ ) parsing table
- 7 Conflicts
  - 7.1 Terminology<sup>[3]</sup>
  - 7.2 LL(1) Conflicts
    - 7.2.1 FIRST/FIRST Conflict

- 7.2.1.1 Special Case: Left Recursion
    - 7.2.2 FIRST/FOLLOW Conflict
  - 7.3 Solutions to LL(1) Conflicts
    - 7.3.1 Left Factoring
    - 7.3.2 Substitution
    - 7.3.3 Left recursion removal<sup>[4]</sup>
- 8 See also
- 9 Notes
- 10 External links

## Overview

For a given context-free grammar, the parser attempts to find the leftmost derivation. Given an example grammar  $G$ :

1.  $S \rightarrow E$
2.  $E \rightarrow (E + E)$
3.  $E \rightarrow i$

the leftmost derivation for  $w = ((i + i) + i)$  is:

$$S \xRightarrow{(1)} E \xRightarrow{(2)} (E + E) \xRightarrow{(2)} ((E + E) + E) \xRightarrow{(3)} ((i + E) + E) \xRightarrow{(3)} ((i + i) + E) \xRightarrow{(3)} ((i + i) + i)$$

Generally, there are multiple possibilities when selecting a rule to expand given (leftmost) non-terminal. In the previous example of the leftmost derivation, in step 2:

$$S \xRightarrow{(1)} E \xRightarrow{(?)} ?$$

We can choose from two rules:

2.  $E \rightarrow (E + E)$
3.  $E \rightarrow i$

To be effective, the parser must be able to make this choice deterministically when possible, without backtracking. For some grammars, it can do this by peeking on the unread input (without reading). In our example, if the parser knows that the next unread symbol is  $($  the only correct rule that can be used is 2.

Generally,  $LL(k)$  parser can look ahead at  $k$  symbols. However, given a grammar, the problem of determining if there exists a  $LL(k)$  parser for some  $k$  that recognizes it is undecidable. For each  $k$ , there is a language that cannot be recognized by  $LL(k)$  parser, but can be by  $LL(k + 1)$ .

We can use the above analysis to give the following formal definition:

Let  $G$  be a context-free grammar and  $k \geq 1$ . We say that  $G$  is  $LL(k)$ , if and only if for any two leftmost derivations:

1.  $S \Rightarrow \dots \Rightarrow wA\alpha \Rightarrow \dots \Rightarrow w\beta\alpha \Rightarrow \dots \Rightarrow wx$
2.  $S \Rightarrow \dots \Rightarrow wA\alpha \Rightarrow \dots \Rightarrow w\gamma\alpha \Rightarrow \dots \Rightarrow wy$

Following condition holds: Prefix of the string  $x$  of length  $k$  equals the prefix of the  $y$  of length  $k$  implies  $\beta = \gamma$ .

In this definition,  $S$  is the starting and  $A$  any non-terminal. The already derived input  $w$ , and yet unread  $x$  and  $y$  are strings of terminals. The greek letters  $\alpha$ ,  $\beta$  and  $\gamma$  represent any string of both terminals and non-terminals (possibly empty). The prefix length corresponds to the lookahead buffer size, and the definition says that this buffer is enough to distinguish between any two derivations of different words.

## Parser

The  $LL(k)$  parser is a deterministic pushdown automaton with the ability to peek on the next  $k$  input symbols without reading. This capability can be emulated by storing the lookahead buffer contents in the finite state space, since both buffer and input alphabet are finite in size. As a result, this does not make the automaton more powerful, but is a convenient abstraction.

The stack alphabet  $\Gamma = N \cup \Sigma$ , where:

- $N$  is the set of non-terminals;
- $\Sigma$  the set of terminal (input) symbols with a special end-of-input (EOI) symbol \$.

The parser stack initially contains the starting symbol above the EOI:  $[S \$]$ . During operation, the parser repeatedly replaces the symbol  $X$  on top of the stack:

- with some  $\alpha$ , if  $X \in N$  and there is a rule  $X \rightarrow \alpha$ ;

- with  $\epsilon$  (in some notations  $\lambda$ ), i.e.  $X$  is popped off the stack, if  $X \in \Sigma$ . In this case, an input symbol  $x$  is read and if  $x \neq X$ , the parser rejects the input.

If the last symbol to be removed from the stack is the EOI, the parsing is successful; the automaton accepts via an empty stack.

The states and the transition function are not explicitly given; they are specified (generated) using a more convenient *parse table* instead. The table provides the following mapping:

- row: top-of-stack symbol  $X$
- column:  $|w| \leq k$  lookahead buffer contents
- cell: rule number for  $X \rightarrow \alpha$  or  $\epsilon$

If the parser cannot perform a valid transition, the input is rejected (empty cells). To make the table more compact, only the non-terminal rows are commonly displayed, since the action is the same for terminals.

## Concrete example

### Set up

To explain an LL(1) parser's workings we will consider the following small LL(1) grammar:

1.  $S \rightarrow F$
2.  $S \rightarrow ( S + F )$
3.  $F \rightarrow a$

and parse the following input:

**( a + a )**

We construct a parsing table for this grammar by expanding all the terminals by column and all nonterminals by row. Later, the expressions are numbered by the position where the columns and rows cross. For example, the terminal '(' and non-terminal 'S' match for expression number 2. The table is as follows:

	(	)	a	+	\$
S	2	-	1	-	-
F	-	-	3	-	-

(Note that there is also a column for the special terminal, represented here as \$, that is used to indicate the end of the input stream.)

## Parsing procedure

In each step, the parser reads the next-available symbol from the input stream, and the top-most symbol from the stack. If the input symbol and the stack-top symbol match, the parser discards them both, leaving only the unmatched symbols in the input stream and on the stack.

Thus, in its first step, the parser reads the input symbol '(' and the stack-top symbol 'S'. The parsing table instruction comes from the column headed by the input symbol '(' and the row headed by the stack-top symbol 'S'; this cell contains '2', which instructs the parser to apply rule (2). The parser has to rewrite 'S' to '( S + F )' on the stack by removing 'S' from stack and pushing '(', 'S', '+', 'F', ')' onto the stack and this writes the rule number 2 to the output. The stack then becomes:

```
[ (, S, +, F, ), $ ]
```

Since the '(' from the input stream did not match the top-most symbol, 'S', from the stack, it was not removed, and remains the next-available input symbol for the following step.

In the second step, the parser removes the '(' from its input stream and from its stack, since they now match. The stack now becomes:

```
[ S, +, F, ), $ ]
```

Now the parser has an 'a' on its input stream and an 'S' as its stack top. The parsing table instructs it to apply rule (1) from the grammar and write the rule number 1 to the output stream. The stack becomes:

```
[ F, +, F, ), $ ]
```

The parser now has an 'a' on its input stream and an 'F' as its stack top. The

parsing table instructs it to apply rule (3) from the grammar and write the rule number 3 to the output stream. The stack becomes:

```
[ a, +, F, ), $ ]
```

In the next two steps the parser reads the '**a**' and '**+**' from the input stream and, since they match the next two items on the stack, also removes them from the stack. This results in:

```
[ F, ), $ ]
```

In the next three steps the parser will replace '**F**' on the stack by '**a**', write the rule number 3 to the output stream and remove the '**a**' and '**)**' from both the stack and the input stream. The parser thus ends with '**\$**' on both its stack and its input stream.

In this case the parser will report that it has accepted the input string and write the following list of rule numbers to the output stream:

```
[ 2, 1, 3, 3 ]
```

This is indeed a list of rules for a leftmost derivation of the input string, which is:

$$S \rightarrow ( S + F ) \rightarrow ( F + F ) \rightarrow ( a + F ) \rightarrow ( a + a )$$

## Parser implementation in C++

Below follows a C++ implementation of a table-based LL parser for the example language:

```
#include <iostream>
#include <map>
#include <stack>

enum Symbols {
    // the symbols:
    // Terminal symbols:
    TS_L_PARENS,    // (
    TS_R_PARENS,    // )
    TS_A,           // a
    TS_PLUS,        // +
    TS_EOS,         // $, in this case corresponds to '\0'
    TS_INVALID,     // invalid token

    // Non-terminal symbols:
    NTS_S,          // S
}
```

```
NTS_F      // F
};

/*
Converts a valid token to the corresponding terminal symbol
*/
enum Symbols lexer(char c)
{
    switch(c)
    {
        case '(': return TS_L_PARENS;
        case ')': return TS_R_PARENS;
        case 'a': return TS_A;
        case '+': return TS_PLUS;
        case '\0': return TS_EOS; // end of stack: the $ terminal symbol
        default: return TS_INVALID;
    }
}

int main(int argc, char **argv)
{
    using namespace std;

    if (argc < 2)
    {
        cout << "usage:\n\tll '(a+a)'" << endl;
        return 0;
    }

    // LL parser table, maps < non-terminal, terminal> pair to action
    map<enum Symbols, map<enum Symbols, int> > table;
    stack<enum Symbols> ss; // symbol stack
    char *p; // input buffer

    // initialize the symbols stack
    ss.push(TS_EOS); // terminal, $
    ss.push(NTS_S); // non-terminal, S

    // initialize the symbol stream cursor
    p = &argv[1][0];

    // setup the parsing table
    table[NTS_S][TS_L_PARENS] = 2;
    table[NTS_S][TS_A] = 1;
    table[NTS_F][TS_A] = 3;

    while(ss.size() > 0)
    {
        if(lexer(*p) == ss.top())
        {
            cout << "Matched symbols: " << lexer(*p) << endl;
            p++;
            ss.pop();
        }
        else
        {
            cout << "Rule " << table[ss.top()][lexer(*p)] << endl;
            switch(table[ss.top()][lexer(*p)])
            {
                case 1: // 1. S → F
                    ss.pop();
                    ss.push(NTS_F); // F
                    break;

                case 2: // 2. S → ( S + F )
```

```
        ss.pop();
        ss.push(TS_R_PARENS);    // )
        ss.push(NTS_F);         // F
        ss.push(TS_PLUS);       // +
        ss.push(NTS_S);         // S
        ss.push(TS_L_PARENS);   // (
        break;

    case 3: // 3. F → a
        ss.pop();
        ss.push(TS_A);          // a
        break;

    default:
        cout << "parsing table defaulted" << endl;
        return 0;
        break;
    }
}

cout << "finished parsing" << endl;

return 0;
}
```

## Parser implementation in Python

```
#All constants are indexed from 0
Term = 0
Rule = 1

# Terminals
T_LPAR = 0
T_RPAR = 1
T_A = 2
T_PLUS = 3
T_END = 4
T_INVALID = 5

# Non-terminals
N_S = 0
N_F = 1

#parse table
table = [[ 1, -1, 0, -1, -1, -1],
          [-1, -1, 2, -1, -1, -1]]

rules = [(Rule,N_F)],
          [(Term,T_LPAR), (Rule,N_S), (Term,T_PLUS), (Rule,N_F), (Term,T_RPAR)],
          [(Term,T_A)]

stack = [(Term,T_END), (Rule,N_S)]

def lexicalAnalysis(inputstring):
    print('Lexical analysis')
    tokens = []
    #cdict = {'+': T_PLUS, '(': T_LPAR, ')': T_RPAR, 'a': T_A}
    #for c in inputstring:
    #    tokens.append(cdict.get(c, T_INVALID))
```



```

#
# in the meantime it has been changed on wikipedia to simple mapping above,
# but the original if-elif-elif-else could be indented to make further distinction
# for multi-character terminals like between '-' and '->' .
for c in inputstring:
    if c == '+': tokens.append(T_PLUS)
    elif c == '(': tokens.append(T_LPAR)
    elif c == ')': tokens.append(T_RPAR)
    elif c == 'a': tokens.append(T_A)
    else: tokens.append(T_INVALID)
tokens.append(T_END)
print(tokens)
return tokens

def syntacticAnalysis(tokens):
    print('Syntactic analysis')
    position = 0
    while len(stack) > 0:
        (stype, svalue) = stack.pop()
        token = tokens[position]
        if stype == Term:
            if svalue == token:
                position += 1
                print('pop', svalue)
                if token == T_END:
                    print('input accepted')
            else:
                print('bad term on input:', token)
                break
        elif stype == Rule:
            print('svalue', svalue, 'token', token)
            rule = table[svalue][token]
            print('rule', rule)
            for r in reversed(rules[rule]):
                stack.append(r)
            print('stack', stack)

inputstring = '(a+a)'
syntacticAnalysis(lexicalAnalysis(inputstring))

```

## Remarks

As can be seen from the example, the parser performs three types of steps depending on whether the top of the stack is a nonterminal, a terminal or the special symbol \$:

- If the top is a nonterminal then it looks up in the parsing table on the basis of this nonterminal and the symbol on the input stream which rule of the grammar it should use to replace it with on the stack. The number of the rule is written to the output stream. If the parsing table indicates that there is no such rule then it reports an error and stops.
- If the top is a terminal then it compares it to the symbol on the input stream and if they are equal they are both removed. If they are not equal the parser reports an error and stops.
- If the top is \$ and on the input stream there is also a \$ then the parser

reports that it has successfully parsed the input, otherwise it reports an error. In both cases the parser will stop.

These steps are repeated until the parser stops, and then it will have either completely parsed the input and written a leftmost derivation to the output stream or it will have reported an error.

## Constructing an LL(1) parsing table

In order to fill the parsing table, we have to establish what grammar rule the parser should choose if it sees a nonterminal  $A$  on the top of its stack and a symbol  $a$  on its input stream. It is easy to see that such a rule should be of the form  $A \rightarrow w$  and that the language corresponding to  $w$  should have at least one string starting with  $a$ . For this purpose we define the *First-set* of  $w$ , written here as  $\mathbf{Fi}(w)$ , as the set of terminals that can be found at the start of some string in  $w$ , plus  $\varepsilon$  if the empty string also belongs to  $w$ . Given a grammar with the rules  $A_1 \rightarrow w_1, \dots, A_n \rightarrow w_n$ , we can compute the  $\mathbf{Fi}(w_i)$  and  $\mathbf{Fi}(A_i)$  for every rule as follows:

1. initialize every  $\mathbf{Fi}(A_i)$  with the empty set
2. set  $\mathbf{Fi}(w_i)$  to  $Fi(w_i)$  for every rule  $A_i \rightarrow w_i$ , where  $Fi$  is defined as follows:
  - $Fi(a w') = \{ a \}$  for every terminal  $a$
  - $Fi(A w') = \mathbf{Fi}(A)$  for every nonterminal  $A$  with  $\varepsilon$  not in  $\mathbf{Fi}(A)$
  - $Fi(A w') = \mathbf{Fi}(A) \setminus \{ \varepsilon \} \cup Fi(w')$  for every nonterminal  $A$  with  $\varepsilon$  in  $\mathbf{Fi}(A)$
  - $Fi(\varepsilon) = \{ \varepsilon \}$
3. add  $\mathbf{Fi}(w_i)$  to  $\mathbf{Fi}(A_i)$  for every rule  $A_i \rightarrow w_i$
4. do steps 2 and 3 until all  $\mathbf{Fi}$  sets stay the same.

Unfortunately, the First-sets are not sufficient to compute the parsing table. This is because a right-hand side  $w$  of a rule might ultimately be rewritten to the empty string. So the parser should also use the rule  $A \rightarrow w$  if  $\varepsilon$  is in  $\mathbf{Fi}(w)$  and it sees on the input stream a symbol that could follow  $A$ . Therefore we also need the *Follow-set* of  $A$ , written as  $\mathbf{Fo}(A)$  here, which is defined as the set of terminals  $a$  such that there is a string of symbols  $\alpha A a \beta$  that can be derived from the start symbol. We use  $\$$  as a special terminal indicating end of input stream and  $S$  as start symbol.

Computing the Follow-sets for the nonterminals in a grammar can be done as follows:

1. initialize  $\mathbf{Fo}(S)$  with  $\{ \$ \}$  and every other  $\mathbf{Fo}(A_i)$  with the empty set

2. if there is a rule of the form  $A_j \rightarrow wA_iw'$ , then
  - if the terminal  $a$  is in  $Fi(w')$ , then add  $a$  to  $\mathbf{Fo}(A_i)$
  - if  $\varepsilon$  is in  $Fi(w')$ , then add  $\mathbf{Fo}(A_j)$  to  $\mathbf{Fo}(A_i)$
  - if  $w'$  has length 0, then add  $\mathbf{Fo}(A_j)$  to  $\mathbf{Fo}(A_i)$
3. repeat step 2 until all  $Fo$  sets stay the same.

Now we can define exactly which rules will be contained where in the parsing table. If  $T[A, a]$  denotes the entry in the table for nonterminal  $A$  and terminal  $a$ , then

$T[A, a]$  contains the rule  $A \rightarrow w$  if and only if  
 $a$  is in  $\mathbf{Fi}(w)$  or  
 $\varepsilon$  is in  $\mathbf{Fi}(w)$  and  $a$  is in  $\mathbf{Fo}(A)$ .

If the table contains at most one rule in every one of its cells, then the parser will always know which rule it has to use and can therefore parse strings without backtracking. It is in precisely this case that the grammar is called an *LL(1) grammar*.

## Constructing an LL( $k$ ) parsing table

Until the mid-1990s, it was widely believed that LL( $k$ ) parsing (for  $k > 1$ ) was impractical, since the parser table would have exponential size in  $k$  in the worst case. This perception changed gradually after the release of the Purdue Compiler Construction Tool Set around 1992, when it was demonstrated that many programming languages can be parsed efficiently by an LL( $k$ ) parser without triggering the worst-case behavior of the parser. Moreover, in certain cases LL parsing is feasible even with unlimited lookahead. By contrast, traditional parser generators like yacc use LALR(1) parser tables to construct a restricted LR parser with a fixed one-token lookahead.

## Conflicts

As described in the introduction, LL(1) parsers recognize languages that have LL(1) grammars, which are a special case of context-free grammars (CFGs); LL(1) parsers cannot recognize all context-free languages. The LL(1) languages are a proper subset of the LR(1) languages which in turn are a proper subset of all context-free languages. In order for a CFG to be an LL(1) grammar, certain conflicts must not arise, which we describe in this section.

## Terminology<sup>[3]</sup>

Let  $A$  be a non-terminal.  $FIRST(A)$  is (defined to be) the set of terminals that can appear in the first position of any string derived from  $A$ .  $FOLLOW(A)$  is the union over  $FIRST(B)$  where  $B$  is any non-terminal that immediately follows  $A$  in the right hand side of a production rule.

## LL(1) Conflicts

There are 2 main types of LL(1) conflicts:

### FIRST/FIRST Conflict

The  $FIRST$  sets of two different grammar rules for the same non-terminal intersect. An example of an LL(1) FIRST/FIRST conflict:

```
S -> E | E 'a'
E -> 'b' | ε
```

$FIRST(E) = \{'b', \epsilon\}$  and  $FIRST(E 'a') = \{'b', 'a'\}$ , so when the table is drawn, there is conflict under terminal 'b' of production rule  $S$ .

### Special Case: Left Recursion

Left recursion will cause a FIRST/FIRST conflict with all alternatives.

```
E -> E '+' term | alt1 | alt2
```

### FIRST/FOLLOW Conflict

The  $FIRST$  and  $FOLLOW$  set of a grammar rule overlap. With an empty string ( $\epsilon$ ) in the  $FIRST$  set it is unknown which alternative to select. An example of an LL(1) conflict:

```
S -> A 'a' 'b'
A -> 'a' | ε
```

The  $FIRST$  set of  $A$  now is  $\{'a', \epsilon\}$  and the  $FOLLOW$  set  $\{'a'\}$ .

## Solutions to LL(1) Conflicts

## Left Factoring

For a general method, see removing left recursion.

A common left-factor is "factored out".

```
A -> X | X Y Z
```

becomes

```
A -> X B  
B -> Y Z | ε
```

Can be applied when two alternatives start with the same symbol like a FIRST/FIRST conflict.

Another example (more complex) using above FIRST/FIRST conflict example:

```
S -> E | E 'a'  
E -> 'b' | ε
```

becomes (merging into a single non-terminal)

```
S -> 'b' | ε | 'b' 'a' | 'a'
```

then through left-factoring, becomes

```
S -> 'b' E | E  
E -> 'a' | ε
```

## Substitution

Substituting a rule into another rule to remove indirect or FIRST/FOLLOW conflicts. Note that this may cause a FIRST/FIRST conflict.

## Left recursion removal<sup>[4]</sup>

A simple example for left recursion removal: The following production rule has left recursion on E

```
E -> E '+' T
```

```
-> T
```

This rule is nothing but list of Ts separated by '+'. In a regular expression form  $T ('+' T)^*$ . So the rule could be rewritten as

```
E -> T Z
Z -> '+' T Z
  -> ε
```

Now there is no left recursion and no conflicts on either of the rules.

However, not all CFGs have an equivalent LL(k)-grammar, e.g.:

```
S -> A | B
A -> 'a' A 'b' | ε
B -> 'a' B 'b' 'b' | ε
```

It can be shown that there does not exist any LL(k)-grammar accepting the language generated by this grammar.

## See also

- Comparison of parser generators
- Parse tree
- Top-down parsing
- Bottom-up parsing

## Notes

1. Rosenkrantz, D. J.; Stearns, R. E. (1970). "Properties of Deterministic Top Down Grammars" (PDF). *Information and Control* **17**: 226–256. doi:10.1016/s0019-9958(70)90446-8.
2. Dick Grune; Criel J.H. Jacobs (29 October 2007). *Parsing Techniques: A Practical Guide*. Springer. pp. 585–. ISBN 978-0-387-68954-8.
3. <http://www.cs.uaf.edu/~cs331/notes/LL.pdf>
4. Modern Compiler Design, Grune, Bal, Jacobs and Langendoen

## External links

- A tutorial on implementing LL(1) parsers in C# (<http://www.itu.dk/people/kfl/parsernotes.pdf>)
- Parsing Simulator (<http://www.supereasyfree.com/software/simulators>)

/compilers/principles-techniques-and-tools/parsing-simulator/parsing-simulator.php) This simulator is used to generate parsing tables LL(1) and to resolve the exercises of the book.

- LL(1) DSL PEG parser (toolkit framework) (<http://expressionscompiler.codeplex.com/>)
- Language theoretic comparison of LL and LR grammars (<http://cs.stackexchange.com/questions/43/language-theoretic-comparison-of-ll-and-lr-grammars>)

Retrieved from "[https://en.wikipedia.org/w/index.php?title=LL\\_parser&oldid=693632794](https://en.wikipedia.org/w/index.php?title=LL_parser&oldid=693632794)"

Categories: Parsing algorithms

---

- This page was last modified on 4 December 2015, at 02:30.
- Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.