

# 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.

## Canonical Finite State Machine Example

Grammar:

$S \rightarrow A\$$

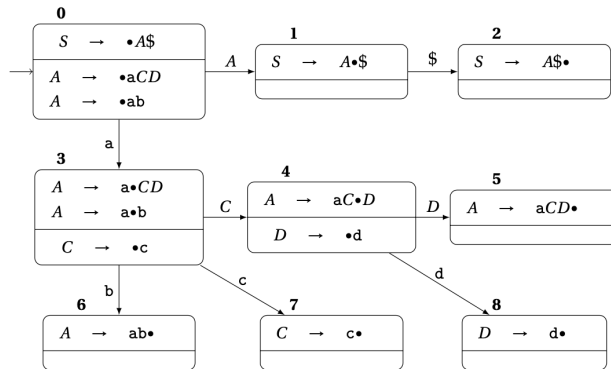
$A \rightarrow aCD$

$A \rightarrow ab$

$C \rightarrow c$

$D \rightarrow d$

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

### 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:  $\mathcal{O}(|w|)$

CFG: polynomial in  $|w|$

Grammar (not necessarily context free): undecidable problem

## Trivia

Q: Give an advantage of hand-built scanners over automata-based scanners? A: They can go beyond regular languages

Q: How can we make sure the sequence 'true' is always deemed a constant and not a string? A: Constants like 'true' should be recognized as a single token

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: Argue that the exact of both type checking and reachability analysis are in fact undecidable. A: The exact versions of both type checking and reachability analysis are undecidable because they would require solving the halting problem, which is proven to be undecidable. For type checking, determining the exact type of an expression in all possible cases would involve analyzing all potential program executions, which is equivalent to solving whether a program halts with certain inputs. Similarly, exact reachability analysis involves determining whether a specific program state can be reached, which also boils down to predicting the behavior of all possible executions, again leading to the halting problem. Since the halting problem is undecidable, both exact type checking and exact reachability analysis are also undecidable.

## 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)**

English language

## Example Grammars

**LL(1) grammar that is not strongly LL(1)**

$S \rightarrow A|B$

$A \rightarrow aA|a$

$B \rightarrow bB|b$

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

$S \rightarrow AB$

$A \rightarrow aAa|a$

$B \rightarrow bBb|b$

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

$S \rightarrow Aa$

$A \rightarrow Bb|Cc$

$B \rightarrow b$

$C \rightarrow c$

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

$S \rightarrow aB|aDc$

$B \rightarrow bBc|c$

$D \rightarrow bc|c$

## Special examples

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

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

Grammar:

$S \rightarrow SS$

$S \rightarrow (S)$

$S \rightarrow \epsilon$