

# Compilers

## Key Definitions and Concepts

**Regular language:** languages that can be recognized by finite automata or described by regular expressions

**CFG:** formal grammar where every production rule has a single non-terminal symbol on the left-hand side and a string of terminals and/or non-terminals on the right-hand side.

**CFL:** if it can be generated by a context-free grammar (CFG)

**LL(k) language:** if there is an LL(k) grammar  $G_L$  that accepts it, i.e.:  $L(G_L) = L$ .

**Maximal munch:** the longest possible string that can be matched by a token is matched.

**Strongly LL(k) Grammar:** A  $CFGG = \langle V, T, P, S \rangle$  is strong LL(k) iff, for all pairs of rules  $A \rightarrow \alpha_1$  and  $A \rightarrow \alpha_2$  in  $P$ , (with  $\alpha_1 \neq \alpha_2$ ):

$First^k(\alpha_1 \cdot Follow^k(A)) \cap First^k(\alpha_2 \cdot Follow^k(A)) = \emptyset$

## Canonical Finite State Machine Example

Grammar:

$S \rightarrow A\$$

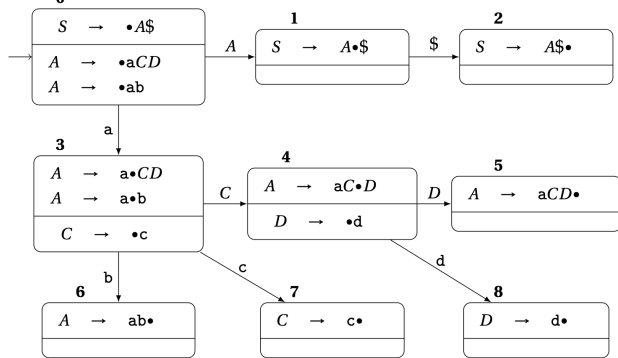
$A \rightarrow aCD$

$A \rightarrow ab$

$C \rightarrow c$

$D \rightarrow d$

CFSM:



The language of this grammar is LR(1).

## First and Follow

### Calculating FIRST

The  $FIRST^k(\alpha)$  set for a non-terminal  $S$  is defined as:  $FIRST^k(\alpha)$  includes terminals that appear as the first  $k$  symbols in any derivation of  $S$ .

To compute  $FIRST^k(\alpha)$ :

- For each production  $S \rightarrow \alpha$ , add the first terminal of  $\alpha$  to  $FIRST^k(\alpha)$ .
- If  $\alpha$  starts with a non-terminal  $A$ , include  $FIRST^k(A)$ .
- Consider  $k$ -length lookahead, ensuring that  $FIRST^k(\alpha)$  captures all terminals up to  $k$  symbols.

### Calculating FOLLOW

The  $FOLLOW^k(\alpha)$  set for a non-terminal  $S$  is defined as:  $FOLLOW^k(\alpha)$  includes terminals that can appear immediately after  $S$  in some string derivable from the start symbol, considering up to  $k$  symbols.

To compute  $FOLLOW^k(\alpha)$ :

- For each production  $A \rightarrow \alpha S \beta$ , add  $FIRST^k(\beta)$  to  $FOLLOW^k(\alpha)$ .
- If  $\beta$  can derive the empty string  $\epsilon$ , add  $FOLLOW^k(A)$  to  $FOLLOW^k(\alpha)$ .
- Ensure that  $FOLLOW^k(\alpha)$  includes all possible terminals appearing after  $S$  within  $k$  symbols.

## Grammars

### Grammar Fundamentals

**Context-Free Grammar (CFG):**  $G = (V, T, P, S)$  where: -  $V$ : finite set of non-terminals -  $T$ : finite set of terminals -  $P$ : finite set of productions  $A \rightarrow \alpha$  where  $A \in V$ ,  $\alpha \in (V \cup T)^*$  -  $S$ : start symbol  $\in V$

**Derivation:**  $\alpha \Rightarrow \beta$  if  $\alpha = \gamma A \delta$ ,  $\beta = \gamma \omega \delta$ , and  $A \rightarrow \omega \in P$

**Language:**  $L(G) = \{w \in T^* \mid S \Rightarrow^* w\}$

**Ambiguous Grammar:** A grammar where some string has multiple parse trees.

### Grammar Transformations

**Left Recursion:**  $A \rightarrow A\alpha \mid \beta$  - **Elimination:** Replace with  $A \rightarrow \beta A'$ ,  $A' \rightarrow \alpha A' \mid \epsilon$

**Left Factoring:**  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$  - **Factor out:**  $A \rightarrow \alpha A'$ ,  $A' \rightarrow \beta_1 \mid \beta_2$

**Grammar Hierarchy:**  $LL(k) \subset LR(k)$ ,  $SLR(k) \subset LALR(k) \subset LR(k)$

### LL(k) Grammar

A CFG is LL(k) if: The grammar can be parsed by a top-down parser with  $k$  lookahead tokens. The parsing table has no ambiguity and is deterministic.

### SLR(k) Grammar

A CFG is SLR(k) if: It is LR(k) (i.e., it can be parsed by an LR(k) parser). The parsing table has no shift/reduce or reduce/reduce conflicts when using the FOLLOW sets to resolve conflicts.

### LALR(k) Grammar

A CFG is LALR(k) if: It is LR(k). The parsing table can be constructed using LALR(k) lookahead by merging states with the same core (LALR(k) states). LALR(k) is a refinement of SLR(k) that reduces the number of states while maintaining LR(k) power.

### LR(k) Grammar

A CFG is LR(k) if: It has no shift/reduce or reduce/reduce conflicts for a given  $k$ . The grammar can be parsed by an LR(k) parser using the complete parsing table.

## Time Complexities

DFA:  $O(|w|)$

CFG: polynomial in  $|w|$

Grammar (not necessarily context free): undecidable problem

## 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  $\Sigma$ . 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  $\Sigma$ . 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  $\Sigma$ . 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.

## 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)**

$L = \{a^n bc^n d \mid n \geq 1\} \cup \{a^n cb^n d \mid n \geq 1\}$

**CFL that is not LR(1)**

English language

## 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:  
 $S \rightarrow SS$   
 $S \rightarrow (S)$   
 $S \rightarrow \epsilon$

# Parsing Grammar Proofs

## LL(k) Proof

A grammar is **LL(k)** if:

1. No left recursion exists.
2. For every non-terminal  $A$  with productions  $A \rightarrow \alpha \mid \beta$ :

$$\text{FIRST}_k(\alpha \cdot \text{FOLLOW}_k(A)) \cap \text{FIRST}_k(\beta \cdot \text{FOLLOW}_k(A)) = \emptyset$$

(If  $\alpha/\beta$  can derive  $\epsilon$ , include  $\text{FOLLOW}_k(A)$ .)

## SLR(k) Proof

A grammar is **SLR(k)** if:

1. Construct SLR(1) states (LR(0) items + FOLLOW for reduce actions).
2. No *shift-reduce* or *reduce-reduce* conflicts in the parsing table.
3. Conflicts tested using FOLLOW sets (not lookaheads).

## LR(k) Proof

A grammar is **LR(k)** if:

1. Construct canonical LR(k) states (with full lookaheads).
2. No conflicts in the parsing table (even with precise lookaheads).

## LALR(k) Proof

A grammar is **LALR(k)** if:

1. Merge LR(1) states with identical *cores* (same productions, different lookaheads).
2. Ensure no conflicts arise after merging (common in practice but weaker than LR).

## Hierarchy

$$\text{SLR} \subset \text{LALR} \subset \text{LR}, \quad \text{LL} \subset \text{CFG (non-overlapping)}$$

**Key Notes:** - For **LL(k)**: Compute  $\text{FIRST}_k$  and  $\text{FOLLOW}_k$  rigorously. - For **SLR**: Conflicts = overlap between shift terminals and FOLLOW of reduced non-terminal. - **LALR** merges LR(1) states; conflicts here imply not LALR but might still be LR. - **Left recursion** invalidates LL(k); **ambiguity** invalidates all.

## Constructing an LR Parsing Table

To determine whether a grammar is LR, SLR, or LALR, an LR parsing table must be constructed. The steps are as follows:

1. **Augment the Grammar:** Add a new start symbol  $S'$  and production  $S' \rightarrow S$ , where  $S$  is the original start symbol. This ensures the parser recognizes when it has fully parsed the input.
2. **Compute LR(0) Items:** Create the set of LR(0) items for the grammar. Each item represents a position within a production (e.g.,  $A \rightarrow \alpha \cdot \beta$ ).
3. **Closure Operation:** For each item  $A \rightarrow \alpha \cdot B\beta$ , add all productions of  $B$  with a dot at the beginning ( $B \rightarrow \cdot \gamma$ ) to the closure. Repeat until no new items can be added.
4. **Goto Function:** For each item set  $I$  and each grammar symbol  $X$  (terminal or non-terminal), compute the *goto*(I, X) by advancing the dot over  $X$  in all items of  $I$  and taking the closure of the result.
5. **Build the Canonical Collection:** Start with the initial state  $I_0 = \text{closure}(\{S' \rightarrow \cdot S\})$ . Use the *goto* function to generate all reachable states. This forms the canonical collection of LR(0) item sets.

### 6. Action and Goto Tables:

- For each state  $I_i$  and terminal  $a$ , if *goto*( $I_i, a$ ) leads to state  $I_j$ , add a **shift** action to the parsing table:  $\text{ACTION}[i, a] = \text{shift } j$ .
- For each state  $I_i$  containing an item  $A \rightarrow \alpha \cdot$  (complete production), add a **reduce** action:  $\text{ACTION}[i, a] = \text{reduce } A \rightarrow \alpha$  for all  $a \in \text{FOLLOW}(A)$  (for SLR). For LR(1) and LALR, use the lookahead symbols associated with the item.
- For the state containing  $S' \rightarrow S \cdot$ , add an **accept** action:  $\text{ACTION}[i, \$] = \text{accept}$ .

### 7. Check for Conflicts:

- **Shift-Reduce Conflict:** Occurs if a state has both a shift and reduce action for the same terminal.
- **Reduce-Reduce Conflict:** Occurs if a state has multiple reduce actions for the same terminal.
- If no conflicts exist, the grammar is:
  - **SLR:** If conflicts are resolved using FOLLOW sets.
  - **LR(1):** If conflicts are resolved using precise lookahead symbols.
  - **LALR:** If the grammar remains conflict-free after merging LR(1) states with identical cores.

8. **Example Construction:** Consider the grammar: Augment the grammar, compute the LR(0) items, build the canonical collection, and populate the parsing table. Check for conflicts to classify the grammar.

**Notes on Classification:** - If the grammar has no conflicts in the LR(0) table, it is **LR(0)**. - If conflicts are resolved using FOLLOW sets, it is **SLR**. - If conflicts are resolved using precise lookaheads, it is **LR(1)**. - If merging LR(1) states resolves conflicts, it is **LALR**.