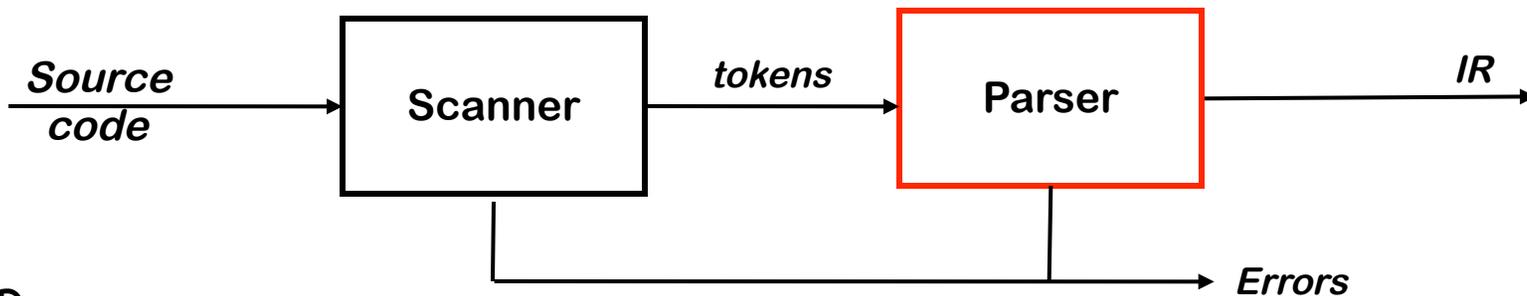




Introduction to Parsing



The Front End



Parser

- Checks the stream of words and their parts of speech (produced by the scanner) for grammatical correctness
- Determines if the input is syntactically well formed
- Guides checking at deeper levels than syntax
- Builds an IR representation of the code

Think of this as the mathematics of diagramming sentences



The Study of Parsing

The process of discovering a *derivation* for some sentence

- Need a mathematical model of syntax — a grammar G
- Need an algorithm for testing membership in $L(G)$
- Need to keep in mind that our goal is building parsers, not studying the mathematics of arbitrary languages

Roadmap

- 1 Context-free grammars and derivations
- 2 Top-down parsing
 - Hand-coded recursive descent parsers
- 3 Bottom-up parsing
 - Generated LR(1) parsers



Specifying Syntax with a Grammar

Context-free syntax is specified with a context-free grammar

$$\begin{aligned} \textit{SheepNoise} &\rightarrow \textit{SheepNoise} \underline{\textit{baa}} \\ &| \underline{\textit{baa}} \end{aligned}$$

This *CFG* defines the set of noises sheep normally make

It is written in a variant of Backus-Naur form

Formally, a grammar is a four tuple, $G = (S, N, T, P)$

- S is the *start symbol* (*set of strings in $L(G)$*)
- N is a set of *non-terminal symbols* (*syntactic variables*)
- T is a set of *terminal symbols* (*words*)
- P is a set of *productions or rewrite rules* ($P: N \rightarrow (N \cup T)^+$)



Deriving Syntax

We can use the *SheepNoise* grammar to create sentences

→ use the productions as *rewriting rules*

Rule	Sentential Form
—	<i>SheepNoise</i>
2	<u>baa</u>

Rule	Sentential Form
—	<i>SheepNoise</i>
1	<i>SheepNoise</i> <u>baa</u>
2	<u>baa</u> <u>baa</u>

Rule	Sentential Form
—	<i>SheepNoise</i>
1	<i>SheepNoise</i> <u>baa</u>
1	<i>SheepNoise</i> <u>baa</u> <u>baa</u>
2	<u>baa</u> <u>baa</u> <u>baa</u>

And so on ...



A More Useful Grammar

To explore the uses of CFGs, we need a more complex grammar

1	<i>Expr</i>	→	<i>Expr Op Expr</i>
2			<u>number</u>
3			<u>id</u>
4	<i>Op</i>	→	+
5			-
6			*
7			/

<i>Rule</i>	<i>Sentential Form</i>
—	<i>Expr</i>
1	<i>Expr Op Expr</i>
2	<id, <u>x</u> > <i>Op Expr</i>
5	<id, <u>x</u> > - <i>Expr</i>
1	<id, <u>x</u> > - <i>Expr Op Expr</i>
2	<id, <u>x</u> > - <num, <u>z</u> > <i>Op Expr</i>
6	<id, <u>x</u> > - <num, <u>z</u> > * <i>Expr</i>
3	<id, <u>x</u> > - <num, <u>z</u> > * <id, <u>y</u> >

- Such a sequence of rewrites is called a *derivation*
- Process of discovering a derivation is called *parsing*

We denote this derivation: $Expr \Rightarrow^* \underline{id} - \underline{num} * \underline{id}$



Derivations

- At each step, we choose a non-terminal to replace
- Different choices can lead to different derivations

Two derivations are of interest

- *Leftmost derivation* — replace leftmost NT at each step
- *Rightmost derivation* — replace rightmost NT at each step

These are the two *systematic* derivations

(We don't care about randomly-ordered derivations!)

The example on the preceding slide was a *leftmost* derivation

- Of course, there is also a *rightmost* derivation
- Interestingly, it turns out to be different



The Two Derivations for $\underline{x} - \underline{2} * \underline{y}$

In both cases, $Expr \Rightarrow^* \underline{id} - \underline{num} * \underline{id}$

- The two derivations produce different parse trees
- The parse trees imply different evaluation orders!

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
3	$\langle id, \underline{x} \rangle Op Expr$
5	$\langle id, \underline{x} \rangle - Expr$
1	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

Leftmost derivation

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
3	$Expr Op \langle id, \underline{y} \rangle$
6	$Expr * \langle id, \underline{y} \rangle$
1	$Expr Op Expr * \langle id, \underline{y} \rangle$
2	$Expr Op \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$
5	$Expr - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

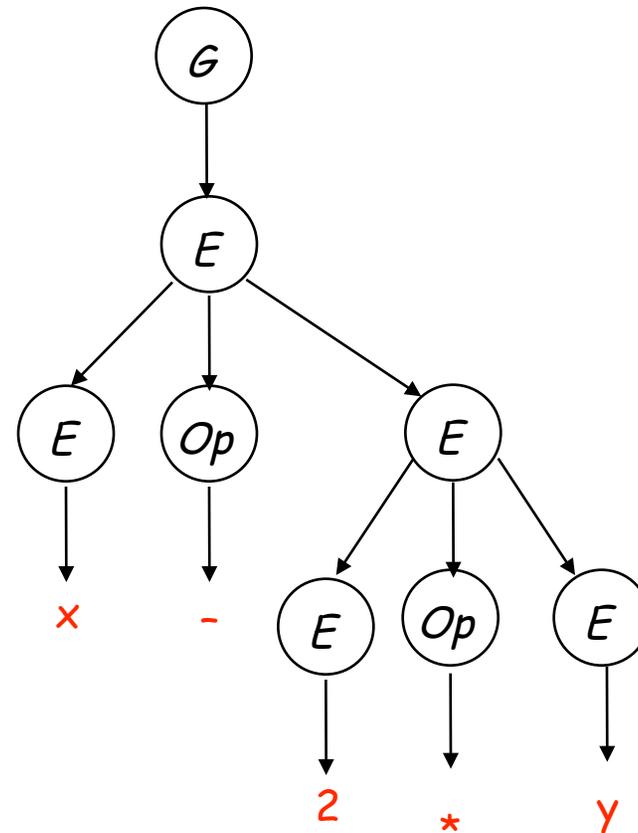
Rightmost derivation



Derivations and Parse Trees

Leftmost derivation

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
3	$\langle id, \underline{x} \rangle Op Expr$
5	$\langle id, \underline{x} \rangle - Expr$
1	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$



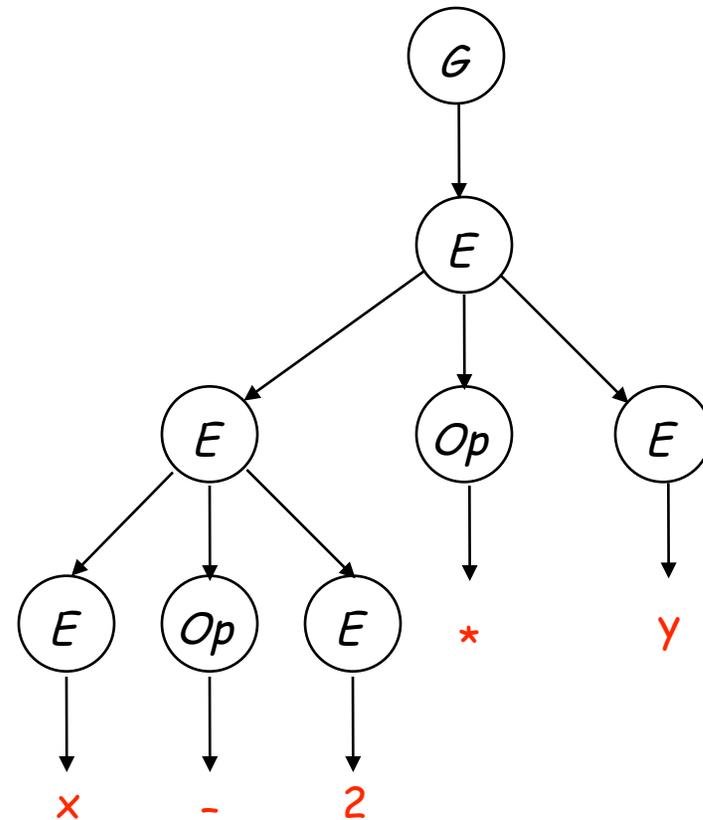
This evaluates as $\underline{x} - (\underline{2} * \underline{y})$



Derivations and Parse Trees

Rightmost derivation

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
3	$Expr Op \langle id, \underline{y} \rangle$
6	$Expr * \langle id, \underline{y} \rangle$
1	$Expr Op Expr * \langle id, \underline{y} \rangle$
2	$Expr Op \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$
5	$Expr - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$



This evaluates as $(x - 2) * y$



Derivations and Precedence

These two derivations point out a problem with the grammar:

It has no notion of precedence, or implied order of evaluation

To add precedence

- Create a non-terminal for each *level of precedence*
- Isolate the corresponding part of the grammar
- Force the parser to recognize high precedence subexpressions first

For algebraic expressions

- Multiplication and division, first *(level one)*
- Subtraction and addition, next *(level two)*



Derivations and Precedence

Adding the standard algebraic precedence produces:

level two	1	Goal	→	Expr
	2	Expr	→	Expr + Term
	3			Expr - Term
	4			Term
level one	5	Term	→	Term * Factor
	6			Term / Factor
	7			Factor
	8	Factor	→	<u>number</u>
	9			<u>id</u>

This grammar is slightly larger

- Takes more rewriting to reach some of the terminal symbols
- Encodes expected precedence
- Produces same parse tree under leftmost & rightmost derivations

Let's see how it parses $x - 2 * y$



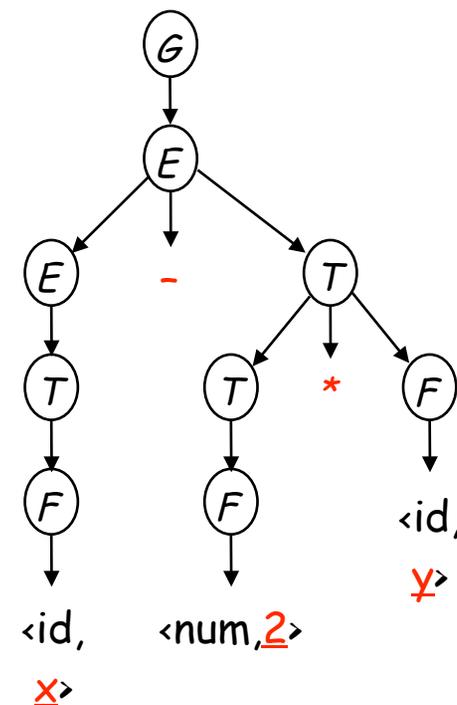
Derivations and Precedence

$$x - 2 * y$$

1	Goal	→	Expr
2	Expr	→	Expr + Term
3			Expr - Term
4			Term
5	Term	→	Term * Factor
6			Term / Factor
7			Factor
8	Factor	→	<u>number</u>
9			<u>id</u>

Rule	Sentential Form
—	Goal
1	Expr
3	Expr - Term
5	Expr - Term * Factor
9	Expr - Term * <id,y>
7	Expr - Factor * <id,y>
8	Expr - <num,2> * <id,y>
4	Term - <num,2> * <id,y>
7	Factor - <num,2> * <id,y>
9	<id,x> - <num,2> * <id,y>

The rightmost derivation



Its parse tree

This produces $x - (2 * y)$, along with an appropriate parse tree.

Both the leftmost and rightmost derivations give the same expression, because the grammar directly encodes the desired precedence.



Ambiguous Grammars

Our original expression grammar had other problems

- This grammar allows multiple leftmost derivations for $x - 2 * y$
- Hard to automate derivation if > 1 choice
- The grammar is *ambiguous*

1	$Expr \rightarrow Expr Op Expr$
2	<u>number</u>
3	<u>id</u>
4	$Op \rightarrow +$
5	-
6	*
7	/

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
①	$Expr Op Expr Op Expr$
3	$\langle id, \underline{x} \rangle Op Expr Op Expr$
5	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

different choice
than the first
time



Two Leftmost Derivations for $x - 2 * y$

The Difference:

- Different productions chosen on the second step
- Both derivations succeed in producing $x - 2 * y$

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
③	$\langle id, \underline{x} \rangle Op Expr$
5	$\langle id, \underline{x} \rangle - Expr$
1	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

Original choice

Rule	Sentential Form
—	$Expr$
1	$Expr Op Expr$
①	$Expr Op Expr Op Expr$
3	$\langle id, \underline{x} \rangle Op Expr Op Expr$
5	$\langle id, \underline{x} \rangle - Expr Op Expr$
2	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle Op Expr$
6	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * Expr$
3	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

New choice



Ambiguous Grammars

Definitions

- If a grammar has more than one leftmost derivation for a single *sentential form*, the grammar is *ambiguous*
- If a grammar has more than one rightmost derivation for a single sentential form, the grammar is *ambiguous*
- The leftmost and rightmost derivations for a sentential form may differ, even in an unambiguous grammar

Classic example — the *if-then-else* problem

$$\begin{aligned} \text{Stmt} \rightarrow & \text{if Expr then Stmt} \\ & | \text{if Expr then Stmt else Stmt} \\ & | \dots \text{ other stmts } \dots \end{aligned}$$

This ambiguity is entirely grammatical in nature

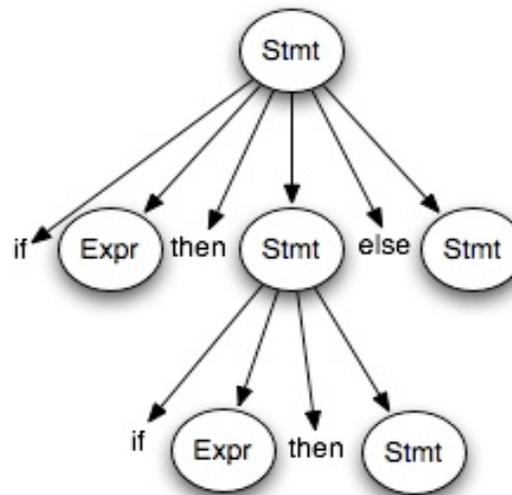


Ambiguity

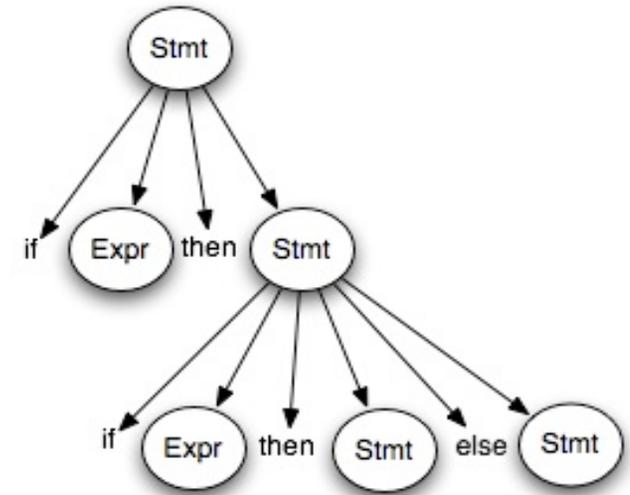
This sentential form has two derivations

if $Expr_1$ then if $Expr_2$ then $Stmt_1$ else $Stmt_2$

- $Stmt \rightarrow$ if $Expr$ then $Stmt$ (1)
- | if $Expr$ then $Stmt$ else $Stmt$ (2)
- | ... other stmts ...



*production 2, then
production 1*



*production 1, then
production 2*



Ambiguity

Removing the ambiguity

- Must rewrite the grammar to avoid generating the problem
- Match each else to innermost unmatched if (*common sense rule*)

With this grammar, the example has only one derivation

```
1 | Statement → if Expr then Statement
2 |           | if Expr then WithElse else Statement
3 |           | Assignment
4 | WithElse → if Expr then WithElse else WithElse
5 |           | Assignment
```

Intuition: binds each else to the innermost if



Ambiguity

if $Expr_1$ then if $Expr_2$ then Assignment₁ else Assignment₂

<i>Rule</i>	<i>Sentential Form</i>
	<i>Statement</i>
1	<u>if</u> $Expr$ <u>then</u> <i>Statement</i>
2	<u>if</u> $Expr$ <u>then</u> <u>if</u> $Expr$ <u>then</u> <i>WithElse</i> <u>else</u> <i>Statement</i>
3	<u>if</u> $Expr$ <u>then</u> <u>if</u> $Expr$ <u>then</u> <i>WithElse</i> <u>else</u> <i>Assignment</i>
5	<u>if</u> $Expr$ <u>then</u> <u>if</u> $Expr$ <u>then</u> <i>Assignment</i> <u>else</u> <i>Assignment</i>

This binds the else controlling Assignment₂ to the inner if



Deeper Ambiguity

Ambiguity usually refers to confusion in the CFG

Overloading can create deeper ambiguity

$$a = f(17)$$

In many Algol-like languages, f could be either a function or a subscripted variable

Disambiguating this one requires context

- Need values of declarations
- Really an issue of *type*, not context-free syntax
- Requires an extra-grammatical solution (not in CFG)
- Must handle these with a different mechanism
 - Step outside grammar rather than use a more complex grammar



Ambiguity - the Final Word

Ambiguity arises from two distinct sources

- Confusion in the context-free syntax (*if-then-else*)
- Confusion that requires context to resolve (*overloading*)

Resolving ambiguity

- To remove context-free ambiguity, rewrite the grammar
- To handle context-sensitive ambiguity takes cooperation
 - Knowledge of declarations, types, ...
 - Accept a superset of $L(G)$ & check it by other means[†]
 - This is a language design problem

Sometimes, the compiler writer accepts an ambiguous grammar

- Parsing techniques that "do the right thing"
- *i.e.*, always select the same derivation

[†]See Chapter 4