

Compilers v2

By Cedric De Vijt

Foundations of Formal Languages

Regular Language: A language that can be expressed using regular expressions or recognized by finite automata (DFA or NFA). Regular languages are closed under union, concatenation, and Kleene star operations. Regular languages are less expressive than context-free languages. **Regular Expression:** A pattern describing a regular language, using symbols like '*' (zero or more), '+' (one or more), '|' (union), and '()' (grouping). **Context-Free Language (CFL):** Languages defined by context-free grammars (CFGs) and recognized by pushdown automata (PDAs). **Context-Free Grammar (CFG):** A grammar consisting of rules of the form $A \rightarrow \alpha$, where A is a nonterminal, and α is a string of terminals and/or nonterminals. CFGs generate context-free languages, which can be recognized by pushdown automata. **Strongly LL(k):** A grammar is strongly LL(k) if, for every nonterminal A and every lookahead string of k tokens w , there is at most one production rule $A \rightarrow \alpha$ that can be applied, determined solely by the first k tokens of the input (the lookahead). This eliminates any need for backtracking or additional context beyond the k -token lookahead.

First and Follow Sets

The **FIRST^k(α)** set for a string α (terminal, non-terminal, or sequence) contains all sequences of up to k terminals that can appear as the first k symbols in any derivation of α . If α derives a string shorter than k , the entire derived string is included.

Steps to Compute FIRST^k(α):

- If α is a terminal a : Add the string a (of length 1) to **FIRST^k(α)**.
- If α is a sequence $\alpha_1 \alpha_2 \dots \alpha_n$:
 - Initialize **FIRST^k(α)** = \emptyset .
 - For $i = 1$ to n :
 - Add to **FIRST^k(α)** all strings of length up to k from **FIRST^k(α_i)** if $\alpha_1, \dots, \alpha_{i-1}$ can all derive the empty string ϵ .
 - For strings shorter than k , concatenate with **FIRST^{k-|s|}($\alpha_{i+1} \dots \alpha_n$)** (where $|s|$ is the length of the string).
 - If α can derive ϵ , include ϵ in **FIRST^k(α)**.
- If α is a non-terminal A :
 - For each production $A \rightarrow \beta$, compute **FIRST^k(β)** and add its strings to **FIRST^k(A)**.
 - Repeat until **FIRST^k(A)** stabilizes (no new strings are added).
- Handle k -length lookahead: Truncate strings longer than k to their first k symbols.

The **FOLLOW^k(A)** set for a non-terminal A contains all sequences of up to k terminals that can appear immediately after A in some derivation from the start symbol. If A appears at the end of a derivation, include the end-of-input marker (e.g., $\$$).

Steps to Compute FOLLOW^k(A):

- Initialize:
 - For the start symbol S , add $\$$ (or a k -length end marker) to **FOLLOW^k(S)**.
 - Set **FOLLOW^k(A)** = \emptyset for all other non-terminals A .
- For each production $B \rightarrow \alpha A \beta$ (where A is a non-terminal):
 - Compute **FIRST^k(β)** and add its strings to **FOLLOW^k(A)**.
 - If β can derive ϵ , add **FOLLOW^k(B)** to **FOLLOW^k(A)**.
- Iterate:
 - Repeat step 2 across all productions until no new strings are added to any **FOLLOW^k(A)**.
- Handle k -length lookahead: Ensure all strings in **FOLLOW^k(A)** are of length up to k , truncating longer strings to their first k symbols.

LR Items

An **LR(k)** item is a production with a dot (\bullet) marking a position in the right-hand side, along with a k -symbol lookahead string. The dot indicates how much of the production has been recognized so far. An item $[A \rightarrow \alpha \bullet \beta, w]$ means we have seen α and expect to see β , with lookahead w .

Types of LR Items:

- Shift item:** $[A \rightarrow \alpha \bullet a \beta, w]$ where a is a terminal
- Reduce item:** $[A \rightarrow \alpha \bullet, w]$ where the dot is at the end
- Goto item:** $[A \rightarrow \alpha \bullet B \beta, w]$ where B is a non-terminal

Steps to Compute CLOSURE(I) for a set of items I :

- Initialize **CLOSURE(I)** = I .
- For each item $[A \rightarrow \alpha \bullet B \beta, w]$ in **CLOSURE(I)** where B is a non-terminal:
 - For each production $B \rightarrow \gamma$:
 - Compute **FIRST^k(βw)** (concatenate β and lookahead w , then take first k symbols).
 - For each string $u \in \text{FIRST}^k(\beta w)$:
 - Add item $[B \rightarrow \bullet \gamma, u]$ to **CLOSURE(I)** if not already present.
- Repeat step 2 until no new items are added to **CLOSURE(I)**.

Steps to Compute GOTO(I, X) for item set I and symbol X :

- Initialize $J = \emptyset$.
- For each item $[A \rightarrow \alpha \bullet X \beta, w]$ in I :
 - Add item $[A \rightarrow \alpha X \bullet \beta, w]$ to J .
- Return **CLOSURE(J)**.

Steps to Construct the LR(k) Automaton:

- Create the initial state $I_0 = \text{CLOSURE}(\{[S' \rightarrow \bullet S, \$^k]\})$ where S' is the augmented start symbol.
- Initialize the set of states $C = \{I_0\}$ and a worklist $W = \{I_0\}$.
- While $W \neq \emptyset$:
 - Remove a state I from W .
 - For each symbol X (terminal or non-terminal) such that **GOTO(I, X) $\neq \emptyset$** :
 - Let $J = \text{GOTO}(I, X)$.
 - If $J \notin C$:
 - Add J to C and W .
 - Add transition $I \xrightarrow{X} J$ to the automaton.
- The resulting automaton defines the LR parsing table.

Parsing Conflicts

Shift-Reduce Conflict

A shift-reduce conflict occurs when the parser cannot decide whether to:

- Shift:** Move the next input symbol onto the stack
 - Reduce:** Apply a production rule to reduce symbols on the stack
- In terms of LR items, this happens when a state contains both:
- A shift item: $[A \rightarrow \alpha \bullet a \beta, w]$ (expecting terminal a)
 - A reduce item: $[B \rightarrow \gamma \bullet, a]$ (ready to reduce with lookahead a)

The conflict manifests in the parsing table as both a shift action and a reduce action for entry $[state, a]$.

Reduce-Reduce Conflict

A reduce-reduce conflict occurs when the parser cannot decide which production to use for reduction. This happens when a state contains multiple reduce items with overlapping lookaheads:

- $[A \rightarrow \alpha \bullet, w]$
- $[B \rightarrow \beta \bullet, w]$

Both items indicate that reduction is possible with the same lookahead w , but the parser cannot determine which production rule to apply.

The conflict manifests in the parsing table as multiple reduce actions for the same entry $[state, w]$.

Resolution

These conflicts indicate that the grammar is not in the respective LR class (LR(k), LALR(k), or SLR(k)). Resolution strategies include:

- Increasing the lookahead length k
- Grammar transformation (left-factoring, eliminating ambiguity)
- Using precedence and associativity rules

Top-Down Parsing (LL)

Remove Unnecessary Productions 1) **Unproductive Symbols:** Variables that cannot derive terminal strings. 2) **Unreachable Symbols:** Symbols not reachable from the start variable. **Procedure:** Remove unproductive symbols first, then unreachable ones.

Left Recursion Elimination 1) **Immediate Left Recursion:** For productions of the form $A \rightarrow A\alpha|\beta$, replace with:

$$A \rightarrow \beta A'$$
$$A' \rightarrow \alpha A' | \epsilon$$

2) **Indirect Left Recursion:** For productions $A \rightarrow B\alpha|\beta$ and $B \rightarrow A\gamma|\delta$, eliminate by substituting B in A 's productions. **Left Factoring** 1) **Common Prefixes:** For productions of the form $A \rightarrow \alpha\beta_1|\alpha\beta_2$, replace with:

$$A \rightarrow \alpha A'$$
$$A' \rightarrow \beta_1|\beta_2$$

2) **Multiple Common Prefixes:** For productions $A \rightarrow \alpha\beta_1|\alpha\beta_2|\alpha\beta_3$, factor out the common prefix:

$$A \rightarrow \alpha A'$$
$$A' \rightarrow \beta_1|\beta_2|\beta_3$$

LL(k) Grammar Requirements 1) **No Left Recursion:** The grammar must not have left recursion. 2) **Disjoint First/Follow Sets:** For any variable A , the sets **FIRST(A)** and **FOLLOW(A)** must be disjoint. 3) **Predictive Parsing Table:** The parsing table must be constructed such that for each variable A and terminal a , there is at most one production $A \rightarrow \alpha$ where $a \in \text{FIRST}(\alpha)$ or $a \in \text{FOLLOW}(A)$. **LL Parsing Table Construction** 1) **Compute FIRST Sets:** For each variable, compute the **FIRST** set, which contains the terminals that can appear at the beginning of any string derived from that variable. 2) **Compute FOLLOW Sets:** For each variable, compute the **FOLLOW** set, which contains the terminals that can appear immediately after that variable in any derivation. 3) **Construct Parsing Table:** For each variable A and terminal a : - If $a \in \text{FIRST}(A)$, add the production $A \rightarrow \alpha$ to the table. - If $\epsilon \in \text{FIRST}(A)$ and $a \in \text{FOLLOW}(A)$, add the production $A \rightarrow \epsilon$. - If $a \notin \text{FIRST}(A)$ and $a \notin \text{FOLLOW}(A)$, the grammar is not LL(k). **LL Parsing Algorithm** 1) **Initialization:** Start with the stack containing the start variable and the input string. 2) **Top of Stack:** If the top of the stack is a terminal, match it with the input symbol. 3) **Variable Handling:** If the top of the stack is a variable A : - Look up the parsing table for A and the current input symbol. - If a production $A \rightarrow \alpha$ is found, replace A on the stack with α . - If no production is found, report an error. 4) **Continue Until Done:** Repeat until the stack is empty or an error occurs.

Grammar and language proof approach

If you have to prove that the language generated by a grammar is LL(k) or LR(k), you can you can alter the grammar to make it easier to prove. The following steps can be used to systematically analyze and classify grammars:

Preliminary checks

- **Well-formed grammar:** Check if the grammar is well-formed, i.e., it adheres to the rules of context-free grammars.
- **Ambiguity detection:** Identify if the grammar is ambiguous by checking for multiple parse trees for the same string.
- **Useless symbol removal:** Eliminate symbols that do not contribute to the language generated by the grammar.

Example Languages

LR(0) language that is not LL(1)
 $L = \{a^n b^n \mid n \geq 0\}$
LR(0) language that is not LL(k) for any k
 $L = \{a^n b^n c^m \mid n, m \geq 1\}$
LL(2) language that is not LL(1)
 $L = \{a^n b c^n d \mid n \geq 1\} \cup \{a^n c b^n d \mid n \geq 1\}$
CFL that is not LR(1)
 $L = \{w w^R \mid w \in \{a, b\}^*\}$ where w^R is the reverse of w

Example Grammars

LL(k) Grammar

LL(1) grammar that is not strongly LL(1)
Theorem 5.4 "All LL(1) grammars are also strong LL(1), i.e. the classes of LL(1) and strong LL(1) grammars coincide."
LL(1) grammar that is not LALR(1)
 $S \rightarrow aX$
 $S \rightarrow Eb$
 $S \rightarrow Fc$
 $X \rightarrow Ec$
 $X \rightarrow Fb$
 $E \rightarrow A$
 $F \rightarrow A$
 $A \rightarrow \epsilon$

LL(2) grammar that is not strongly LL(2)
 $S \rightarrow aAa$
 $S \rightarrow bABa$
 $A \rightarrow b$
 $A \rightarrow \epsilon$
 $B \rightarrow b$
 $B \rightarrow c$

LL(2) and SLR(1) grammar but neither LL(1) nor LR(0) $S \rightarrow ab$
 $S \rightarrow ac$
 $S \rightarrow a$

LALR(k) Grammar

LALR(1) grammar that is not SLR(1)
 $S \rightarrow Aa$
 $S \rightarrow bAc$
 $S \rightarrow dc$
 $A \rightarrow d$

LR(k) Grammar

LR(1) grammar that is not LALR(1) and not LL(1)
 $S \rightarrow aAd$
 $S \rightarrow bBd$
 $S \rightarrow aBe$
 $S \rightarrow bAe$
 $A \rightarrow c$
 $B \rightarrow c$

LR(k + 1) grammar that is not LR(k) and not LL(k + 1)
 $S \rightarrow Ab^k c$
 $S \rightarrow Bb^k d$
 $A \rightarrow a$
 $B \rightarrow a$

Context Free Grammar

CFG that is not LR(k)
 $S \rightarrow aAc$
 $A \rightarrow bAb$
 $A \rightarrow b$

CFG that is not LR(0)
 $S \rightarrow a$
 $S \rightarrow abS$

SLR(k) Grammar

SLR(1) grammar that is not LR(0)
 $S \rightarrow Exp$
 $Exp \rightarrow Exp + Prod$
 $Exp \rightarrow Prod$
 $Prod \rightarrow Prod * Atom$
 $Prod \rightarrow Atom$
 $Atom \rightarrow Id$
 $Atom \rightarrow (Exp)$

Special examples

LL(1) language with a non-LL(1) grammar $L = \{w \mid w \text{ is a string of balanced parentheses}\}$
Grammar:
 $S \rightarrow SS$
 $S \rightarrow (S)$
 $S \rightarrow \epsilon$

Trivia

Q: Advantage of automata-based scanners over hand-built scanners: They require less coding effort. **Q:** Advantage of hand-built scanners over automata-based scanners: They can go beyond regular languages. **Q:** If you are implementing an LL(k) parser based on a given grammar, which version of recursion should you favor when manipulating or rewriting the grammar? Right recursion, LL parsers work top-down, predicting which production to use by looking ahead at the next k tokens. Left recursion causes problems for LL parsers because it leads to infinite recursion during parsing. **Q:** LLVM intermediate representation language: True: The number of registers is unlimited, The code has to be in SSA: single static assignment, Jumps can only target the start of a basic block. False: Basic blocks can end without a terminator **Q:** Argue that the exact versions of type checking and reachability analysis are in fact undecidable. Exact type checking and reachability analysis are undecidable because they reduce to the halting problem. For type checking, determining if a program's type is exactly correct requires predicting all possible execution paths, which may not terminate, akin to deciding if a Turing machine halts. Similarly, exact reachability analysis involves determining whether a program state is reachable, which also requires solving the halting problem, as it involves predicting if a computation leads to a specific state. Since the halting problem is undecidable, both tasks are undecidable in their exact forms. **Q:** Context free grammar for a language L on Σ. For all words $w \in \sigma^*$, we can use the grammar to determine whether $w \in L$ in time: polynomial in $|w|$; CYK (Cocke-Younger-Kasami) algorithm: Runs in $O(n^3)$ time where $n = |w|$, assuming the grammar is in Chomsky Normal Form. **Q:** Suppose we are given a grammar (not necessarily context free!) for a language L on Σ. For all words $w \in \sigma^*$, we can use the grammar to determine whether $w \in L$ in time: nope, that is an undecidable problem: For a general grammar (not necessarily context-free), the problem of determining whether a word $w \in L$ is undecidable. This is because general grammars correspond to Turing machines, and the membership problem for languages generated by unrestricted grammars is equivalent to the halting problem, which is undecidable. There is no algorithm that can decide, for all words $w \in \Sigma^*$, whether $w \in L$ in a finite amount of time. **Q:** Suppose we are given a deterministic finite automaton with state set Q for a language L on Σ. For all words $w \in \sigma^*$, we can use the automaton to determine whether $w \in L$ in time: $O(|w|)$: A deterministic finite automaton (DFA) processes an input word w by transitioning between states for each symbol in w . Since the state set Q is fixed, each symbol is processed in constant time, $O(1)$, by looking up the transition in the DFA's transition table. For a word of length $|w|$, the DFA makes $|w|$ transitions, resulting in a total time complexity of $O(|w|)$. **Q:** How can we make sure the sequence "true" is always deemed a contant and not a string? Define "true" as a reserved keyword or boolean literal in the lexer rules, distinct from string literals (e.g., quoted strings like "true"). **Q:** What is the longest match principle? A: The principle states that when there is a choice between several possible matches during tokenization, the lexer should choose the longest possible match that forms a valid token. **Q:** Give an arithmetic-expression tree such that the minimal number of registers required to evaluate it is exactly 7: We need to ensure that the computation's register usage peaks at 7, meaning at least one point in the evaluation requires 7 registers to hold intermediate results, and no evaluation requires more. $(((((a + b) + c) + d) + e) + f) + g)$