

Compiler Construction

Lecture 7: Syntax Analysis III ($LL(1)$ Parsing)

Thomas Noll

Lehrstuhl für Informatik 2
(Software Modeling and Verification)



noll@cs.rwth-aachen.de

<http://moves.rwth-aachen.de/teaching/ss-14/cc14/>

Summer Semester 2014

- 1 Recap: $LL(k)$ Grammars
- 2 Characterization of $LL(1)$
- 3 Computing Lookahead Sets
- 4 Deterministic Top-Down Parsing

$LL(k)$: reading of input from Left to right with k -lookahead, computing a Leftmost analysis

Definition ($LL(k)$ grammar)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$ and $k \in \mathbb{N}$. Then G has the $LL(k)$ property (notation: $G \in LL(k)$) if for all leftmost derivations of the form

$$S \Rightarrow_I^* wA\alpha \left\{ \begin{array}{l} \Rightarrow_I w\beta\alpha \Rightarrow_I^* wx \\ \Rightarrow_I w\gamma\alpha \Rightarrow_I^* wy \end{array} \right.$$

such that $\beta \neq \gamma$, it follows that $\text{first}_k(x) \neq \text{first}_k(y)$
(i.e., different productions must not yield the same lookahead).

Motivation:

- $k = 1$ sufficient to resolve nondeterminism in “most” practical applications
- Implementation of $LL(k)$ parsers for $k > 1$ rather involved (cf. ANTLR [ANother Tool for Language Recognition; formerly PCCTS] at <http://www.antlr.org/>)

Abbreviations: $\text{fi} := \text{first}_1$, $\text{fo} := \text{follow}_1$, $\Sigma_\varepsilon := \Sigma \cup \{\varepsilon\}$

Corollary

- ① For every $\alpha \in X^*$,

$$\text{fi}(\alpha) = \{a \in \Sigma \mid \text{ex. } w \in \Sigma^* : \alpha \Rightarrow^* aw\} \cup \{\varepsilon \mid \alpha \Rightarrow^* \varepsilon\} \subseteq \Sigma_\varepsilon$$

- ② For every $A \in N$,

$$\text{fo}(A) = \{x \in \text{fi}(\alpha) \mid \text{ex. } w \in \Sigma^*, \alpha \in X^* : S \Rightarrow_i^* wA\alpha\} \subseteq \Sigma_\varepsilon.$$

Definition (Lookahead set)

Given $\pi = A \rightarrow \beta \in P$,

$$\text{la}(\pi) := \text{fi}(\beta \cdot \text{fo}(A)) \subseteq \Sigma_\varepsilon$$

is called the **lookahead set** of π (where $\text{fi}(\Gamma) := \bigcup_{\gamma \in \Gamma} \text{fi}(\gamma)$).

Corollary

① For all $a \in \Sigma$,

$$a \in \text{la}(A \rightarrow \beta) \text{ iff } a \in \text{fi}(\beta) \text{ or } (\beta \Rightarrow^* \varepsilon \text{ and } a \in \text{fo}(A))$$

② $\varepsilon \in \text{la}(A \rightarrow \beta) \text{ iff } \beta \Rightarrow^* \varepsilon \text{ and } \varepsilon \in \text{fo}(A)$

- 1 Recap: $LL(k)$ Grammars
- 2 Characterization of $LL(1)$
- 3 Computing Lookahead Sets
- 4 Deterministic Top-Down Parsing

Characterization of $LL(1)$

Theorem 7.1 (Characterization of $LL(1)$)

$G \in LL(1)$ iff for all pairs of rules $A \rightarrow \beta \mid \gamma \in P$ (where $\beta \neq \gamma$):

$$la(A \rightarrow \beta) \cap la(A \rightarrow \gamma) = \emptyset.$$

Proof.

on the board



Remark: the above theorem generally does not hold if $k > 1$ (cf. exercises)

- 1 Recap: $LL(k)$ Grammars
- 2 Characterization of $LL(1)$
- 3 Computing Lookahead Sets
- 4 Deterministic Top-Down Parsing

Computing Lookahead Sets I

(see Waite/Goos: *Compiler Construction*, p. 164f)

Lemma 7.2 (Computation of fi/fo)

The sets $\text{fi}(\alpha) \subseteq \Sigma_\varepsilon$ (for $\alpha \in X^*$) and $\text{fo}(A) \subseteq \Sigma_\varepsilon$ (for $A \in N$) are the least sets such that:

① $\text{fi}(Y)$ for $Y \in X$:

- $Y \in \Sigma \implies \text{fi}(Y) = \{Y\}$
- $Y \rightarrow A_1 \dots A_k Z \alpha \in P, k \in \mathbb{N}, Z \in X, \varepsilon \in \text{fi}(A_1) \cap \dots \cap \text{fi}(A_k), a \in \text{fi}(Z) \implies a \in \text{fi}(Y)$
- $Y \rightarrow A_1 \dots A_k \in P, k \in \mathbb{N}, \varepsilon \in \text{fi}(A_1) \cap \dots \cap \text{fi}(A_k) \implies \varepsilon \in \text{fi}(Y)$

② $\text{fi}(Y_1 \dots Y_n)$ for $n \in \mathbb{N}, Y_i \in X$:

- $\varepsilon \in \text{fi}(Y_1) \cap \dots \cap \text{fi}(Y_{k-1}), a \in \text{fi}(Y_k), k \in [n] \implies a \in \text{fi}(Y_1 \dots Y_n)$
- $\varepsilon \in \text{fi}(Y_1) \cap \dots \cap \text{fi}(Y_n) \implies \varepsilon \in \text{fi}(Y_1 \dots Y_n)$

③ $\text{fo}(A)$ for $A \in N$:

- $\varepsilon \in \text{fo}(S)$
- $A \rightarrow \alpha B \beta \in P, a \in \text{fi}(\beta) \implies a \in \text{fo}(B)$
- $A \rightarrow \alpha B \beta \in P, \varepsilon \in \text{fi}(\beta), x \in \text{fo}(A) \implies x \in \text{fo}(B)$

Computing Lookahead Sets II

Corollary 7.3

- ① $A \rightarrow a\alpha \in P \implies a \in \text{fi}(A)$
- ② $A \rightarrow B\alpha \in P, a \in \text{fi}(B) \implies a \in \text{fi}(A)$
- ③ $A \rightarrow \varepsilon \in P \implies \varepsilon \in \text{fi}(A)$
- ④ $\text{fi}(\varepsilon) = \{\varepsilon\}$
- ⑤ $a \in \text{fi}(A) \implies a \in \text{fi}(A\alpha)$
- ⑥ $A \rightarrow \alpha B \in P, x \in \text{fo}(A) \implies x \in \text{fo}(B)$

Example 7.4

Grammar for
arithmetic expressions
(cf. Example 5.10):

$$G_{AE} : \begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid a \mid b \end{array}$$

- $F \rightarrow a \in P \implies a \in \text{fi}(F)$
- $T \rightarrow F \in P, a \in \text{fi}(F) \implies a \in \text{fi}(T)$
- $a \in \text{fi}(T) \implies \text{la}(T \rightarrow T * F) = \text{fi}(T * F \cdot \text{fo}(T)) \ni a$
- $a \in \text{fi}(F) \implies \text{la}(T \rightarrow F) = \text{fi}(F \cdot \text{fo}(T)) \ni a$
- $\implies a \in \text{la}(T \rightarrow T * F) \cap \text{la}(T \rightarrow F) \neq \emptyset$
- $\implies G_{AE} \notin LL(1)$

Fixing the Problem

(general methods later)

Example 7.5 (continuing Example 7.4)

Restructuring (such that $L(G'_{AE}) = L(G_{AE})$):

$$\begin{aligned} G'_{AE} : \quad E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow (E) \mid a \mid b \end{aligned}$$

| $A \in N$ | $\text{fi}(A)$ | $\text{fo}(A)$ |
|-----------|----------------------|----------------------------|
| E | $\{(, a, b\}$ | $\{\varepsilon,)\}$ |
| E' | $\{+, \varepsilon\}$ | $\{\varepsilon,)\}$ |
| T | $\{(, a, b\}$ | $\{+, \varepsilon,)\}$ |
| T' | $\{*, \varepsilon\}$ | $\{+, \varepsilon,)\}$ |
| F | $\{(, a, b\}$ | $\{*, +, \varepsilon,)\}$ |

| | |
|------------------------------|--|
| $A \rightarrow \beta \in P$ | $\text{la}(A \rightarrow \beta) = \text{fi}(\beta \cdot \text{fo}(A))$ |
| $E \rightarrow TE'$ | $\{(, a, b\}$ |
| $E' \rightarrow +TE'$ | $\{+\}$ |
| $E' \rightarrow \varepsilon$ | $\{\varepsilon,)\}$ |
| $T \rightarrow FT'$ | $\{(, a, b\}$ |
| $T' \rightarrow *FT'$ | $\{*\}$ |
| $T' \rightarrow \varepsilon$ | $\{+, \varepsilon,)\}$ |
| $F \rightarrow (E)$ | $\{((\}$ |
| $F \rightarrow a$ | $\{a\}$ |
| $F \rightarrow b$ | $\{b\}$ |

$\implies G'_{AE} \in LL(1)$

- 1 Recap: $LL(k)$ Grammars
- 2 Characterization of $LL(1)$
- 3 Computing Lookahead Sets
- 4 Deterministic Top-Down Parsing

Deterministic Top-Down Parsing

Approach: given $G \in CFG_{\Sigma}$,

- ① Verify that $G \in LL(1)$ by computing the lookahead sets and checking alternatives for disjointness
- ② Start with nondeterministic top-down parsing automaton $NTA(G)$
- ③ Use **1-symbol lookahead** to control the choice of expanding productions:

- $(aw, A\alpha, z) \vdash (aw, \beta\alpha, zi)$
if $\pi_i = A \rightarrow \beta$ and $a \in la(\pi_i)$

- $(\varepsilon, A\alpha, z) \vdash (\varepsilon, \beta\alpha, zi)$
if $\pi_i = A \rightarrow \beta$ and $\varepsilon \in la(\pi_i)$

- [matching steps as before: $(aw, a\alpha, z) \vdash (w, \alpha, z)$]

\implies **deterministic top-down parsing automaton $DTA(G)$**

Remarks:

- $DTA(G)$ is actually **not a pushdown automaton** (a is read but not consumed). But: can be simulated using the finite control.
- Advantage of using lookahead is **twofold**:
 - Removal of nondeterminism
 - Earlier detection of syntax errors
(in configurations $(aw, A\alpha, z)$ where $a \notin \bigcup_{A \rightarrow \beta \in P} la(A \rightarrow \beta)$)

The Deterministic Top-Down Automaton I

Definition 7.6 (Deterministic top-down parsing automaton)

Let $G = \langle N, \Sigma, P, S \rangle \in LL(1)$. The **deterministic top-down parsing automaton** of G , DTA(G), is defined by the following components.

- Input alphabet Σ , pushdown alphabet X , output alphabet $[p]$
- Configurations $\Sigma^* \times X^* \times [p]^*$, initial configuration (w, S, ε) , final configurations $\{\varepsilon\} \times \{\varepsilon\} \times [p]^*$ (as NTA(G))
- Action function

$\text{act} : \Sigma_\varepsilon \times X_\varepsilon \rightarrow \{(\alpha, i) \mid \pi_i = A \rightarrow \alpha\} \cup \{\text{pop}, \text{accept}, \text{error}\}$
with $\text{act}(x, A) := (\alpha, i)$ if $\pi_i = A \rightarrow \alpha$ and $x \in \text{la}(\pi_i)$
 $\text{act}(a, a) := \text{pop}$
 $\text{act}(\varepsilon, \varepsilon) := \text{accept}$
 $\text{act}(x, y) := \text{error} \quad \text{otherwise}$

- Transitions for $x \in \Sigma_\varepsilon$, $w \in \Sigma^*$, $Y \in X$, $\beta \in X^*$, and $z \in [p]^*$:

$$(xw, Y\beta, z) \vdash \begin{cases} (xw, \alpha\beta, zi) & \text{if } \text{act}(x, Y) = (\alpha, i) \\ (w, \beta, z) & \text{if } \text{act}(x, Y) = \text{pop} \end{cases}$$

The Deterministic Top-Down Automaton II

Example 7.7 (cf. Example 7.5)

$$\begin{aligned}
 G'_{AE} : \quad E &\rightarrow TE' & (1) \\
 E' &\rightarrow +TE' \mid \varepsilon & (2, 3) \\
 T &\rightarrow FT' & (4) \\
 T' &\rightarrow *FT' \mid \varepsilon & (5, 6) \\
 F &\rightarrow (E) \mid a \mid b & (7, 8, 9)
 \end{aligned}$$

| | |
|------------------------------|----------------------------------|
| $A \rightarrow \beta \in P$ | $\text{la}(A \rightarrow \beta)$ |
| $E \rightarrow TE'$ | $\{., a, b\}$ |
| $E' \rightarrow +TE'$ | $\{+\}$ |
| $E' \rightarrow \varepsilon$ | $\{\varepsilon, .\}$ |
| $T \rightarrow FT'$ | $\{., a, b\}$ |
| $T' \rightarrow *FT'$ | $\{*\}$ |
| $T' \rightarrow \varepsilon$ | $\{+, \varepsilon, .\}$ |
| $F \rightarrow (E)$ | $\{(.\)}$ |
| $F \rightarrow a$ | $\{a\}$ |
| $F \rightarrow b$ | $\{b\}$ |

$\text{act} : \Sigma_\varepsilon \times X_\varepsilon \rightarrow \{(\alpha, i) \mid \pi_i = A \rightarrow \alpha\} \cup \{\text{pop}, \text{accept}, \text{error}\}$ (empty = error)

| act | E | E' | T | T' | F | a | b | () | * | + | ε |
|---------------|------------|--------------------|------------|--------------------|------------|--------|---|-----|---|---|---------------|
| a | $(TE', 1)$ | | $(FT', 4)$ | | $(a, 8)$ | pop | | | | | |
| b | $(TE', 1)$ | | $(FT', 4)$ | | $(b, 9)$ | pop | | | | | |
| (| $(TE', 1)$ | | $(FT', 4)$ | | $((E), 7)$ | pop | | | | | |
|) | | $(\varepsilon, 3)$ | | $(\varepsilon, 6)$ | | pop | | | | | |
| * | | | | $(*FT', 5)$ | | pop | | | | | |
| + | | $(+TE', 2)$ | | $(\varepsilon, 6)$ | | pop | | | | | |
| ε | | $(\varepsilon, 3)$ | | $(\varepsilon, 6)$ | | accept | | | | | |

The Deterministic Top-Down Automaton III

Example 7.7 (continued)

| act | E | E' | T | T' | F | a | b | () | * | + | ϵ |
|------------|------------|-----------------|------------|-----------------|------------|--------|---|-----|---|---|------------|
| a | $(TE', 1)$ | | $(FT', 4)$ | | $(a, 8)$ | pop | | | | | |
| b | $(TE', 1)$ | | $(FT', 4)$ | | $(b, 9)$ | pop | | | | | |
| (| $(TE', 1)$ | | $(FT', 4)$ | | $((E), 7)$ | pop | | | | | |
|) | | $(\epsilon, 3)$ | | $(\epsilon, 6)$ | | pop | | | | | |
| * | | | | $(*FT', 5)$ | | pop | | | | | |
| + | | $(+TE', 2)$ | | $(\epsilon, 6)$ | | pop | | | | | |
| ϵ | | $(\epsilon, 3)$ | | $(\epsilon, 6)$ | | accept | | | | | |

Leftmost analysis of $(a)*b$:

| | |
|---|--|
| $\vdash ((a)*b, E , \epsilon)$ | $\vdash (\epsilon *b, E') T'E' , 1471486)$ |
| $\vdash ((a)*b, TE' , 1)$ | $\vdash (\epsilon *b,) T'E' , 14714863)$ |
| $\vdash ((a)*b, FT'E' , 14)$ | $\vdash (*b, T'E' , 14714863)$ |
| $\vdash ((a)*b, (E) T'E' , 147)$ | $\vdash (*b, *FT'E' , 147148635)$ |
| $\vdash (a)*b, E) T'E' , 147)$ | $\vdash (b, FT'E' , 147148635)$ |
| $\vdash (a)*b, TE') T'E' , 1471)$ | $\vdash (b, bT'E' , 1471486359)$ |
| $\vdash (a)*b, FT'E') T'E' , 14714)$ | $\vdash (\epsilon, T'E' , 1471486359)$ |
| $\vdash (a)*b, aT'E') T'E' , 147148)$ | $\vdash (\epsilon, E' , 14714863596)$ |
| $\vdash ()*b, T'E') T'E' , 147148)$ | $\vdash (\epsilon, \epsilon , 147148635963)$ |