

Compilers

Calculate First and Follow

First $^k(\alpha)$

1. Initialization: For every terminal a , set $First^k(a) = a$ (since a terminal generates only itself). Initialize $First^k(A)$ to an empty set for every variable A .
2. Iterative Update: Repeat until no changes occur: For each variable A with a production rule $A \rightarrow X_1 X_2 \dots X_n$: Compute $First^k(X_1), First^k(X_2), \dots, First^k(X_n)$. Concatenate elements from $First^k(X_1), First^k(X_2), \dots, First^k(X_n)$ and truncate the results to length k . Update $First^k(A)$ with these truncated elements if they are not already present.
3. Completion: The process stops when no $First^k$ sets are updated in an iteration, meaning all sets have stabilized.

Follow $^k(\alpha)$

1. Initialization: Initialize $Follow^k(X)$ to an empty set for each variable X .
2. Iterative Update: Repeat until no changes occur: For each rule $A \rightarrow B$: Compute $First^k()$ (where $$ is the part of the right-hand side following B). Add $First^k()$ elements to $Follow^k(B)$, completing them with $Follow^k(A)$ if needed to ensure they reach length k . Update $Follow^k(B)$ with these completed elements if they are not already present.
3. Completion: The process stops when no $Follow^k$ sets are updated in an iteration, meaning all sets have stabilized.

Grammars

LL(k)

1. Identify Production Rules for Each Non-Terminal
Initialization: For each non-terminal A , list all production rules of the form $A \rightarrow \alpha_1$ and $A \rightarrow \alpha_2$, where α_1 and α_2 are sequences of symbols (terminals and/or non-terminals). Compare Productions: For each pair of distinct production rules $A \rightarrow \alpha_1$ and $A \rightarrow \alpha_2$: Note that the possible right-hand sides are α_1 and α_2 .
2. Analyze Derivations with Lookahead
Derive Sentential Forms: For each pair $(A \rightarrow \alpha_1, A \rightarrow \alpha_2)$, consider the derivations of the form: $S \Rightarrow^* wA\gamma \Rightarrow w\alpha_1\gamma \Rightarrow^* wx_1$
 $S \Rightarrow^* wA\gamma \Rightarrow w\alpha_2\gamma \Rightarrow^* wx_2$ Here, w is a prefix of terminals, and γ is a sequence of symbols. Compute Lookahead Sets: Determine the lookahead sets $First_k(x_1)$ and $First_k(x_2)$ where x_1 and x_2 are terminal strings derived from $\alpha_1\gamma$ and $\alpha_2\gamma$, respectively.

3. Check for Disambiguation

Lookahead Comparison: For each pair $(A \rightarrow \alpha_1, A \rightarrow \alpha_2)$, compare the lookahead sets: Check if $First_k(x_1)$ equals $First_k(x_2)$. Verify Consistency: If $First_k(x_1) = First_k(x_2)$, ensure that α_1 and α_2 are identical: If $\alpha_1 \neq \alpha_2$ when the lookahead sets are equal, the grammar is not LL(k).

4. Conclude LL(k) Status

If all pairs of distinct production rules for every non-terminal either have different lookahead sets or identical right-hand sides, the grammar is LL(k). If any pair of rules has the same lookahead sets but different right-hand sides, the grammar is not LL(k).

Strong LL(k)

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

LR(k)

The main limitation to the bottom-up parser we have just defined is that it is non-deterministic. There are two sources of non-determinism in this parser:

Reduce-Reduce conflict: such conflicts occur when the top of the stack is the handle to two different rules, and the parser cannot decide which rule to reduce;

Shift-Reduce conflicts: such conflicts occur when the top of the stack constitutes a handle to some rule, but the parser cannot determine whether it should continue shifting or whether it should reduce now.

Time Complexities

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

CFG: polynomial in $|w|$

Grammar (not necessarily context free): undecidable problem

Undecidability of Type Checking and Reachability Analysis

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.

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

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)

$S \rightarrow aB|aDc$

$B \rightarrow bBc|c$

$D \rightarrow bc|c$