# Automata Theory :: LL Parsing

Jörg Endrullis

Vrije Universiteit Amsterdam

# Top-down Parsing

**Top-down parsing** tries to derive the input word from the starting variable $S$.

**Simple leftmost strategy:**

- Always expand the leftmost variable $A$.
  (Replace $A$ by $u$ if there is a rule $A \rightarrow u$.)

- Backtrack when a mismatch with the input string is found.
  (Then try another rule.)

Disadvantage: backtracking is expensive and difficult.

# LL Parsing

## LL parsing

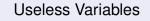Parsing **top-down** with a **leftmost** strategy.

Backtracking is not allowed.

LL parsing does not work for every context-free grammar.

Starting point is a context-free grammar $G = (V, T, S, P)$:
- without useless variables
- $\lambda$-productions and unit productions are allowed (elimination often increases the size of the grammar)

Steps of LL parsing:
- Construct sets First($A$) and Follow($A$) for every variable $A$.
- Construct a parsing table.
- Parse the input word using the parsing table.

# Useless Variables

# Removal of Useless Variables

A variable $A$ is **useless** for a context-free grammar if there exists <span style="color:red">no</span> derivation of the form

$$S \Rightarrow^* uAv \Rightarrow^+ w \qquad \text{with } w \in T^*.$$

Removing production rules that contain a useless variable from a grammar does not change the generated language.

$$S \rightarrow aSb \mid BC \mid \lambda \qquad A \rightarrow Sb \qquad B \rightarrow a \qquad C \rightarrow C$$

Which variables are useless?

- $A$ because there is no derivation $S \Rightarrow^* uAv$
- $C$ because there is no derivation $C \Rightarrow^* w$ with $w \in T^*$
- $B$ because $B$ can be reached only together with $C$

The resulting grammar is <span style="color:red">$S \rightarrow aSb \mid \lambda$</span>.

# Removal of Useless Variables

## Question

How to determine useless variables of a context-free grammar?

## Construction

A variable $A$ is called **productive** if $A \Rightarrow^+ w$ with $w \in T^*$.

We determine all productive variables:

- If $A \rightarrow y$ is a rule and all variables in $y$ are productive, then $A$ is productive.

Remove all rules that contain a non-productive variable.

We determine all **reachable** variables as follows:

- $S$ is reachable.

- If $A \rightarrow y$ and $A$ is reachable, then so are all variables in $y$.

Remove all rules that contain a non-reachable variable.

A variable is useless if it is not in one of the remaining rules.

## Removal of Useless Variables

$$S \rightarrow aSb \mid BC \mid \lambda \qquad A \rightarrow Sb \qquad B \rightarrow a \qquad C \rightarrow C$$

Which variables are non-productive?

- $C$ is not productive

We remove all rules containing non-productive variables:

$$S \rightarrow aSb \mid \lambda \qquad A \rightarrow Sb \qquad B \rightarrow a$$

Which variables are reachable from $S$?

- only $S$ is reachable

We remove all rules containing non-reachable variables:

$$S \rightarrow aSb \mid \lambda$$

Hence only $S$ is useful, the variables $A$, $B$, $C$ are not useful.

First($A$)

# First($A$)

We consider the first terminal letters derivable from a word:

$$\text{First}(w) = \{\, a \in T \mid w \Rightarrow^* a \dots \,\} \cup \{\, \lambda \mid w \Rightarrow^* \lambda \,\}$$

### Algorithm

Let PreFirst($w$) be the smallest set such that:

- $w \in$ PreFirst($w$)

- $a \in$ PreFirst($w$) if $av \in$ PreFirst($w$)

- $B \in$ PreFirst($w$) if $Bv \in$ PreFirst($w$)

- $v \in$ PreFirst($w$) if $Bv \in$ PreFirst($w$) and $B$ erasable

- $v \in$ PreFirst($w$) for every $A \in$ PreFirst($w$) and rule $A \to v$

Then First($w$) consists of

- all terminal letters $a \in T$ such that $a \in$ PreFirst($w$), and

- $\lambda$ if $w = A_1 A_2 \dots A_n$ for erasable variables $A_1, \dots, A_n$.

## Exercise

$$S \rightarrow AAc \qquad A \rightarrow Ba \mid \lambda \qquad B \rightarrow Ab \mid d$$

The erasable variables ($V \Rightarrow^+ \lambda$) are: $A$ .

We determine PreFirst($A$), PreFirst($B$) and PreFirst($S$):

$\text{PreFirst}(A) = \{\, A, \underbrace{Ba}_{\text{from } A}, \underbrace{\lambda}_{\text{from } A}, \underbrace{B}_{\text{from } Ba}, \underbrace{Ab}_{\text{from } B}, \underbrace{d}_{\text{from } B}, \underbrace{b}_{\text{from } Ab} \,\}$

$\text{PreFirst}(B) = \{\, B, \underbrace{Ab}_{\text{from } B}, \underbrace{d}_{\text{from } B}, \underbrace{b}_{\text{from } Ab}, \underbrace{A}_{\text{from } Ab} \,\} \cup \text{PreFirst}(A)$

$\qquad\qquad = \{\, A, Ba, \lambda, B, Ab, d, b \,\}$

$\text{PreFirst}(S) = \{\, S, \underbrace{AAc}_{\text{from } S}, \underbrace{Ac}_{\text{from } AAc}, \underbrace{c}_{\text{from } Ac}, \underbrace{A}_{\text{from } AAc} \,\} \cup \text{PreFirst}(A)$

$\qquad\qquad = \{\, S, AAc, Ac, c, A, Ba, \lambda, B, Ab, d, b \,\}$

Thus we get

$\text{First}(A) = \{\, b, d, \lambda \,\} \quad \text{First}(B) = \{\, b, d \,\} \quad \text{First}(S) = \{\, b, c, d \,\}$

Follow($A$)

# Follow($A$)

The sets First($A$) are not yet sufficient for 'predictive' parsing, if there are derivations $A \Rightarrow^+ \lambda$.

We consider the terminal letters that can follow a variable:

$$\text{Follow}(A) = \{ a \in T \mid S \Rightarrow^* \ldots Aa \ldots \}$$

Intuition: $a \in \text{Follow}(A)$ if $A$ can be followed by $a$ in a derivation.

We use $\$$ as a special '**end of word**' symbol.

### Algorithm

- Follow($S$) $\supseteq \{ \$ \}$
- Follow($A$) $\supseteq$ First($w$) $\setminus \{ \lambda \}$ for every rule $B \rightarrow vAw$
- Follow($A$) $\supseteq$ Follow($B$) for rules $B \rightarrow vAw$ with $\lambda \in$ First($w$)

## Example

- Follow($S$) $\supseteq$ { \$ }
- Follow($A$) $\supseteq$ First($w$) $\setminus$ { $\lambda$ } for every rule $B \rightarrow vAw$
- Follow($A$) $\supseteq$ Follow($B$) for rules $B \rightarrow vAw$ with $\lambda \in$ First($w$)

If $C \rightarrow AB$, then:

- First($B$) $\subseteq$ Follow($A$)
  Example: $C \Rightarrow AB \Rightarrow^* Aaw$ if $B \rightarrow aw$

- Follow($C$) $\subseteq$ Follow($B$)
  Example: $S \Rightarrow Ca \Rightarrow ABa$ if $S \rightarrow Ca$

- Follow($C$) $\subseteq$ Follow($A$) if $B \Rightarrow^* \lambda$
  Example: $S \Rightarrow Ca \Rightarrow ABa \Rightarrow Aa$ if $S \rightarrow Ca$ and $B \rightarrow \lambda$

## Exercise

- Follow($S$) $\supseteq \{ \$ \}$
- Follow($A$) $\supseteq$ First($w$) $\setminus \{ \lambda \}$ for every rule $B \rightarrow vAw$
- Follow($A$) $\supseteq$ Follow($B$) for rules $B \rightarrow vAw$ with $\lambda \in$ First($w$)

$$S \rightarrow Dc \qquad A \rightarrow Ba \mid \lambda$$
$$D \rightarrow AA \qquad B \rightarrow Ab \mid d$$

We have

$$\text{First}(S) = \{ b, c, d \} \qquad \text{First}(A) = \{ \lambda, b, d \}$$
$$\text{First}(D) = \{ \lambda, b, d \} \qquad \text{First}(B) = \{ b, d \}$$

Determine Follow($S$), Follow($D$), Follow($A$), Follow($B$):

Follow($S$) $\supseteq \{ \$ \}$

Follow($D$) $\supseteq \{ c \}$

Follow($A$) $\supseteq$ (First($A$) $\setminus \{ \lambda \}$) $\cup \{ b \} \cup$ Follow($D$) $\supseteq \{ b, c, d \}$

Follow($B$) $\supseteq \{ a \}$

Parser Tables

# Parser Tables

The **parser table** for a context-free grammar is a table with

- columns indexed by terminals $T \cup \{\$\}$,
- rows indexed by variables $V$,

At place $[a \in T \cup \{\$\}, B \in V]$ it contains rules $B \to u$ for which

- $a \in \text{First}(u)$, or (never the case for $a = \$$)
- $\lambda \in \text{First}(u)$ and $a \in \text{Follow}(B)$.

$$S \to aSb \mid \lambda$$

We have

- $\text{First}(aSb) = \{a\}$, $\text{First}(\lambda) = \{\lambda\}$, $\text{First}(S) = \{\lambda, a\}$
- $\text{Follow}(S) = \{b, \$\}$,

Thus the parser table is:

|  | $a$ | $b$ | $\$$ |
|---|---|---|---|
| $S$ | $S \to aSb$ | $S \to \lambda$ | $S \to \lambda$ |

# LL(1) Grammars and Parsing

A grammar is LL(1) if its parser table contains in ever cell $[a \in T \cup \{\$\}, B \in V]$ at most one production rule.

An LL(1) parser reads from Left to right, performs a Leftmost derivation, and looks always at 1 symbol of the input.

Given an LL(1)-grammar and parsing table.

To parse $a_1 \cdots a_n$, we start with $\langle S\$, a_1 \cdots a_n\$ \rangle$.

From a state $\langle v, w \rangle$ we can do the following steps:

- $\langle av', aw' \rangle$ becomes $\langle v', w' \rangle$
- $\langle Bv', aw' \rangle$ becomes $\langle uv', aw' \rangle$ if $B \to u$ at position $[a, B]$
- $\langle Bv', \$ \rangle$ becomes $\langle v', \$ \rangle$ if $B \to u$ at position $[\$, B]$
- $\langle \$, \$ \rangle$ results in **accept**
- In all other cases, $\langle v, w \rangle$ results in **reject**!

# Example

$$S \rightarrow aSb \mid \lambda$$

The parser table is:

|   | a | b | $ |
|---|---|---|---|
| S | $S \rightarrow aSb$ | $S \rightarrow \lambda$ | $S \rightarrow \lambda$ |

$\langle S\$, ab\$ \rangle \rightarrow \langle aSb\$, ab\$ \rangle \rightarrow \langle Sb\$, b\$ \rangle \rightarrow \langle b\$, b\$ \rangle$
$\rightarrow \langle \$, \$ \rangle$ **accept**

$\langle S\$, abb\$ \rangle \rightarrow \langle aSb\$, abb\$ \rangle \rightarrow \langle Sb\$, bb\$ \rangle \rightarrow \langle b\$, bb\$ \rangle$
$\rightarrow \langle \$, b\$ \rangle$ **reject**

$\langle S\$, aab\$ \rangle \rightarrow \langle aSb\$, aab\$ \rangle \rightarrow \langle Sb\$, ab\$ \rangle \rightarrow \langle aSbb\$, ab\$ \rangle$
$\rightarrow \langle Sbb\$, b\$ \rangle \rightarrow \langle bb\$, b\$ \rangle \rightarrow \langle b\$, \$ \rangle$ **reject**

JavaCC (Java Compiler Compiler) automatically generates a parser from an LL(1) grammar.

# LL($k$) Grammars

# LL($k$) Grammars

The class of LL($1$) grammars is often to restrictive in practice.

LL($1$) parsers looks at 1 symbol to decide which rule to use.

An LL($k$) parser looks $k$ symbols ahead to choose the rule.

The parser table is constructed with $k$ symbols look-ahead.

A grammar is LL($k$) if this table has in every cell $\leq 1$ rule.

LL($k$) is strictly contained in LL($k + 1$).

**Disadvantage**: size of the parser table grows exponential in $k$.

# Exercises

Can ambiguous grammars be LL($k$) for some $k \geq 1$?

Is the following grammar LL($k$) for some $k \geq 1$?
$$S \rightarrow aSa \mid \lambda$$

# Left Factorisation

# Left Factorisation

**Left factorisation**: rewrite rules $A \to uv \mid uw$ $(u \neq \lambda)$ into

$$A \to uB \qquad \text{and} \qquad B \to v \mid w$$

where $B$ is a fresh variable.

The grammar $S \to ab \mid ac$ is **not** LL(1):

|     | $a$              | $b$ | $c$ | $\$$ |
| --- | ---------------- | --- | --- | --- |
| $S$ | $S \to ab$       |     |     |     |
|     | $S \to ac$       |     |     |     |

After left factorisation we get: $\quad S \to aA \quad A \to b \mid c$

The grammar $S \to aA$, $A \to b \mid c$ **is** LL(1):

|     | $a$        | $b$        | $c$        | $\$$ |
| --- | ---------- | ---------- | ---------- | --- |
| $S$ | $S \to aA$ |            |            |     |
| $A$ |            | $A \to b$  | $A \to c$  |     |