



Lexical Analysis – Part II: Constructing a Scanner from Regular Expressions

Lecture 6

CS434

Spring 2005

Department of Computer Science

University of Alabama

Joel Jones

Copyright 2003, Keith D. Cooper, Ken Kennedy & Linda Torczon, all rights reserved.
Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.

What can be so hard?

(Fortran 66/77)



```
INTEGERFUNCTIONA
PARAMETER(A=6,B=2)
IMPLICIT CHARACTER*(A-B)(A-B)
INTEGER FORMAT(10), IF(10), DO9E1
100 FORMAT(4H)=(3)
200 FORMAT(4 )=(3)
    DO9E1=1
    DO9E1=1,2
    9 IF(X)=1
      IF(X)H=1
      IF(X)300,200
300 CONTINUE
    END
C THIS IS A "COMMENT CARD"
$ FILE(1)
    END
```

How does a compiler scan this?

- First pass finds & inserts blanks
- Can add extra words or tags to create a scannable language
- Second pass is normal scanner

Example due to Dr. F.K. Zadeck



Building Faster Scanners from the DFA

Table-driven recognizers waste effort

- Read (& classify) the next character
- Find the next state
- Assign to the state variable
- Trip through case logic in action()
- Branch back to the top

We can do better

- Encode state & actions in the code
- Do transition tests locally
- Generate ugly, spaghetti-like code
- Takes (many) fewer operations per input character

```
char ← next character;  
state ← s0 ;  
call action(state,char);  
while (char ≠ eof)  
    state ← δ(state,char);  
    call action(state,char);  
    char ← next character;
```

```
if T(state) = final then  
    report acceptance;  
else  
    report failure;
```



Building Faster Scanners from the DFA

A direct-coded recognizer for r Digit Digit^{*}
goto s_0 ;

s_0 : word $\leftarrow \emptyset$;

- char \leftarrow next character;
- if (char = 'r')
- then goto s_1 ;
- else goto s_e ;

s_1 : word \leftarrow word + char;

- char \leftarrow next character;
- if ('0' \leq char \leq '9')
- then goto s_2 ;
- else goto s_e ;

s_2 : word \leftarrow word + char;

char \leftarrow next character;

if ('0' \leq char \leq '9')

 then goto s_2 ;

 else if (char = eof)

 then report

 success;

 else goto s_e ;

s_e : print error message;

 return failure;

- Many fewer operations per character
- Almost no memory operations
- Even faster with careful use of fall-through cases



Building Faster Scanners

Hashing keywords versus encoding them directly

- Some (**well-known**) compilers recognize keywords as identifiers and check them in a hash table
- Encoding keywords in the DFA is a better idea
 - $O(1)$ cost per transition
 - Avoids hash lookup on each identifier

It is hard to beat a well-implemented DFA scanner



Building Scanners

The point

- All this technology lets us automate scanner construction
- Implementer writes down the regular expressions
- Scanner generator builds NFA, DFA, minimal DFA, and then writes out the (table-driven or direct-coded) code
- This reliably produces fast, robust scanners

For most modern language features, this works

- You should think twice before introducing a feature that defeats a DFA-based scanner
- The ones we've seen (e.g., insignificant blanks, non-reserved keywords) have not proven particularly useful or long lasting



Some Points of Disagreement with EAC

- Table-driven scanners are not fast
 - EAC doesn't say they are slow; it says you can do better
- Faster code can be generated by embedding scanner in code
 - This was shown for both LR-style parsers and for scanners in the 1980s
- Hashed lookup of keywords is slow
 - EAC doesn't say it is slow. It says that the effort can be folded into the scanner so that it has no extra cost. Compilers like *GCC* use hash lookup. A word must fail in the lookup to be classified as an identifier. With collisions in the table, this can add up. At any rate, the cost is unneeded, since the DFA can do it for $O(1)$ cost per character.



Building Faster Scanners from the DFA

Table-driven recognizers waste a lot of effort

- index** • Read (& classify) the next character
- index** • Find the next state
- Assign to the state variable
- index** • Trip through case logic in action()
- Branch back to the top

We can do better

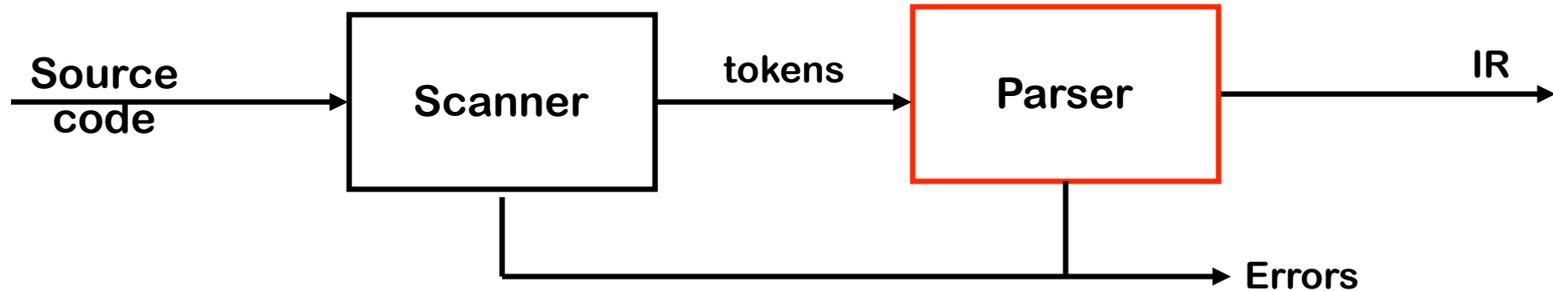
- Encode state & actions in the code
- Do transition tests locally
- Generate ugly, spaghetti-like code
- Takes (many) fewer operations per input character

```
char ← next character;
state ← s0 ;
call action(state,char);
while (char ≠ eof)
    state ← δ(state,char);
    call action(state,char);
    char ← next character;
```

```
if T(state) = final then
    report acceptance;
else
    report failure;
```



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

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

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
—	<i>SheepNoise</i>
2	<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>

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

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

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

Leftmost derivation

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

Rightmost derivation

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

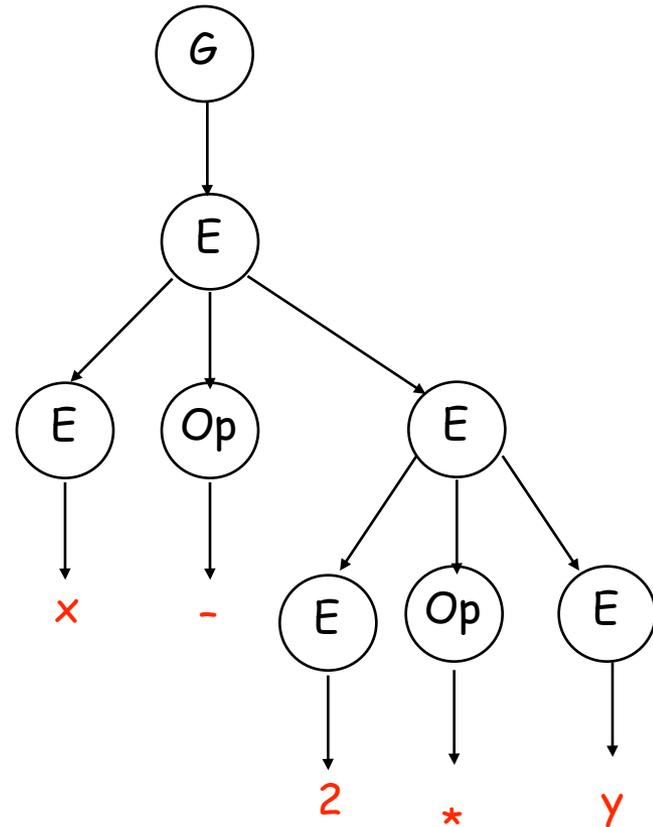
- The two derivations produce different parse trees
- The parse trees imply different evaluation orders!

Derivations and Parse Trees



Leftmost derivation

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



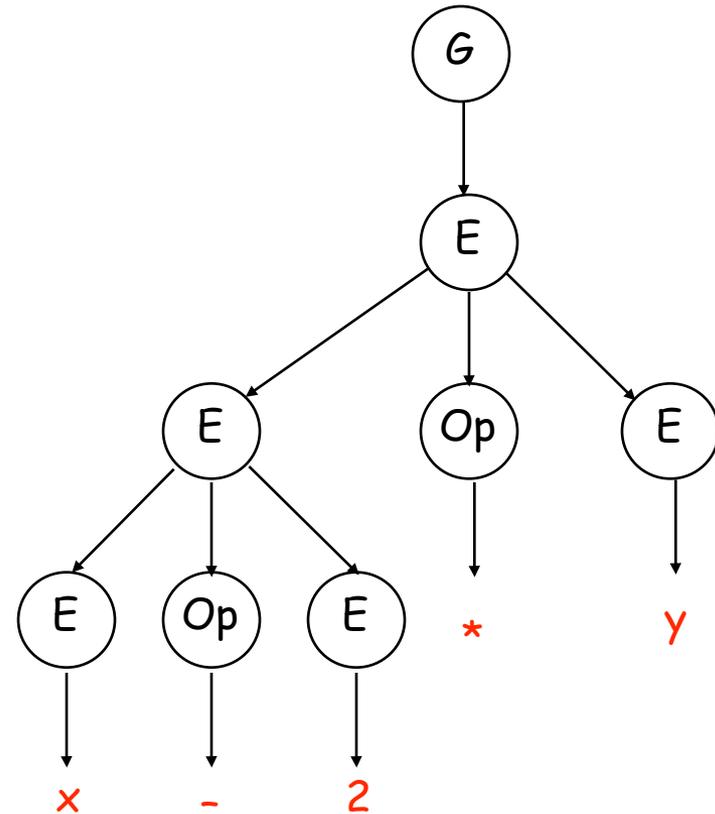
This evaluates as $x - (2 * y)$



Derivations and Parse Trees

Rightmost derivation

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



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	<i>Goal</i>	→	<i>Expr</i>
	2	<i>Expr</i>	→	<i>Expr + Term</i>
	3			<i>Expr - Term</i>
level one	4			<i>Term</i>
	5	<i>Term</i>	→	<i>Term * Factor</i>
	6			<i>Term / Factor</i>
	7			<i>Factor</i>
	8	<i>Factor</i>	→	<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

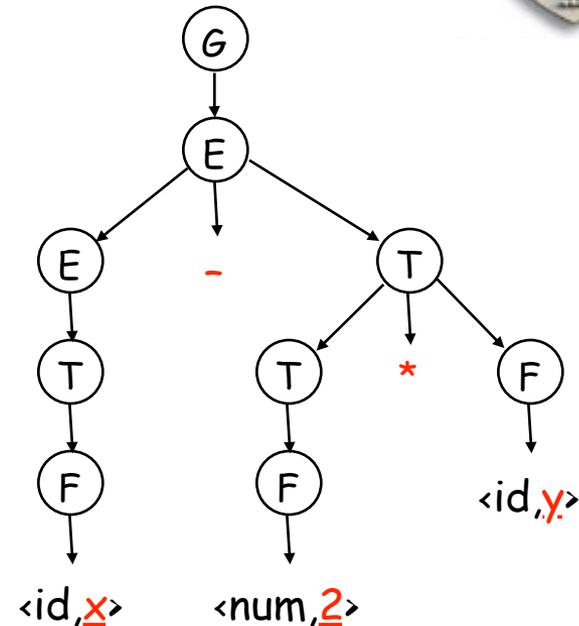
Let's see how it parses $x - 2 * y$



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

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

Original choice

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

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

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

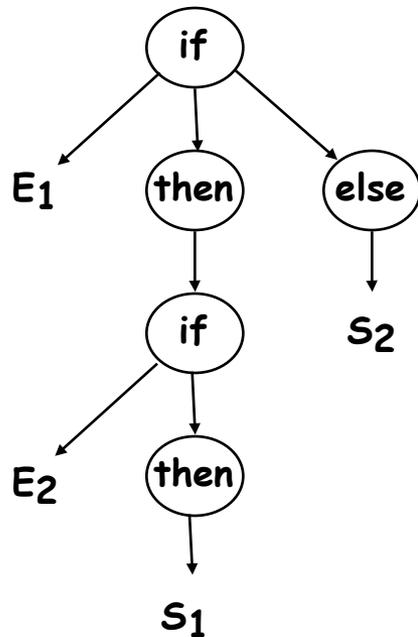
This ambiguity is entirely grammatical in nature

Ambiguity

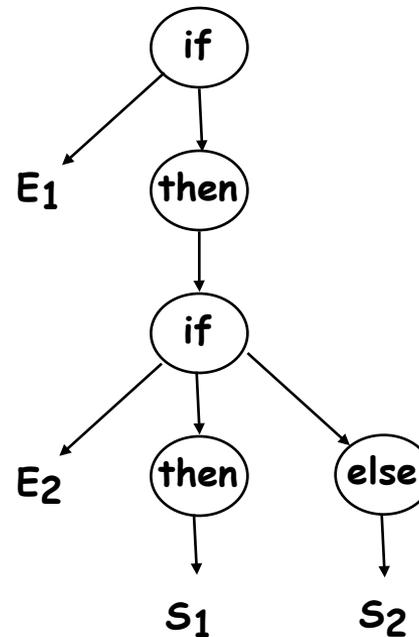


This sentential form has two derivations

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 → if Expr then WithElse else WithElse
4 |         | OtherStmt
5 | NoElse → if Expr then Stmt
6 |         | if Expr then WithElse else 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