

---

# EECS 6083

## Intro to Parsing

### Context Free Grammars

Based on slides from text web site: Copyright 2003, Keith D. Cooper, Ken Kennedy & Linda Torczon, all rights reserved.

# Parsing

---

sequence of tokens  $\xrightarrow{\text{parser}}$  syntax tree

- ◆ **Check the syntax (structure) of a program and create a tree representation of the program**
- ◆ **Programming languages have non-regular constructs.**
  - ◆ **Nesting**
  - ◆ **Recursion**
- ◆ **Context-Free Grammars are used to express the syntax for programming languages**

# Syntax vs. Semantics

---

- ◆ **Syntax** – structure of the program, expressed with grammatical rules.
  - ◆ **Typically a Context Free Grammar (CFG) in Backus-Naur form (BNF) or Extended Backus-Naur form (EBNF).**
- ◆ **Semantics** – meaning of the program, expressed with descriptive text or with *inference rules*
- ◆ **Consider the following English sentences:**

I    fed    the elephant.  
subject    verb    object

I    fed    the brick wall,  
subject    verb    object

*Both are syntactically correct, but the second is not semantically reasonable.*

# Context Free Grammars

---

- ◆ **Comprised of:**
  - ◆ **a set of tokens or *terminal* symbols**
  - ◆ **a set of *non-terminal* symbols**
  - ◆ **a set of rules or *productions* which express the legal relationships between the symbols**
  - ◆ **A *start* or *goal* symbol**

- ◆ **Example:**

- (1)  $\text{expr} \rightarrow \text{expr} - \text{digit}$
- (2)  $\text{expr} \rightarrow \text{expr} + \text{digit}$
- (3)  $\text{expr} \rightarrow \text{digit}$
- (4)  $\text{digit} \rightarrow \mathbf{0|1|2|\dots|9}$

*Terminals: -,+, 0,1,2,...,9*

*Nonterminals: expr, digit*

*Start symbol: expr*

# Some Example CFGs

---

- ◆ **Palindromes over the alphabet {a, b, c}:
  - ◆ (a palindrome is a word that has the same spelling backwards as forwards)
  - ◆ aabcbaa**

# Some Example CFGs (continued)

---

- ◆ **Palindromes over the alphabet {a, b, c}:**
  - ◆ (a palindrome is a word that has the same spelling backwards as forwards)
  - ◆ abba, c, abbcbcbbba
- ◆ **CFG for Palindromes**

**$S \rightarrow aSa$**

**$S \rightarrow bSb$**       *terminal symbols: {a, b, c}*

**$S \rightarrow cSc$**       *non-terminal symbols: {S}*

**$S \rightarrow a$**       *Goal symbol: S*

**$S \rightarrow b$**

**$S \rightarrow c$**

**$S \rightarrow \epsilon$**

# Some Example CFGs (continued)

---

- ◆ **Balanced Parenthesis and Square Brackets**
  - ◆ **E.g.** ( [ [ ] ( ( ) [ ( ) ] [ ] ) ] )

# Some Example CFGs (continued)

---

- ◆ **Balanced Parenthesis and Square Brackets**

- ◆ **E.g.**  $( [ [ ] ( ( ) [ ( ) ] [ ] ) ] )$

- ◆ **The CFG:**

- $B \rightarrow (B)$

- $| [B]$

- $| BB$

- $| \varepsilon$



# Checking for correct Syntax

---

- ◆ **Given a grammar for a language and a program how do you know if the syntax of the program is legal?**
- ◆ **A legal program can be *derived* from the start symbol of the grammar.**

# Deriving a string

---

- The derivation begins with the start symbol
- At each step of a derivation the right hand side of a grammar rule is used to replace a non-terminal symbol.
- Continue replacing non-terminals until only terminal symbols remain

(1) **expr** -> **expr - digit**

(2) **expr** ->**expr + digit**

(3) **expr** -> **digit**

(4) **digit** -> **0|1|2|...|9**

*Rule (1)*                      *Rule (4)*                      *Rule (2)*  
 $\text{expr} \Rightarrow \text{expr} - \text{digit} \Rightarrow \text{expr} - \mathbf{2} \Rightarrow \text{expr} + \text{digit} - 2$

*Rule (4)*                      *Rule (3)*                      *Rule (4)*  
 $\Rightarrow \text{expr} + \mathbf{8} - \mathbf{2} \Rightarrow \text{digit} + \mathbf{8} - \mathbf{2} \Rightarrow 3 + 8 - 2$

Example Input:

3 + 8 - 2

# Rightmost and leftmost derivations

---

- ◆ In a *rightmost derivation* the rightmost non-terminal is replaced at each step.
  - ◆  $\text{expr} \Rightarrow \text{expr} - \text{digit} \Rightarrow \text{expr} - 2 \Rightarrow \text{expr} + \text{digit} - 2 \Rightarrow \text{expr} + 8 - 2 \Rightarrow \text{digit} + 8 - 2 \Rightarrow 3 + 8 - 2$
  - ◆ corresponds to a postorder numbering in reverse of the internal nodes of the parse tree
- ◆ In a *leftmost derivation* the leftmost non-terminal is replaced at each step.
  - ◆  $\text{expr} \Rightarrow \text{expr} - \text{digit} \Rightarrow \text{expr} + \text{digit} - \text{digit} \Rightarrow \text{digit} + \text{digit} - \text{digit} \Rightarrow 3 + \text{digit} - \text{digit} \Rightarrow 3 + 8 - \text{digit} \Rightarrow 3 + 8 - 2$
  - ◆ corresponds to a preorder numbering of the nodes of a parse tree.

# Parse tree

(1)  $\text{expr} \rightarrow \text{expr} - \text{digit}$

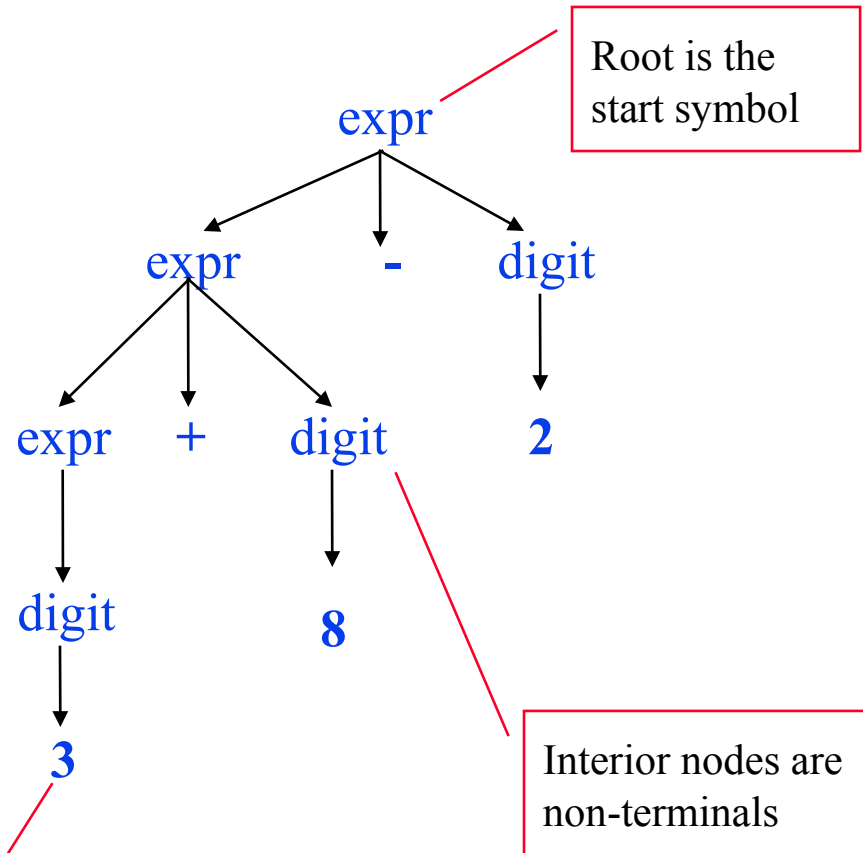
(2)  $\text{expr} \rightarrow \text{expr} + \text{digit}$

(3)  $\text{expr} \rightarrow \text{digit}$

(4)  $\text{digit} \rightarrow 0|1|2|\dots|9$

Example Input:

3 + 8 - 2



# 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>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> >

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

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

# The Two Derivations for $x - 2 * y$

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	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \langle \text{id}, \underline{y} \rangle$

*Leftmost derivation*

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

*Rightmost derivation*

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

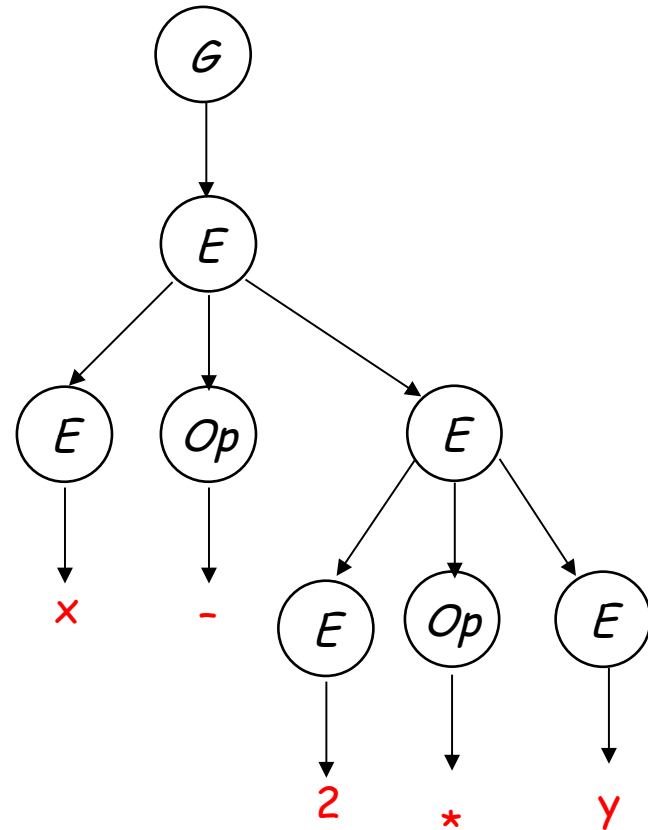
- ◆ The two derivations produce different parse trees
- ◆ The parse trees imply different evaluation orders!

# Derivations and Parse Trees

## Leftmost derivation

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

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

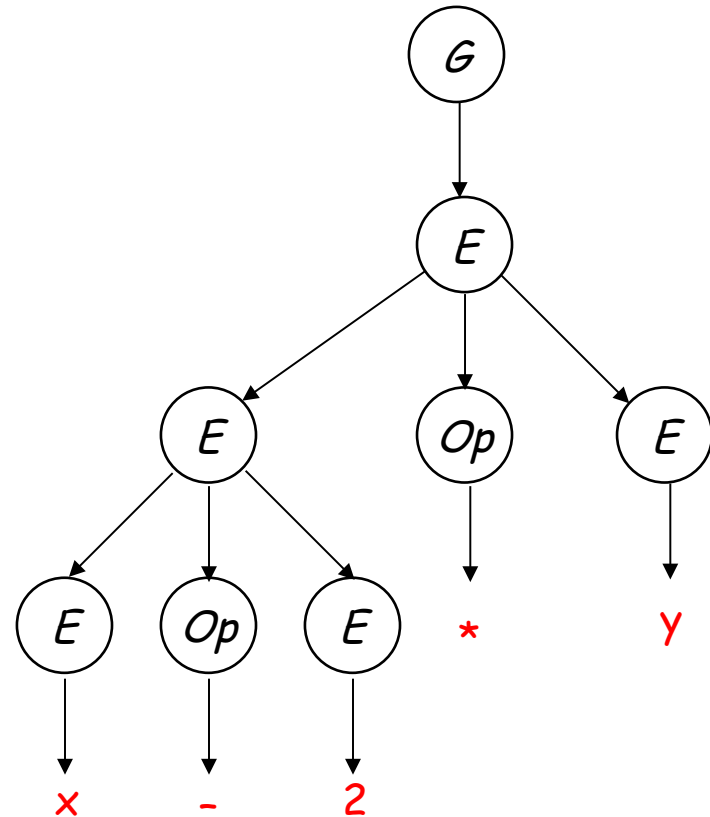


# Derivations and Parse Trees

## 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:

	1	<i>Goal</i>	→	<i>Expr</i>
level two	2	<i>Expr</i>	→	<i>Expr + Term</i>
	3			<i>Expr - Term</i>
	4			<i>Term</i>
level one	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

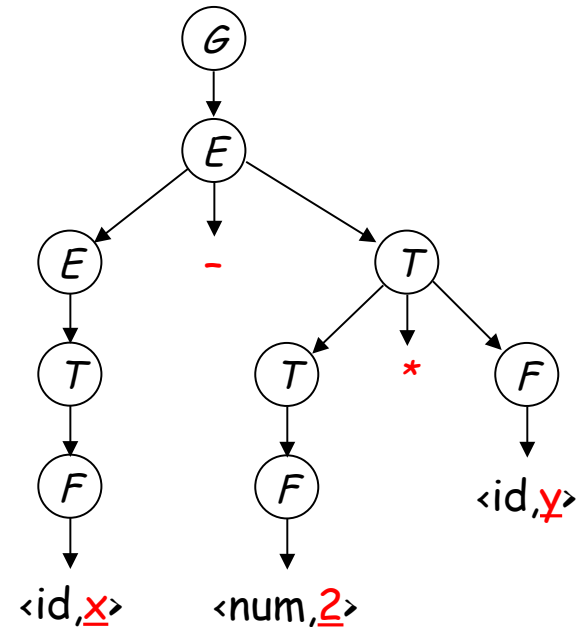
- Takes more rewriting to reach some of the terminal symbols
- Encodes expected precedence
- Produces same parse tree under leftmost & rightmost derivations

*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	$Expr \rightarrow$	$Expr Op Expr$
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	$\langle id, \underline{x} \rangle - \langle num, \underline{2} \rangle * \langle id, \underline{y} \rangle$

- ◆ This grammar allows multiple leftmost derivations for  $\underline{x} - \underline{2} * \underline{y}$
- ◆ 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

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

*This ambiguity is entirely grammatical in nature*

# Ambiguity: Dangling Else example

---

*stmt* → if *expr* then *stmt*

| if *expr* then *stmt* else *stmt*

| other

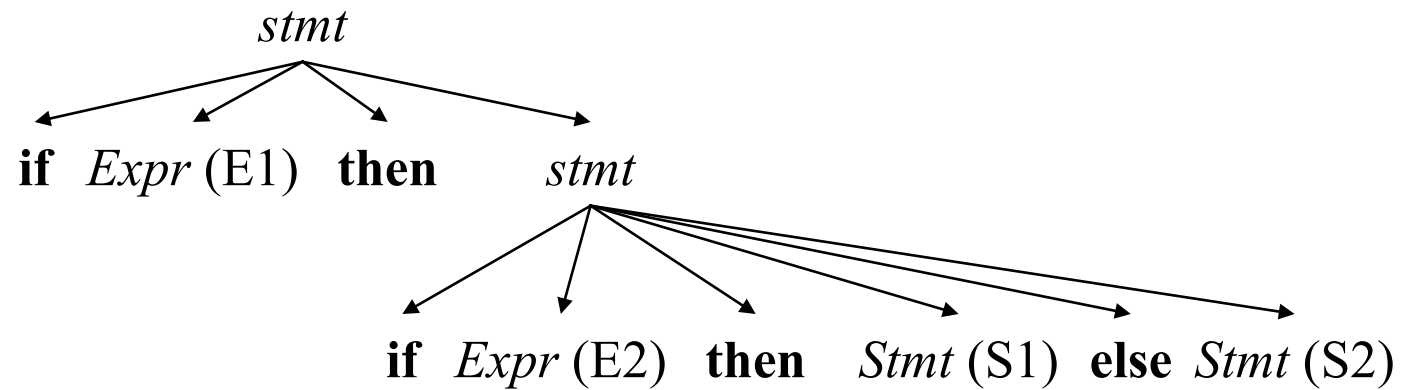
Two parse trees for the legal sentence:

if E1 then if E2 then S1 else S2

# Ambiguity: Dangling Else example (continued)

---

**if E1 then ( if E2 then S1 else S2 )**

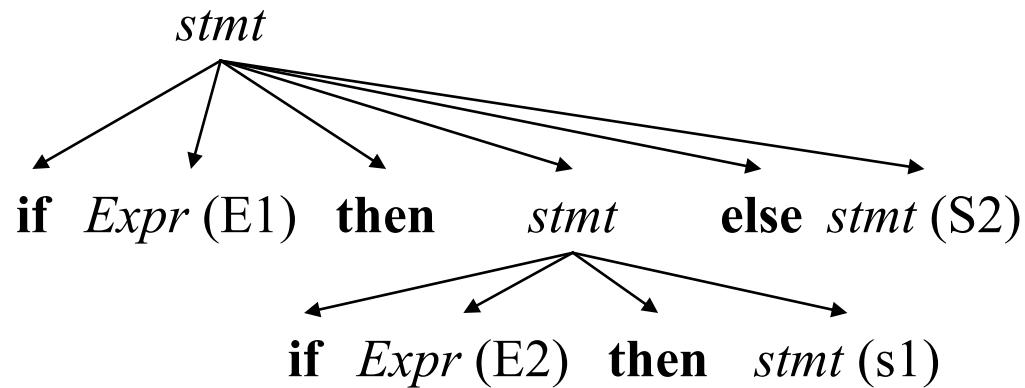




# Ambiguity: Dangling Else example (continued)

---

if E1 then ( if E2 then S1 ) else S2



# 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<sup>†</sup>
  - ◆ 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