



SWEN430 - Compiler Engineering

Lecture 4 - Parsing II

David J. Pearce

*School of Engineering and Computer Science
Victoria University of Wellington*



Language Theory

Language Theory

- Context-Free Grammar (CFG) is a set of rules of the form

$$N \longrightarrow \alpha_1 \mid \dots \mid \alpha_2$$

Here, N is a *non-terminal*, and $\alpha_1, \dots, \alpha_2$ are strings of *terminals* and non-terminals.

$$E \longrightarrow N \mid V \mid (E + E)$$

$$N \longrightarrow [0 - 9]^+$$

$$V \longrightarrow [a - z A - Z _]^+ [a - z A - Z 0 - 9 _]^*$$

- Valid words in this language include:

1, (2 + 3), ((_x + 3) + 103282)

- Invalid words in this language include:

0x, 2 + x, ((_x + 3))

Predictive Parsing

- **Predictive parsers** choose production based on *next input symbol(s)*. They are nice and easy to implement

$$\begin{aligned} S &\longrightarrow \boxed{a} S \mid \boxed{b} T && (1, 2) \\ T &\longrightarrow \boxed{c} \mid \boxed{b} T && (3, 4) \end{aligned}$$

- One symbol of lookahead required above, so it's an LL(1) grammar

$$\begin{aligned} S &\longrightarrow \boxed{a} \boxed{a} S \mid \boxed{a} \boxed{b} T && (1, 2) \\ T &\longrightarrow \boxed{c} \mid \boxed{b} T && (3, 4) \end{aligned}$$

- Above requires two symbols of lookahead, so it's LL(2)

LL(1) — first() sets

first(γ) set

Let γ be a sequence of terminal and non-terminal symbols. Then $first(\gamma)$ is the set of all terminal symbols which begin a string derived from γ .

$$\begin{aligned} S &\longrightarrow T \boxed{a} \mid U \boxed{b} && (1, 2) \\ T &\longrightarrow \boxed{d} T \mid \boxed{e} && (3, 4) \\ U &\longrightarrow \boxed{c} U \mid \boxed{f} && (5, 6) \end{aligned}$$

$$first(T) = \{ \boxed{d}, \boxed{e} \}$$

$$first(U) = \{ \boxed{c}, \boxed{f} \}$$

- What does this tell us about production rules 1 and 2 ?

LL(1) — Condition 1

Choice Condition

Consider an *LL(1)* Context-Free Grammar. For all productions $N \rightarrow \alpha$ and $N \rightarrow \beta$, it must hold that $first(\alpha) \cap first(\beta) = \emptyset$.

- For example:

$$(i) \quad S \rightarrow \boxed{a} S \mid \boxed{a}$$

$$(ii) \quad \begin{array}{l} S \rightarrow T \boxed{a} \mid U \boxed{b} \\ T \rightarrow \boxed{c} T \mid \boxed{d} \\ U \rightarrow \boxed{c} U \mid \boxed{e} \end{array}$$

$$(iii) \quad \begin{array}{l} S \rightarrow T \boxed{a} \mid U \boxed{c} \\ T \rightarrow \boxed{d} T \mid \boxed{e} \\ U \rightarrow \boxed{c} U \mid \epsilon \end{array}$$

- What are $first()$ sets here? (note: ϵ is *empty token*)

LL(1) — follow() sets

follow(γ) set

Let γ be a non-terminal symbol. Then $follow(\gamma)$ is the set of all terminal symbols which may follow a string derived from γ .

- For example:

$$\begin{aligned} S &\longrightarrow T \boxed{a} \mid U \boxed{c} && (1, 2) \\ T &\longrightarrow \boxed{d} T \mid \boxed{e} && (3, 4) \\ U &\longrightarrow \boxed{c} U \mid \epsilon && (5, 6) \end{aligned}$$

$$follow(U) = \{ \boxed{c} \}$$

$$follow(T) = \{ \boxed{a} \}$$

- What does this tell us about production rules 5 and 6 ?
- Do we need *follow* sets if there are no ϵ productions ?

LL(1) — Condition 2

ϵ Condition

Consider an $LL(1)$ Context-Free Grammar. For all non-terminal symbols N where $N \longrightarrow \epsilon$, it must hold that $first(N) \cap follow(N) = \emptyset$.

- For example:

$$(i) \quad \begin{array}{l} S \longrightarrow \boxed{a} \mid T \boxed{b} \\ T \longrightarrow \boxed{b} \mid \epsilon \end{array}$$

$$(ii) \quad \begin{array}{l} S \longrightarrow \boxed{(} T \mid \epsilon \\ T \longrightarrow S \boxed{)} S \end{array}$$

- What are the follow() sets for these grammars?

Left factoring

Left factoring

Suppose we have productions $N \longrightarrow x \alpha$ and $N \longrightarrow x \beta$, where x is either a terminal or non-terminal. Then, we can rewrite this into $N \longrightarrow x M$ and $M \longrightarrow \alpha \mid \beta$.

- Example:

$$List \longrightarrow () \mid (ListBody)$$
$$ListBody \longrightarrow ListElt \mid ListElt , ListBody$$
$$ListElt \longrightarrow N \mid List$$

Becomes $[N=List, x=(, \alpha=), \beta=ListBody]$:

$$List \longrightarrow (RestOfList$$
$$RestOfList \longrightarrow) \mid ListBody)$$

Eliminating Left Recursion

- Consider this grammar:

$$E \longrightarrow E \boxed{+} T \mid T$$
$$T \longrightarrow [\boxed{0} - \boxed{9}]^+$$

This grammar is **left-recursive**. Why is it not $LL(1)$?

- Can (sometimes) eliminate left-recursion by using **right-recursion** instead. Consider this alternative to the above:

$$E \longrightarrow T E'$$
$$E' \longrightarrow \boxed{+} T E' \mid \epsilon$$
$$T \longrightarrow [\boxed{0} - \boxed{9}]^+$$

Is this equivalent to the left-recursive grammar?



Parser Generators

Parser Generators

Yacc

Yet Another Compiler-Compiler generates *bottom up* parses for LALR parser grammars. It was originally developed for UNIX by Stephen C. Johnson.

JavaCC

The **Java Compiler Compiler** generates top-down recursive descent parsers for $LL(k)$ grammars.

ANTLR

Another Tool for Language Recognition generates bottom up parsers for $LL(*)$ grammars. It was originally developed by Terrance Parr (see `antlr.org`).

Example of Java Grammar in ANTLR

```
grammar Expressions;  
  
expr: expr op=('*' | '/') expr  
     | expr op=('+' | '-') expr  
     | INT  
     | IDENTIFIER  
     ;  
  
INT: [0-9]+;  
IDENTIFIER: [a-zA-Z][a-zA-Z0-9]+;  
WS: [\t\r\n]+ -> skip;
```

- ANTLR grammar for simple expressions