

## Regular Expressions

- In **formal languages**, *regular expressions* are used to define **strings of symbols**. The set of all possible strings that a regular expression defines is called the formal language of the regular expression.
- When defining the syntax of programming languages the symbols of interest are characters, e.g. ASCII or UNICODE characters.
- In programming languages, regular expressions and grammars are used to define the syntax of the programming language. **Regular expressions** are used to define the basic syntactical elements – identifiers, keywords, numbers, strings<sup>1</sup>, operators, brackets etc. **Grammars** are used to define how these basic elements can be composed into complex syntactical elements – expressions, statements, declarations, functions, programs etc.
- In a compiler, the **lexical analyser**<sup>2</sup> transforms a stream of characters (i.e. symbols) into a stream of tokens<sup>3</sup> (i.e. objects).
- Although lexical analysers can be written by hand, a more professional approach is to use a tool<sup>4</sup> that reads a list of regular expressions that define the basic syntactical elements and generates the lexical analyser which can be linked into the rest of the compiler<sup>5</sup>.
- Lexical analyser generators are implemented by converting regular expressions into very fast **finite automata**<sup>6</sup>, which are formal way for describing certain algorithms, e.g. pattern matching.

Symbols	
Is <b>A</b> a symbol	Yes
Is <b>B</b> a symbol	Yes
Is <b>C</b> a symbol	Yes
Is <b>£</b> a symbol	Yes
Is <b>♥</b> a symbol	Yes
Is <b>☺</b> a symbol	Yes
Strings	
Is <b>AA</b> a symbol	No
Is <b>AA</b> a string	Yes
Is <b>A</b> a string	Yes
Is <b>ABBA</b> a string	Yes
Is <b>ε</b> the <i>empty string</i>	Yes, <b>ε</b> is a string with no symbols
Regular Expressions - Symbol	
Is <b>A</b> a symbol	No
Is <b>A</b> a string	No
Is <b>A</b> a regular expression	Yes
Is <b>A</b> a regular expression	No
Will <b>A</b> be recognised by <b>A</b>	Yes
Will <b>a</b> be recognised by <b>A</b>	No

<sup>1</sup> Not to be confused with strings of symbols, which are a formal concept.

<sup>2</sup> A lexical analyser is also known as a **scanner**.

<sup>3</sup> When we study grammars, we will refer to tokens as terminals or terminal symbols.

<sup>4</sup> The classic example is **LEX** and its modern version FLEX that outputs lexical analysers in C.

<sup>5</sup> Modern tools like **ANTLR** often handle all syntax generating both the lexical analyser and the parser from the regular expressions and the grammar.

<sup>6</sup> Students are highly encouraged to read more on this in textbooks and also last year's slides.

Will <b>B</b> be recognised by <b>A</b>	No
Can <b>A</b> generate <b>A</b>	Yes
Can <b>A</b> generate <b>B</b>	No
Is $\epsilon$ a regular expression	Yes
Will $\epsilon$ be recognised by <b>E</b>	Yes
<b>Regular Expressions - Concatenation</b>	
Is <b>AA</b> a regular expression	Yes
Will <b>A</b> be recognised by <b>AA</b>	No
Will <b>AA</b> be recognised by <b>AA</b>	Yes
Will <b>AAA</b> be recognised by <b>AA</b>	No
Can <b>AA</b> generate <b>A</b>	No
Can <b>AA</b> generate <b>AA</b>	Yes
Can <b>AA</b> generate <b>AAA</b>	No
Is <b>WHILE</b> a regular expression	Yes, will match <b>WHILE</b>
Is <b>ABBA</b> = <b>A B B A</b>	Yes, spaces are ignored unless escaped (see later)
Is <b>23</b> a regular expression	Yes, will match <b>23</b>
Is <b>&lt;&gt;</b> a regular expression	Yes, will match <b>&lt;&gt;</b>
Will <b>A</b> be recognised by <b>AB</b>	No
Will <b>B</b> be recognised by <b>AB</b>	No
Will <b>AB</b> be recognised by <b>AB</b>	Yes
Will <b>BA</b> be recognised by <b>AB</b>	No
Is <b>AB</b> = <b>BA</b>	No
Will <b>A</b> be recognised by <b>A<math>\epsilon</math></b>	Yes
Will <b>A</b> be recognised by <b><math>\epsilon</math>A</b>	Yes
Can <b>A<math>\epsilon</math></b> generate <b>A</b>	Yes
Can <b><math>\epsilon</math>A</b> generate <b>A</b>	Yes
<b>Regular Expressions - Repetition</b>	
Is <b>A*</b> a regular expression	Yes, will match zero or more <b>A</b> 's
Will <b>A</b> be recognised by <b>A*</b>	Yes
Will <b>AA</b> be recognised by <b>A*</b>	Yes
Will <b>AAA</b> be recognised by <b>A*</b>	Yes
Can <b>A*</b> generate <b>A</b>	Yes
Can <b>A*</b> generate <b>AA</b>	Yes
Can <b>A*</b> generate <b>AAA</b>	Yes
Can <b>A*</b> generate $\epsilon$	Yes
Can <b>A*</b> generate <b>A*</b>	No
Is <b>*</b> a regular expression	No
Will <b>*</b> be recognised by <b>*</b>	No, <b>*</b> is illegal by itself
Will <b>*</b> be recognised by <b>'*'</b>	Yes, quotes are often used to escape special regex characters
Will <b>*</b> be recognised by <b>\*</b>	Yes, some authors/systems use backslash to escape

<p>Will <b>**</b> be recognised by <b>'**'</b></p> <p>Will <b>**</b> be recognised by <b>'*'*'</b></p> <p>Will <b>\</b> be recognised by <b>'\'</b></p>	<p>1 character</p> <p>Yes, matches two asterisks or <b>\*\*</b></p> <p>Yes, matches zero or more asterisks or by <b>\**</b></p> <p>Yes, or by <b>\\</b></p>
<p>Will <b>A</b> be recognised by <b>AA*</b></p> <p>Will <b>A</b> be recognised by <b>A*A</b></p> <p>Can <b>AA*</b> generate <b>A</b></p> <p>Can <b>AA*</b> generate <math>\epsilon</math></p> <p>Can <b>A*A</b> generate <math>\epsilon</math></p> <p>Will <b>AABBB</b> be recognised by <b>A*B*</b></p> <p>Will <b>BBB</b> be recognised by <b>A*B*</b></p> <p>Will <b>A*B*</b> be recognised by <b>A*B*</b></p> <p>Will <b>A*B*</b> be recognised by <b>'A*B*'</b></p>	<p>Yes</p> <p>Yes</p> <p>Yes</p> <p>No</p> <p>No</p> <p>Yes</p> <p>Yes</p> <p>No</p> <p>Yes, or by <b>A\*B\*</b></p>
<b>Regular Expressions – Alternation</b>	
<p>Will <b>A</b> be recognised by <b>A B</b></p> <p>Will <b>B</b> be recognised by <b>A B</b></p> <p>Can <b>A B</b> generate <b>A</b></p> <p>Can <b>A B</b> generate <b>B</b></p> <p>Can <b>A B</b> generate <b>AB</b></p> <p>Is <b>A B</b> = <b>B A</b></p> <p>Will <b>A</b> be recognised by <b>A <math>\epsilon</math></b></p> <p>Will <b>B</b> be recognised by <b><math>\epsilon</math> A</b></p> <p>Can <b>A <math>\epsilon</math></b> generate <b>A</b></p> <p>Can <b>A <math>\epsilon</math></b> generate <math>\epsilon</math></p>	<p>Yes</p> <p>Yes</p> <p>Yes</p> <p>Yes</p> <p>No</p> <p>Yes</p> <p>Yes</p> <p>No</p> <p>Yes</p> <p>Yes</p>
<b>Regular Expressions – Grouping</b>	
<p>Will <b>AB</b> be recognised by <b>A B*</b></p> <p>Will <b>AB</b> be recognised by <b>(A B)*</b></p> <p>Will <b>ABBA</b> be recognised by <b>(A B)*</b></p> <p>Will <b>AABBA</b> be recognised by <b>(A B)*</b></p> <p>Will <b>(A B)*</b> b.r.b. <b>'(A B)*'</b></p>	<p>No</p> <p>Yes</p> <p>Yes</p> <p>Yes</p> <p>Yes, or by <b>\(A B\) \*</b></p>
<b>Regular Expressions – Precedence (lowest)</b>	
Repetition (highest), concatenation, alternation	
<p>Is <b>AB*</b> = <b>(AB)*</b></p> <p>Is <b>AB*</b> = <b>A(B*)</b></p> <p>Is <b>A*B*</b> = <b>(A*B)*</b></p> <p>Is <b>A*B*</b> = <b>(A*)(B*)</b></p> <p>Is <b>A**</b> = <b>(A*)*</b> = <b>A*</b></p>	<p>No</p> <p>Yes</p> <p>No</p> <p>Yes</p> <p>Yes</p>
<p>Is <b>A B*</b> = <b>(A B)*</b></p> <p>Is <b>A B*</b> = <b>A (B*)</b></p>	<p>No</p> <p>Yes</p>
<p>Is <b>AB BA</b> = <b>BA AB</b></p>	<p>Yes</p>

Is $AB^*BA^* = (ABBA)^*$	No
Is $AB^*BA^* = A(B^*)B(A^*)$	Yes
Is $AB^{**} = (AB^*)^*$	No
Is $AB^{**} = A((B^*)^*)$	Yes
<b>Regular Expressions – Shortcuts</b>	
Is $A? = A  \epsilon$	Yes, 0 or 1
Is $A+ = AA^*$	Yes, 1 or more
Is $[ABC] = A B C$	Yes, any symbols in the set
Is $[A-C] = A B C$	Yes, range of symbols
Is $[0-9] = 0 1 2 3 4 5 6 7 8 9$	Yes, range of symbols
Is $[0-2A-C] = 0 1 2 A B C$	Yes, two ranges of symbols
$\square$ denotes any permissible symbol	Any symbol <sup>7</sup> in the alphabet (see below)
$[^A]$ denotes ...	Any permissible symbol except A
$[^AB]$ denotes ...	Any permissible symbol except A and B

## Formal Languages

Formal Language	
The set of all possible strings that can be generated by a regular expression <b>regex</b> is called the <i>formal language</i> of the regular expression and is written $L(\text{regex})$ .	Formal languages (sets of strings) defined by regular expressions are called <i>regular languages</i> .
Is $L(A) = \{A\}$	Yes
$L(\epsilon) = \{\epsilon\}$	Yes
$L(\emptyset) = \{\}$	Yes
$L(A B) = \{A, B\}$	Yes
$L(AB) = \{AB\}$	Yes
$L((A B)C) = \{AC, BC\}$	Yes
$L(A^*) = \{\epsilon, A, AA, AAA, \dots\}$	Yes
Alphabet	
Is $\{A, B, C\}$ an alphabet	Yes, the alphabet is the set of all permitted symbols
Is A a symbol of the alphabet $\{A, B, C\}$	Yes
Is Z a symbol of the alphabet $\{A, B, C\}$	No
Alphabets are typically denoted by $\Sigma$	Yes
Can $\epsilon$ be part of an alphabet	No
For $\Sigma = \{A, B\}$ , is $\Sigma^* = \{\epsilon, A, B, AA, AB, BA, BB, AAA, \dots\}$	Yes, $\Sigma^*$ is the <i>closure</i> of the alphabet $\Sigma$ and is the set of all possible strings of symbols from $\Sigma$ including $\epsilon$
Is $\Sigma^*$ a formal language	Yes, $\Sigma^*$ is the formal language of the alphabet $\Sigma$
For the alphabet $\Sigma = \{A, B\}$ , how many strings in $\Sigma^*$ have length 2, 3 and n	$2^2, 2^3, 2^n$

<sup>7</sup> In most tools, . (dot) stands for any character other an end-of-line character.

Exercises	
What does $(0 1)^+$ recognise	Binary numbers
What does $(0 1)^*0$ recognise	Even binary numbers
What does $[0-9]^+$ recognise	Unsigned decimal numbers
What does $(+ -)[0-9]^+$ recognise	Signed decimal numbers
What does $(+ -)?[0-9]^+$ recognise	Optionally signed decimal numbers
Give a regular expression for:	
(i) a letter	$[a-zA-Z]$
(ii) an identifier	$[a-zA-Z][a-zA-Z0-9_]*$
(iii) the keywords <b>if</b> , <b>while</b> , and <b>do</b>	$if   while   do$
(iv) comment of the form $\{ \dots \}$	$\{ [^}]^* \}$
(v) comment of the form $/* \dots */$	$'/*' [^']* ('*' ([^*/] *)?)* '*/'$
For the alphabet $\Sigma = \{A, B, C\}$ give a regular expression for all strings that:	
(i) contain exactly one <b>B</b>	$(A C)^* B (A C)^*$
(ii) contain at most one <b>B</b>	$(A C)^* B? (A C)^*$
(iii) do not contain 2 consecutive <b>B</b> 's	$(A C BA BC)^* B?$
For the alphabet $\Sigma = \{A, B\}$ give a regular expression for all strings:	
(i) beginning with a <b>B</b> followed by zero or more <b>A</b> 's	$BA^*$
(ii) in which every <b>A</b> is immediately followed by at least two <b>B</b> 's	$(B ABB)^*$
(iii) containing exactly two <b>B</b> 's	$A^* B A^* B A^*$
(iv) containing two consecutive <b>A</b> 's or two consecutive <b>B</b> 's	$(A B)^* (AA BB) (A B)^*$
(v) an even number of <b>A</b> 's and an even number of <b>B</b> 's	$(AA   BB   (AB BA) (AA BB)^* (AB BA) )^*$

## Tokens

- Lexical analysers typically provide a function that returns a token (an object) that corresponds to the next basic syntactic element. There are several common types of token.
- **Identifier tokens** are usually classified into keyword identifiers and non-keyword identifiers.
- **Keyword identifiers** e.g. `class`, `package`, `while` have special meaning in the programming language and are represented by their own token e.g. CLASS token, PACKAGE token, WHILE token. There are small number in most languages, e.g. <50 keywords.
- **Non-keyword identifiers** e.g. `year`, `quad09`. These are defined by the programmer and normally represented by a general identifier token, e.g. IDENT plus a string attribute holding the actual name. For example, `year` could be mapped to the token IDENT("year").
- Lexical analysers need to be able to quickly determine if an identifier is a keyword. A fast technique for this is to use a **perfect hash function**. More generally, lexical analysis is one of the simplest parts of a compiler, a bad-implementation can lead to very poor performance. Profiling a compiler often reveals the compiler can spend more than 50% of its time in the lexical analyser.
- **Literal tokens** (i.e. *constants*) are usually classified into unsigned integers, reals and strings.
- **Unsigned integers** e.g. `123` are represented by an INTEGER token, plus an integer attribute holding the actual value e.g. INTEGER(`123`). What about negative integers?
- **Unsigned reals** e.g. `123.45` are represented by a REAL token for reals, plus a floating-point attribute holding the actual value e.g. REAL(`123.45`).
- **Strings** e.g. `"hello world"` are represented by a STRING token for strings, plus a string attribute holding the actual value e.g. STRING(`"hello world"`)
- We need to take care that the programming language that the compiler is implemented in correctly represents values of the compiled language, e.g. Unicode strings using a compiler that only supports ASCII strings, big integers using a compiler that only supports short integers.
- In addition to identifiers and literals, there are several other tokens types:
- **Operators and other punctuation** e.g. `+`, `:`, `<=`, `(` are normally represented by their own token e.g. PLUS, COLON, LESSEQUAL, LPAREN
- **Whitespace** e.g. spaces, tabs, newlines are *not* normally returned as tokens (unless they occur in a string literal). Whitespace is normally needed to separate adjacent identifiers, literals e.g. in `public static void`. They can be represented by regular expressions however.
- **Comments** are also not normally returned as tokens.
- What extra information would be useful for the lexical analyser to record in returned tokens? Hint: think about for reporting errors later in the compilation or at runtime.

## Disambiguation Rules

- If a string can be matched by more than one regular expression we need rules to resolve the ambiguity. A workable set of disambiguation rules is:
  - (1) match the longest matching string, otherwise
  - (2) assume that the list of regular expressions defined for the programming language have textual precedence such that the earliest written regular expression takes precedence.
- For example, if we would write the regular expression for keywords before the regular expression for identifiers, the string `doughnut` would return an IDENTIFIER token, while the string `do` would return a DO token.

## Context Free Grammars

- Regular expressions are not sufficiently expressive to define the more complex syntactic structures of programming languages, for example, nested constructs. For this, we need grammars.
- There are many classes of grammars. An important classification is due to Chomsky, who identified 4 classes: **regular grammars**, **context free grammars**, **context-sensitive grammars** and **unrestricted grammars** - succeeding grammar is a superset of the preceding grammar. We'll consider the second of these, since the syntax of many programming languages can be described by a context-free grammar.
- In formal languages, a **context-free grammar**  $G$  consists of:
  - (1) a set of **terminal symbols**. When defining programming languages, terminal symbols correspond to the tokens returned by the lexical analyser, that is, they correspond to basic syntactical elements, like identifiers, keywords, numbers, operators and other punctuation characters.
  - (2) a set of **non-terminal symbols** (disjoint from the set of terminal symbols). Non-terminal symbols correspond to grammar rules<sup>8</sup>. This set includes a **start symbol** for the grammar, which is either explicitly identified or corresponds to the non-terminal symbol on the left-hand side of the first rule. We use the latter approach for the start symbol.
  - (3) a set of **rules** of the form:  $N \rightarrow \alpha$  where  $N$  is a non-terminal symbol, and  $\alpha$  is a string of terminal and non-terminal symbols.
- The formal language corresponding to a grammar is the set of all strings<sup>9</sup> of terminals that can be produced from the start symbol using the rules of the grammar.
- In these notes, we will write terminal symbols either in (1) lowercase italics, e.g. *id*, *while*, *num* or (2) or in quotes for operators and other punctuation characters/character sequences, e.g. '+', '<=', or (3)  $\epsilon$  for the empty string.
- Non-terminal symbols will be written as 1st-letter uppercase identifiers, e.g. *Program*, *Expression*

```
Program    → Block
Block      → Statement Block
Block      →  $\epsilon$ 
Statement  → id '=' Expression
Statement  → read id
Statement  → write Expression
Expression → Expression Operator Expression
Expression → '(' Expression ')'
Expression → num
Operator   → '+'
Operator   → '-'
Operator   → '*'
Operator   → '/'
```

- In the example above, the terminals are *id*, *num*, *read*, *write*, +, \* Non-terminals are *Program*, *Block*, *Statement*, *Expression* and *Op*. The start symbol is *Program*.
- For conciseness, vertical bars can be used to group alternative rules<sup>10</sup> for the same non-term terminal, e.g.

```
Program    → Block
Block      → Statement Block |  $\epsilon$ 
Statement  → id '=' Expression | read id | write Expression
Expression → Expression Operator Expression | '(' Expression ')' | num
Operator   → '+' | '-' | '*' | '/'
```

<sup>8</sup> Rules are also known as **productions**.

<sup>9</sup> Some textbooks use the term **sentence** instead of string.

<sup>10</sup> Note: a rule with  $N$  alternatives are really  $N$  rules not one.

- The rule  $\text{Expression} \rightarrow \text{Expression Operator } \underline{\text{num}}$  is *left-recursive*.
- The rule  $\text{Expression} \rightarrow \underline{\text{num}} \text{ Operator Expression}$  is *right-recursive*.

## Derivations

- A **derivation** is a sequence of substitutions of rules occurring on the right hand side of grammar. One substitution in each derivation step ( $\Rightarrow$ ). Example:  

$$\begin{aligned} \text{Expression} &\Rightarrow \text{Expression Operator Expression}^{11} \\ &\Rightarrow \text{Expression '+' Expression} \\ &\Rightarrow \underline{\text{num}} \text{ '+' Expression} \\ &\Rightarrow \underline{\text{num}} \text{ '+' } \underline{\text{num}} \end{aligned}$$
- Note: each num token in the derivation would include the actual value of the number in the input, e.g. if the input was 45+35, the first num token would be num(45) while the second num token would be num(35). We will sometimes just use 45 or 35 to emphasise concrete terminals.
- The **formal language** corresponding to a grammar is the set of all strings of terminals that can be derived from the start symbol using derivations i.e.  $\{ \text{string} \mid \text{startSymbol} \Rightarrow^* \text{string} \}$  where  $\Rightarrow^*$  is a derivation.
- For the grammar corresponding to the start symbol **Program** above, each string in the formal language corresponds to a **syntactically legal program**.

Grammar	Formal Language
$E \rightarrow E 'a' \mid 'a'$	$\{ a^n, n \geq 1 \}$
$E \rightarrow '(' E ')' \mid 'a'$	$\{ (a)^n, n \geq 0 \}$
$E \rightarrow '(' E )$	$\{ \}$ - no strings can be derived, no base case

- There can be more than one derivation for the same terminal string. For example for num '+' num there are 5 other derivations including the following two:
- If at each step, the leftmost rule is substituted, the derivation is called a **leftmost-derivation**, e.g.:  

$$\begin{aligned} \text{Expression} &\Rightarrow \text{Expression Operator Expression} \\ &\Rightarrow \underline{\text{num}} \text{ Operator Expression} \\ &\Rightarrow \underline{\text{num}} \text{ '+' Expression} \\ &\Rightarrow \underline{\text{num}} \text{ '+' } \underline{\text{num}} \end{aligned}$$
- If at each step, the rightmost rule is substituted, the derivation is called a **rightmost-derivation**, e.g.:  

$$\begin{aligned} \text{Expression} &\Rightarrow \text{Expression Operator Expression} \\ &\Rightarrow \text{Expression Operator } \underline{\text{num}} \\ &\Rightarrow \text{Expression '+' } \underline{\text{num}} \end{aligned}$$

Derivations	Quotes around 1 and 9 omitted for clarity
For $E \rightarrow E '9' \mid '1'$ give derivations for: 1 19 199 91	$E \Rightarrow 1$ $E \Rightarrow E 9 \Rightarrow 1 9$ $E \Rightarrow E 9 \Rightarrow E 9 9 \Rightarrow E 1 9 9$ not possible
For $E \rightarrow '9' E \mid '1'$ give derivations for: 1 19 199 91	$E \Rightarrow 1$ not possible not possible $E \Rightarrow 9 E \Rightarrow 9 1$

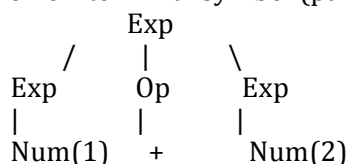
<sup>11</sup> Each right-hand-side is called a *sentential form*. The final sentential form of terminals is called the *yield* of the derivation.



For $E \rightarrow '1' E \mid '9'$ give derivations for: 1 19 199 91	not possible $E \Rightarrow 1 E \Rightarrow 1 9$ not possible not possible
For $E \rightarrow E '1' \mid '9'$ give derivations for: 1 19 199 91	not possible not possible not possible $E \Rightarrow E 1 \Rightarrow 9 1$

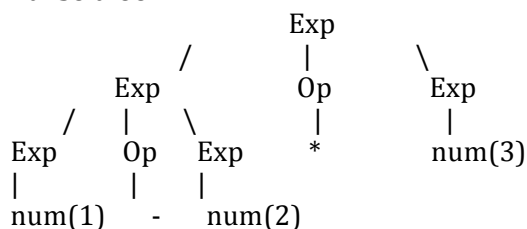
## Parse Trees and Ambiguous Grammars

- The **parse tree** for a derivation is constructed by connecting each symbol (child node) in a derivation to the non-terminal symbol (parent node) from which it was derived. For the 1+2 above, we have:

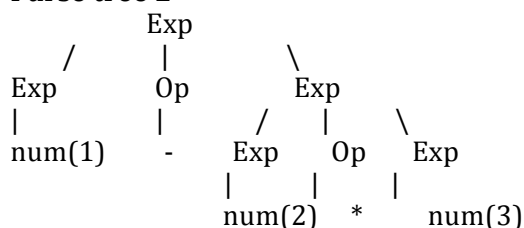


- A context free grammar is **ambiguous** if for at least one string in the language there are two possible derivations that correspond to distinct parse trees. For example, for the string 1-2\*3, has two different parse trees:

**Parse tree 1**



**Parse tree 2**



- An alternative way of stating this, is that a grammar is ambiguous if there is exists a string in the language that has multiple left-most derivations (or multiple right-most derivations).
- Ambiguity is a property of grammars not formal languages.

Ambiguous Grammars	
For each of the following show that the grammar is ambiguous by giving two leftmost derivations for the same string: Quotes around terminal symbols are omitted for clarity.	
$E \rightarrow '9'$ $F \rightarrow '9'$	$E \Rightarrow 9$ $F \Rightarrow 9$
$E \rightarrow E \mid \varepsilon$	$E \Rightarrow \varepsilon$

	$E \Rightarrow E \Rightarrow \epsilon$
$E \rightarrow E '-' E \mid '9'$	$E \Rightarrow E-E \Rightarrow 9-E \Rightarrow 9-E-E \Rightarrow 9-9-E \Rightarrow 9-9-9$ $E \Rightarrow E-E \Rightarrow E-E-E \Rightarrow 9-E-E \Rightarrow 9-9-E \Rightarrow 9-9-9$
$E \rightarrow E E \mid '<' E '>' \mid \epsilon$ for the string $<<<>> <<>>$	$E \Rightarrow EE \Rightarrow <EE>E \Rightarrow <\epsilon E>E \Rightarrow <\epsilon \epsilon>E \Rightarrow <\epsilon \epsilon>\epsilon$ $E \Rightarrow <EE> \Rightarrow <E> \Rightarrow <>$

- **Why is ambiguity undesirable?**

- Although the language for an ambiguous grammar is completely determined, the different parse trees/derivations might have different semantics.
- For example, should  $1-2*3$  mean  $(1-2)*3=-3$  or  $1-(2*3)=-5$ ? Similarly, do we want  $1-2-3$  to mean  $(1-2)-3=-4$  or  $1-(2-3)=0$ ?
- A nice practical solution, when it's available, is to have disambiguating rules. For example, we could handle the first ambiguity with a disambiguating rule that defines the *precedence* of operators, while the second ambiguity could be handled by a disambiguating rule that defines the *associativity* of operators.

- An alternative approach is to rewrite the grammar to remove the ambiguity. For example, by writing extra rules, **Expression** can be rewritten as:

```

Expression → Expression AddOp Term | Term
AddOp      → '+' | '-'
Term       → Term MulOp Factor | Factor
MulOp      → '*' | '/'
Factor     → '(' Expression ')' | num

```

- This formulation is unambiguous.
- Addition and subtraction are grouped in Expression and closer to the root of the parse tree so will be of lower precedence than multiplication and division.
- Using **Expression AddOp Term** makes addition and subtraction *left-associative* (**Term AddOp Expression** would make them *right-associative*). Similarly for multiplication and division.
- For expressions, which approach is better? Although the second approach is purer in the sense that no disambiguating rules are required, it leads to grammars that are harder to understand (for humans) and more complex parse trees (for the compiler writer to deal with). So *a system that supports and resolves some ambiguous grammars is desirable*.
- Another classic example of an ambiguous grammar is the **dangling-else** found in many programming languages:

```

IfStatement → if Expression then Statement
IfStatement → if Expression then Statement else Statement

```

- For **if A then if B then X else Y** this can generate 2 distinct derivations where **Y** binds either with **A** or with **B**. The different derivations have different meaning and would generate different code..

- We can remove this ambiguity in favour of **Y** associating with **B** by re-writing the grammar as:

```

IfStatement → if Expression then Statement
IfStatement → if Expression then WithElse else Statement
WithElse    → if Expression then WithElse else WithElse

```

- Now each **else** binds with the nearest preceding unmatched **if**.
- An alternative, quite nice approach is to change the syntax to the **if-then-elseif-endif** form.
- **Note1:** if ambiguity doesn't affect the semantics then there is no need to remove the ambiguity. For example for **Seq**  $\rightarrow$  **Seq** **;** **Seq** | **Statement** and **Statement**  $\rightarrow$  **id** it doesn't matter what the parse tree is.
- **Note2:** there is no algorithm for removing ambiguity in context-free grammars. In fact there is no algorithm for determining whether a grammar is ambiguous!
- **Note3:** certain context-free formal languages are **inherently ambiguous** - every grammar for the language is ambiguous. For example the formal language  $\{0^a 1^b 2^c \mid a=b \text{ or } b=c\}$  is inherently ambiguous.

## Extended Backus-Naur Form (EBNF)

- As we have seen, the rules for a context free grammar take the following general form:  $N \rightarrow \alpha$  where  $N$  is a non-terminal and  $\alpha$  is a possibly empty sequence of terminals and non-terminals. Note:  $N$  can be defined more than once.
- In the computer science community this basic form is often referred to as **Backus-Naur**<sup>12</sup> Form (BNF) after two of the early pioneers of programming language design.
- Because Repetitive and Optional constructs are very common in programming languages many authors/systems use an Extended form of BNF for writing context free grammars - EBNF. EBNF doesn't make the grammar any more powerful however.

- Repetition:  $\{\alpha\}$       0 or more occurrences of  $\alpha$  §  
Optional:     $[\alpha]$       0 or 1 occurrence of  $\alpha$   
Choice:       $\alpha|\beta$        $\alpha$  or  $\beta$ .  
Grouping:     $(\alpha)$        $\alpha$ . Useful for grouping alternatives within a rule e.g.  $(\alpha|\beta)$
- Examples:

BNF:	$Seq \rightarrow Stmt \text{ ';' } Seq$	EBNF:	$Seq \rightarrow Stmt \{ \text{ ';' } Stmt \}$	left recursive or
	$Seq \rightarrow Stmt$		$Seq \rightarrow \{ Stmt \text{ ';' } \} Stmt$	right-recursive
BNF:	$Exp \rightarrow Exp \text{ AddOp } Term$	EBNF:	$Exp \rightarrow Term \{ \text{ AddOp } Term \}$	left associative or
	$Exp \rightarrow Term$		$Exp \rightarrow Term [ \text{ AddOp } Exp ]$	right-associative

- For the expression example above, an EBNF version is:

```
Expression → Term { ('+' | '-' ) Term }
Term       → Factor { ('*' | '/' ) Factor }
Factor     → '(' Expression ')' | num
```

## Efficiency of Parsers for Context Free Grammars

- Given that most programming languages can be described by a context-free grammar, why are they not implemented in practice?
- The primary reason is that the **worst-case time complexity** to parse using a context-free grammar is  $O(n^3)$ . This is unacceptable for all but the smallest of programs.
- Fortunately, there are subsets of context-free grammars that can be parsed in  $O(n)$  time and are able to describe the syntax of most programming languages. The two subsets that we will consider are **LL(1)** grammars and **LR(1)**<sup>13</sup> grammars. LR(1) grammars are a superset of LL(1) grammars.
- The main issue to resolve is to ensure that the context-free grammar is amenable to LL(1) or LR(1) parsing

### LL(1) Grammars

- A context-free grammar where it is possible to determine which rule to use in a leftmost-derivation by scanning the input left-to-right (beginning-to-end) and by looking at the next token only is called LL(1). LL( $k$ ) parsers have  $k$ -token lookahead.
- LL parsers construct a parse-tree from the root (from the start symbol) down to the leaf nodes (the terminal symbols), and know at each step what to do next.
- There are two main types of LL parser: recursive descent, and LL parse tables / pushdown automaton. We'll only consider recursive descent parsers.

<sup>12</sup> John Backus and Peter Naur who used BNF to define the syntax of Algol 60. John Backus also designed Fortran.

<sup>13</sup> and its more common variant: **LALR(1)**.

- Recursive descent parsers can be written by hand or generated by a parser-generator tool like **ANTLR**.
- Like ambiguity removal for a context-free grammar, there is no automatic method of ensuring that a context-free grammar is an LL(1) grammar. We need to rewrite the grammar if any LL(1) issues arise.

## Recursive Descent Parser

- We'll first outline a method to translate a set of EBNF rules into code.
  - Our recursive descent parsers will consist of:
  - a set of *parse functions*, one for each non-terminal **A** e.g. **A()** Note: the parse function could also return an AST for the rule.
  - the *current input token* e.g. global variable **token**
  - a *token match and advance to next token* function, e.g.

```
def match(expected):
    if token==expected: token=lex.get_token()
    else error("unexpected token encountered...")
```

- A non-terminal **N** will be translated to **N()**
- A terminal **T** will be translated to **match(T)**
- Concatenation **A B** becomes

```
A(); B();
```

- Alternation **A | B** becomes

```
if token in FIRST(A):    A()
elif token in FIRST(B):  B()
```

- Repetition **{A}** becomes

```
while token in FIRST(A):    A()
```

- Optional **[A]** becomes

```
if token in FIRST(A): A()
```

- For the translations above, use **match(A)** instead of **A()** if **A** is a terminal.
- **FIRST(A)** corresponds to the set of all tokens that could possibly start an **A**. If A is empty, then include tokens that could follow A also.
- **FOLLOW(A)** corresponds to the set of all tokens that could possibly follow an **A** anywhere in the grammar.
- For alternation, if the FIRST sets of the alternatives are not disjoint, e.g. if there is a token in both FIRST(A) and FIRST(B) then the grammar is not LL(1).

- In our code, token types will be given in uppercase, for example **IF** for the if token.

Example:

```
Statement      → IfStatement | BeginStatement | PrintStatement
IfStatement    → if Expression then Statement [else Statement]
BeginStatement → begin Statement { ';' Statement } end
PrintStatement → print Expression
```

```
def Statement():
    if token == IF:      IfStatement()
    elif token == BEGIN: BeginStatement()
    elif token == PRINT: PrintStatement()

def IfStatement():
    match(IF)
    Expression()
    match(THEN)
    Statement()
    if token == ELSE:
        match(ELSE)
        Statement()
    else error()

def BeginStatement():
    match(BEGIN)
    Statement()
    while token == SEMICOLON:
        match(SEMICOLON)
        Statement()
    match(END)

def PrintStatement():
    match(PRINT)
    Expression()
```

•

### • Parsing and Building an Abstract Syntax Tree

- The example above will recognise a legal statement but will not produce a parse tree nor an abstract syntax tree that captures the syntax of the program for use by the semantic analysis or code generation phases of the compiler.
- We can adapt our approach by returning an AST node from each parse function.
- We do this by defining the following classes:

```
class StatementAST extends AST:
    Position p # useful for error handling

class IfAST extends StatementAST:
    ExprAST e
    StatementAST thenS, elseS

class BeginAST extends StatementAST:
    StatementAST statlist[] # list

class PrintAST extends StatementAST:
    ExprAST e
```

- We can now implement the parser for our example as follows:

```
def Statement():
```

```

if token == IF:    return IfStatement()
elif token == BEGIN: return BeginStatement()
elif token == PRINT: return PrintStatement()
else error()

def IfStatement():
    match(IF)
    e = Expression()
    match(THEN)
    thenS = Statement()
    if token == ELSE:
        match(ELSE)
        elseS = Statement()
    else:
        elseS = None
    return IfAST(e,thenS,elseS)

def BeginStatement():
    match(BEGIN)
    statlist = []
    statlist.append(Statement())
    while token == SEMICOLON:
        match(SEMICOLON)
        statlist.append(Statement())
    match(END)
    return BeginAST(statlist)

def PrintStatement():
    match(PRINT)
    return PrintAST(Expression())

```

## Expression calculator

- The approach to translating EBNF into parsing functions can also be used for other tasks. For example for calculating expressions with correct precedence and associativity we could use:

```

def Expression:
    result = Term()
    while token in {PLUS, MINUS}:
        op=token; match(op); right=Term()
        if op==PLUS: result += right else result -= right
    match(token);
    return result

def Term:
    result = Factor()
    while token in {MULT, DIVIDE}:
        op=token; match(op); right=Factor()
        if op==MULT: result *= right
        else if right !=0: result /= right
        else error("Divide by zero")
    match(token);
    return result

def Factor:
    if token == LPAREN:
        match(LPAREN); result=Expression(); match(RPAREN)
    elif token == NUM: result=token.value; match(token);
    else error("(" or number expected.")
    return result

```

## CFG to LL(1)

- In order to employ top-down parsing we need to ensure that our grammar is LL(1). Unfortunately transforming a non-LL(1) context-free grammar to LL(1) cannot be fully automated and often requires care and a little ingenuity. In particular we need to ensure that the transformed grammar has the same semantics as the original grammar and if possible retains the readability of the original.
- The following 3 transformations are commonly tried: (1) Left Factorisation, (2) Substitution, (3) Left Recursion Removal. Left recursion removal and substitution are normally carried out first then left factorisation.

- (1) **Left Factorisation**. If two or more alternatives of a rule share a common prefix we can factor out the common prefix using Left Factorisation.

In EBNF:  $A \rightarrow B C \mid B D$  can be changed to  $A \rightarrow B (C \mid D)$

$A \rightarrow B C \mid B$  can be changed to  $A \rightarrow B [C]$

In BNF:  $A \rightarrow B C \mid B D$  can be changed to  $A \rightarrow B X$  and  $X \rightarrow C \mid D$

$A \rightarrow B C \mid B$  can be changed to  $A \rightarrow B X$  and  $X \rightarrow C \mid \epsilon$

If we don't do this the parser won't know whether to choose the first alternative or the second.

- (2) **Substitution** involves replacing a non-terminal  $N$  on the right-hand side of a rule by each of the alternatives of  $N$ . Substitution is useful when conflicts are "indirect" in order to make the conflict "direct" and hopefully amenable to left factoring. Substitution is also used to write clearer grammars.
- Consider the following example:

```
ForStatement → for ControlVar ':=' Expr Direction Expr do Statement
ControlVar   → id
Direction    → to | downto
```

- Although this grammar is LL(1) it seems sensible to remove both **ControlVar** and **Direction** using substitution. The **ControlVar** rule seems to be defined purely as commentary - it serves no grammatical role otherwise. **Direction** is defined as a rule because the above grammar is given in BNF which doesn't support grouping of alternatives. By using substitution and EBNF we obtain a simpler rule:

```
ForStatement → for id ':=' Expr (to | downto) Expr do Statement
```

- Now consider the following example:

```
Statement → Assignment | Call | pass
Assignment → id ':=' Expr
Call       → id '(' Expr ')
```

- The grammar is not LL(1) because id can start both **Assignment** and **Call** and these are alternatives of **Statement**. In addition the grammar isn't in a form suitable for left factorisation. We can however bypass these issues by substituting **Assignment** and **Call** in **Statement** and then left factorising **Statement**:

```
Statement → id (':=' Expr | '(' Expr ')) | pass
```

- Although this form of the grammar is LL(1) it merges/obscures the clean separation in the original grammar. An LL(1) parser will need to handle this by re-introducing ASTnodes and parse functions for the alternatives and passing id to them in order to preserve the original form in later passes.

- (3) **Left Recursion Removal.** For example  $\text{Expr} \rightarrow \text{Expr Op Expr}$  cannot be handled by an LL parser because we would recurse forever as soon as we enter the parse function. We need to eliminate left-recursion.
- We'll cover the common case of direct left recursion removal, which leaves the grammar largely intact. To eliminate a direct left recursive rule we re-write the rule using right-repetition (in BNF we can do the rewrite using a new auxiliary rule and right-recursion) – ANTLR 4 does this for you! Handling of other indirect recursive cases is more involved and often results in incomprehensible LL(1) grammars.
- In EBNF:  $A \rightarrow X \mid A Y$  can be changed to  $A \rightarrow X \{Y\}$   
More generally we have:  

$$A \rightarrow X_1 \mid X_2 \mid \dots \mid X_N \mid A Y_1 \mid A Y_2 \mid \dots \mid A Y_M \text{ can be changed to}$$

$$A \rightarrow (X_1 \mid X_2 \mid \dots \mid X_N) \{Y_1 \mid Y_2 \mid \dots \mid Y_M\}$$
- **LL(2) Parser.** When developing an LL(1) grammar and a hand-written parser, a practical solution to developing clean abstract syntax trees and nice code is to just peek-ahead and look at another token. This is equivalent to an LL(2) parser. **ANTLR4** is capable of emitting LL(k) parsers.

### A word about handling compilation errors

- The response of a compiler to errors is an important factor in its usefulness.
- Some compilers stop after producing a single error message.
- Most attempt some form of recovery and are capable of detecting further errors.
- Some compilers even attempt to correct syntactic errors in order to continue with semantic analysis and even code generation.
- The following are good guidelines:
  - Produce informative messages for the user.
  - Report the position of an error as accurately as possible.
  - If the compiler is part of an Integrated Development Environment (IDE), highlight the region of each error in an editor window.
  - If error recovery is attempted, skip as little as possible in order to parse as much of the remaining code as possible, avoid infinite looping in the recovery process, and avoid spurious error messages being generated as a result of the recovery.
  - If error correction is attempted, ensure that the corrected program has the same syntax tree and semantics as the original (assuming we know what that is!).



## LR(1) Grammars

- The second class of practical (i.e. expressive and fast) parsers are those for LR(1) grammars.
- LR(1) grammars are a superset of LL(1) grammars and allow us to write clearer grammars, for example, they allow left recursion in rules.
- Unlike LL(1) parsers, LR(1) parsers are too complex to write by hand. We need to use a LR parser-generator tool like **bison** (GNU version of yacc), Sable CC, CUP, Beaver etc
- Most LR parser generators support a subset of LR(1), called **LALR(1)** which is more memory efficient, but very close in expressiveness to LR(1). We will assume that LR tools support LR(1) grammars even though many actually support LALR(1) grammars. We won't look at how LR parser generators are implemented <sup>14</sup>.
- LR(1) builds the parse tree bottom up, from the leaves (non-terminals) to the root (start symbol).
- They achieve this by **reducing** child nodes corresponding to the RHS (**handle**) of a rule to a parent node corresponding to the left-hand side of the rule. A stack is used to maintain the current state of the parse.
- For example, given the input **9+4** and the grammar

```

Expression → Expression '+' Term | Term
Term       → Term '*' Factor | Factor
Factor     → '(' Expression ')' | num

```

- an LR parser will proceed to build the parse tree as follows

	Stack Action	Stack	Tree Action	Parse Tree
1	<b>Shift</b> <sup>15</sup> 9 onto the stack.	9	Create a leaf node for 9	9
2	Reduce 9 i.e. recognise that 9 is a <i>handle</i> (corresponds to a RHS for a rule) and replace 9 by Factor on stack.	Factor	Create a Factor parent node for 9 node	F   9
3	Reduce Factor, i.e. recognise that <b>Factor</b> is a handle and replace Factor by Term on stack	Term	Create a parent <b>Term</b> node for Factor	T   F   9
4	Shift <sup>16</sup> + onto the stack and a	Term+	Create a leaf node for +	T   +   F   9
5	Shift 4 onto the stack	Term+4	Create a leaf node for 4	T   +   4   F   9
6	Reduce 4 to Factor	Term+Factor	Create a parent Factor	T   +   F 

<sup>14</sup> Interested students are highly encouraged to study the implementation of LR parsers in textbooks and also last year's slides.

<sup>15</sup> Traditional LR parsing term for *push*.

<sup>16</sup> Why not reduce Term to Factor?

			node for 4.	$\begin{array}{c} F \quad 4 \\   \\ 9 \end{array}$
7	Reduce <sup>17</sup> Term+Factor to Term	Term	Create a parent Term node for Term, +, Factor child nodes	$\begin{array}{ccccc} & & T & & \\ & / &   & \backslash & \\ T & & + & & F \\   & & & &   \\ F & & & & 4 \\   & & & & \\ 9 & & & & \end{array}$
8	Reduce Term to Expression	Expression	Create a parent Expression node for Term child node	$\begin{array}{ccccc} & & E & & \\ & &   & & \\ & & T & & \\ & / &   & \backslash & \\ T & & + & & F \\   & & & &   \\ F & & & & 4 \\   & & & & \\ 9 & & & & \end{array}$

---

<sup>17</sup> Again why not reduce Factor to Term?