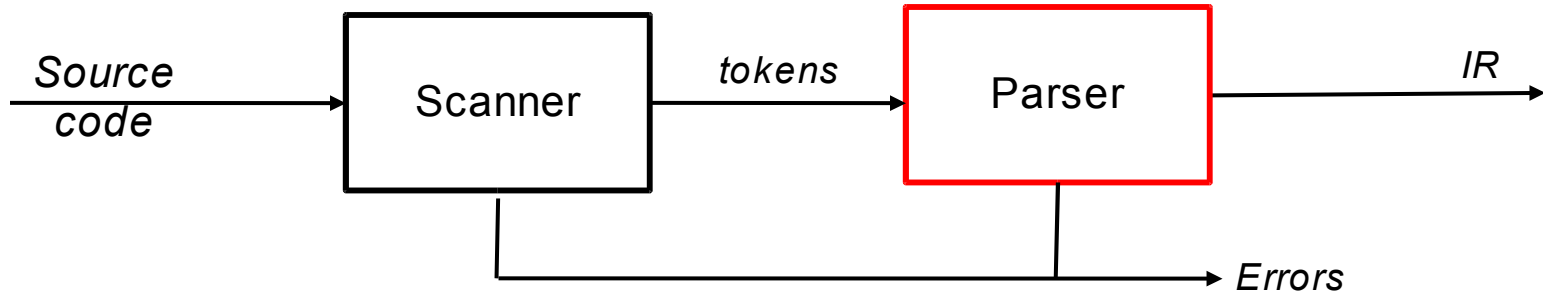

Introduction to Parsing

Copyright 2003, Keith D. Cooper, Ken Kennedy & Linda Torczon, all rights reserved.

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

$$\text{SheepNoise} \rightarrow \text{SheepNoise } \underline{\text{baa}}$$
$$\quad \quad \quad | \quad \underline{\text{baa}}$$

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 ($P : N \rightarrow (N \cup T)^+$)

Example due to Dr. Scott K. Warren

Deriving Syntax

We can use the SheepNoise grammar to create sentences

- use the productions as rewriting rules

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

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

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

And so on ...

While it is cute, this example quickly runs out of intellectual steam ...

A More Useful Grammar

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

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

Rule	Sentential Form
—	Expr
1	Expr Op Expr
2	<id, <u>x</u> > Op Expr
5	<id, <u>x</u> > - Expr
1	<id, <u>x</u> > - Expr Op Expr
2	<id, <u>x</u> > - <num, <u>z</u> > Op Expr
6	<id, <u>x</u> > - <num, <u>z</u> > * Expr
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: $\text{Expr} \Rightarrow^* \underline{\text{id}} - \underline{\text{num}} * \underline{\text{id}}$

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 $\underline{x} - \underline{2} * \underline{y}$

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

Leftmost derivation

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

Rightmost
derivation

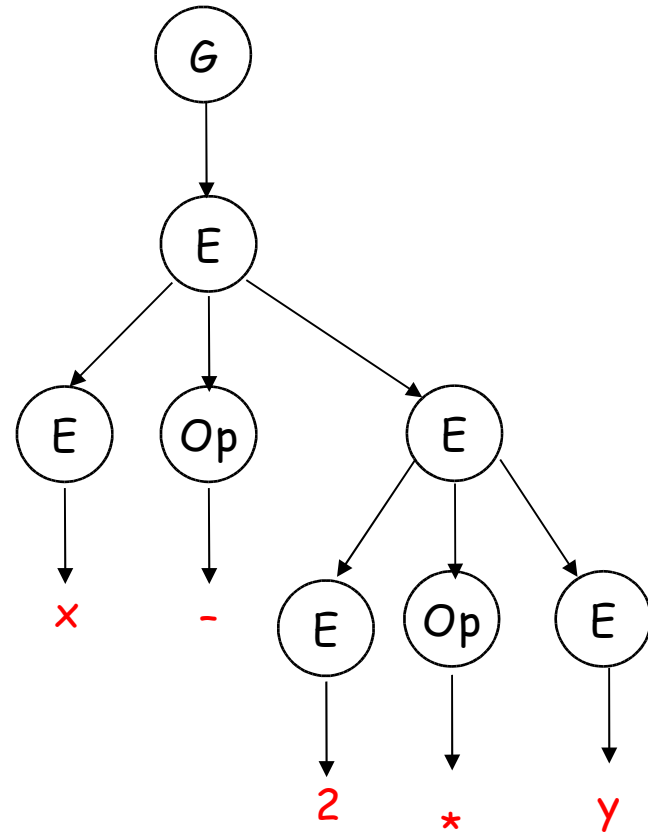
In both cases, $\text{Expr} \Rightarrow^* \underline{\text{id}} - \underline{\text{num}} * \underline{\text{id}}$

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

Derivations and Parse Trees

Leftmost derivation

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

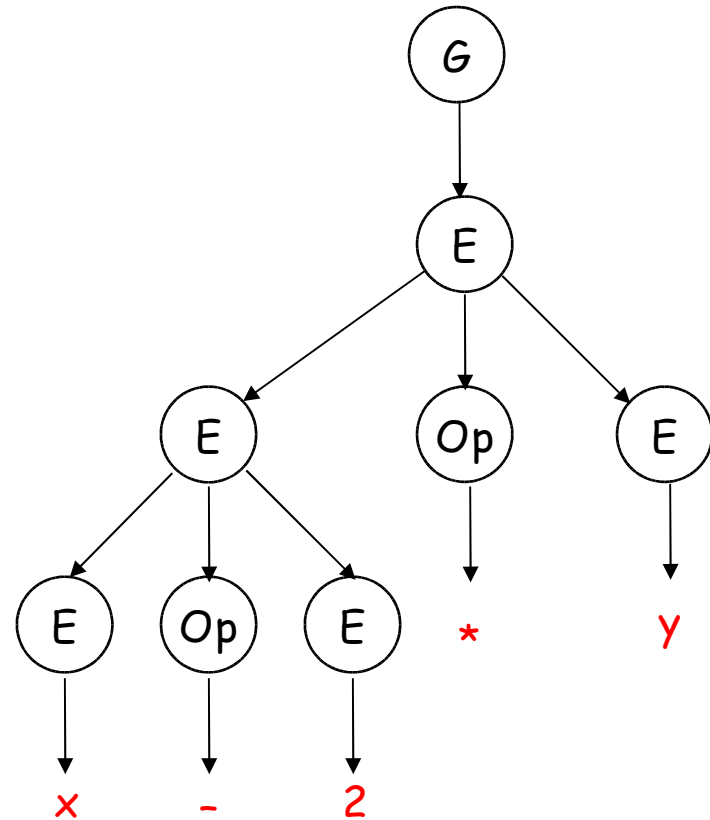


This evaluates as $\underline{x} - (\underline{2} * \underline{y})$

Derivations and Parse Trees

Rightmost derivation

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



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

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

This grammar is slightly larger

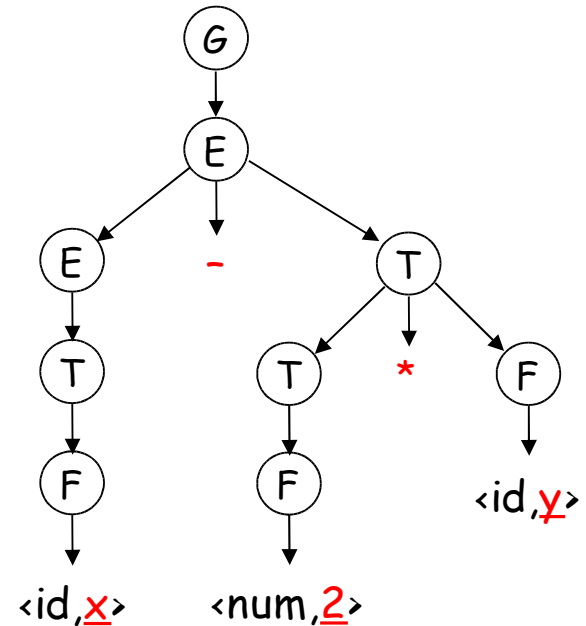
- Takes more rewriting to reach some of the terminal symbols
- Encodes expected precedence
- Produces same parse tree under leftmost & rightmost derivations

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

Derivations and Precedence

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,z> * <id,y>
4	Term - <num,z> * <id,y>
7	Factor - <num,z> * <id,y>
9	<id,x> - <num,z> * <id,y>

The rightmost derivation



Its parse tree

This produces $x - (z * y)$, 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

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

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

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

different choice than
the first time

Two Leftmost Derivations for $x - 2 * y$

The Difference:

➤ Different productions chosen on the second step

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

Original choice

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

New choice

Both derivations succeed in producing $x - 2 * y$

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

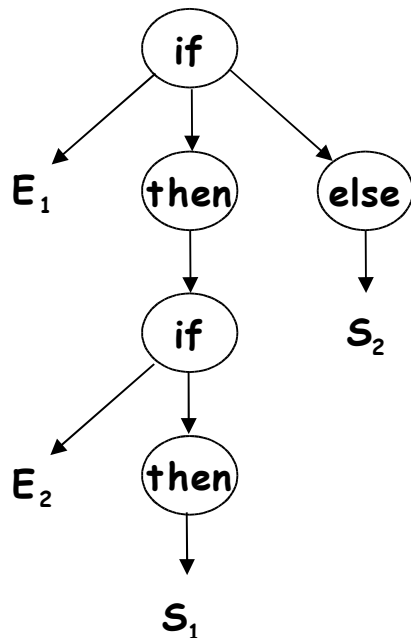
$$\begin{array}{l} \text{Stmt} \rightarrow \text{if Expr then Stmt} \\ \quad | \text{if Expr then Stmt else Stmt} \\ \quad | \dots \text{ other stmts } \dots \end{array}$$

This ambiguity is entirely grammatical in nature

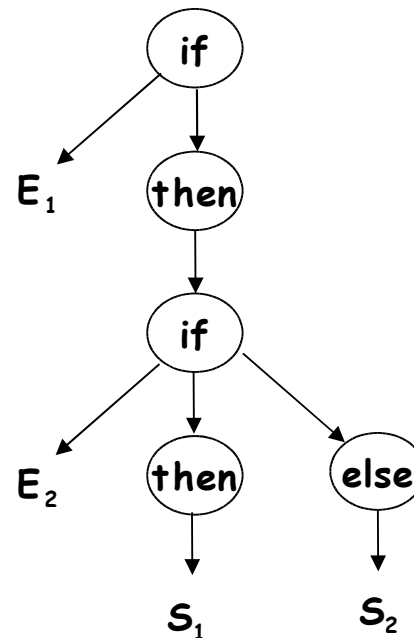
Ambiguity

This sentential form has two derivations

if Expr₁ then if Expr₂ then Stmt₁ else Stmt₂



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)

1	Stmt	→	WithElse
2			NoElse
3	WithElse	→	<u>if</u> Expr <u>then</u> WithElse <u>else</u> WithElse
4			OtherStmt
5	NoElse	→	<u>if</u> Expr <u>then</u> Stmt
6			<u>if</u> Expr <u>then</u> WithElse <u>else</u> NoElse

Intuition: a NoElse always has no else on its last cascaded else if statement

With this grammar, the example has only one derivation

Ambiguity

if Expr₁ then if Expr₂ then Stmt₁ else Stmt₂

Rule	Sentential Form
—	Stmt
2	NoElse
5	<u>if</u> Expr <u>then</u> Stmt
?	<u>if</u> E ₁ <u>then</u> Stmt
1	<u>if</u> E ₁ <u>then</u> WithElse
3	<u>if</u> E ₁ <u>then</u> <u>if</u> Expr <u>then</u> WithElse <u>else</u> WithElse
?	<u>if</u> E ₁ <u>then</u> <u>if</u> E ₂ <u>then</u> WithElse <u>else</u> WithElse
4	<u>if</u> E ₁ <u>then</u> <u>if</u> E ₂ <u>then</u> S ₁ <u>else</u> WithElse
4	<u>if</u> E ₁ <u>then</u> <u>if</u> E ₂ <u>then</u> S ₁ <u>else</u> S ₂

This binds the else controlling S₂ 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