

## review

## The main work of lexical analysis

Classify program substrings according to token class

```

    if (i == j)
        Z = 0;
    else
        Z = 1;

\tif (i == j)\n\t\tz = 0;\n\telse\n\t\tz = 1;

```

## Token class

Identifier  
 keyword  
 number  
 whitespace  
 special symbol  
 and so on

## Token

**<token class, lexeme>**

(IF, "if") (PLUS, "+") (ID, "x")

Example: foo=42

## Summary

- An implementation must do two things:
  1. Recognize substrings corresponding to tokens
    - The *lexemes*
  2. Identify the token class of each lexeme

## Questions: Regular Languages

**What's the function of Regular expressions?**

Specify the lexical structure of programming languages

**lexical structure=token classes**

writing patterns to match a specific sequence of characters (also known as a string)

## Questions:Regular Languages

- lexical structure=token classes

**Identifies:** sequences of letters and digits beginning with a letter

**Numbers:** numeric constants, such as 42, 3.14

**Reserve words:** Fixed strings of characters that have special meaning in the language

**Special symbols:** include arithmetic operations, assignment, equality and so on

## Multiple Choices

Choose the regular languages that are equivalent to the given regular language:  
 $(0 \mid 1)^* 1 (0 \mid 1)^*$

$\Sigma = \{0, 1\}$

☐ A  $(01 \mid 11)^* (0 \mid 1)^*$

☐ B  $(0 \mid 1)^* (10 \mid 11 \mid 1)(0 \mid 1)^*$

☐ C  $(1 \mid 0)^* 1 (1 \mid 0)^*$

☐ D  $(0 \mid 1)^* (0 \mid 1)(0 \mid 1)^*$

## Lexical Specification

Extensions of **regular expression** to construct a full lexical specification on the programming language

## Notations for regular expressions

- At least one:  $A^+$   $\equiv AA^*$
- Union:  $A \mid B$   $\equiv A + B$
- Option:  $A?$   $\equiv A + \varepsilon$
- Range:  $'a' + 'b' + \dots + 'z'$   $\equiv [a-z]$
- Excluded range:  
 $\text{complement of } [a-z] \equiv [^a-z]$

## Partitioning a string into tokens

- Write a rexp for the lexemes of each token class
  - Number =  $\text{digit}^+$
  - Keyword =  $'\text{if}' + '\text{else}' + \dots$
  - Identifier =  $\text{letter}(\text{letter} + \text{digit})^*$
  - OpenPar =  $'('$
- Construct  $R$ , matching all lexemes for all tokens

$R = \text{Keyword} + \text{Identifier} + \text{Number} + \dots$   
 $= R_1 + R_2 + \dots$

## Partitioning a string into tokens

- Let input be  $x_1 \dots x_n$   
 For  $1 \leq i \leq n$  check  
 $x_1 \dots x_i \in L(R)$
- If success, then we know that  
 $x_1 \dots x_i \in L(R_j)$  for some  $j$
- Remove  $x_1 \dots x_i$  from input and go to (3)

## Questions: Finite Automata

### What's the function of finite automata?

Represent an abstracted version of a program

Implementation model for regular expression

Regular expression=specification

Finite automata=implementation

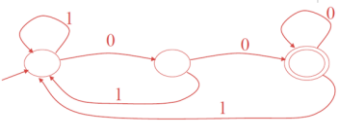
Example

## Single Choice

### In-Video Quiz

### Finite Automata

Select the regular language that denotes the same language as this finite automaton



- ☐ A  $(0 + 1)^*$ 
☒ B  $(0 + 1)^*00$ 
☐ C  $1^* + (01)^* + (001)^* + (000^*1)^*$

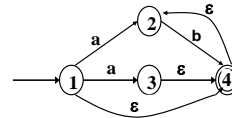
## Questions Finite Automata

### What's the difference between DFA and NFA?

- **Deterministic Finite Automata (DFA)**
  - One transition per input per state
  - No  $\epsilon$ -moves
- **Nondeterministic Finite Automata (NFA)**
  - Can have multiple transitions for one input in a given state
  - Can have  $\epsilon$ -moves

## Questions Finite Automata

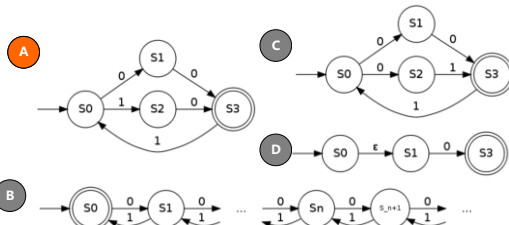
### Give the paths for recognizing string **abb**



## Single Choice

### Quiz 1 Question 11

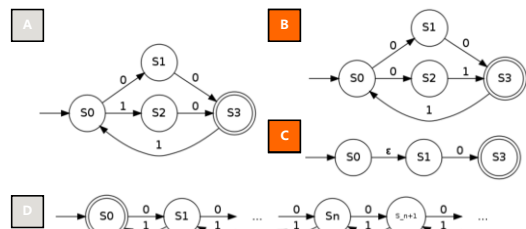
Which of the following automata are DFAs?



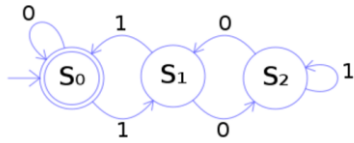
## Multiple Choices

### Quiz 1 Question 12

Which of the following automata are NFAs?



## Making Useful Programs with DFAs



What strings will return "True" when passed through this DFA?  
 this automaton will return "True" for any binary number that represents a multiple of 3.

## Chapter 3 Context-Free Grammars and Parsing

### Study Goals:

Understand:

Context-free grammar, Derivation, Reduction, Parse tree, Abstract syntax tree

Know:

Hierarchy of grammar, Ambiguous Grammars

### Main Content of Parsing Study

Specification of syntax structure :

[Context-free grammar](#)

Method of turning grammar rules into code for parsing

[Top-down parsing method](#)

[Bottom-up parsing method](#)

	Scanning	Parsing
<b>Task</b>	determining the structure of tokens	determining the syntax or structure of a program
<b>Describing Tools</b>	regular expression	context-free grammar
<b>Algorithmic Method</b>	represent by DFA	top-down parsing bottom-up parsing
<b>Result Data Structure</b>	liner structure	parser tree or syntax tree, they are recursive

### 3.1 The Parsing Process

### 3.2 Context-Free Grammars

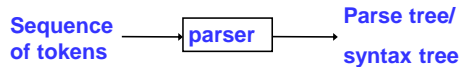
### 3.3 Parse Trees and Abstract Syntax Tree

### 3.4 Ambiguity

### 3.1 The Parsing Process

#### Task of the parser

- Determine the syntactic structure of a program from the tokens produced by the scanner
- Either explicitly or implicitly, construct a parse tree or syntax tree that represent this structure.



### Interface of the parser

#### Input:

The parser calls a scanner procedure to fetch the next token from the input as it is needed during the parsing process

#### Output:

An explicit or implicit syntax tree needs to be constructed. Each node of the syntax tree includes the attributes needed for the remainder of the compilation process

### Error Handling of Parser

#### Error Recover

Report meaningful error messages and resume the parsing as close to the actual error as possible (to find as many errors as possible)

#### Error repair

Infer a possible corrected code version from the incorrect version presented to it (This is usually done only in simple cases)



### 3.2 Context-Free Grammars

#### Function

A context-free grammar is a specification for the syntactic structure of a programming language

It is similar to regular expressions except that a context-free grammar involves recursive rules.

#### Example

context-free grammar for integer arithmetic expression

$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number}$

$\text{op} \rightarrow + \mid - \mid *$

regular expression for number

$\text{number} = \text{digit digit}^*$

$\text{digit} = 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

#### 3.2.1 Definition of Context-free Grammar

#### Definition

A context-free grammar  $G = (V_T, V_N, P, S)$ :

1.  $V_T$  is a set of **terminals**
2.  $V_N$  is a set of **nonterminals**,  $V_N \cap V_T = \emptyset$
3.  $P$  is a set of **productions**, or grammar rules, of the form  $A \rightarrow \alpha$ , where  $A \in V_N$  and  $\alpha \in (V_N \cup V_T)^*$
4.  $S$  is a start symbol,  $S \in V_N$

### Example

The grammar of simple arithmetic expressions  
 $G=(V_T, V_N, P, S)$

- $V_T = \{\text{num}, +, -, *, /, (, )\}$      $\text{exp} \rightarrow \text{exp op exp}$
- $V_N = \{\text{exp}, \text{op}\}$      $\text{exp} \rightarrow (\text{exp})$
- $\text{exp}$  is the start symbol     $\text{exp} \rightarrow \text{num}$
- productions are:     $\text{op} \rightarrow +$
- $\text{op} \rightarrow -$
- $\text{op} \rightarrow *$
- $\text{op} \rightarrow /$     ➡

### Explanation

1.  $V_T$  are the basic symbols from which strings are formed. Terminals are tokens
2.  $V_N$  are names for structures that denote sets of strings
3. The set of strings that start symbol "**S**" denotes is the language defined by the grammar
4. A production defines a structure whose name is to the left of the arrow. The layout of the structure is defined by the right of the arrow

5. " $\rightarrow$ " shows the left of the arrow cannot simply be replaced by its definition, as a result of the recursive nature of the definition
6. The form of productions ( $A \rightarrow \alpha$ ) is called Backus-Naur form (or BNF)

### Notation Conventions

With the following conventions, we can only write productions of a grammar

Unless otherwise stated, the left side of the first production is the start symbol  
 Using lower-case letters to represent terminals

Using upper-case letters or name with  $\langle \dots \rangle$  to represent nonterminals

If  $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$  are all productions with  $A$  on the left, we may write  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$

### Example

Grammar for simple arithmetic expressions can be written as:

$E \rightarrow E O E \mid (E) \mid \text{num}$

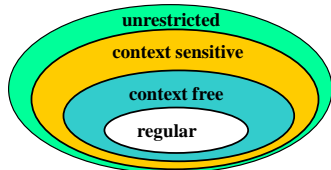
$O \rightarrow + \mid - \mid * \mid /$

### Chomsky hierarchy of grammar

There are four kinds of grammars:

1. Unrestricted grammar (type 0)
2. Context sensitive grammar (type 1)
3. Context free grammar (type 2)
4. Regular grammar (type 3)

The Relationship of these four grammars:



Kind	Production	Explanation
unrestricted (type 0)	$\alpha \rightarrow \beta \quad V = V_N \cup V_T$ $\alpha \in V^+, \beta \in V^*$	there is no restriction to production
context sensitive (type 1)	$\alpha A \gamma \rightarrow \alpha \beta \gamma, \alpha, \gamma \in V^*$ $A \in V_N, \beta \in V^+$	A may be replaced only if $\alpha$ occurs before A and $\gamma$ occurs after A
context free (type 2)	$A \rightarrow \beta, A \in V_N, \beta \in V^*$	A may be replaced anywhere, regardless of where A occurs
regular (type 3)	$A \rightarrow aB \text{ or } A \rightarrow a, A, B \in V_N, a \in V_T$	equivalent to regular expression



### 3.2.2 Derivation and Reduction

#### Function of derivation and reduction

Context-free grammar rules determine the set of syntactically legal strings of tokens

For example: Corresponding to grammar

$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number}$

$\text{op} \rightarrow + \mid - \mid *$

$(34-3)*42$  is a legal string

$34-3*42$  is not a legal string

Grammar rules determine the legal strings of tokens by means of derivation or reduction

#### A derivation step and a reduction step

$A \rightarrow \beta$  is a production of  $G$ , if there are strings  $v$  and  $w : v = \alpha A \gamma, w = \alpha \beta \gamma$ , where  $\alpha, \gamma \in (V_N \cup V_T)^*$ , we say there is a **derivation step** from  $v$  to  $w$ , or a **reduction step** from  $w$  to  $v$ , written as  $v \Rightarrow w$

- A **derivation step** is a replacement of a nonterminal by the right-hand side of the production
- A **reduction step** is a replacement of the right-hand side of production by the nonterminal on the left

#### Example

Grammar  $G$ :  $S \rightarrow 0S1, S \rightarrow 01$

A derivation step

$0S1 \Rightarrow 00S11 \quad (S \rightarrow 0S1)$

$00S11 \Rightarrow 000S111 \quad (S \rightarrow 0S1)$

$000S111 \Rightarrow 00001111 \quad (S \rightarrow 01)$

$S \Rightarrow 0S1 \quad (S \rightarrow 0S1)$

#### Derivation and Reduction

- The closure of  $\Rightarrow$ ,  $\alpha \Rightarrow^* \beta$

$\alpha \Rightarrow^* \beta$  if and only if there is a sequence of 0 or more derivation

steps ( $n \geq 0$ ),  $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_{n-1} \Rightarrow \alpha_n$ ,

such that  $\alpha = \alpha_1$  and  $\beta = \alpha_n$

- $S \Rightarrow^* w$ , where  $w \in V_T^*$  and  $S$  is the start symbol of  $G$  is called a **derivation** from  $S$  to  $w$  or a **reduction** from  $w$  to  $S$

**Example**

$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{num}$

$\text{op} \rightarrow + \mid - \mid *$

a derivation for  $(34-3)*42$  is:

- (1)  $\text{exp} \Rightarrow \text{exp op exp}$  ( $\text{exp} \rightarrow \text{exp op exp}$ )
- (2)  $\Rightarrow \text{exp op num}$  ( $\text{exp} \rightarrow \text{num}$ )
- (3)  $\Rightarrow \text{exp} * \text{num}$  ( $\text{op} \rightarrow *$ )
- (4)  $\Rightarrow (\text{exp}) * \text{num}$  ( $\text{exp} \rightarrow (\text{exp})$ )
- (5)  $\Rightarrow (\text{exp op exp}) * \text{num}$  ( $\text{exp} \rightarrow \text{exp op exp}$ )
- (6)  $\Rightarrow (\text{exp op num}) * \text{num}$  ( $\text{exp} \rightarrow \text{num}$ )
- (7)  $\Rightarrow (\text{exp-num}) * \text{num}$  ( $\text{op} \rightarrow -$ )
- (8)  $\Rightarrow (\text{num-num}) * \text{num}$  ( $\text{exp} \rightarrow \text{num}$ )

**3.2.3 The Language Defined by a Grammar**

A **sentential form** of G

**S** is the start symbol of G, if  $S \Rightarrow^* \alpha$ ,  $\alpha \in (V_N \cup V_T)^*$ ,  $\alpha$  is a sentential form of G

A **sentence** of G

**w** is a sentential form of G, if **w** contains only terminals, then **w** is a sentence of G

**Example**

G:  $S \rightarrow 0S1$ ,  $S \rightarrow 01$

$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000S111 \Rightarrow 00001111$

- $S, 0S1, 00S11, 000S111, 00001111$  are all sentential forms of G
- $00001111$  is a sentence of G

**The language defined by G, written as  $L(G)$** 

$L(G) = \{w \in V_T^* \mid \text{there exists a derivation } S \Rightarrow^* w \text{ of G}\}$

that is,  $L(G)$  is the set of sentences derivable from S

Example:

G:  $S \rightarrow 0S1$ ,  $S \rightarrow 01$

$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 0^3S1^3 \Rightarrow \dots \Rightarrow 0^{n-1}S1^{n-1}$

$\Rightarrow 0^n1^n$

$L(G) = \{0^n1^n \mid n \geq 1\}$

**The Relation between Grammar G and  $L(G)$** 

- Given G,  $L(G)$  can be obtained by derivation

Example

G:  $E \rightarrow E+T \mid T$      $T \rightarrow T \times F \mid F$      $F \rightarrow (E) \mid a$

$E \Rightarrow E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow a+T \Rightarrow a+T \times F$   
 $\Rightarrow a+F \times F \Rightarrow a+a \times F \Rightarrow a+a \times a$

$L(G)$  are arithmetic expressions consisted of **a, +,  $\times$ , ( and )**

- Given the description of **L**, we can design grammar for **L**

Example

A language consists of 0 and 1, every string of the language has the same number of 0 and 1

Grammar for L is:

$A \rightarrow 0B \mid 1C$

$B \rightarrow 1 \mid 1A \mid 0BB$

$C \rightarrow 0 \mid 0A \mid 1CC$





### 3.3 Parse Trees and Abstract Syntax Trees

#### 3.3.1 Parse Trees

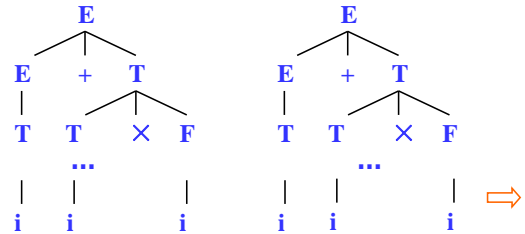
Function of parse trees

- A parse tree is a useful representation of the structure of a string of tokens
- Parse trees represent derivations visually

□ Example:  $E \rightarrow E+T \mid T$      $T \rightarrow T \times F \mid F$      $F \rightarrow (E) \mid i$   
the derivation of "i+i×i" and the parse tree

$E \Rightarrow E+T \Rightarrow E+T \times F \Rightarrow T+T \times F \Rightarrow \dots \Rightarrow i+i \times i$

$E \Rightarrow E+T \Rightarrow T+T \Rightarrow T+T \times F \Rightarrow \dots \Rightarrow i+i \times i$



Explanation:

- A parse tree of a string corresponds in general to many derivations of the string
- Parse trees abstract the essential features of derivations while factoring out superficial difference in order
- Derivations do not uniquely represent the structure of the strings they construct, while parse trees do

#### Definition of Parse tree

A parse tree over the grammar  $G$  is a rooted labeled tree with the following properties:

1. The root node is labeled with the start symbol  $S$
2. Each leaf node is labeled with a terminal or with  $\epsilon$
3. Each nonleaf node is labeled with a nonterminal
4. If a node with label  $A \in V_N$  has  $n$  children with labels  $X_1, X_2, \dots, X_n$  (which may be terminals or nonterminals), then

⇐  $A \rightarrow X_1 X_2 \dots X_n \in P$

#### Leftmost and Rightmost Derivation

Leftmost derivation

- A derivation in which the leftmost nonterminal is replaced at each step in the derivation
- It corresponds to a preorder traversal of the parse tree

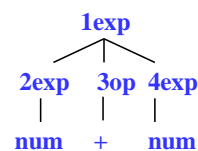
#### Example

$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{num}$

$\text{op} \rightarrow + \mid - \mid *$

❖ The leftmost derivation of "num+num" is:

$\text{exp} \Rightarrow \text{exp op exp} \Rightarrow \text{num op exp}$   
 $\Rightarrow \text{num} + \text{exp} \Rightarrow \text{num} + \text{num}$

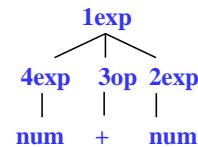


### Rightmost derivation

- A derivation in which the rightmost nonterminal is replaced at each step in the derivation
- It corresponds to the **reverse** of a postorder traversal of the parse tree

❖ The rightmost derivation of “num+num” is:

exp => exp op exp => exp op num  
=> exp + num => num + num



- Leftmost and rightmost derivation are unique for the string they construct
- They are uniquely associated with the parse tree

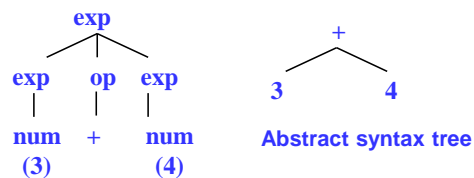
### Relation of Parse tree and Derivation

- Each derivation gives rise to a parse tree
- Many derivations may give rise to the same parse tree
- Each parse tree has a unique leftmost and rightmost derivation that give rise to it
- Parse trees uniquely express the structure of syntax, as do leftmost and rightmost derivations, but not other derivations in general

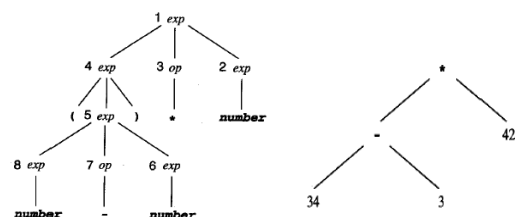
## 3.3.2 Abstract Syntax Trees

### The need of abstract syntax tree

- A parse tree contains much more information than is absolutely necessary for a compiler to produce executable code
- For example



The parse tree and abstract syntax tree for expression (34-3)\*42



### Definition of Abstract Syntax trees

Abstract syntax trees represent abstractions of the actual string of tokens. Nevertheless they contain all the information needed for translation (represent precisely the semantic content of the string).

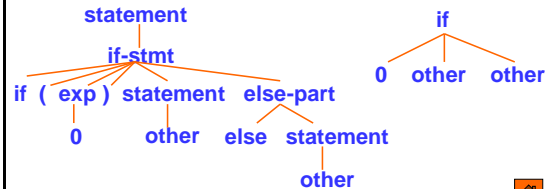
In more efficient form than parse trees

A parser will go through all the steps represented by a parse tree, but will usually only construct an abstract syntax tree

### Example Grammar for statement

$\langle \text{statement} \rangle \rightarrow \langle \text{if-stmt} \rangle \mid \text{other}$   
 $\langle \text{if-stmt} \rangle \rightarrow \text{if } (\langle \text{exp} \rangle) \langle \text{statement} \rangle \langle \text{else-part} \rangle$   
 $\langle \text{else-part} \rangle \rightarrow \text{else } \langle \text{statement} \rangle \mid \epsilon$   
 $\langle \text{exp} \rangle \rightarrow 0 \mid 1$

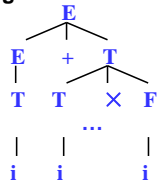
The parse tree and abstract syntax tree for string "if (0) other else other"



### 3.4 Ambiguity

In general, a string of tokens has one parse tree, which corresponds to more than one derivations of the string

Parse tree of string "i+i×i"



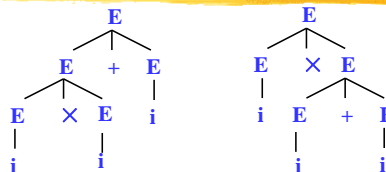
$E \Rightarrow E+T \Rightarrow E+T \times F \Rightarrow T+T \times F \Rightarrow \dots \Rightarrow i+i \times i$

$E \Rightarrow E+T \Rightarrow T+T \Rightarrow T+T \times F \Rightarrow \dots \Rightarrow i+i \times i$

It is possible for some grammars to permit a string to have more than one parse tree (or leftmost/rightmost derivation)

For example:

Integer arithmetic grammar:  $E \rightarrow E+E \mid E \times E \mid (E) \mid i$   
String "i×i+i" has two different parse trees:



Corresponding to two leftmost derivations:

1:  $E \Rightarrow E+E \Rightarrow E \times E+E \Rightarrow i \times E+E \Rightarrow i \times i+E \Rightarrow i \times i+i$

2:  $E \Rightarrow E \times E \Rightarrow i \times E \Rightarrow i \times E+E \Rightarrow i \times i+E \Rightarrow i \times i+i$

### Ambiguity

A grammar  $G$  is ambiguous if there exists a string  $w \in L(G)$  such that  $w$  has two distinct parse trees (or leftmost /rightmost derivations)

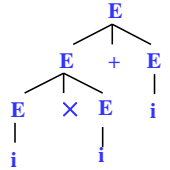
## How to deal with ambiguity

Two basic methods are used to deal with ambiguities

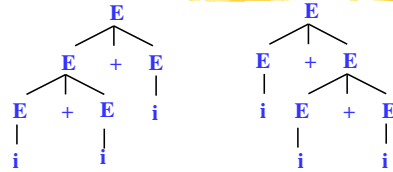
1. State a rule that specifies in each ambiguous case which of the parse trees is the correct one

Disambiguating rule:

- > **MUL** has a higher precedence than **ADD**
- > **MUL** and **ADD** are left associative



Integer arithmetic grammar:  $E \rightarrow E+E | E \times E | (E) | i$   
String "i+i+i" has two different parse trees:



Base on disambiguating rule, the first one is correct

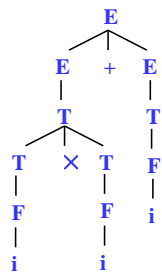
2. Change the grammar into a form that forces the construction of the correct parse tree

For example:  $E \rightarrow E+E | E \times E | (E) | i$

- > **Add precedence**  
Group the operators into groups of equal precedence, and for each precedence we must write a different rule:

$E \rightarrow E+E | T \quad T \rightarrow T \times T | F \quad F \rightarrow (E) | i$

"i×i+i" is not ambiguous, the leftmost derivation for it is:  
 $E \Rightarrow E+E \Rightarrow T+E \Rightarrow T \times T+E \Rightarrow \dots \Rightarrow i \times i+i$



but "i+i+i" is still ambiguous, which has two leftmost derivations:

$E \Rightarrow E+E \Rightarrow E+E+E \Rightarrow \dots$

$E \Rightarrow E+E \Rightarrow T+E \Rightarrow F+E \Rightarrow i+E \Rightarrow \dots$

- > **Add associativity**: a left recursive rule makes its operators associate on the left, while a right recursive rule makes them associate on the right

$E \rightarrow E+T | T \quad T \rightarrow T \times F | F \quad F \rightarrow (E) | i$

"i+i+i" is not ambiguous, the leftmost derivation for it is:

$E \Rightarrow E+T \Rightarrow E+T+T \Rightarrow \dots \Rightarrow i+i+i$

