

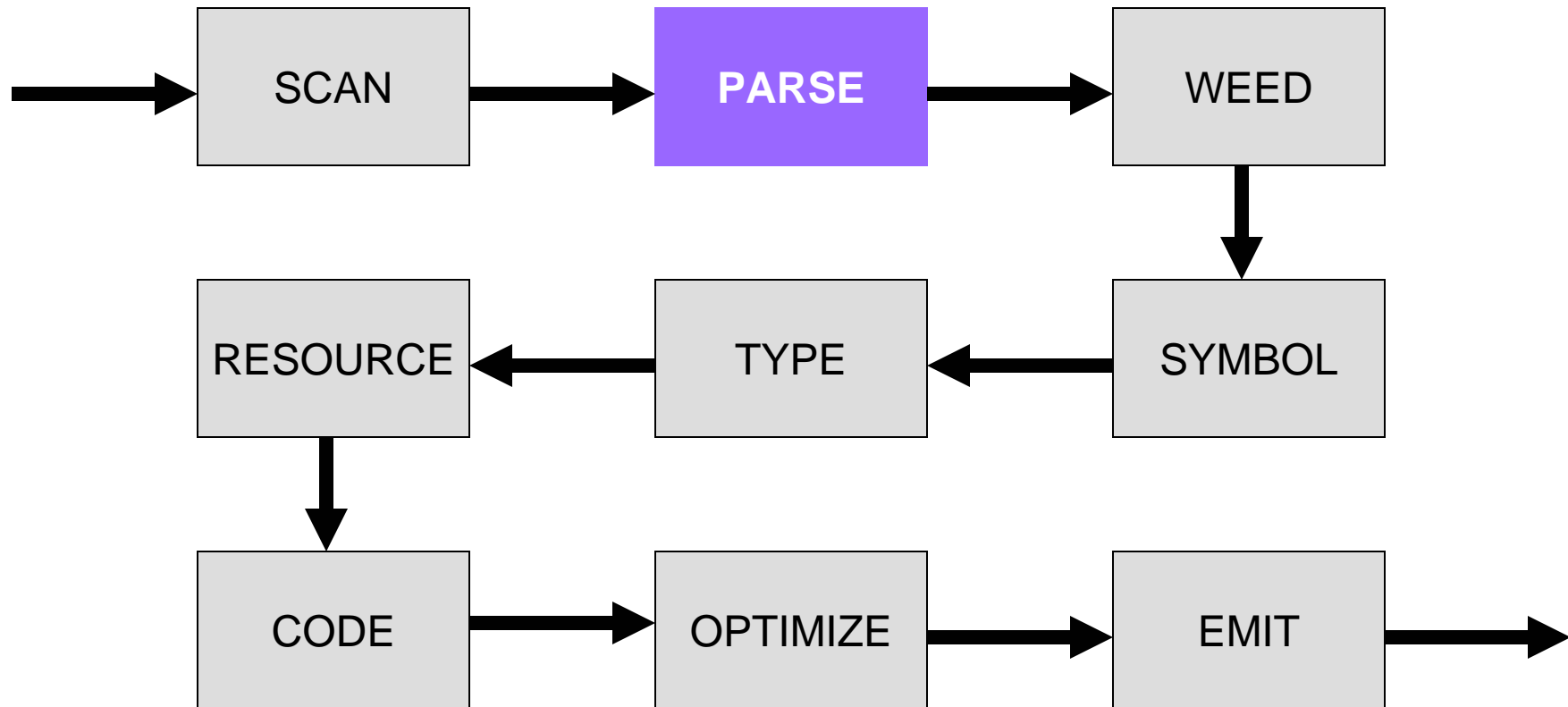


Compiler

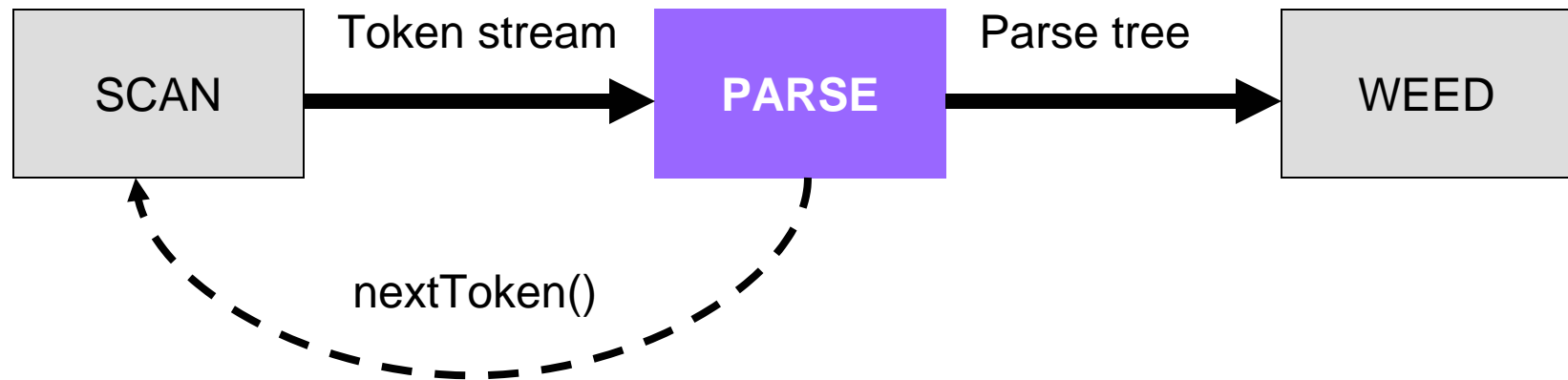
Parsing

© Copyright 2005, Matthew B. Dwyer and Robby. The syllabus and lectures for this course are copyrighted materials and may not be used in other course settings outside University of Nebraska-Lincoln and Kansas State University in their current form or modified form without express written permission of one of the copyright holders. During this course, students are prohibited from selling notes to or being paid for taking notes by any person or commercial firm without the express written permission of one of the copyright holders.

Compiler Architecture



Compiler Architecture





The Functionality of the Parser

- **Input:** sequence of tokens from lexer
- **Output:** parse tree of the program
 - parse tree is generated if the input is a legal program
 - if input is an illegal program, syntax errors are issued
- **Note:** Instead of parse tree, some parsers produce:
 - abstract syntax tree (AST) + symbol table
- In the following, we'll assume that parse tree is generated.



Comparison with Scanner

<i>Phase</i>	<i>Input</i>	<i>Output</i>
Scanner	String of characters	String of tokens
Parser	String of tokens	Parse tree

Example

- The program:

$x * y + z$

- Input to parser:

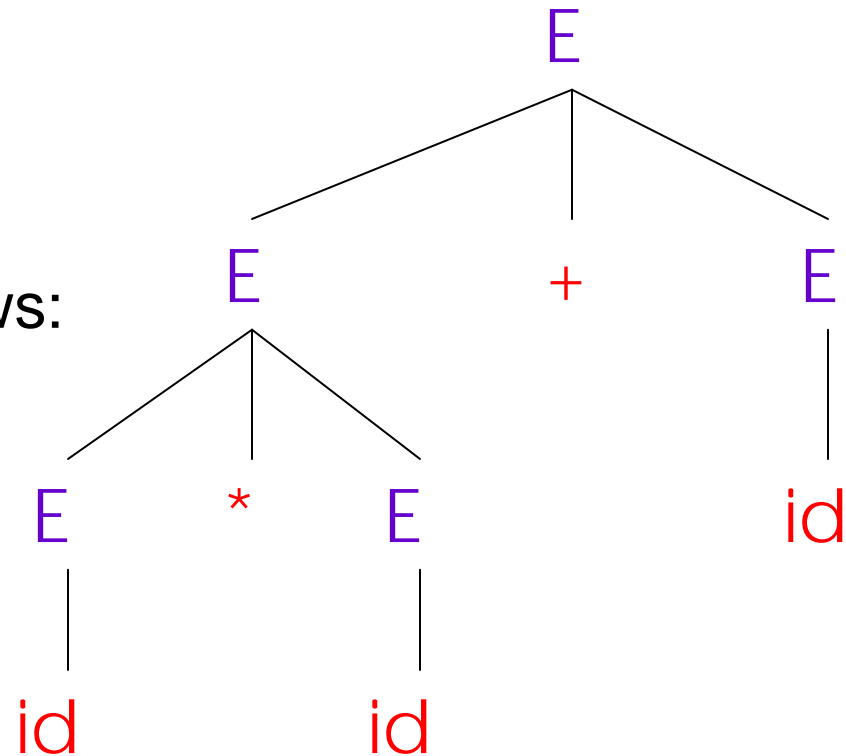
ID TIMES ID PLUS ID

we'll write tokens as follows:

id * id + id

- Output of parser:

the parse tree





Can we use REs?

Write an automaton that accepts strings

$\{ "a", "(a)", "((a))", "(((a)))" \}$

$\{ "a", "(a)", "((a))", "(((a)))", \dots "(^ka)^k" \}$

where a^k means $\underbrace{a a a \dots a}_{k \text{ times}}$

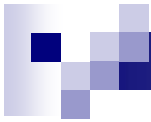


Can we use REs?

What programs are generated by?

digit+ (("+" | "-" | "" | "/") digit+)**

What important properties does this RE fail to express?



Parser : Overview

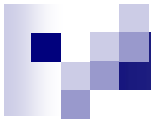
Parser plays two important roles

Recognizer: not all strings of tokens are programs

- must distinguish between valid and invalid strings of tokens

Translator: must expose program structure

- e.g., associativity and precedence
- hence must return the parse tree



Parser : Overview

Similar requirements as for scanner, we need:

A language for describing valid strings of tokens

- ☐ To function like REs in scanner
- ☐ Context-free Grammars (CFG)

A method for distinguishing valid from invalid strings of tokens

- ☐ To function like the DFA and algorithm in scanner
- ☐ Push-down Automata (PDA)



Context-free Grammar

In English:

- An integer is an arithmetic expression.
- If exp_1 and exp_2 are arithmetic expressions, then so are the following:

$\text{exp}_1 - \text{exp}_2$

$\text{exp}_1 / \text{exp}_2$

(exp_1)

As a CFG:

$E \rightarrow \text{int} | \text{t}$

$E \rightarrow E - E$

$E \rightarrow E / E$

$E \rightarrow (E)$



Reading the CFG

- The grammar has five *terminal* symbols: `intlit`, `-`, `/`, `(`, `)`
 - terminals are tokens returned by the scanner.
- The grammar has one *non-terminal* symbol: `E`
 - non-terminals describe valid (sub)sequences of tokens
- The grammar has four *productions* or *rules*
 - each of the form: $E \rightarrow \alpha$
 - left-hand side is a single non-terminal
 - right-hand side (α) is either
 - a sequence of one or more terminals and/or non-terminals
 - ε (an empty production)



Example, revisited

A more compact way to write previous grammar:

$$E \rightarrow \text{intlit} \mid E - E \mid E / E \mid (E)$$

or

$$\begin{array}{l} E \rightarrow \text{intlit} \\ \mid E - E \\ \mid E / E \\ \mid (E) \end{array}$$



A formal definition of CFGs

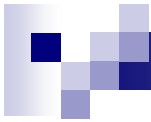
- A CFG consists of

- A set of *terminals* T
- A set of *non-terminals* N
- A *start symbol* S (a non-terminal)
- A set of *productions*:

$$X \rightarrow Y_1 Y_2 \dots Y_n$$

where $X \in N$ and $Y_i \in T \cup N \cup \{\epsilon\}$

- It is called context-free because X can always be substituted regardless where X occurs



Notational Conventions

- In these lecture notes
 - Non-terminals are written upper-case
 - Terminals are written lower-case
 - The start symbol is the left-hand side of the first production



The *Language* of a CFG

The language defined by a CFG is the set of terminal strings that can be *derived* from the start symbol of the grammar.

Derivation: Read productions as rules:

$$X \rightarrow Y_1 \dots Y_n$$

Means X can be replaced by $Y_1 \dots Y_n$



Derivation

*A sequence of grammar rule applications
leading from the start symbol*

Basic Idea

1. Begin with a string consisting of the start symbol **S**
2. Replace any non-terminal **X** in the string by a the right-hand side of **some** production

$$X \rightarrow Y_1 .. Y_n$$

3. Repeat (2) until there are no non-terminals in the string



Derivation: an example

CFG:

$E \rightarrow \text{id}$

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

derivation:

E

$\rightarrow E + E$

$\rightarrow E * E + E$

$\rightarrow \text{id} * E + E$

$\rightarrow \text{id} * \text{id} + E$

$\rightarrow \text{id} * \text{id} + \text{id}$

Is string $\text{id} * \text{id} + \text{id}$ in the
language defined by the grammar?



The Language of a CFG

More formally, write

$$X_1 \dots X_i \dots X_n \rightarrow X_1 \dots X_{i-1} Y_1 \dots Y_m X_{i+1} \dots X_n$$

if there is a production

$$X_i \rightarrow Y_1 \dots Y_m$$

A sentential form is a sequence of terminals and non-terminals



The Language of a CFG (Cont.)

Write

$$X_1 \dots X_n \xrightarrow{*} Y_1 \dots Y_m$$

if

$$X_1 \dots X_n \rightarrow \dots \rightarrow \dots \rightarrow Y_1 \dots Y_m$$

in 0 or more steps



The Language of a CFG

Let G be a context-free grammar with start symbol S . Then the language of G is:

$$\{a_1 \dots a_n \mid S \xrightarrow{*} a_1 \dots a_n \wedge a_i \in T\}$$



Examples

Strings of balanced parentheses

$$\{(^i)^i \mid i \geq 0\}$$

The grammar:

$$\begin{array}{l} S \rightarrow (S) \\ S \rightarrow \varepsilon \end{array}$$

*same
as*

$$\begin{array}{l} S \rightarrow (S) \\ \quad \mid \varepsilon \end{array}$$



Arithmetic Example

Simple arithmetic expressions:

$$E \rightarrow id \mid E + E \mid E * E \mid (E)$$

Some elements of the language:

id

id + id

(id)

id * id

(id) * id

id * (id)



Notes on CFGs

Membership in a language is “yes” or “no”

- we also need parse tree of the input!
- furthermore, we must handle errors gracefully

Need an “implementation” of CFG’s,

- i.e. the parser

Form of the grammar is important for
generating a parser



Derivations and Parse Trees

A *derivation* is a sequence of productions

$$S \rightarrow \dots \rightarrow \dots \rightarrow \dots$$

A derivation can be drawn as a tree

- Start symbol is the tree's root
- For a production $X \rightarrow Y_1 \dots Y_n$
add children $Y_1 \dots Y_n$ to node X



Derivation Example

- Grammar

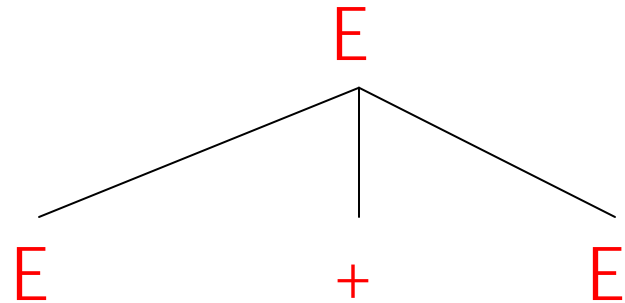
$$E \rightarrow id \mid E + E \mid E * E \mid (E)$$

- String

$$id * id + id$$

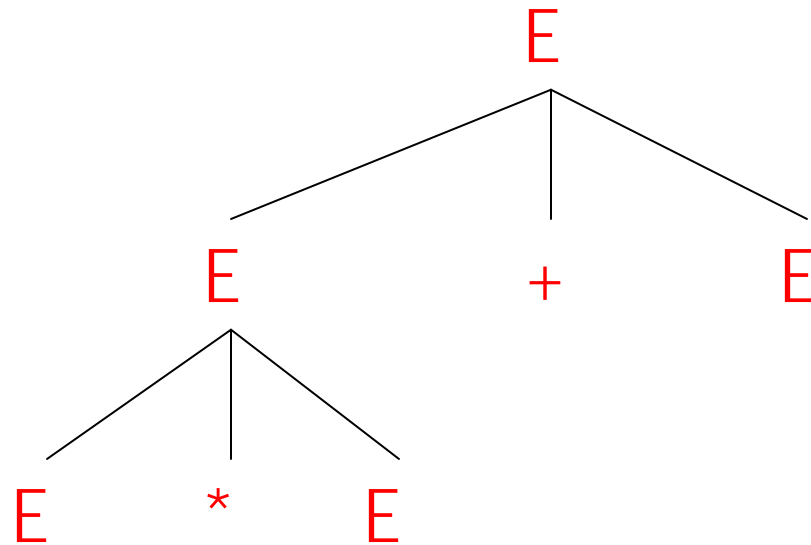
Derivation Example (Cont.)

E
 $\rightarrow E + E$



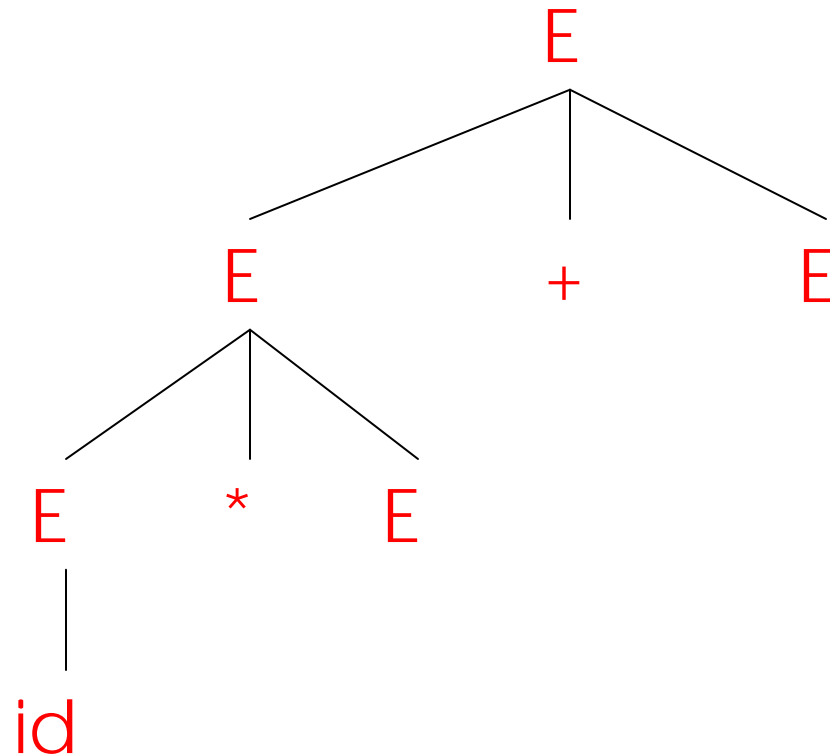
Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E * E + E$



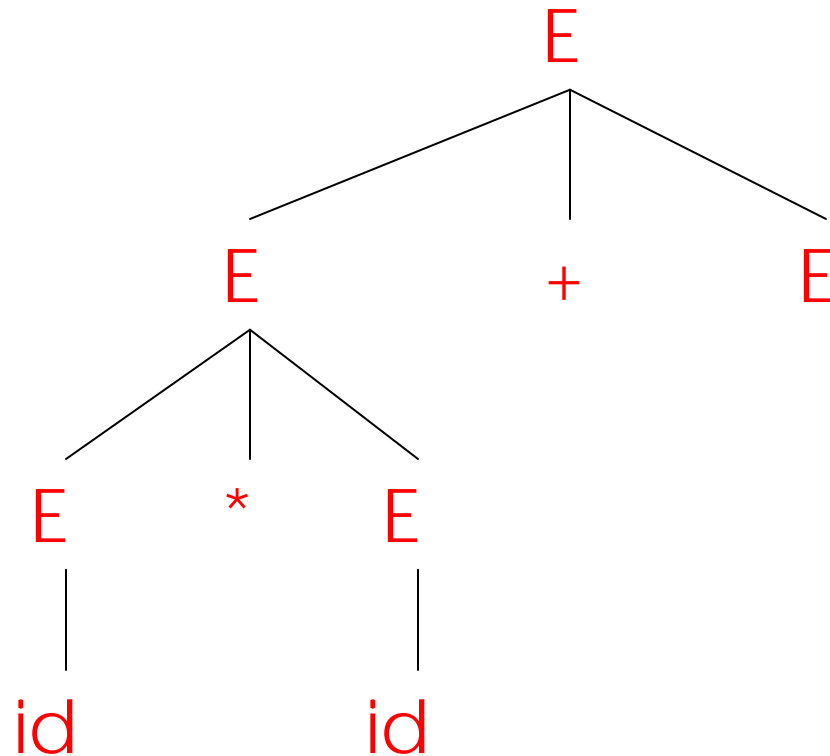
Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E * E + E$
 $\rightarrow id * E + E$



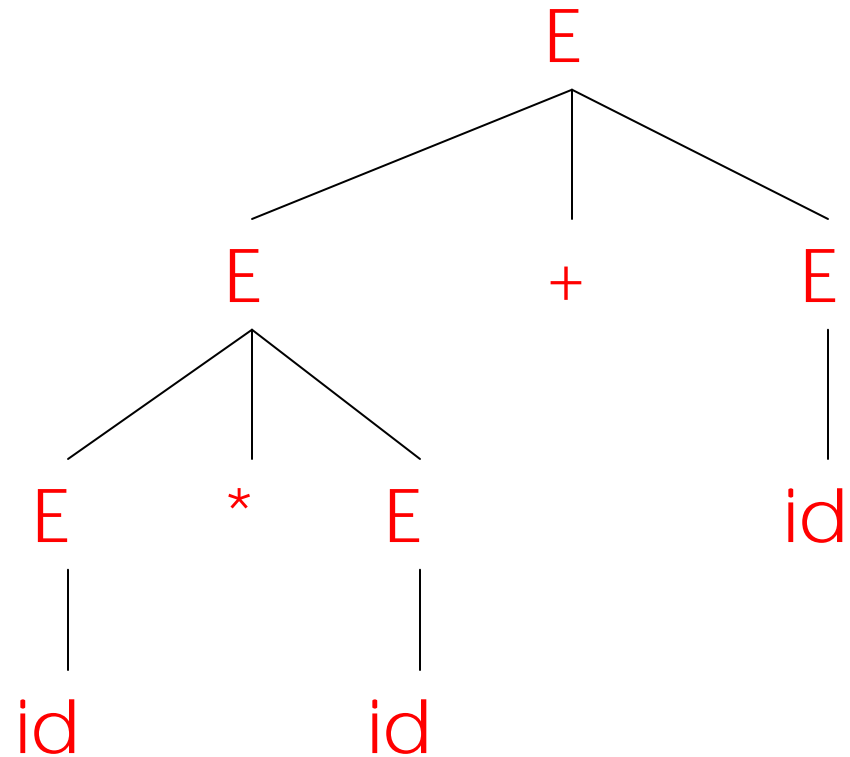
Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E * E + E$
 $\rightarrow id * E + E$
 $\rightarrow id * id + E$



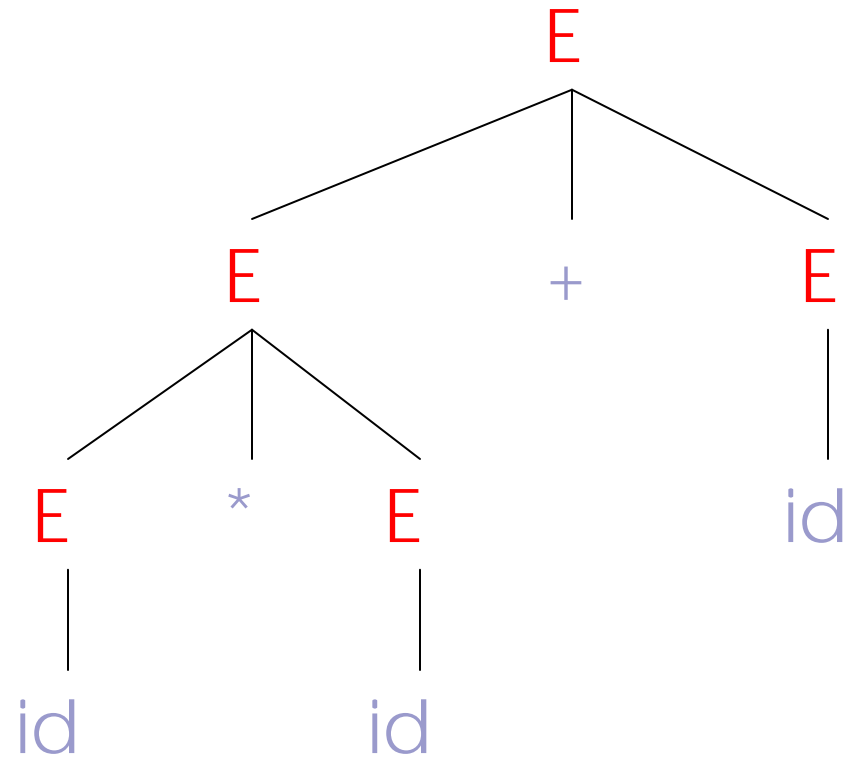
Derivation Example (Cont.)

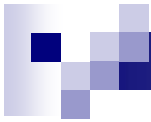
E
 $\rightarrow E + E$
 $\rightarrow E * E + E$
 $\rightarrow id * E + E$
 $\rightarrow id * id + E$
 $\rightarrow id * id + id$



Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E * E + E$
 $\rightarrow id * E + E$
 $\rightarrow id * id + E$
 $\rightarrow id * id + id$



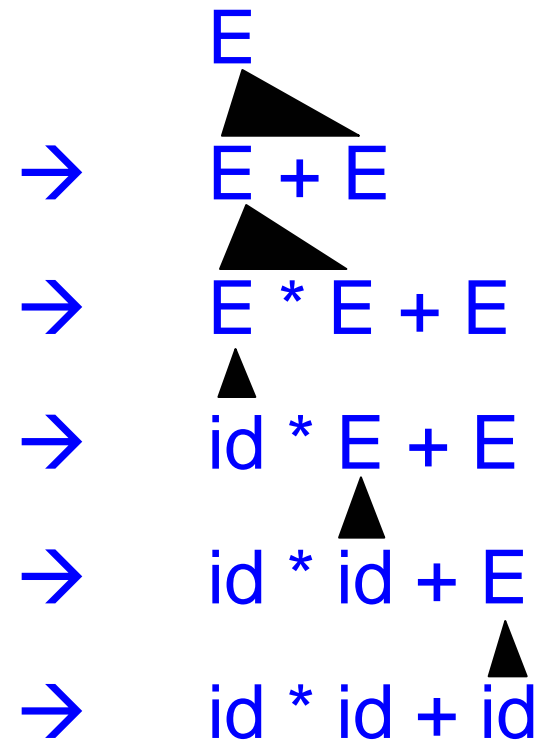


Notes on Derivations

- A parse tree has
 - Terminals at the leaves
 - Non-terminals at the interior nodes
- An in-order traversal of the leaves is the original input
- The parse tree shows the association of operations, the input string does not

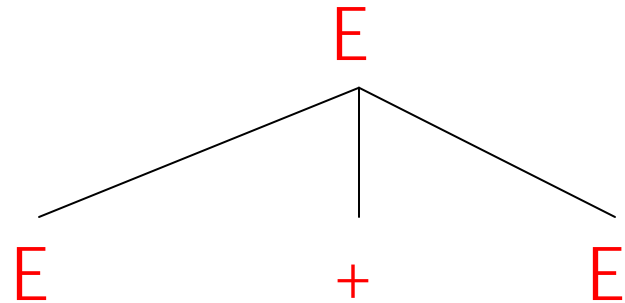
Left-most and Right-most Derivations

- The example is a *left-most* derivation
 - At each step, replace the left-most non-terminal
- There is an equivalent notion of a *right-most* derivation



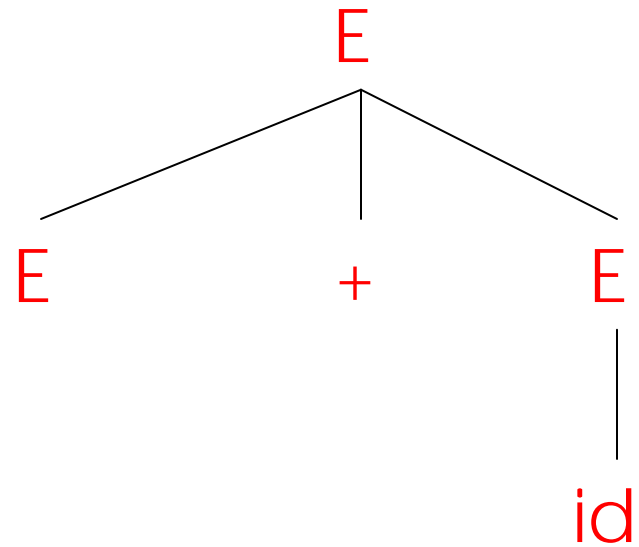
Derivation Example (Cont.)

E
 $\rightarrow E + E$



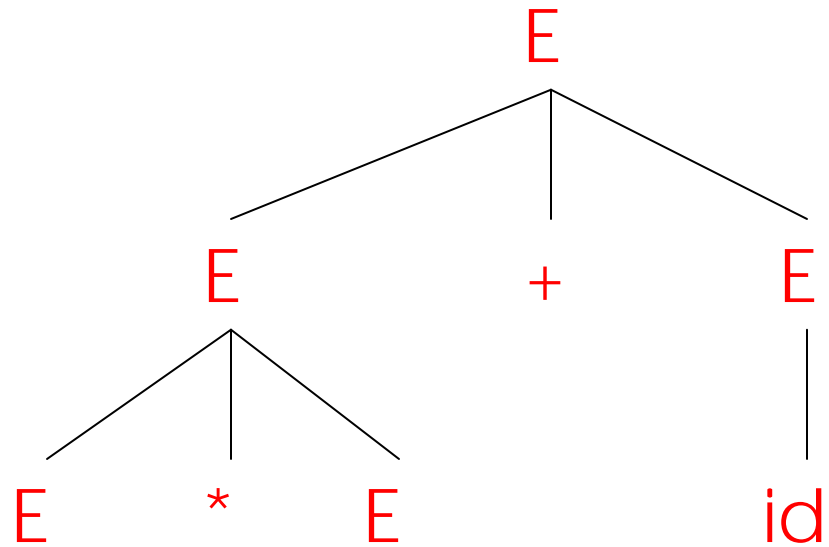
Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E + id$



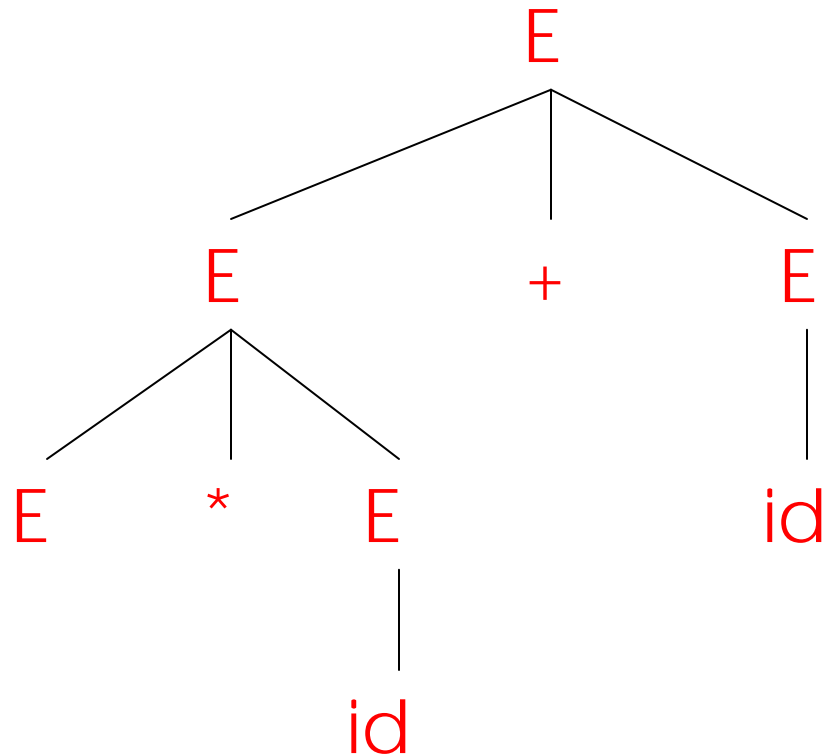
Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E + id$
 $\rightarrow E * E + id$



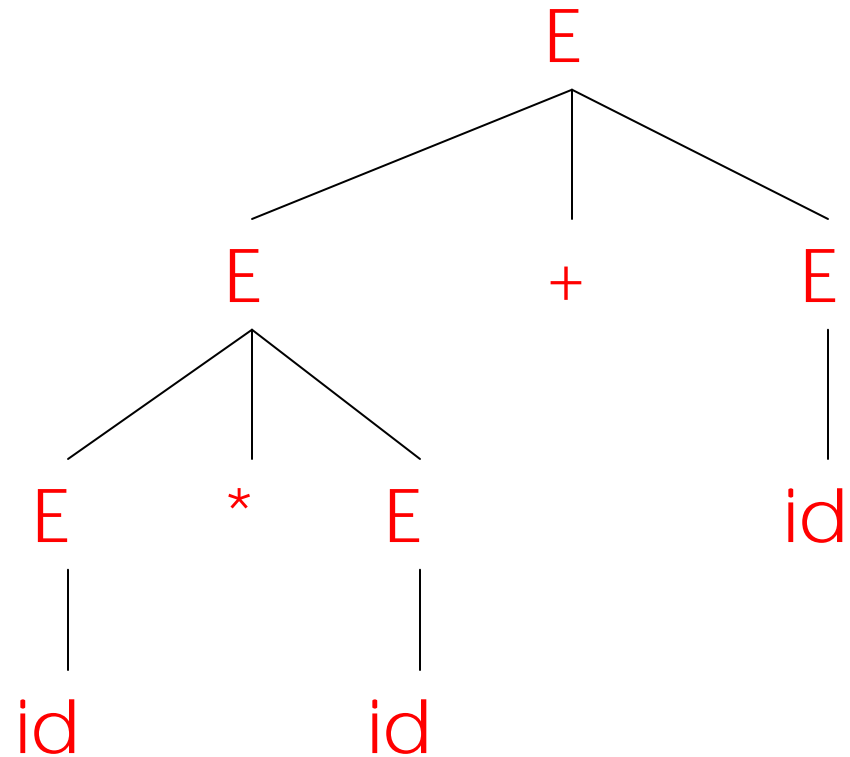
Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E + id$
 $\rightarrow E * E + id$
 $\rightarrow E * id + id$



Derivation Example (Cont.)

E
 $\rightarrow E + E$
 $\rightarrow E + id$
 $\rightarrow E * E + id$
 $\rightarrow E * id + id$
 $\rightarrow id * id + id$





Derivations and Parse Trees

- Note that right-most and left-most derivations have the same parse tree
- The difference is the order in which branches are added

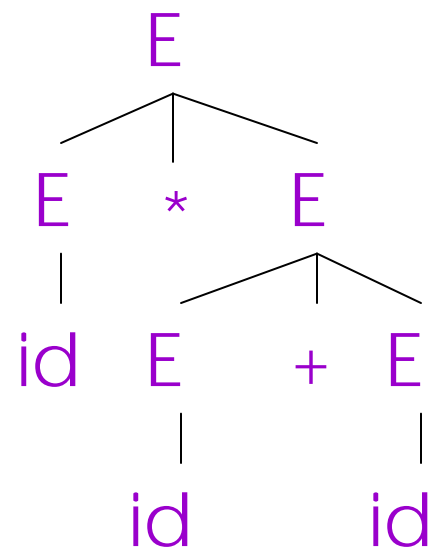
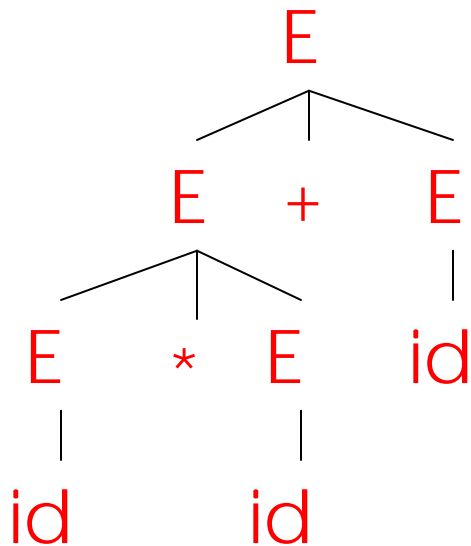


Summary of Derivations

- We are not just interested in whether $s \in L(G)$
 - we need a parse tree for s
- A derivation defines a parse tree
 - but one parse tree may have many derivations
- Left-most and right-most derivations are important in parser implementation

Multiple Parse Trees

The string “id * id + id” has two parse trees





For You To Do

For each of the two parse trees, find the corresponding **left**-most derivation

For each of the two parse trees, find the corresponding **right**-most derivation



Ambiguity

- A grammar is *ambiguous* if for some string (the following three conditions are equivalent)
 - it has more than one parse tree
 - if there is more than one right-most derivation
 - if there is more than one left-most derivation
- *Ambiguity* is **BAD**
 - Want operator precedence enforced!

Resolving Ambiguity

- Rewrite the grammar
 - use a different nonterminal for each precedence level
 - start with the lowest precedence (MINUS)

$$E \rightarrow E - E \mid E / E \mid (E) \mid id$$

rewrite to

$$E \rightarrow E - E \mid T$$
$$T \rightarrow T / T \mid F$$
$$F \rightarrow id \mid (E)$$



For You To Do

Attempt to construct a parse tree for
“id – id / id”
that shows the *wrong* precedence.



Associativity

- The grammar captures operator precedence, but it still does not capture the proper meaning
- Fails to express that both subtraction and division are *left* associative;
$$5-3-2 = ((5-3)-2) \text{ and } \textit{not} (5-(3-2))$$



For You To Do

Draw two parse trees for the expression
5-3-2

using the grammar given above;

- ☐ one that correctly groups 5-3
- ☐ one that incorrectly groups 3-2



Recursion

A grammar is *recursive* in nonterminal X if $X \xrightarrow{+} \dots X \dots$

- in one or more steps, X derives a sequence of symbols that includes an X

A grammar is *left recursive* in X if $X \xrightarrow{+} X \dots$

- in one or more steps, X derives a sequence of symbols that *starts* with an X

A grammar is *right recursive* in X if $X \xrightarrow{+} \dots X$

- in one or more steps, X derives a sequence of symbols that *ends* with an X



Resolving Associativity

- The grammar given above is both left and right recursive in nonterminals E and T
- To correctly expresses operator associativity:
 - For left associativity, use left recursion.
 - For right associativity, use right recursion.
- Here's the correct grammar:

$$\begin{aligned} E &\rightarrow E - T \mid T \\ T &\rightarrow T / F \mid F \\ F &\rightarrow \text{id} \mid (E) \end{aligned}$$



Ambiguity: The Dangling Else

Consider the grammar

$$\begin{array}{l} E \rightarrow \text{if } E \text{ then } E \\ \quad | \text{if } E \text{ then } E \text{ else } E \\ \quad | \text{print} \end{array}$$

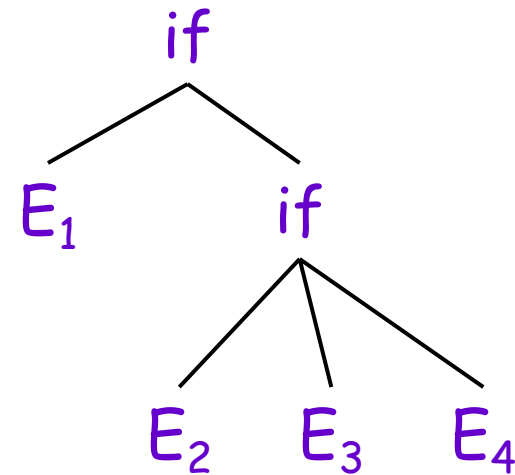
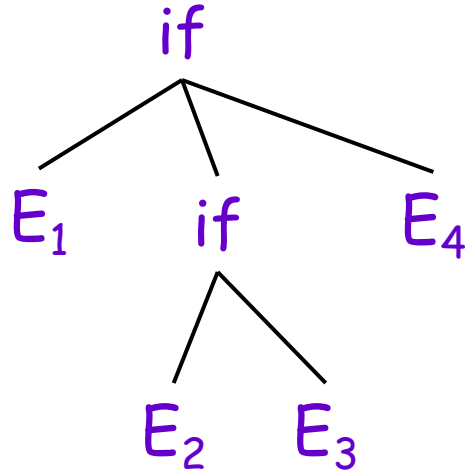
This grammar is also ambiguous

The Dangling Else: Example

The expression

if E_1 then if E_2 then E_3 else E_4

has two parse trees





The Dangling Else: A Fix

else matches the closest unmatched **then**

We can describe this in the grammar

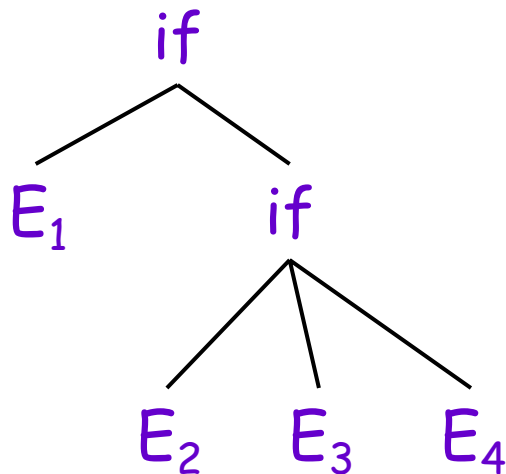
$$\begin{array}{ll} E & \rightarrow \text{MIF} \quad /* \text{ all then are matched } */ \\ & | \text{ UIF} \quad /* \text{ some then are unmatched } */ \end{array}$$

$$\begin{array}{l} \text{MIF} \rightarrow \text{if } E \text{ then MIF else MIF} \\ \quad | \text{ print} \end{array}$$

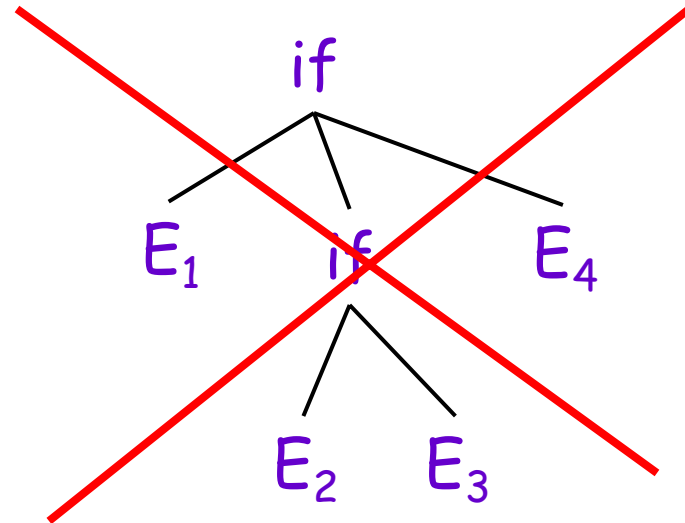
$$\begin{array}{l} \text{UIF} \rightarrow \text{if } E \text{ then } E \\ \quad | \text{ if } E \text{ then MIF else UIF} \end{array}$$

The Dangling Else: Example Revisited

The expression $\text{if } E_1 \text{ then if } E_2 \text{ then } E_3 \text{ else } E_4$



- A valid parse tree (for a **UIF**)



- Not valid because the **then** expression is not a **MIF**



Precedence and Associativity Declarations

- Instead of rewriting the grammar
 - Use the more natural (ambiguous) grammar
 - Along with disambiguating declarations
- Most parser generators allow precedence and associativity declarations to disambiguate grammars



Extended Backus-Naur Form (EBNF)

- StaticJava and ExtendedStaticJava concrete syntax grammar are written in EBNF
- BNF is a meta-syntax for CFG
$$\langle \text{symbol} \rangle ::= \langle \text{symbol}_1 \rangle \dots \langle \text{symbol}_n \rangle$$
 - **symbols** that appear in the **left-hand** side are **non-terminals** (otherwise, **terminals**)
 - use **|** for **choice**
 - whitespace is **usually** assumed in between
- EBNF adds *****, **?**, **+**, ... similar to regular expression



Notes on StaticJava EBNF Grammar

- The concrete syntax is optimized for reading
 - it is ambiguous
 - when in doubt, remember that it is a strict subset of Java (thus, it is a Java program)
- Does it have the dangling-else problem?
 - why?



ANTLR Grammar for StaticJava

Preview

EBNF

<rulename> ::=
 <alternative₁> |
 <alternative₂> |
 ... |
 <alternative_n>

ANTLR

rulename :
 alternative₁ |
 alternative₂ |
 ... |
 alternative_n ;

...more ANTLR after we learn LL(k) parsing next time!