

Compiler Construction

Lecture 6: Syntax Analysis II ($LL(k)$ Grammars)

Thomas Noll

Lehrstuhl für Informatik 2
(Software Modeling and Verification)

RWTH Aachen University

noll@cs.rwth-aachen.de

<http://www-i2.informatik.rwth-aachen.de/i2/cc12/>

Summer Semester 2012

1 Repetition: Nondeterministic Top-Down Parsing

2 Correctness of NTA(G)

3 Adding Lookahead

4 $LL(k)$ Grammars

5 Follow Sets

6 $LL(1)$ Grammars

Approach:

- ① Given $G \in CFG_{\Sigma}$, construct a **nondeterministic pushdown automaton** (PDA) which accepts $L(G)$ and which additionally computes corresponding leftmost derivations (similar to the proof of " $L(CFG_{\Sigma}) \subseteq L(PDA_{\Sigma})$ ")
 - input alphabet: Σ
 - pushdown alphabet: X
 - output alphabet: $[p]$
 - state set: not required
- ② Remove nondeterminism by allowing **lookahead** on the input:
 $G \in LL(k)$ iff $L(G)$ recognizable by deterministic PDA with lookahead of k symbols

Definition (Nondeterministic top-down parsing automaton)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$. The **nondeterministic top-down parsing automaton** of G , $NTA(G)$, is defined by the following components.

- **Input alphabet:** Σ
- **Pushdown alphabet:** X
- **Output alphabet:** $[p]$
- **Configurations:** $\Sigma^* \times X^* \times [p]^*$ (top of pushdown to the left)
- **Transitions** for $w \in \Sigma^*$, $\alpha \in X^*$, and $z \in [p]^*$:
 - expansion steps: if $\pi_i = A \rightarrow \beta$, then $(w, A\alpha, z) \vdash (w, \beta\alpha, zi)$
 - matching steps: for every $a \in \Sigma$, $(aw, a\alpha, z) \vdash (w, \alpha, z)$
- **Initial configuration** for $w \in \Sigma^*$: (w, S, ε)
- **Final configurations**: $\{\varepsilon\} \times \{\varepsilon\} \times [p]^*$

Remark: $NTA(G)$ is nondeterministic iff G contains $A \rightarrow \beta \mid \gamma$

1 Repetition: Nondeterministic Top-Down Parsing

2 Correctness of $NTA(G)$

3 Adding Lookahead

4 $LL(k)$ Grammars

5 Follow Sets

6 $LL(1)$ Grammars

Theorem 6.1 (Correctness of NTA(G))

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$ and NTA(G) as before. Then, for every $w \in \Sigma^*$ and $z \in [p]^*$,

$(w, S, \varepsilon) \vdash^* (\varepsilon, \varepsilon, z)$ iff z is a leftmost analysis of w

Theorem 6.1 (Correctness of NTA(G))

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$ and NTA(G) as before. Then, for every $w \in \Sigma^*$ and $z \in [p]^*$,

$(w, S, \varepsilon) \vdash^* (\varepsilon, \varepsilon, z)$ iff z is a leftmost analysis of w

Proof.

\implies (soundness): see Exercise 2

\impliedby (completeness): on the board



- 1 Repetition: Nondeterministic Top-Down Parsing
- 2 Correctness of $NTA(G)$
- 3 Adding Lookahead
- 4 $LL(k)$ Grammars
- 5 Follow Sets
- 6 $LL(1)$ Grammars

Adding Lookahead

Goal: resolve nondeterminism of $\text{NTA}(G)$ by supporting **lookahead of $k \in \mathbb{N}$ symbols** on the input
⇒ determination of expanding A -production by next k symbols

Adding Lookahead

Goal: resolve nondeterminism of $NTA(G)$ by supporting **lookahead of $k \in \mathbb{N}$ symbols** on the input
⇒ determination of expanding A -production by next k symbols

Definition 6.2 (first_k set)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$, $\alpha \in X^*$, and $k \in \mathbb{N}$. Then the **first_k set** of α , $\text{first}_k(\alpha) \subseteq \Sigma^*$, is given by

$$\begin{aligned} \text{first}_k(\alpha) := & \{v \in \Sigma^k \mid \text{ex. } w \in \Sigma^* \text{ such that } \alpha \Rightarrow^* vw\} \cup \\ & \{v \in \Sigma^{<k} \mid \alpha \Rightarrow^* v\} \end{aligned}$$

Adding Lookahead

Goal: resolve nondeterminism of $NTA(G)$ by supporting **lookahead of $k \in \mathbb{N}$ symbols** on the input
⇒ determination of expanding A -production by next k symbols

Definition 6.2 (first_k set)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$, $\alpha \in X^*$, and $k \in \mathbb{N}$. Then the **first_k set** of α , $\text{first}_k(\alpha) \subseteq \Sigma^*$, is given by

$$\begin{aligned}\text{first}_k(\alpha) := & \{v \in \Sigma^k \mid \text{ex. } w \in \Sigma^* \text{ such that } \alpha \Rightarrow^* vw\} \cup \\ & \{v \in \Sigma^{<k} \mid \alpha \Rightarrow^* v\}\end{aligned}$$

Remark: $\text{first}_k(\alpha)$ is effectively computable

Adding Lookahead

Goal: resolve nondeterminism of $NTA(G)$ by supporting **lookahead of $k \in \mathbb{N}$ symbols** on the input
⇒ determination of expanding A -production by next k symbols

Definition 6.2 (first_k set)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$, $\alpha \in X^*$, and $k \in \mathbb{N}$. Then the **first_k set** of α , $\text{first}_k(\alpha) \subseteq \Sigma^*$, is given by

$$\begin{aligned} \text{first}_k(\alpha) := & \{v \in \Sigma^k \mid \text{ex. } w \in \Sigma^* \text{ such that } \alpha \Rightarrow^* vw\} \cup \\ & \{v \in \Sigma^{<k} \mid \alpha \Rightarrow^* v\} \end{aligned}$$

Remark: $\text{first}_k(\alpha)$ is effectively computable

Example 6.3 (first_k set)

Let $G : S \rightarrow aSb \mid \varepsilon$.

① $\text{first}_1(ab) = \{a\} = \text{first}_2(a)$

Adding Lookahead

Goal: resolve nondeterminism of $NTA(G)$ by supporting **lookahead of $k \in \mathbb{N}$ symbols** on the input
⇒ determination of expanding A -production by next k symbols

Definition 6.2 (first_k set)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$, $\alpha \in X^*$, and $k \in \mathbb{N}$. Then the **first_k set** of α , $\text{first}_k(\alpha) \subseteq \Sigma^*$, is given by

$$\begin{aligned}\text{first}_k(\alpha) := \{v \in \Sigma^k \mid \text{ex. } w \in \Sigma^* \text{ such that } \alpha \Rightarrow^* vw\} \cup \\ \{v \in \Sigma^{<k} \mid \alpha \Rightarrow^* v\}\end{aligned}$$

Remark: $\text{first}_k(\alpha)$ is effectively computable

Example 6.3 (first_k set)

Let $G : S \rightarrow aSb \mid \varepsilon$.

- 1 $\text{first}_1(ab) = \{a\} = \text{first}_2(a)$
- 2 $\text{first}_3(S) = \{\varepsilon, ab, aab, aaa\}$

Adding Lookahead

Goal: resolve nondeterminism of $NTA(G)$ by supporting **lookahead** of $k \in \mathbb{N}$ **symbols** on the input
⇒ determination of expanding A -production by next k symbols

Definition 6.2 (first_k set)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$, $\alpha \in X^*$, and $k \in \mathbb{N}$. Then the **first_k set** of α , $\text{first}_k(\alpha) \subseteq \Sigma^*$, is given by

$$\begin{aligned}\text{first}_k(\alpha) := \{v \in \Sigma^k \mid \text{ex. } w \in \Sigma^* \text{ such that } \alpha \Rightarrow^* vw\} \cup \\ \{v \in \Sigma^{<k} \mid \alpha \Rightarrow^* v\}\end{aligned}$$

Remark: $\text{first}_k(\alpha)$ is effectively computable

Example 6.3 (first_k set)

Let $G : S \rightarrow aSb \mid \varepsilon$.

- 1 $\text{first}_1(ab) = \{a\} = \text{first}_2(a)$
- 2 $\text{first}_3(S) = \{\varepsilon, ab, aab, aaa\}$
- 3 $\text{first}_3(Sa) = \{a, aba, aab, aaa\}$

- 1 Repetition: Nondeterministic Top-Down Parsing
- 2 Correctness of $NTA(G)$
- 3 Adding Lookahead
- 4 $LL(k)$ Grammars
- 5 Follow Sets
- 6 $LL(1)$ Grammars

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

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

Definition 6.4 ($LL(k)$ 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).

Remarks:

- If $G \in LL(k)$, then the leftmost derivation step for $wA\alpha$ in

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

is determined by the next k symbols following w .

Remarks:

- If $G \in LL(k)$, then the leftmost derivation step for $wA\alpha$ in

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

is determined by the next k symbols following w .

- Corresponding computations of NTA(G):

$$(wx, S, \varepsilon) \vdash^* (x, A\alpha, z) \stackrel{(*)}{\vdash} (x, \beta\alpha, zi) \vdash^* (\varepsilon, \varepsilon, ziz')$$

$$(wy, S, \varepsilon) \vdash^* (y, A\alpha, z) \stackrel{(*)}{\vdash} (y, \gamma\alpha, zj) \vdash^* (\varepsilon, \varepsilon, zjz'')$$

where $\pi_i = A \rightarrow \beta$ and $\pi_j = A \rightarrow \gamma$

- Deterministic decision in $(*)$ possible if $\text{first}_k(x) \neq \text{first}_k(y)$

Remarks:

- If $G \in LL(k)$, then the leftmost derivation step for $wA\alpha$ in

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

is determined by the next k symbols following w .

- Corresponding computations of NTA(G):

$$(wx, S, \varepsilon) \vdash^* (x, A\alpha, z) \stackrel{(*)}{\vdash} (x, \beta\alpha, zi) \vdash^* (\varepsilon, \varepsilon, ziz')$$

$$(wy, S, \varepsilon) \vdash^* (y, A\alpha, z) \stackrel{(*)}{\vdash} (y, \gamma\alpha, zj) \vdash^* (\varepsilon, \varepsilon, zjz'')$$

where $\pi_i = A \rightarrow \beta$ and $\pi_j = A \rightarrow \gamma$

- Deterministic decision in $(*)$ possible if $\text{first}_k(x) \neq \text{first}_k(y)$
- **Problem:** how to determine the A -production from the lookahead (potentially infinitely many derivations $\beta\alpha \Rightarrow_I^* x$ / $\gamma\alpha \Rightarrow_I^* y$)?

Lemma 6.5 (Characterization of LL(k))

$G \in \text{LL}(k)$ iff for all leftmost derivations of the form

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

such that $\beta \neq \gamma$, it follows that $\text{first}_k(\beta\alpha) \cap \text{first}_k(\gamma\alpha) = \emptyset$.

Lemma 6.5 (Characterization of LL(k))

$G \in \text{LL}(k)$ iff for all leftmost derivations of the form

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

such that $\beta \neq \gamma$, it follows that $\text{first}_k(\beta\alpha) \cap \text{first}_k(\gamma\alpha) = \emptyset$.

Proof.

omitted

□

Lemma 6.5 (Characterization of LL(k))

$G \in LL(k)$ iff for all leftmost derivations of the form

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

such that $\beta \neq \gamma$, it follows that $\text{first}_k(\beta\alpha) \cap \text{first}_k(\gamma\alpha) = \emptyset$.

Proof.

omitted

□

Remarks:

- If $G \in LL(k)$, then the A -production is determined by the lookahead sets $\text{first}_k(\beta\alpha)$ (for every $A \rightarrow \beta \in P$).
- **Problem:** still infinitely many right contexts α to be considered (if β [or γ] “too short”, i.e., $\text{first}_k(\beta\alpha) \neq \text{first}_k(\beta)$).
- **Idea:** α derives to “everything that follows A ”

- 1 Repetition: Nondeterministic Top-Down Parsing
- 2 Correctness of $NTA(G)$
- 3 Adding Lookahead
- 4 $LL(k)$ Grammars
- 5 Follow Sets
- 6 $LL(1)$ Grammars

Goal: determine all possible lookaheads from production alone
(by combining all possible right contexts)

Goal: determine all possible lookaheads from production alone
(by combining all possible right contexts)

Definition 6.6 (follow $_k$ set)

Let $G = \langle N, \Sigma, P, S \rangle \in CFG_{\Sigma}$, $A \in N$, and $k \in \mathbb{N}$. Then the follow_k set of A , $\text{follow}_k(A) \subseteq \Sigma^*$, is given by

$\text{follow}_k(A) := \{v \in \text{first}_k(\alpha) \mid \text{ex. } w \in \Sigma^*, \alpha \in X^* \text{ such that } S \Rightarrow_I^* wA\alpha\}.$

- 1 Repetition: Nondeterministic Top-Down Parsing
- 2 Correctness of $NTA(G)$
- 3 Adding Lookahead
- 4 $LL(k)$ Grammars
- 5 Follow Sets
- 6 $LL(1)$ Grammars

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/>)

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 6.7

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 6.8 (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 $\text{fi}(\Gamma) := \bigcup_{\gamma \in \Gamma} \text{fi}(\gamma)$).

Definition 6.8 (Lookahead set)

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

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

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

Corollary 6.9

1 For all $a \in \Sigma$,

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

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