

Parsers, Context-free Grammars & Top-down Parsing

MICHAEL WOLLOWSKI

Purpose of Parsers

A parser takes the input from the lexical analyzer,

Determines whether the input is a valid sentence in the source language and

Produces a parse tree for further processing

Limitations of Regular Expressions

REs are not sufficient to capture the structure of a language.

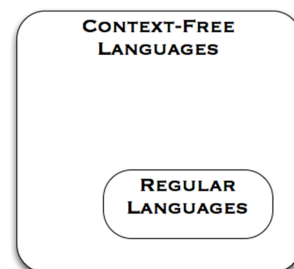
Consider blocks in Java:

```
{  
  stmt  
  stmt  
  {  
    stmt  
  }  
}
```

Why can this structure not be specified by a regular expressions?

Context-Free Grammars

Context-free grammars are more expressive than regular expressions



Context-Free Grammars

A *context-free grammar* G is a quadruple (T, NT, S, P) where:

- T is the set of terminal symbols, or words, in the language $L(G)$. Terminal symbols correspond to syntactic categories returned by the scanner.
- NT is the set of nonterminal symbols that appear in the productions of G . Non-terminals are variables to provide abstraction and structure in the productions.
- S is a nonterminal designated as the *start symbol* of the grammar. S represents the set of sentences in $L(G)$.
- P is the set of productions or rewrite rules in G . Each rule in P has the form $NT \rightarrow (T \cup NT)^*$. It replaces a single nonterminal with a string of one or more grammar symbols.

Context-Free Grammars

Here is a sample grammar:

1	$Expr$	\rightarrow	(\underline{Expr})
2			$Expr Op$ name
3			name
4	Op	\rightarrow	+
5			-
6			\times
7			\div

Derivations

A derivation is the process of showing that a sentence is a member of the language defined by a grammar.

It begins with the start symbol, *Expr*

In the example to the right, we derive $(a + b) * c$

Rule	Sentential Form
	<i>Expr</i>
2	<i>Expr Op name</i>
6	<i>Expr × name</i>
1	(<i>Expr</i>) × name
2	(<i>Expr Op name</i>) × name
4	(<i>Expr + name</i>) × name
3	(name + name) × name

Rightmost Derivation of $(a + b) \times c$

Leftmost vs Rightmost Derivation

In a rightmost derivation, we attempt to use a rule which replaces the rightmost nonterminal.

In a leftmost derivation, we attempt to rewrite the leftmost nonterminal.

Example:

Rule	Sentential Form
	<i>Expr</i>
2	<i>Expr Op name</i>
1	(<i>Expr</i>) <i>Op name</i>
2	(<i>Expr Op name</i>) <i>Op name</i>
3	(name <i>Op name</i>) <i>Op name</i>
4	(name + name) <i>Op name</i>
6	(name + name) × name

Leftmost Derivation of $(a + b) \times c$

Leftmost vs Rightmost Derivation

While the derivations are different, the parse trees are the same.

Rule	Sentential Form
	<i>Expr</i>
2	<i>Expr Op name</i>
1	<u>(Expr) Op name</u>
2	<u>(Expr Op name) Op name</u>
3	<u>(name Op name) Op name</u>
4	<u>(name + name) Op name</u>
6	<u>(name + name) × name</u>

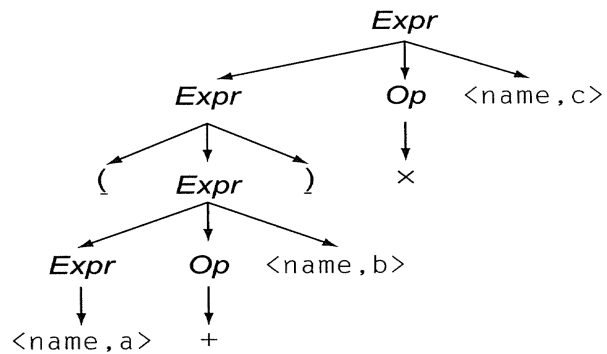
Leftmost Derivation of $(a + b) \times c$

Rule	Sentential Form
	<i>Expr</i>
2	<i>Expr Op name</i>
6	<i>Expr × name</i>
1	<u>(Expr) × name</u>
2	<u>(Expr Op name) × name</u>
4	<u>(Expr + name) × name</u>
3	<u>(name + name) × name</u>

Rightmost Derivation of $(a + b) \times c$

Derivations

The corresponding parse tree looks like this:



Ambiguity

A grammar G is ambiguous, if some sentence in $L(G)$ has more than one rightmost (or leftmost) derivation.

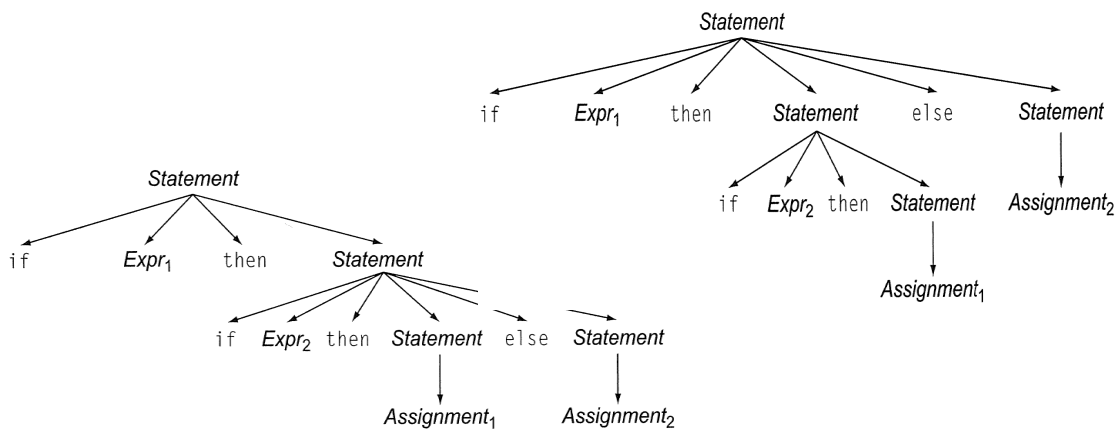
Consider the following grammar:

$$\begin{aligned} \text{Statement} \rightarrow & \text{if Expr then Statement else Statement} \mid \\ & \text{if Expr then Statement} \mid \\ & \text{Assignment} \end{aligned}$$

The following code fragment has two distinct derivations:

`if Expr1 then if Expr2 then Assignment1 else Assignment2`

Ambiguity Example



Removing Ambiguity

The general rule is to match each **else** with the closest previously unmatched **then**.

Rewrite the grammar so that a statement appearing between a **then** and **else** must be matched.

In other words, it must not end with an unmatched **then** followed by any statement

Removing Ambiguity

The original grammar:

Statement → if *Expr* then *Statement* else *Statement* |
if *Expr* then *Statement* |
Other

Ambiguity removed:

Statement → *Matched_statement* | *Unmatched_statement*
Matched_statement → if *Expr* then *Matched_Statement* else *Matched_Statement* |
Other
Unmatched_statement → if *Expr* then *Statement* |
if *Expr* then *Matched_Statement* else *Unmatched_Statement* |
Other

Precedence

Consider the following grammar:

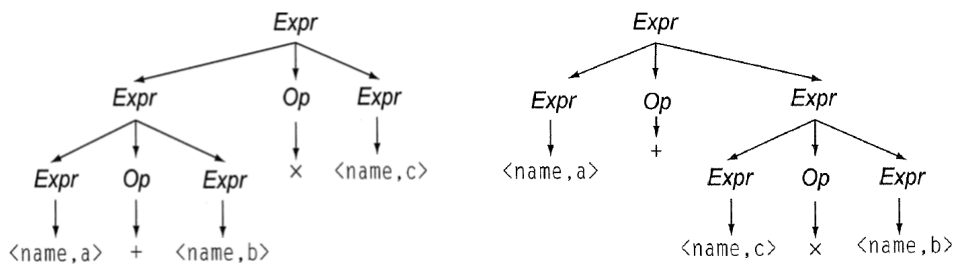
1	$Expr$	\rightarrow	$(Expr)$
2		$ $	$Expr Op Expr$
3		$ $	name
4	Op	\rightarrow	$+$
5		$ $	$-$
6		$ $	\times
7		$ $	\div

Derive the expression: $a + b * c$

Precedence

We can obtain two different parse trees.

One obeys the precedence of the arithmetic operators.



Precedence

We eliminate the ambiguity that leads to potentially incorrect parse trees by modifying the grammar.

In essence, we move operators with higher precedence further down in the grammar.

0	<i>Goal</i>	→	<i>Expr</i>
1	<i>Expr</i>	→	<i>Expr</i> + <i>Term</i>
2			<i>Expr</i> - <i>Term</i>
3			<i>Term</i>
4	<i>Term</i>	→	<i>Term</i> × <i>Factor</i>
5			<i>Term</i> ÷ <i>Factor</i>
6			<i>Factor</i>
7	<i>Factor</i>	→	(<i>Expr</i>)
8			num
9			name

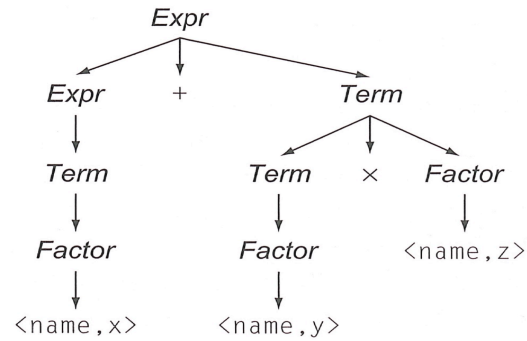
Precedence In-class Assignment

Using the revised grammar, derive $a + b * c$, again

0	<i>Goal</i>	→	<i>Expr</i>
1	<i>Expr</i>	→	<i>Expr</i> + <i>Term</i>
2			<i>Expr</i> - <i>Term</i>
3			<i>Term</i>
4	<i>Term</i>	→	<i>Term</i> × <i>Factor</i>
5			<i>Term</i> ÷ <i>Factor</i>
6			<i>Factor</i>
7	<i>Factor</i>	→	(<i>Expr</i>)
8			num
9			name

Precedence

The resulting parse tree:



19

Top-Down Parsers

A *top-down* parser, as the name suggests, begins with the start symbol and systematically applies the rules of the CFG, until there are no more non-terminal symbols left.

By contrast, a *bottom-up* parser, begins with the tokens returned by the lexical analyzer and using the rules of the CFG, replaces them with non-terminals until the start symbol is reached.

Graphically speaking, a top-down parser builds the parse tree from the root, while a bottom-up parser begins with the leaves, combining partial parse trees.

20

Oracles

Given an unmodified CFG, the challenge of a parser is to decide which rule to pick so as to avoid back-tracking.

If the parser were to be all knowing, it would always select the right rule, as shown:

Rule	Sentential Form	Input
	<i>Expr</i>	↑ name + name x name
1	<i>Expr</i> + <i>Term</i>	↑ name + name x name
3	<i>Term</i> + <i>Term</i>	↑ name + name x name
6	<i>Factor</i> + <i>Term</i>	↑ name + name x name
9	name + <i>Term</i>	↑ name + name x name
→	name + <i>Term</i>	name ↑ + name x name
→	name + <i>Term</i>	name + ↑ name x name
4	name + <i>Term</i> x <i>Factor</i>	name + ↑ name x name
6	name + <i>Factor</i> x <i>Factor</i>	name + ↑ name x name
9	name + name x <i>Factor</i>	name + ↑ name x name
→	name + name x <i>Factor</i>	name + name ↑ x name
→	name + name x <i>Factor</i>	name + name x ↑ name
9	name + name x name	name + name x ↑ name
→	name + name x name	name + name x name ↑

Classes of Context-free Grammars

There are no oracles!

However, we can modify the grammar so as to make oracles unnecessary.

Classes of Context-free Grammars

LR(1) grammars include a large subset of the unambiguous CFGs.

LR(1) grammars can be parsed, bottom-up, in a linear scan from left to right, looking at most one word ahead of the current input symbol. The “R” stands for a right-most derivation.

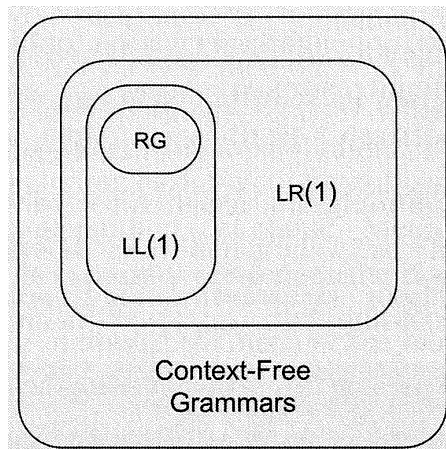
The widespread availability of tools that derive parsers from LR(1) grammars has made LR(1) parsers “everyone’s” favorite grammar.

LL(1) grammars are an important subset of the LR(1) grammars.

LL(1) grammars can be parsed, top-down, in a linear scan from left to right, looking at most one word ahead of the current input symbol. The “L” stands for a left-most derivation.

Many programming languages can be written in an LL(1) grammar.

Classes of Context-free Grammars



Top-Down Parsing: Eliminating Left-Recursion

We will focus on Top-down parsing for now.

Consider the following excerpt from our grammar:

Suppose we wish to derive: **a + b * c**

The grammar is recursive.

Computers being as systematic as they are, our compiler may decide to always use the first rule, leading to:

$$Expr \rightarrow Expr + Expr + Expr + \dots$$

$$\begin{array}{l} Expr \rightarrow Expr + Term \\ | Expr - Term \\ | Term \end{array}$$

Top-Down Parsing: Eliminating Left-Recursion

We are going to re-write the productions that are left-recursive to be right-recursive.

This requires the introduction of additional non-terminals.

The general pattern of this transformation is as follows:

$$A \rightarrow A\alpha \mid \beta$$

is transformed to:

$$A \rightarrow \beta R$$

$$R \rightarrow \alpha R \mid \epsilon$$

Top-Down Parsing: Eliminating Left-Recursion

Result of eliminating left-recursion from the productions for *Expr* and *Term*

Notice the right-recursive nature of the grammar

Original		Transformed	
<i>Expr</i>	\rightarrow <i>Expr</i> + <i>Term</i>	<i>Expr</i>	\rightarrow <i>Term</i> <i>Expr'</i>
	<i>Expr</i> - <i>Term</i>	<i>Expr'</i>	\rightarrow + <i>Term</i> <i>Expr'</i>
	<i>Term</i>		- <i>Term</i> <i>Expr'</i>
			ϵ
<i>Term</i>	\rightarrow <i>Term</i> \times <i>Factor</i>	<i>Term</i>	\rightarrow <i>Factor</i> <i>Term'</i>
	<i>Term</i> \div <i>Factor</i>	<i>Term'</i>	\rightarrow \times <i>Factor</i> <i>Term'</i>
	<i>Factor</i>		\div <i>Factor</i> <i>Term'</i>
			ϵ

Top-Down Parsing: Eliminating Left-Recursion

Visually, this works out to be like so (courtesy of the dragon book)

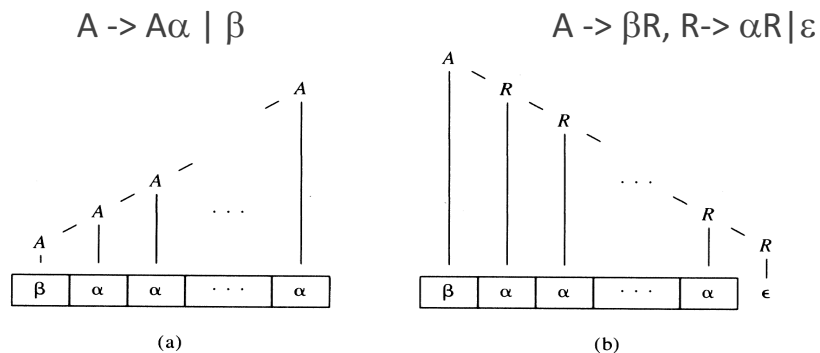


Fig. 2.18. Left- and right-recursive ways of generating a string.

Eliminating Left-Recursion: In-class Assignment

Using the revised grammar, evaluate: $a + b * c$

0	$Goal \rightarrow Expr$	6	$Term' \rightarrow \times Factor Term'$
1	$Expr \rightarrow Term Expr'$	7	$ \div Factor Term'$
2	$Expr' \rightarrow + Term Expr'$	8	$ \epsilon$
3	$ - Term Expr'$	9	$Factor \rightarrow (Expr)$
4	$ \epsilon$	10	$ num$
5	$Term \rightarrow Factor Term'$	11	$ name$