

Compilers v2

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 $\text{FIRST}^k(\alpha)$ 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 $\text{FIRST}^k(\alpha)$:

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

The $\text{FOLLOW}^k(A)$ 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 $\text{FOLLOW}^k(A)$:

- Initialize:
 - For the start symbol S , add \$ (or a k -length end marker) to $\text{FOLLOW}^k(S)$.
 - Set $\text{FOLLOW}^k(A) = \emptyset$ for all other non-terminals A .
- For each production $B \rightarrow \alpha A \beta$ (where A is a non-terminal):
 - Compute $\text{FIRST}^k(\beta)$ and add its strings to $\text{FOLLOW}^k(A)$.
 - If β can derive ϵ , add $\text{FOLLOW}^k(B)$ to $\text{FOLLOW}^k(A)$.
- Iterate:
 - Repeat step 2 across all productions until no new strings are added to any $\text{FOLLOW}^k(A)$.
- Handle k -length lookahead: Ensure all strings in $\text{FOLLOW}^k(A)$ 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 (•) 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 $\text{CLOSURE}(I) = I$.
- For each item $[A \rightarrow \alpha \bullet B\beta, w]$ in $\text{CLOSURE}(I)$ where B is a non-terminal:
 - For each production $B \rightarrow \gamma$:
 - Compute $\text{FIRST}^k(\beta w)$ (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 $\text{CLOSURE}(I)$ if not already present.
- Repeat step 2 until no new items are added to $\text{CLOSURE}(I)$.

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

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

LR(k) Automaton Construction

- Augment the Grammar:** Add a new start symbol S' and production $S' \rightarrow S\$,$ where S is the original start symbol and \$ is the k -symbol end-of-input marker.
- Initialize the First State:** Compute the initial state $I_0 = \text{CLOSURE}(\{[S' \rightarrow \bullet S \$]\})$, where the dot (•) indicates the current position in the production.
- Compute CLOSURE:** For a set of items I :
 - Start with $\text{CLOSURE}(I) = I$.
 - For each item $[A \rightarrow \alpha \bullet B\beta, w]$ in $\text{CLOSURE}(I)$, where B is a non-terminal:
 - For each production $B \rightarrow \gamma$, compute $\text{FIRST}^k(\beta w)$.
 - For each $u \in \text{FIRST}^k(\beta w)$, add $[B \rightarrow \bullet \gamma, u]$ to $\text{CLOSURE}(I)$ if not present.
 - Repeat until no new items are added.
- Compute GOTO:** For a state I and symbol X (terminal or nonterminal):
 - Initialize $J = \emptyset$.
 - For each item $[A \rightarrow \alpha \bullet X\beta, w]$ in I , add $[A \rightarrow \alpha X \bullet \beta, w]$ to J .
 - Return $\text{CLOSURE}(J)$.
- Build the State Machine:**
 - Initialize the set of states $C = \{I_0\}$ and worklist $W = \{I_0\}$.
 - While $W \neq \emptyset$:
 - Remove a state I from W .
 - For each symbol X where $\text{GOTO}(I, X) \neq \emptyset$:
 - 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.

6. Construct the Parsing Table:

- For each state $I \in C$:
 - If $[A \rightarrow \alpha \bullet a\beta, w] \in I$ and $\text{GOTO}(I, a) = J$, add shift action to state J for terminal a .
 - If $[A \rightarrow \alpha \bullet, w] \in I$, add reduce action by $A \rightarrow \alpha$ for lookahead w .
 - If $[S' \rightarrow S \bullet \$^k] \in I$, add accept action for $\k .
- Check for conflicts: Shift-reduce or reduce-reduce conflicts indicate the grammar is not LR(k).

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

Grammar Proofs

Note: If a language described by a grammar is specified, the grammar can be altered to achieve a certain parser type while preserving the language generated.

Preliminary Checks

Before testing for specific parser types, perform these essential checks:

1. Well-formed grammar verification

- Ensure the CFG has no useless symbols (nonterminals that cannot derive any terminal string or cannot be reached from the start symbol).
- Remove useless symbols:
 - Identify *productive* symbols (those deriving a terminal string).
 - Identify *reachable* symbols (those derivable from the start symbol).
 - Eliminate all other symbols.
- Verify the grammar has at least one terminal production and the start symbol is productive and reachable.

2. Ambiguity check

- An ambiguous grammar cannot be LL(k) for any k , as LL(k) requires deterministic parsing.

- LR(k), SLR(k), and LALR(k) can handle some ambiguous grammars if the parsing table is conflict-free.
- Testing for ambiguity is undecidable in general; heuristically, attempt to construct multiple leftmost (for LL(k)) or rightmost (for LR-based) derivations for the same string, or rely on parsing table conflicts.

3. Left recursion elimination

- For LL(k): Grammar must be free of direct or indirect left recursion to avoid infinite recursion in top-down parsing.
- Apply left recursion elimination (introducing new nonterminals to replace left-recursive productions) while preserving the language.
- Detect indirect left recursion by computing the transitive closure of the “derives” relation.
- For LR-based parsers (SLR, LALR, LR): Left recursion is acceptable.

LL(k) Parser Check

Definition: A grammar is LL(k) if it can be parsed top-down with k -lookahead symbols without backtracking, using a deterministic parsing table.

Algorithm Steps

1. Compute FIRST_k and FOLLOW_k sets

- See the dedicated section on First and Follow Sets for detailed computation steps.

2. Construct LL(k) parsing table

- For each production $A \rightarrow \alpha$, compute lookahead set:

$$\text{Lookahead}(A \rightarrow \alpha) = \text{FIRST}_k(\alpha) \cup (\text{FOLLOW}_k(A) \text{ if } \alpha \text{ is nullable}).$$

- Place production in table at $[A, t]$ for each terminal string t in the lookahead set.

3. Check for conflicts

- Grammar is LL(k) if no cell $[A, t]$ contains multiple productions.
- If conflicts arise, consider left factoring, increasing k , or rewriting the grammar.

LR(k) Parser Check

Definition: A grammar is LR(k) if it can be parsed bottom-up with k -lookahead symbols using a deterministic LR parsing table based on canonical LR(k) items.

Algorithm Steps

1. Construct LR(k) item sets

- Augment grammar with $S' \rightarrow S$.
- See the LR Items section for detailed item formats and CLOSURE/GOTO operations.
- Build canonical LR(k) DFA using closure and goto functions.

2. Build parsing table

- **Shift:** If $[A \rightarrow \alpha \cdot a\beta, w]$ and $\text{goto}(I, a) = J$, add shift to state J .
- **Reduce:** If $[A \rightarrow \alpha \cdot, w]$, add reduce by $A \rightarrow \alpha$ for lookahead w .
- **Accept:** If $[S' \rightarrow S \cdot, \epsilon]$, accept on end-of-input.
- Empty cells indicate error actions.

3. Check for conflicts

- Grammar is LR(k) if no shift-reduce or reduce-reduce conflicts exist.
- If conflicts occur, test higher k or check for inherent ambiguity.

4. Test incrementally

- Start with $k = 0$ (LR(0) uses empty lookaheads $w = \epsilon$).
- Increase k if needed, noting exponential growth in table size.

SLR(k) Parser Check

Definition: SLR(k) is a simplified LR(k) using FOLLOW_k sets for lookaheads instead of precise LR(k) lookaheads, making it less powerful but easier to compute.

Algorithm Steps

1. Build LR(0) DFA

- See LR Items section for DFA construction using LR(0) items.

2. Compute FOLLOW_k sets

- See the dedicated section on First and Follow Sets for computation steps.

3. Construct SLR(k) parsing table

- Shift and goto actions same as LR(0).
- For item $[A \rightarrow \alpha \cdot]$ in state I , add reduce by $A \rightarrow \alpha$ for each $t \in \text{FOLLOW}_k(A)$.
- Note: FOLLOW_k may introduce conflicts not present in LR(k).

4. Check for conflicts and test incrementally

- Grammar is SLR(k) if table has no conflicts.
- If conflicts arise, consider LALR(k) or LR(k).

LALR(k) Parser Check

Definition: LALR(k) is a compromise between SLR(k) and LR(k), using LR(0) DFA but with lookaheads computed by merging LR(k) states.

Algorithm Steps

1. Build LR(0) DFA

- See LR Items section for DFA construction.

2. Compute LALR(k) lookaheads

- Use spontaneous generation and propagation of lookaheads through the LR(0) DFA.
- Merge lookaheads for items in the same LR(0) state.

3. Construct LALR(k) parsing table

- Shift and goto actions identical to LR(0).
- For item $[A \rightarrow \alpha \cdot, w]$, add reduce by $A \rightarrow \alpha$ for lookahead w .

4. Check for conflicts and test incrementally

- Grammar is LALR(k) if table has no conflicts.
- If conflicts arise, test LR(k).

Parser Type Relationships and Optimization

- $\text{LR}(k) \supseteq \text{LALR}(k) \supseteq \text{SLR}(k)$.
- $\text{LL}(k) \not\subseteq \text{LR}(k)$ and $\text{LR}(k) \not\subseteq \text{LL}(k)$.
- If a grammar is not SLR(k), it may still be LALR(k) or LR(k), as LALR(k) and LR(k) use more precise lookaheads.
- If a grammar is LALR(k), it is also LR(k); if it is SLR(k), it is also LALR(k) and LR(k).

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)

$$L = \{w w^R \mid w \in \{a, b\}^*\} \text{ where } w^R \text{ is the reverse of } w$$

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)

$$\begin{aligned} S &\rightarrow aX \\ S &\rightarrow Eb \\ S &\rightarrow Fc \\ X &\rightarrow Ec \\ X &\rightarrow Fb \\ E &\rightarrow A \\ F &\rightarrow A \\ A &\rightarrow \epsilon \end{aligned}$$

LL(2) grammar that is not strongly LL(2)

$$\begin{aligned} S &\rightarrow aAa \\ S &\rightarrow bABa \\ A &\rightarrow b \\ A &\rightarrow \epsilon \\ B &\rightarrow b \\ B &\rightarrow c \end{aligned}$$

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

$$\begin{aligned} S &\rightarrow ac \\ S &\rightarrow a \end{aligned}$$

LALR(k) Grammar

LALR(1) grammar that is not SLR(1)

$$\begin{aligned} S &\rightarrow Aa \\ S &\rightarrow bAc \\ S &\rightarrow dc \\ A &\rightarrow d \end{aligned}$$

LR(k) Grammar

LR(1) grammar that is not LALR(1) and not LL(1)

$$\begin{aligned} S &\rightarrow aAd \\ S &\rightarrow bBd \\ S &\rightarrow aBe \\ S &\rightarrow bAe \\ A &\rightarrow c \\ B &\rightarrow c \end{aligned}$$

LR(k + 1) grammar that is not LR(k) and not LL(k + 1)

$$\begin{aligned} S &\rightarrow Ab^k c \\ S &\rightarrow Bb^k d \\ A &\rightarrow a \\ B &\rightarrow a \end{aligned}$$

Context Free Grammar

CFG that is not LR(k)

$$\begin{aligned} S &\rightarrow aAc \\ A &\rightarrow bAb \\ A &\rightarrow b \end{aligned}$$

CFG that is not LR(0)

$$\begin{aligned} S &\rightarrow a \\ S &\rightarrow abS \end{aligned}$$

SLR(k) Grammar

SLR(1) grammar that is not LR(0)

$$\begin{aligned} S &\rightarrow Exp \\ Exp &\rightarrow Exp + Prod \\ Exp &\rightarrow Prod \\ Prod &\rightarrow Prod * Atom \\ Prod &\rightarrow Atom \\ Atom &\rightarrow Id \\ Atom &\rightarrow (Exp) \end{aligned}$$

Special examples

LL(1) language with a non-LL(1) grammar

$$L = \{w \mid w \text{ is a string of balanced parentheses}\}$$

Grammar:

$$\begin{aligned} S &\rightarrow SS \\ S &\rightarrow (S) \\ S &\rightarrow \epsilon \end{aligned}$$

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 constant 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. $(((((a + b) + c) + d) + e) + f) + g)$