

Grammars and parsing* with C# 2.0

Peter Sestoft^{†‡} and Ken Friis Larsen[‡]
sestoft@dina.kvl.dk ken@friislarsen.net

2006-08-31

*Based on earlier versions for Standard ML and for Java.

[†]Department of Mathematics and Physics, Royal Veterinary and Agricultural University, Denmark.

[‡]IT University of Copenhagen, Denmark.

Contents

1	Grammars and Parsing	3
2	Grammars	4
2.1	Grammar notation	4
2.2	Derivation	4
3	Parsing theory	6
3.1	Parsing: reconstruction of a derivation tree	6
3.2	A more machine-oriented view of parsing	8
3.3	Left factorization	9
3.4	Left recursive nonterminals	10
3.5	First-sets, follow-sets and selection sets	11
3.6	Summary of parsing theory	15
4	Parser construction in C#	16
4.1	C# representation of input strings	16
4.2	Constructing parsing methods in C#	17
4.3	Parsing methods follow the derivation tree	20
4.4	Summary of parser construction	20
5	Scanners	21
5.1	Character classification methods	21
5.2	Scanning names	22
5.3	Distinguishing names from keywords	23
5.4	Scanning floating-point numerals	24
5.5	Reading string and text files with C#	26
5.6	Summary of scanners	26
6	Parsers with attributes	27
6.1	Constructing attributed parsers	27
6.2	Building representations of input	30
6.3	Summary of parsers with attributes	33
7	A larger example: arithmetic expressions	34
7.1	A grammar for arithmetic expressions	34
7.2	The parser constructed from the grammar	35
7.3	A scanner for arithmetic expressions	36
7.4	Evaluating arithmetic expressions	37
8	Some background	38
8.1	History and notation	38
8.2	Extended Backus-Naur Form	38
8.3	Classes of languages	39
8.4	Further reading	39
9	Exercises	40
	References	43

1 Grammars and Parsing

Often the input to a program is given as a text, but internally in the program it is better represented more abstractly: by a C# object, for instance. The program must read the input text, check that it is well-formed, and convert it to the internal form. This is particularly challenging when the input is in ‘free format’, with no restrictions on the layout.

For example, think of a program for doing symbolic differentiation of mathematical expressions. It needs to read an expression involving arithmetic operators, variables, parentheses, etc. It must check that the parentheses match, it should allow any number of blanks around operators, and so on, and must build a suitable internal representation of the expression. Doing this without a systematic approach is very hard.

Example 1 This text file describes the probable states of a slightly defective gas gauge in a car, given the state of the car’s battery and its gas tank:

```
probability(GasGauge | BatteryPower, Gas)
{
    (0, 0): 100.0, 0.0;
    (0, 1): 100.0, 0.0;
    (1, 0): 100.0, 0.0;
    (1, 1): 0.1, 99.9;
}
```

These lecture notes explain how to create programs that can read an input text file such as the above, check that its format is correct, and build an internal representation (an array or a list) of the data in the input file. Here we shall not be concerned with the meaning¹ of these data. □

Thus we provide simple tools to perform these tasks:

- systematic *description* of the structure of input data, and
- systematic *construction* of programs for reading and checking the input, and for converting it to internal form.

The input descriptions are called *grammars*, and the programs for reading input are called *parsers*. We explain grammars and the construction of parsers in C#. The methods shown here are essentially independent of C#, and can be used with suitable modifications in any language that has recursive procedures (Java, Ada, C, ML, Modula, Pascal, etc.)

The order of presentation is as follows. First we introduce grammars, then we explain parsing, formulate some requirements on grammars, and show how to construct a parser skeleton from a grammar which satisfies the requirements. These parsers usually read sequences of symbols instead of raw texts. So-called *scanners* are introduced to convert texts to symbol sequences. Then we show how to extend the parsers to build an internal representation of the input while reading and checking it.

Throughout we illustrate the techniques using a *very* simple language of arithmetic expressions. At the end of the notes, we apply the techniques to parse and evaluate more realistic arithmetic expressions, such as $3.1*(7.6-9.6/-3.2)+(2.0)$.

When reading these notes, keep in mind that although it may look ‘theoretical’ at places, the goal is to provide a *practically* useful tool.

¹The lines (0, 0) and (0, 1) say that if the battery is completely uncharged (0) and the tank is empty (0) or non-empty (1), then the meter will indicate Empty with probability 100%. The line (1, 0) says that if the battery is charged (1) and the tank is empty (0), then the gas gauge will indicate Empty with probability 100% also. Finally, the line (1, 1) says that even when the battery is charged (1) and the tank is non-empty (1), the gas gauge will (erroneously) indicate Empty with probability 0.1% and Nonempty with probability 99.9%.

2 Grammars

2.1 Grammar notation

A *grammar* G is a set of rules for combining symbols to a well-formed text. The symbols that can appear in a text are called *terminal symbols*. The combinations of terminal symbols are described using *grammar rules* and *nonterminal symbols*. Nonterminal symbols cannot appear in the final texts; their only role is to help generating texts: strings of terminal symbols.

A *grammar rule* has the form $A = f_1 \mid \dots \mid f_n$ where the A on the left hand side is the nonterminal symbol defined by the rule, and the f_i on the right hand side show the legal ways of deriving a text from the nonterminal A .

Each *alternative* f is a *sequence* $e_1 \dots e_m$ of symbols. We write Λ for the empty sequence (that is, when $m = 0$).

A *symbol* is either a *nonterminal symbol* A defined by some grammar rule, or a *terminal symbol* " c " which stands for c .

The *starting symbol* S is one of the nonterminal symbols. The well-formed texts are precisely those derivable from the starting symbols.

The grammar notation is summarized in Figure 1.

A *grammar* $G = (T, N, R, S)$ has a set T of terminals, a set N of nonterminals, a set R of rules, and a starting symbol $S \in N$.

A *rule* has form $A = f_1 \mid \dots \mid f_n$, where $A \in N$ is a nonterminal, each alternative f_i is a sequence, and $n \geq 1$.

A *sequence* has form $e_1 \dots e_m$, where each e_j is a symbol in $T \cup N$, and $m \geq 0$. When $m = 0$, the sequence is empty and is written Λ .

Figure 1: Grammar notation

Example 2 Simple arithmetic expressions of arbitrary length built from the subtraction operator '-' and the numerals 0 and 1 can be described by the following grammar:

$$\begin{aligned} E &= T \text{ "-" } E \mid T . \\ T &= \text{"0"} \mid \text{"1"} . \end{aligned}$$

The grammar has terminal symbols $T = \{\text{"-"}, \text{"0"}, \text{"1"}\}$, nonterminal symbols $N = \{E, T\}$, two rules in R with two alternatives each, and starting symbol E . Usually the starting symbol is listed first. □

2.2 Derivation

The grammar rule $T = \text{"0"} \mid \text{"1"}$ above says that we may derive either the string "0" or the string "1" from the nonterminal T , by replacing or substituting either "0" or "1" for T . These *derivations* are written $T \Rightarrow \text{"0"}$ and $T \Rightarrow \text{"1"}$.

Similarly, from nonterminal E we can derive T , for instance. From T we could derive "0", for example, which shows that from E we can derive "0", written $E \Rightarrow \text{"0"}$.

Choosing the other alternative for E we might get the derivation

$$\begin{aligned} E &\Rightarrow T \text{ "-" } E \\ &\Rightarrow \text{"0"} \text{ "-" } E \\ &\Rightarrow \text{"0"} \text{ "-" } T \\ &\Rightarrow \text{"0"} \text{ "-" } \text{"1"} \end{aligned}$$

In each step of a derivation we replace a nonterminal with one of the alternatives on the right hand side of its rule. A derivation can be shown as a tree; see Figure 2.

Every internal node in the tree is labelled by a nonterminal, such as E. The sequence of children of an internal node, such as T, "-", E, represents an alternative from the corresponding grammar rule.

A leaf of the tree is labelled by a terminal symbol, such as "-". Taking the leaves in sequence from left to right gives the string derived from the symbol at the root of the tree: "0" "-" "1".

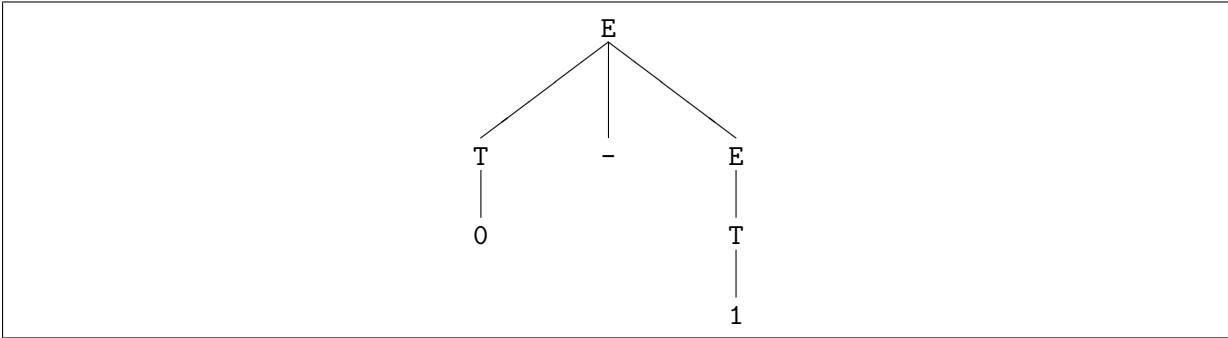


Figure 2: A derivation tree

One can think of a grammar G as a generator of strings of terminal symbols. Let T^* be the set of all strings of symbols from T , including the empty string Λ . When A is a nonterminal, the set of strings derivable from A is called $L(A)$:

$$L(A) = \{ w \in T^* \mid A \implies w \}$$

When grammar G has starting symbol S , the language generated by G is $L(G) = L(S)$. Grammars are useful because they are finite and compact descriptions of usually infinite languages.

In the example above we have $L(E) = \{0, 1, 0-0, 0-1, 1-0, 1-1, 0-0-0, \dots\}$, namely, the set of well-formed texts according to the grammar. As shown here, the quotes around strings of terminals are often left out.

Example 3 In mathematics, a rather liberal notation is used for writing down polynomials in x , such as $x^3 - 2x^2$. The following grammar describes such polynomials:

- Poly = Term
 - | Plusminus Term
 - | Poly Plusminus Term .
- Term = Natnum "x" Exponent
 - | Natnum
 - | "x" Exponent .
- Exponent = "^" Natnum
 - | Λ .
- Plusminus = "+" | "-" .

Assume that Natnum stands for any natural number 0, 1, 2, ...

Check that the following strings are derivable: "0", "-0", "2x + 5", "x^3 - 2x^2", and that the following strings are not derivable: "2xx", "+-1", "5 7", "x^3 + - 2x^2". \square

3 Parsing theory

We have seen that a grammar can be used to derive symbol strings having a certain structure. When a program reads an input file, the problem is the opposite: given an input text and a grammar, we want to see whether the text *could* have been generated by the grammar. Moreover, *when* this is the case, we want to know *how* it was generated, that is, which grammar rules and which alternatives were used in the derivation. This process is called *parsing* or *syntax analysis*.

For the class of grammars defined in Figure 1 it is always possible to reconstruct a correct derivation. In the method below we shall further restrict the grammars so that there is a simple and efficient way to perform the reconstruction.

This section explains a simple parsing principle. Section 4 explains how to construct C# parser programs working according to this principle.

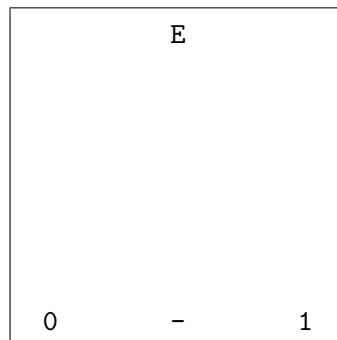
3.1 Parsing: reconstruction of a derivation tree

An attempt to reconstruct the derivation of a given string is called *parsing*. In these notes, we perform the reconstruction by working from the starting symbol down towards the given string. This method is called *top-down parsing*.

Consider again the grammar in Example 2:

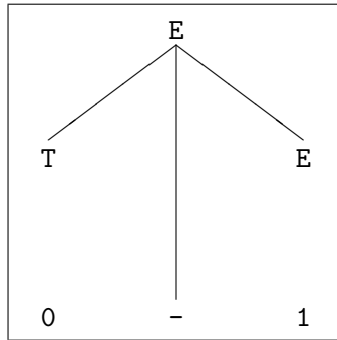
$$\begin{aligned} E &= T \text{ "-" } E \mid T . \\ T &= \text{"0"} \mid \text{"1"} . \end{aligned}$$

Let us reconstruct a derivation of the string "0" "-" "1" from the starting symbol E. We will do it by reconstructing the derivation tree, and therefore draw a box, write the starting symbol E at the top, and write the given input string "0" "-" "1" at the bottom of the box:



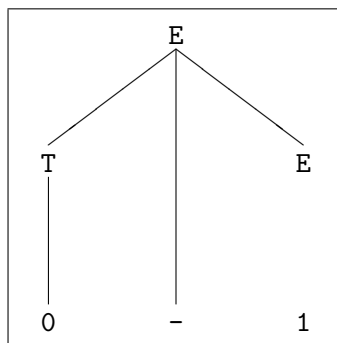
(a)

Our task is to find a derivation tree which connects E with the string at the bottom. We start from the top, and must derive something from E. According to the grammar, there are two possibilities, $E \implies T \text{ "-" } E$ and $E \implies T$. Only the first alternative is useful because the string, which involves a minus sign, could never be derived from T. So we extend the tree with the branches T, "-", and E, as shown in box (b):



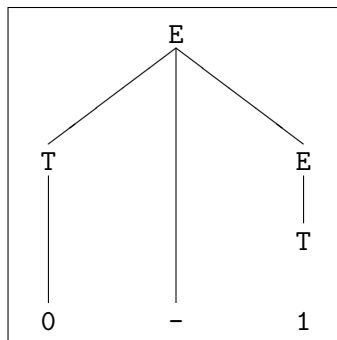
(b)

The next task is to derive the string "0" from T; luckily the grammar allows $T \Rightarrow "0"$, so we can extend the tree with the branch from T to "0" as shown in box (c):



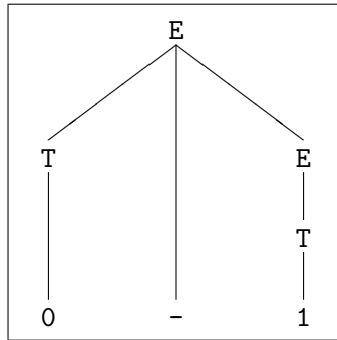
(c)

Next we must see how the remaining input symbol "1" can be derived from E. The $E \Rightarrow T$ alternative is reasonable, so we extend the tree with a branch from E to T, as shown in box (d):



(d)

Finally, we must derive "1" from T, but again there is a rule $T \Rightarrow "1"$, so we extend the tree with a branch from T to "1", as shown in box (e):



(e)

The parsing is complete: given the input string "0" "-" "1" we have constructed a derivation tree for it. When a derivation tree is the result of parsing, it is usually called a *parse tree*.

The derivation tree tells us two things. First, the input string *is* derivable from the starting symbol E. Second, we know at least one *way* it can be derived.

In the reconstruction, we worked from the top (E) and downwards; thus *top-down* parsing. Also note that in each step we extended the tree at the *leftmost* nonterminal.

3.2 A more machine-oriented view of parsing

We now consider another way to explain top-down parsing, more suited for programming than the trees shown above. We solve the same problem once more: can the string "0" "-" "1" be derived from E using the grammar in Example 2?

Previously we wrote down the string, wrote the nonterminal E above it, and reconstructed a derivation tree connecting the two. Now we write the string to the left, and the nonterminal E to the right:

"0" "-" "1" E

This corresponds to the situation in box (a). In general there is a string of remaining input symbols on the left and a sequence of remaining grammar symbols (nonterminals or terminals) on the right. This situation can be read as an equation "0" "-" "1" = E between the two sides. Parsing solves the equation in a number of steps. Each step rewrites the leftmost nonterminal on the right hand side, until the input string has been derived. Whenever the same symbol is at the head of both sides, we can cancel it. This is much like cancellation in algebra, where $x + y = x + z$ can be reduced to $y = z$ by cancelling x . The parsing is successful when both sides are empty, that is, Λ .

Returning to our task, we must rewrite E. There are two possibilities, $E \implies T \text{ "-" } E$ and $E \implies T$. It is easy to see for the human reader that "0" "-" "1" can be derived only from the first alternative, because of the "-" symbol. We now rewrite E to T "-" E and have the configuration

"0" "-" "1" T "-" E

This corresponds to the situation in box (b). Since $T \implies "0"$, we can get

"0" "-" "1" "0" "-" E

corresponding to the situation in box (c). Now we can cancel "0" and then "-" in both columns, so we need only see how the remaining input symbol "1" can be derived from E. Choosing the $E \implies T$ alternative and then $T \implies "1"$, we get in turn:

"1"	E
"1"	T
"1"	"1"

The two latter lines correspond to the situations in box (d) and (e). Now we can cancel "1" on both sides, leaving the empty string Λ on both sides, so the parsing process was successful. The complete sequence of parsing steps was:

"0" "-" "1"	E
"0" "-" "1"	T "-" E
"0" "-" "1"	"0" "-" E
"1"	E
"1"	T
"1"	"1"
Λ	Λ

We want to mechanize the parsing process by writing a program to perform it, but there is one problem. To decide which alternative of **E** to use (in the first parsing step), we had to look ahead in the input string to find the symbol "-". This lookahead is complicated to do in a program.

If our parsing program could choose the alternative by looking only at the *first* symbol of the remaining input, then the program would be simpler and more efficient.

3.3 Left factorization

The problem is with rules such as $E = T \text{ "-" } E \mid T$, where both alternatives start with the same symbol, T. We would like to *factorize* the right hand side into 'T ("-" E | Λ)', pulling the T outside a parenthesis, so to speak, and thus postponing the choice between the alternatives until after T has been parsed.

However, our grammar notation does not allow such parenthesized grammar fragments. To solve this problem we introduce a new nonterminal **Eopt** defined by $Eopt = \text{"-"} E \mid \Lambda$, and redefine E as $E = T Eopt$. Thus **Eopt** represents the parenthesized grammar fragment above.

Moreover, in Section 6 below it will prove useful to replace **E** in the **Eopt** rule with its only alternative T **Eopt**.

Example 4 Left factorization of the Example 2 grammar therefore gives

E	= T Eopt .
Eopt	= "-" T Eopt Λ .
T	= "0" "1" .

The set of strings derivable from **E** is the same as in Example 2, but the derivations will be different. □

Now the derivation of "0" "-" "1" (in fact, any derivation) must begin with $E \implies T Eopt$, and we need to see how "0" "-" "1" can be derived from T **Eopt**. Since $T \implies \text{"0"}$, we can cancel the "0" and only need to see how the remaining input "-" "1" can be derived from **Eopt**. There are two alternatives, $Eopt \implies \Lambda$ and $Eopt \implies \text{"-"} T Eopt$.

Since Λ can derive only the empty string, whereas the other alternative can derive strings starting with "-", we choose the latter. We now must see how "-" "1" can be derived from "-" T **Eopt**. The "-" is cancelled, and we must see how "1" can be derived from T **Eopt**. Now $T \implies \text{"1"}$, we cancel the "1", and we are left with the empty input. Clearly the empty input can be derived from **Eopt** only by its first alternative, Λ .

The parsing steps for "0" "-" "1" with the left factorized grammar of Example 4 are:

"0"	"-"	"1"	E
"0"	"-"	"1"	T Eopt
"0"	"-"	"1"	"0" Eopt
	"-"	"1"	"-" T Eopt
		"1"	T Eopt
		"1"	"1" Eopt
		Λ	Eopt
		Λ	Λ

Notice that now we can always choose between the alternatives by looking only at the first symbol of the remaining input. The corresponding derivation tree is shown in Figure 3.

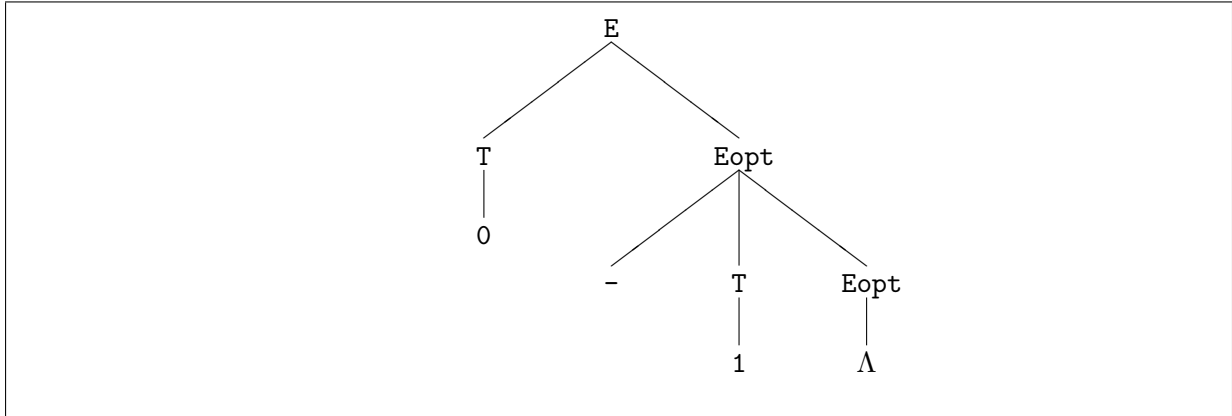


Figure 3: Derivation tree for the left factorized grammar

3.4 Left recursive nonterminals

There is another type of grammar rules we want to avoid. Consider the grammar

$$\begin{aligned} E &= E \text{ "-" } T \mid T . \\ T &= \text{"0"} \mid \text{"1"} . \end{aligned}$$

Some reflection (or experimentation) shows that it generates the same strings as the grammar from Example 2. However, E is *left recursive*: there is a derivation $E \Rightarrow E \dots$ from E to a symbol string that begins with E . It is even *self left recursive*: there is an alternative for E that begins with E itself. This means that the grammar is no good for top-down parsing, since we cannot choose between the alternatives for E by looking only at the first input symbol (in fact, not even by looking at any bounded number of symbols at the beginning of the string).

Left factorization is not possible for the above grammar, since the alternatives begin with different nonterminals. The only solution is to change the grammar to one that is not left recursive. Fortunately, this is always possible. In the present case, the original Example 2 grammar is a good solution.

In general, consider a grammar in which nonterminal A is self left recursive:

$$A = A g_1 \mid \dots \mid A g_m \mid f_1 \mid \dots \mid f_n .$$

The g_i and f_j stand for sequences of grammar symbols (possibly Λ). We require that $m, n \geq 1$, and that no f_j can derive a string beginning with A , so the only left recursion is through the first m alternatives.

Observe that every string derived from A must begin with an f_j , and continue with zero or more g_i 's. Therefore we can construct the following equivalent grammar where A is not self left recursive:

$$\begin{aligned} A &= f_1 A_{opt} \mid \dots \mid f_n A_{opt} \ . \\ A_{opt} &= g_1 A_{opt} \mid \dots \mid g_m A_{opt} \mid \Lambda \ . \end{aligned}$$

The role of the new nonterminal A_{opt} is to derive sequences of zero or more g_i 's.

The new grammar produced by this transformation usually is not left recursive, and it generates the same strings as the original one (namely, an f_j followed by zero or more g_i 's). The transformation sometimes produces a new grammar which is again left recursive. In that case, one must apply (more) cleverness to find a non left recursive grammar.

3.5 First-sets, follow-sets and selection sets

Consider a rule $A = f_1 \mid f_2$, and assume we have an input string ' $\tau \dots$ ' whose first input symbol is τ . We want to decide whether this string could be derived from A . Moreover, we want to choose between the alternatives f_1 and f_2 by looking only at the first input symbol.

There are two ways it might make sense to choose f_1 . First, if we can derive a string *starting* with τ from f_1 , then choosing f_1 might be sensible. Secondly, if we can derive the empty string Λ from f_1 , and we can derive a string starting with τ from something *following* whatever is derived from A , then choosing f_1 might be sensible.

To make the choice between f_1 and f_2 simple, we shall *require* that for a given input symbol τ , it makes sense to choose f_1 , or f_2 , or none of them, but it must never make sense to choose both. Accordingly, the parser chooses f_1 , or f_2 , or rejects the input as wrong. We now make this idea more precise.

The set of terminal symbols that can begin a string derivable from f is called its *first-set* and is written $First(f)$. The set of symbols that can follow a nonterminal A is called its *follow-set* and is written $Follow(A)$.

The *selection set* for an alternative f_i of a nonterminal $A = f_1 \mid \dots \mid f_n$ is $First(f_i)$ if f_i cannot derive the empty string Λ , and $First(f_i) \cup Follow(A)$ if f_i can derive Λ :

$$Select(f_i) = \begin{cases} First(f_i) \cup Follow(A) & \text{if } f_i \implies \Lambda \\ First(f_i) & \text{otherwise} \end{cases}$$

Intuitively, the selection set $Select(f_i)$ is the set of input symbols for which it is sensible to choose f_i . Why? It makes sense to choose f_i only if the first input symbol can be derived from f_i , or if the input symbol can follow A , and A can derive Λ via f_i .

How can we compute $First(f)$? If f is Λ , we have $First(\Lambda) = \{\}$ because the empty string does not start with any symbol.

If f is a terminal symbol " c ", we have $First("c") = \{c\}$ because the only string derivable is " c ", which begins with c .

If f is a nonterminal A whose rule is $A = f_1 \mid \dots \mid f_n$, then the set of strings derivable is the union of those derivable from the alternatives f_i . Therefore $First(A)$ is the union of the first-sets of the alternatives.

If f is a sequence $e_1 e_2 \dots e_m$, the set of strings derivable is the concatenation of strings derivable from the elements. Thus $First(f)$ includes $First(e_1)$. Moreover, if e_1 can derive Λ , then every string derivable from $e_2 \dots e_m$ is derivable also from $e_1 e_2 \dots e_m$. Therefore when e_1 can derive Λ , $First(f)$ includes $First(e_2 \dots e_m)$ also.

The computation of $First(f)$ is summarized in Figure 4.

$First(\Lambda)$	$=$	$\{\}$	
$First("c")$	$=$	$\{c\}$	for terminal "c"
$First(A)$	$=$	$First(f_1) \cup \dots \cup First(f_n)$	for nonterminal A where A is defined by $A = f_1 \mid \dots \mid f_n$
$First(e_1 e_2 \dots e_m)$	$=$	$\begin{cases} First(e_1) \cup First(e_2 \dots e_m) & \text{if } e_1 \Rightarrow \Lambda \\ First(e_1) & \text{otherwise} \end{cases}$	

Figure 4: Computation of first-sets

How can we compute the follow-set $Follow(A)$ of a nonterminal A? Assume that A appears in the rule $B = \dots \mid \dots A f \mid \dots$ for nonterminal B, where f is a string of grammar symbols (possibly Λ). Then $Follow(A)$ must include everything that f can begin with, and, if f can derive Λ , then also everything that can follow B. This is expressed by Figure 5.

The follow-set $Follow(A)$ of nonterminal A is the least (smallest) set of terminal symbols satisfying for every rule $B = \dots \mid \dots A f \mid \dots$, that			
$Follow(A)$	\supseteq	$\begin{cases} First(f) \cup Follow(B) & \text{if } f \Rightarrow \Lambda \\ First(f) & \text{otherwise} \end{cases}$	

Figure 5: Computation of follow-sets

With these definitions, the requirement on grammars for parser construction is that the selection sets of distinct alternatives f_i and f_j are disjoint: $Select(f_i) \cap Select(f_j) = \{\}$. Then a given input symbol c can belong to the selection set of at most one alternative, so the input symbol determines which alternative to choose.

However, in practice we shall use the more easily checkable sufficient requirements given in Figure 6.

Every grammar rule must have one of the forms			
Form 0:	$A = f_1$		
Form 1:	$A = f_1 \mid \dots \mid f_n$	$n \geq 2$	
Form 2:	$A = f_1 \mid \dots \mid f_n \mid \Lambda$	$n \geq 1$	
For rules of form 1 or 2 we require:			
<ul style="list-style-type: none"> • For distinct f_i and f_j we must have $First(f_i) \cap First(f_j) = \{\}$. • No f_i can derive Λ. • In rules of form 2, we must have $First(f_i) \cap Follow(A) = \{\}$ for all f_i. 			

Figure 6: Sufficient requirements on grammar for parsing

The requirements in Figure 6 imply that the grammar does not contain a left recursive nonterminal (unless the nonterminal is unreachable from the starting symbol, and therefore irrelevant).

Looking again at the left factorization example, we see that it does not satisfy the first requirement in Figure 6.

Example 5 Clearly $First("0") = \{0\}$ and $First("1") = \{1\}$, so in the grammar from Example 2

$$\begin{aligned} E &= T \text{"-"} E \mid T . \\ T &= "0" \mid "1" . \end{aligned}$$

we have

$$\begin{aligned} First(T) &= First("0") \cup First("1") = \{0, 1\} \\ First(T \text{"-"} E) &= First(T) = \{0, 1\} \end{aligned}$$

The rule $E = T \text{"-"} E \mid T$ is of form 1 and does not satisfy the requirement on first-sets in Figure 6, since the first-sets of the alternatives are not disjoint; they are identical. This problem occurs whenever two alternatives begin with the same symbol. \square

Example 6 Let us compute $Follow(T)$ for the grammar shown above. Consulting Figure 5, we see that $Follow(T)$ is the smallest set of terminal symbols which satisfies the two inequalities

$$\begin{aligned} Follow(T) \supseteq First(\text{"-"} E) &= First(\text{"-"}) = \{-\} \\ Follow(T) \supseteq Follow(E) & \end{aligned}$$

The first inequality is caused by the alternative $E = T \text{"-"} E \mid \dots$, and the second one by the alternative $E = \dots \mid T$. In the latter case, the ϵ following T is the empty string Λ .

But what is $Follow(E)$? It is the empty set $\{\}$, since $Follow(E)$ is defined to be the least set which satisfies

$$Follow(E) \supseteq Follow(E)$$

because there is a rule $E = T \text{"-"} E \mid \dots$. Any set satisfies this inequality. In particular the empty set does, and this is clearly the least such set.

Using this fact, we see that $Follow(T) = \{-\}$. \square

Example 7 In the left factorized Example 4 grammar

$$\begin{aligned} E &= T Eopt . \\ Eopt &= \text{"-"} T Eopt \mid \Lambda . \\ T &= "0" \mid "1" . \end{aligned}$$

the $Eopt$ rule has form 2, and we have for the alternatives of $Eopt$:

$$\begin{aligned} First(\text{"-"} T Eopt) &= First(\text{"-"}) = \{-\} \\ First(\Lambda) &= \{\} \end{aligned}$$

Reasoning as for $Follow(E)$ in the previous example, we also find that $Follow(Eopt) = \{\}$.

The first-sets $\{\}$ and $\{-\}$ of the two alternatives are disjoint, and $First(\text{"-"} T Eopt) \cap Follow(Eopt) = \{\}$, so the E rule satisfies the grammar requirements. The selection sets for the two alternatives are $\{\text{"-"}\}$ and $\{\}$. This shows how to choose between the alternatives of E : if the first input symbol is "-" , then choose the first alternative ($\text{"-"} T Eopt$), and if the input is empty, then choose the second alternative (Λ). \square

Now we know about first-sets, consider again the left recursive rule $E = E \text{"-"} T \mid T$ from Section 3.4. It has form 1, and for the first alternative we have

$$\begin{aligned} First(E \text{"-"} T) &= First(E) \\ &= First(E \text{"-"} T) \cup First(T) \\ &\supseteq First(T) \end{aligned}$$

Since $First(T) = \{0, 1\}$ is not empty, the first-sets of the alternatives $E \text{ "-" } T$ and T are not disjoint, and therefore the requirements of Figure 6 are not satisfied.

Example 8 The following grammar describes more realistic arithmetic expressions:

$$\begin{aligned} E &= E \text{ "+" } T \mid E \text{ "-" } T \mid T \text{ .} \\ T &= T \text{ "*" } F \mid T \text{ "/" } F \mid F \text{ .} \\ F &= \text{Real} \mid \text{"(" } E \text{ ")} \text{ .} \end{aligned}$$

Here E stands for expression, T for term, and F for factor. So an expression is the sum or difference of an expression and a term, or just a term. A term is the product or quotient of a term and a factor, or just a factor. A factor is a constant number, or an expression surrounded by parentheses.

The rules for E and T must be transformed to remove left recursion as explained in Section 3.4:

$$\begin{aligned} E &= T \text{ Eopt} \text{ .} \\ \text{Eopt} &= \text{"+" } T \text{ Eopt} \mid \text{"-" } T \text{ Eopt} \mid \Lambda \text{ .} \\ T &= F \text{ Topt} \text{ .} \\ \text{Topt} &= \text{"*" } F \text{ Topt} \mid \text{"/" } F \text{ Topt} \mid \Lambda \text{ .} \\ F &= \text{Real} \mid \text{"(" } E \text{ ")} \text{ .} \end{aligned}$$

Now we must check the grammar requirements. First we compute the follow-sets.

To determine $Follow(E)$ we list the requirements imposed by Figure 5, by considering all right hand side occurrences of E . There is only one, in the F rule:

$$\begin{aligned} Follow(E) &\supseteq First(\text{"})") \\ &= \{ \text{"})" \} \end{aligned}$$

Now $Follow(E)$ is the smallest set satisfying this requirement, so

$$Follow(E) = \{ \text{"})" \}$$

To determine $Follow(\text{Eopt})$ we similarly find the requirements

$$\begin{aligned} Follow(\text{Eopt}) &\supseteq Follow(E) \\ Follow(\text{Eopt}) &\supseteq Follow(\text{Eopt}) \end{aligned}$$

Again, $Follow(\text{Eopt})$ is the least set satisfying these requirements, so we conclude that

$$Follow(\text{Eopt}) = \{ \text{"})" \}$$

To determine $Follow(T)$ we note the sole requirement

$$\begin{aligned} Follow(T) &\supseteq First(\text{Eopt}) \cup Follow(\text{Eopt}) \\ &= \{ \text{"+"}, \text{"-" } \} \cup Follow(\text{Eopt}) \end{aligned}$$

for which the smallest solution is

$$Follow(T) = \{ \text{"+"}, \text{"-"}, \text{"})" \}$$

To determine $Follow(\text{Topt})$ we note the requirements

$$\begin{aligned} Follow(\text{Topt}) &\supseteq Follow(T) \\ Follow(\text{Topt}) &\supseteq Follow(\text{Topt}) \end{aligned}$$

for which the smallest solution is

$$\begin{aligned}
Follow(\text{Topt}) &= Follow(\text{T}) \\
&= \{ "+", "-", ")" \}
\end{aligned}$$

Now let us check the grammar requirements from Figure 6:

- The rules for **E** and **T** are of type 0 and therefore OK.
- The rule for **F** is of type 1 and OK because the first-sets $\{\text{Real}\}$ and $\{ "(" \}$ are disjoint.
- The rule for **Eopt** is of type 2 and OK because the first-sets $\{ "+", "-" \}$ and the follow-set $Follow(\text{Eopt}) = \{ ")" \}$ are disjoint.
- The rule for **Topt** is of type 2 and OK because the first-sets $\{ "*", "/" \}$ and the follow-set $Follow(\text{Topt}) = \{ "+", "-", ")" \}$ are disjoint.

Thus the transformed grammar satisfies the requirements. □

3.6 Summary of parsing theory

We have shown informally how top-down parsing works. We defined the concepts of first-set and follow-set. Using these concepts, we formulated a sufficient requirement on grammars for parser construction. For a grammar to satisfy this requirement, it must have no two alternatives starting with the same symbol, and no left recursive rules.

4 Parser construction in C#

We now show a systematic way to write a *parser skeleton* (in C#) for input described by a grammar satisfying the requirements in Figure 6. The parser skeleton checks that the input is well-formed, but does not build an internal representation of it; this will be done in Section 6.

4.1 C# representation of input strings

The raw input (from, say, a text file) is a stream of characters. For parsing, we need to turn this into a stream of tokens, where a *token* is the internal representation of a terminal symbol.

In C# we can represent a token stream by an object `ts` of type `IEnumerator<Token>` where `Token` is the type that represent a single token. Then `ts.Current` holds the current token, and the method `ts.MoveNext()` reads the next token in the stream. To signal an error in the input we shall throw an exception of type `ApplicationException`. The parser must access the input through the token stream `ts` only. We shall make a slight abuse of `IEnumerator<Token>` and demand that a valid token stream always ends with the special token `EOF` (end-of-file).

The type `Token` must have a field (or property) `kind` of some enum type `Kind` that specifies which kind it is. The use of an enum type allows us to test the kind of a token using C#'s `switch` statement. A simple token corresponding to a terminal symbol such as `"-"` may be represented by a token whose `kind` field is set to an enum value such as `SUB`.

However, some tokens occur in families. For instance, the terminal symbol `Real` may stand for all real numbers. Clearly, we cannot have an enum value for each real number, so terminal symbols belonging to this family will be represented by a token with the `kind` field set to the enum value `NUM` together with a field (or property) `nval` of type `double` in the token object.

Thus, we shall represent tokens with: an enum declaration, `Kind`, and a structure declaration, `Token`. Example given:

```
enum Kind { EOF, NUM, SUB, ... }
struct Token {
    public readonly Kind kind;
    public readonly double nval;
    private Token(Kind k) { kind = k; nval = 0; }
    private Token(double n) { kind = Kind.NUM; nval = n;}
    public override string ToString() {
        if ( kind == Kind.NUM ) return "NUM("+nval+")";
        else return kind.ToString();
    }
    // Factory methods
    static public Token FromKind(Kind k) { return new Token(k); }
    static public Token FromDouble(double d) { return new Token(d); }
}
```

To use the name `TokenStream`, rather than `System.Collection.Generic.IEnumerator` we shall use the following `using`-declaration when we work with token streams:

```
using TokenStream = System.Collections.Generic.IEnumerator<Token>;
```

We return to the subject of how to transform a stream of characters into a token stream in Section 5.

4.2 Constructing parsing methods in C#

A parser for a grammar G satisfying the requirements of Figure 6 can be constructed systematically from the grammar rules. The parser will consist of a set of mutually recursive *parsing methods*, one for each nonterminal in the grammar.

The parsing method corresponding to nonterminal A is called A also. It tries to find a string derivable from A at the beginning of the current token stream. If it succeeds, then it just returns, possibly after having read more tokens from the token stream. If it fails, then it throws an exception of class `ApplicationException`.

Grammar The parser for a grammar $G = (T, N, R, S)$ has the form

```
void A1 (TokenStream ts) { ... }
void A2 (TokenStream ts) { ... }
...
void Ak (TokenStream ts) { ... }

void Parse(TokenStream ts) {
    ts.MoveNext();
    S(ts);
    if (ts.Current.kind != Kind.EOF)
        throw new ApplicationException("Expected end of file");
    return;
}
```

where $\{A_1, \dots, A_k\} = N$ is the set of nonterminals and S is the starting symbol. The main method is `Parse`; it checks that no input remains after parsing. If parsing succeeds and no input remains, it just returns; otherwise it throws an exception.

Rule of form 0 The parsing method for a rule of form $A = f_1$ is

```
void A(TokenStream ts) {
    parse code for f1
    return;
}
```

Rule of form 1 The parsing method for a rule of form $A = f_1 \mid \dots \mid f_n$ is

```
void A(TokenStream ts) {
    switch(ts.Current.kind) {
    case t11 ... case t1a1:
        parse code for alternative f1
        return;
    ...
    case tn1 ... case tnan:
        parse code for alternative fn
        return;
    default:
        throw new ApplicationException("Expected t11 or ... or tnan");
    }
}
```

where $\{t_{i1}, \dots, t_{ia_i}\} = First(f_i)$ is the first-set of alternative f_i , for $i = 1, \dots, n$.

Rule of form 2 The parsing method for a rule of form $A = f_1 \mid \dots \mid f_n \mid \Lambda$ is

```
void A(TokenStream ts) {
    switch(ts.Current.kind) {
        case t11 ... case t1a1:
            parse code for alternative f1
            return;
        ...
        case tn1 ... case tnan:
            parse code for alternative fn
            return;
        default:
            return;
    }
}
```

where $\{t_{i1}, \dots, t_{ia_i}\} = First(f_i)$ as above.

Sequence The parse code for an alternative f which is a sequence $e_1 e_2 \dots e_m$ is

```
 $\mathcal{P}(e_1)$ 
 $\mathcal{P}(e_2)$ 
...
 $\mathcal{P}(e_m)$ 
```

where the parse code $\mathcal{P}(e_i)$ for each symbol e_i is defined below. Note that when the sequence is empty (that is, $m = 0$), the parse code is empty, too.

Nonterminal The parse code $\mathcal{P}(A)$ for a nonterminal A is a call $A(ts)$ to its parsing method.

Terminal The parse code $\mathcal{P}("c")$ for a terminal "c" depends on its position e_j in the sequence $e_1 \dots e_m$.

If the terminal is *not* the first symbol e_1 , then we must check that c is the current symbol in the token stream ts , and, if so, read the next token:

```
if (ts.Current.kind != Kind.c)
    throw new ApplicationException("Expected c");
ts.MoveNext();
```

If the terminal *is* the first symbol e_1 , then this check has already been made by the **switch** code for alternatives, so we just need to read the next token;

```
ts.MoveNext();
```

Note that in the parser construction for grammar rules of form 1 and 2, the grammar requirements ensure that the first-sets are disjoint, so the **case**-alternatives are all distinct. Also, for rules of form 2, the grammar requirements ensure that every first-set is disjoint from the follow-set of A , so an f_i alternative is never wrongly chosen instead of the **default** alternative (for Λ), which occurs last.

Example 9 Applying this construction method to the left factorized grammar from Example 4 gives the parser below. The token stream is provided by an object ts of type `TokenStream`.

```

enum Kind { EOF, SUB, ZERO, ONE }
struct Token {
    public readonly Kind kind;
    public Token(Kind k) { kind = k; }
    static public Token FromKind(Kind k) { return new Token(k); }
}

class Example9 {
    void E(TokenStream ts) {
        T(ts); Eopt(ts); return;
    }
    void Eopt(TokenStream ts) {
        switch (ts.Current.kind) {
            case Kind.SUB:
                ts.MoveNext(); T(ts); Eopt(ts); return;
            default:
                return;
        }
    }
    void T(TokenStream ts) {
        switch (ts.Current.kind) {
            case Kind.ZERO:
                ts.MoveNext(); return;
            case Kind.ONE:
                ts.MoveNext(); return;
            default:
                throw new ApplicationException("Expected 0 or 1");
        }
    }
    public void Parse(TokenStream ts) {
        ts.MoveNext();
        E(ts);
        if (ts.Current.kind != Kind.EOF)
            throw new ApplicationException("Expected end of file, but got:"
                +ts.Current.kind);
        return;
    }
}

```

To demonstrate the construction method we have followed it mindlessly in this example. Parts of the program may be improved, but it is advisable to postpone such improvements until you have acquired more practice with parser construction.

We can test the parser by using a List to represent a token stream:

```

class ListTestExample9 {
    public static void Main(string[] args) {
        List<Token> toks = new List<Token>();
        toks.Add(Token.FromKind(Kind.ZERO));
        toks.Add(Token.FromKind(Kind.SUB));
        toks.Add(Token.FromKind(Kind.ONE));
        toks.Add(Token.FromKind(Kind.EOF));
        Example9 p = new Example9();
        p.Parse(toks.GetEnumerator());
    }
}

```

□

4.3 Parsing methods follow the derivation tree

It is interesting to consider how the token stream $\langle '0', '- ', '1', \text{EOF} \rangle$, which corresponds to the input "0" "-" "1", is parsed by the parsing methods above. It turns out that the sequence of calls closely follows the derivation tree:

```
E calls T with  $\langle '0', '- ', '1', \text{EOF} \rangle$ 
  T finds a '0' and returns to E; now  $\langle '- ', '1', \text{EOF} \rangle$  remains
E calls Eopt with  $\langle '- ', '1', \text{EOF} \rangle$ 
  Eopt finds a '- '; now  $\langle '1', \text{EOF} \rangle$  remains
  Eopt calls T with  $\langle '1', \text{EOF} \rangle$ 
    T finds a '1' and returns to Eopt; now  $\langle \text{EOF} \rangle$  remains
  Eopt calls a second instance of Eopt with  $\langle \text{EOF} \rangle$ 
    Eopt finds EOF and returns to the first Eopt; now  $\langle \text{EOF} \rangle$  remains
  Eopt returns to E
```

If we draw this sequence of calls as a tree, we get precisely the derivation tree in Figure 3. The parsing methods walk through the derivation tree from left to right. This is no coincidence, but a result of the systematic construction shown above: the derivation tree corresponds to a particular string and a particular grammar, and the parser was constructed systematically from this grammar.

4.4 Summary of parser construction

We have shown a systematic way to construct a parser skeleton from a grammar satisfying the requirements in Figure 6. The parser skeleton just checks that the input, which is a stream of terminal symbols, follows the grammar. In Section 6 we show how to make the parser return more information, such as an internal representation of the input.

The parser in Example 9 can be found in file `Example9.cs`.

5 Scanners

In the parser above the input was represented as a token stream. As explained in the previous section, real input files are character streams, but it is inconvenient to work with bare characters. Therefore parsing of text files is usually divided into two phases.

In the *first* phase, the character stream is converted to a stream of tokens, and lay-out information (such as blanks) in the input text is removed. This is called *scanning* or *lexical analysis* and is explained in this section.

In the *second* phase, the stream of tokens is parsed as described in the previous sections.

The division into two phases gives a convenient way to allow any number of blanks *between* numerals and names without allowing blanks *inside* numerals and names. By a *blank* we mean a space character, a tabulation character, or a newline. The scanner decides what is a numeral, what is a name, and so on, and throws away all extra blanks. Then the parser never sees a blank: only numerals, names, and so on.

Although a scanner could be constructed systematically from a grammar for terminal symbols, we will not do that here. As mentioned in the previous section we represent a token stream as an object of class `TokenStream` (that is, `IEnumerator<Token>`). Thus, a scanner will implement the `IScanner` interface:

```
public interface IScanner {
    TokenStream Scan(TextReader reader);
}
```

Our scanners will not deal directly with files, but will instead take an object `reader` of type `System.IO.TextReader`. This allows us to reuse the scanners (and thereby parsers) for input not coming from files.

Example 10 The scanner below could be used with the parser in Example 9. It ignores all blanks, and considers all characters other than '-', '0' and '1' as errors.

```
class ZeroOneScan : IScanner {
    public TokenStream Scan(TextReader reader) {
        while ( reader.Peek() != -1 ) {
            if ( Char.IsWhiteSpace((char) reader.Peek()) ) reader.Read();
            else
                switch(reader.Read()) {
                    case '-': yield return Token.FromKind(Kind.SUB); break;
                    case '0': yield return Token.FromKind(Kind.ZERO); break;
                    case '1': yield return Token.FromKind(Kind.ONE); break;
                    default: throw new ApplicationException("Illegal character");
                }
        }
        yield return Token.FromKind(Kind.EOF);
    }
}
```

□

5.1 Character classification methods

In Example 10 we used the library method `Char.IsWhiteSpace` to decide whether an input character is a blank. The `Char` struct contains other useful character classification methods, for example:

<code>IsDigit</code>	Is the given character a decimal digit.
<code>IsLower</code>	Is the given character a lower case letter.
<code>IsUpper</code>	Is the given character a upper case letter.
<code>IsLetter</code>	Is the given character a letter.
<code>IsLetterOrDigit</code>	Is the given character a letter or a digit.

5.2 Scanning names

Suppose we need to scan and parse an input language (such as a programming language) which contains *names* of variables, procedures, or similar. Names are also called *identifiers*. A name in this sense is typically a nonempty sequence of letters and digits, beginning with a letter, and containing no blanks. Names can be described by this grammar:

```

Name      = Letter Letdigs .
Letdigs   = Letter Letdigs | Digit Letdigs | Λ .
Letter    = "A" | ... | "Z" | "a" | ... | "z" .
Digit     = "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9" .

```

A scanning method `ScanName` for names is shown below. It is called from the main scanner method `Scan` when a letter is met.

```

static string ScanName(TextReader reader) {
    StringBuilder sb = new StringBuilder();
    while ( Char.IsLetterOrDigit((char) reader.Peek() )
           sb.Append( (char) reader.Read() );
    return sb.ToString();
}

```

It accumulates the characters of the name using the `StringBuilder sb` until it meets a character which is neither a letter nor a digit, or until the input ends. It returns the collected string.

The name-scanner is invoked by adding an extra conditional branch in the `Scan` method. Futhermore, the `Token` struct must be extended with an extra field `val` for holding the string representing the name, similar we need an extra enum value in the enum type `Kind`.

```

enum Kind { ..., NAME }
struct Token {
    public readonly string sval;
    private Token(string s) { kind = Kind.NAME; sval = s; ... }
    private Token(Kind k) { kind = k; sval = null; ... }
    static public Token FromString(string s) { return new Token(s); }
    ...
}
class Scanner : IScanner {
    ...
    static string ScanName(TextReader reader) { ... }
    public IEnumerator<Token> Scan(TextReader reader) {
        while ( reader.Peek() != -1 ) {
            if ( Char.IsWhiteSpace((char) reader.Peek()) ) reader.Read();
            else if ( Char.IsLetter((char) reader.Peek()) )
                yield return Token.FromString(ScanName());
            else
                ...
        }
        yield return Token.FromKind(Kind.EOF);
    }
}

```

5.3 Distinguishing names from keywords

Most (programming) languages contain so-called *keywords* or *reserved names* which are sequences of letters that cannot be used as names. For instance, C# and Java have the keywords ‘class’, ‘interface’, ‘while’, and so on.

Keywords are represented by other terminal symbols than names. For instance, if ‘class’, ‘interface’, and ‘while’ are keywords, then the corresponding terminal symbols may be represented as tokens with kind CLASS, INTERFACE, and WHILE (assuming that Kind contains these enum values), instead of a token with kind set to NAME and sval set to "class" etc.

A scanner distinguishes names from keywords as follows. Whenever something that looks like a name has been found by the scanner, it is compared to a list of keywords. It is classified as a keyword if it is in the list, otherwise it is classified as a name. For example, the following extension of the scanner above will distinguish keywords from names:

```

enum Kind { ..., NAME, CLASS, INTERFACE, WHILE }
struct Token { ... }
class Scanner : IScanner {
    ...
    static string ScanName(TextReader reader) { ... }
    public IEnumerator<Token> Scan(TextReader reader) {
        while ( reader.Peek() != -1 ) {
            if ( Char.IsWhiteSpace((char) reader.Peek()) ) {
                reader.Read();
            } else if ( Char.IsLetter((char) reader.Peek()) ) {
                string sval = ScanName(reader);
                switch ( sval ) {
                    case "class":
                        yield return Token.FromKind(Kind.CLASS); break;
                    case "interface":
                        yield return Token.FromKind(Kind.INTERFACE); break;
                    case "while":
                        yield return Token.FromKind(Kind.WHILE); break;
                    default:
                        yield return Token.FromString(sval); break;
                }
            } else
                ...
        }
        yield return Token.FromKind(Kind.EOF);
    }
}

```

If the set of keywords is large, one may use more efficient means to find the token corresponding to a given string. For instance, the hash-table Dictionary class from System.Collections.Generic may be useful.

5.4 Scanning floating-point numerals

A numeral is a string of characters that represents a number, such as "3.1414". Floating-point numerals can be described by this grammar:

```

Real    = Digits "." Digits .
Digits = Digit | Digit Digits .
Digit  = "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9" .

```

The scanning method `ScanReal` shown in the following is called as `ScanReal(reader)` from the main scanning method when a digit `c` has been met by the scanner.


```

static double ScanReal(TextReader reader) {
    double n = reader.Read() - '0';
    while ( Char.IsDigit((char) reader.Peek()) ) {
        n = 10 * n + reader.Read() - '0';
    }
    if ( reader.Peek() == '.' ) {
        reader.Read();
        return ScanFrac(n, 0.1, reader);
    } else {
        return n;
    }
}

static double ScanFrac(double n, double wt, TextReader reader) {
    while( Char.IsDigit((char) reader.Peek()) ) {
        n += wt * (reader.Read() - '0');
        wt /= 10.0;
    }
    return n;
}

```

The `ScanReal` method accumulates the value `n` of the digits before the point until a point is met, or another non-digit is met, or the input ends. In the first case, there must be a digit after the point. If there is, then method `ScanFrac` is called. Method `ScanFrac` accumulates the value of the digits after the point until a non-digit is met. Then it returns the scanned number.

The scanner method for reals is invoked by adding an extra conditional branch in the `Scan` method. Furthermore, the `Token` struct must be extended with an extra field `nval` for holding the double representing the number, similarly we need an extra enum value in the enum type `Kind` (just like we did for names). We shall only show the extension of the `Scan` method

```

class Scanner : IScanner {
    ...
    public IEnumerator<Token> Scan(TextReader reader) {
        while ( reader.Peek() != -1 ) {
            if ( Char.IsWhiteSpace((char) reader.Peek()) ) {
                reader.Read();
            } else if ( Char.IsDigit((char) reader.Peek()) ) {
                yield return Token.FromDouble(ScanReal(reader));
            } else {
                ...
            }
        }
        yield return Token.FromKind(Kind.EOF);
    }
    ...
}

```

5.5 Reading string and text files with C#

Scanners and parsers can be tested with short strings like this:

```
class StringTestExample9 {
    public static void Main(string[] args) {
        IScanner scanner = new ZeroOneScan();
        Example9 parser = new Example9();
        TextReader r = new StringReader("0 - 1");
        parser.Parse(scanner.Scan(r));
    }
}
```

For larger input texts this is impractical. Hence, we use the following code to read file `example1.zo` and check that it is well-formed:

```
class TestExample9 {
    public static void Main(string[] args) {
        IScanner scanner = new ZeroOneScan();
        Example9 parser = new Example9();
        using(TextReader r = File.OpenText("example1.zo")) {
            parser.Parse(scanner.Scan(r));
        }
    }
}
```

5.6 Summary of scanners

Scanning is the first phase in the parsing of a text. It turns the string of characters into a stream of tokens, which is then read by a parser in the second phase.

6 Parsers with attributes

So far a parser just checks that an input string can be generated by the grammar: only the *form* or *syntax* of the input is handled. Of course we usually want to know more about the input, so we extend the parsers to return a representation of the input.

For this, every parsing method must return an additional result. Some parsing methods take an additional parameter too. The additional parameters and results are called *attributes*.

6.1 Constructing attributed parsers

We still use parser skeletons constructed as in Section 4.2, but we add code to handle the attributes. So far a parsing method `A` has had signature

```
void A(TokenStream ts) { ... }
```

but from now on its type will be

```
OutValue A(InValue inval, TokenStream ts) { ... }
```

where `InValue` and `OutValue` are the types of the attributes. We call `A` an *attributed parser*. The `inval` argument is called an *inherited* attribute and the return value is called a *synthesized* attribute. One may decorate parse trees with attribute values. An inherited attribute is sent down the tree as an additional argument to a parsing method, and a synthesized attribute is sent up the tree as an additional result from a parsing method.

Some parsing methods do not take any inherited attributes, but most attributed parsing methods return a synthesized attribute. An attributed parsing method `A` either returns the synthesized attribute, or raises an exception.

One cannot say in general how to turn a parser skeleton into an attributed parser. What extensions and changes are required depends on the kind of information we need about the parsed input. Below we consider a typical example: simple expressions.

The main method `Parse` of a parser now either returns a result or raises an exception. For example, if the type `OutValue` is really `double`, then `Parse` is written like this:

```
public double Parse(TokenStream ts) {
    ts.MoveNext();
    double result = E(ts);
    switch( ts.Current.kind ) {
    case Kind.EOF: return result;
    default: throw new ApplicationException("Parse error: "+ts.Current);
    }
}
```

where `E` is the starting symbol of the grammar. In the following examples and exercises, method `Parse` will always have this form, and is therefore not shown.

Arithmetic expressions are usually evaluated from left to right. One also says that the arithmetic operators, such as `'-'`, *associate to the left*, that is, group to the left.

Example 11 Recall the parser skeleton for arithmetic expressions in Example 9. We extend it so that every parsing method returns a synthesized attribute which is the value of the expression parsed by that method. To evaluate from left to right, we also extend method `Eopt` with an inherited attribute `inval` which at any point is the value of the expression parsed so far. When a `T` is parsed in the `SUB` branch of method `Eopt`, its value `tv` is subtracted from `inval`, and the result is passed to `Eopt` in the recursive call.

```

class EvalParser {
    double E(TokenStream ts) {
        double tv = T(ts);
        double ev = Eopt(tv, ts);
        return ev;
    }
    double Eopt(double inval, TokenStream ts) {
        switch( ts.Current.kind ) {
        case Kind.SUB:
            ts.MoveNext();
            double tv = T(ts);
            double ev = Eopt(inval - tv, ts);
            return ev;
        default: return inval;
        }
    }
    double T(TokenStream ts) {
        switch (ts.Current.kind) {
        case Kind.ZERO: ts.MoveNext(); return 0;
        case Kind.ONE: ts.MoveNext(); return 1;
        default:
            throw new ApplicationException("Expected 0 or 1");
        }
    }
    public double Parse(TokenStream ts) { ... }
}

```

Method E first calls T. Method T parses a "0" or a "1", and returns the double 0 or 1, which gets bound to tv. This value is passed to Eopt as an inherited attribute inval. In method Eopt there are two possibilities: *either* it parses "-" T Eopt in which case it calls T again, subtracts the new T-value tv from inval, and passes the difference to Eopt in the recursive call, *or* it parses Λ , in which case it just returns inval.

At any point, inval is the value of the expression parsed so far. When the input is empty, inval is the value of the entire expression. □

Figure 7 shows the attribute values when parsing the input string "0" "-" "1" "-" "1" and evaluating from left to right, as done by the parser above. The inherited attribute inval is shown to the left of the lines, and the synthesized attributes are shown to the right.

A typical use of the attributed parser is

```

class MainClass {
    public static void Main(string[] args) {
        IScanner scanner = new ZeroOneScan();
        EvalParser parser = new EvalParser();
        string s = "0-1-1";
        TextReader r = new StringReader(s);
        Console.WriteLine("{0} = {1}", s, parser.Parse(scanner.Scan(r)));
    }
}

```

where ZeroOneScan is the scanner class defined in Example 10. Executing the program will print 0-1-1 = -2 on the console.

Some of the previous parsing methods could be simplified. For instance, the method E could be simplified to:

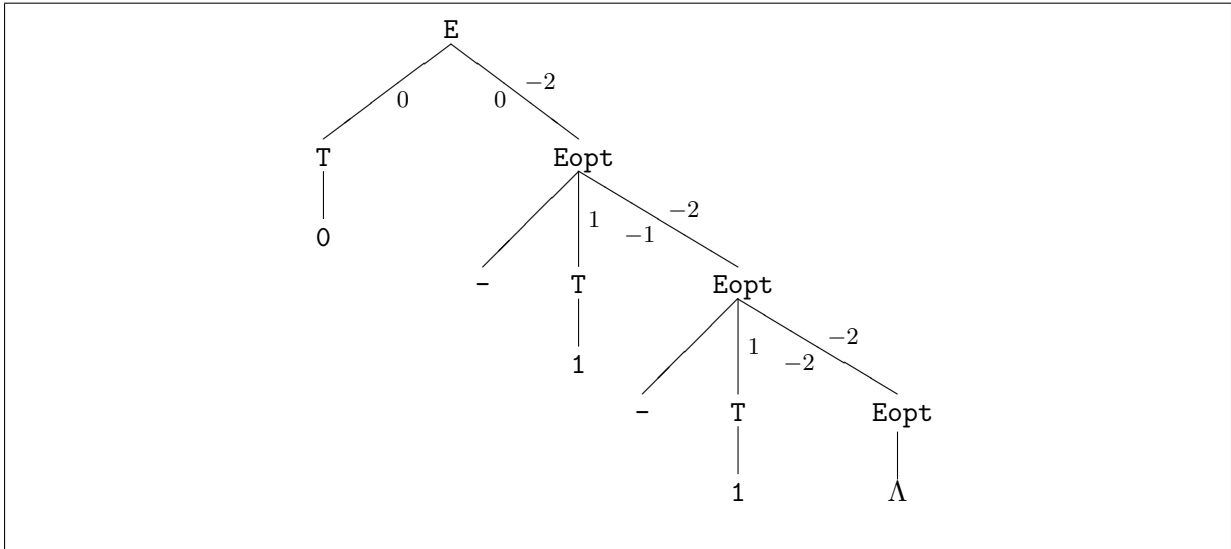


Figure 7: Parse tree with attributes for left-to-right evaluation

```
double E(TokenStream ts)
{ return Eopt(T(ts), ts); }
```

Such simplifications are best done after the parser has been written. Their effect on execution time is limited, so they are mostly of cosmetic value. Also, simplifications require a good understanding of expression evaluation order in C#; otherwise one may introduce subtle errors.

Above we defined left-to-right evaluation of arithmetic expressions, which is usual in programming languages. What if we had a bizarre desire to evaluate from right to left (as in the programming language APL)? This can be done with a small change to the attributed parser from Example 11.

Example 12 The attributed parser in Example 11 can be changed to evaluate the expression from right to left as follows:

```

class RightToLeftEvalParser {
    double E(TokenStream ts) {
        double tv = T(ts);
        double ev = Eopt(tv, ts);
        return ev;
    }
    double Eopt(double inval, TokenStream ts) {
        switch( ts.Current.kind ) {
        case Kind.SUB:
            ts.MoveNext();
            double tv = T(ts);
            double ev = Eopt(tv, ts);
            return inval - ev;
        default: return inval;
        }
    }
    double T(TokenStream ts) {
        switch (ts.Current.kind) {
        case Kind.ZERO: ts.MoveNext(); return 0;
        case Kind.ONE: ts.MoveNext(); return 1;
        default:
            throw new ApplicationException("Expected 0 or 1");
        }
    }
    public double Parse(TokenStream ts) { ... }
}

```

The only change is in the SUB branch of the Eopt method. The subtraction is now done *after* the recursive call to Eopt.

At any point, *inval* is the value (0 or 1) of the last T parsed. The value *ev* of the remaining expression is subtracted from *inval* *after* the recursive call to Eopt. No subtractions are done until the entire expression has been parsed; and then they are done from right to left. □

Figure 8 shows the attribute values when parsing the input string "0" "-" "1" "-" "1", and evaluating from right to left, as done by the parser above. The inherited attribute *inval* is shown to the left of the lines, and the synthesized attributes are shown to the right.

6.2 Building representations of input

An important application of attributed parsers is to build representations of the input that has been read by the parser. Such representations are often called abstract syntax trees. An *abstract syntax tree* is a representation of a text which shows the structure of the text and leaves out irrelevant information, such as the number of blanks between symbols.

A flexible and general way to represent abstract syntax trees in C# is to use classes and subclasses. For instance, to represent simple expressions as defined in Example 2, one may define an abstract class Expr of expressions, and concrete classes Zero, One, and Minus, corresponding to each of the three kinds of expressions:

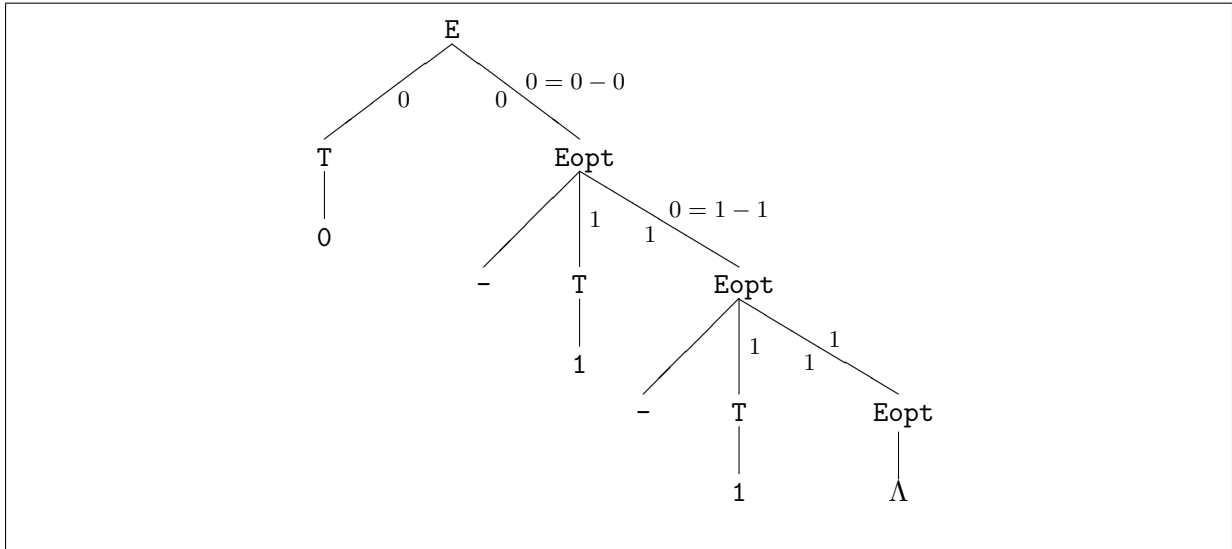


Figure 8: Parse tree with attributes for right-to-left evaluation

```

abstract class Expr { }

class Zero : Expr {
  public override string ToString()
  { return "0"; }
}

class One : Expr {
  public override string ToString()
  { return "1"; }
}

class Minus : Expr {
  Expr E1, E2;
  public Minus(Expr e1, Expr e2)
  { E1 = e1; E2 = e2; }
  public override string ToString()
  { return String.Format("{0}-{1}", E1, E2); }
}

```

These class declarations say: there is an abstract concept `Expr` of expression. There are concrete expression concepts `Zero` and `One`, corresponding to `0` and `1`. There is a concrete expression concept `Minus`, which consists of two subexpressions `E1` and `E2`. Any object belonging to class `Expr`, or one of its subclasses, knows how to convert itself to a `string` (by overriding the `ToString` method from class `Object`).

Thus the expression `0-1` can be represented as `new Minus(new Zero(), new One())`. The expression `0-1-1` can be represented either as `new Minus(new Minus(new Zero(), new One()), new One())` or as `new Minus(new Zero(), new Minus(new One(), new One()))`. The first representation corresponds to a left-to-right reading, and the second one corresponds to a right-to-left reading.

Let us make an attributed parser which builds the representation corresponding to a left-to-right reading of simple arithmetic expressions. Such a parser will be very similar to the parser

for left-to-right evaluation in Example 11.

Example 13 This parser builds abstract syntax trees for simple arithmetic expressions.

```
class AstParser {
    Expr E(TokenStream ts) {
        Expr tv = T(ts);
        Expr ev = Eopt(tv, ts);
        return ev;
    }
    Expr Eopt(Expr inval, TokenStream ts) {
        switch (ts.Current.kind) {
            case Kind.SUB:
                ts.MoveNext();
                Expr tv = T(ts);
                Expr ev = Eopt(new Minus(inval, tv), ts);
                return ev;
            default: return inval;
        }
    }
    Expr T(TokenStream ts) {
        switch (ts.Current.kind) {
            case Kind.ZERO: ts.MoveNext(); return new Zero();
            case Kind.ONE: ts.MoveNext(); return new One();
            default:
                throw new ApplicationException("Expected 0 or 1");
        }
    }
    public Expr Parse(TokenStream ts) { ... }
}
```

Instead of returning an integer (0 or 1), method `T` now returns the representation of an expression: an object of class `Zero` or `One`. Instead of subtracting one number from another, returning a number, method `Eopt` now builds and returns a representation of an expression (in the `SUB` branch).

Since the new representation is built before `Eopt` is called recursively to parse the rest of the expression, the representation is built from left to right as in Example 11. At any point, `inval` is the representation of the expression parsed so far. \square

The new parsing methods return a representation of the parsed expression rather than its value. Hence they have return type `Expr`, not `double`. A typical application of the attributed parser in Example 13 would look like this:

```
class MainClass {
    public static void Main(string[] args) {
        IScanner scanner = new ZeroOneScan();
        EvalParser parser = new EvalParser();
        string s = "0-1-1";
        TextReader r = new StringReader(s);
        Console.WriteLine("{0} = {1}", s, parser.Parse(scanner.Scan(r)));
    }
}
```

Calling method `parser.Parse` will scan and parse the string and build a representation of its contents, as an object of class `Expr`. Printing this object will invoke its `ToString()` method to

convert the object to a string.

6.3 Summary of parsers with attributes

To make parsing methods return information about the input, we add new components to their results and (possibly) to their arguments. Different ways of handling the new results and arguments give different effects, such as left-to-right or right-to-left evaluation. Looking at parse trees is helpful for understanding attribute evaluation.

An abstract syntax tree is a representation of a text without unnecessary detail. Parsers can be extended with attributes to construct the abstract syntax tree for a text while parsing it.

7 A larger example: arithmetic expressions

We now consider arithmetic expressions such as $4.0+5.0*7.0$ and $(20.0-5.0)/3.0$, which are found in almost all programming languages, and show how to scan, parse, and evaluate them.

7.1 A grammar for arithmetic expressions

Here is a first attempt at a grammar for arithmetic expressions:

```
E    = E "+" E
      | E "-" E
      | E "*" E
      | E "/" E
      | Real
      | "(" E ")" .
```

This grammar does not satisfy the grammar requirements, but could easily be transformed to do so. However, the grammar does not express the structure of arithmetic expressions very well. In arithmetics, the multiplication and division operators bind more strongly than addition and subtraction. Thus $4.0+5.0*7.0$ should be thought of as $4.0+(5.0*7.0)$, giving 39, and not as $(4.0+5.0)*7.0$, giving 63. We say that multiplication and division have higher *precedence* than addition and subtraction.

A subexpression which is a numeral or a parenthesized expression is called a *factor*. A subexpression involving only multiplications and divisions of factors is called a *term*. An expression is a sequence of additions or subtractions of terms.

Then the precedence can be expressed as follows: Factors must be evaluated first, and terms must be evaluated before additions and subtractions.

To ensure that terms are parsed as units, we introduce a separate nonterminal T for them, and similarly for factors F. This gives the following grammar for arithmetic expressions:

```
E    = E "+" T | E "-" T | T .
T    = T "*" F | T "/" F | F .
F    = Real | "(" E ")" .
```

The rule for E generates strings of form $T \text{ "+" } T \text{ "-" } \dots \text{ "+" } T$ with one or more T's separated by additions and subtractions. Similarly, the rule for T generates $F \text{ "*" } F \text{ "/" } \dots \text{ "*" } F$ with one or more F's separated by multiplications and divisions. Note that *Real* stands for a class of terminal symbols: the real numerals.

To avoid the left recursive rules, we transform the E and T rules as described in Section 3.4. We obtain the following grammar:

```
E    = T Eopt .
Eopt = "+" T Eopt | "-" T Eopt |  $\Lambda$  .
T    = F Topt .
Topt = "*" F Topt | "/" F Topt |  $\Lambda$  .
F    = Real | "(" E ")" .
```

This grammar satisfies the requirements in Figure 6, as argued in Example 8.

7.2 The parser constructed from the grammar

The terminal symbols of the grammar are the operators '+', '-', '*', '/', the parentheses '(' and ')', and the real numerals:

```
enum Kind { EOF, PLUS, MINUS, MULT, DIV, LPAR, RPAR, NUM }
struct Token {
    public readonly Kind kind;
    public readonly double nval;
    private Token(Kind k) { kind = k; nval = 0;}
    private Token(double n) { kind = Kind.NUM; nval = n;}
    public override string ToString() {
        if ( kind == Kind.NUM ) return "NUM("+nval+")";
        else return kind.ToString();
    }
    static public Token FromKind(Kind k) { return new Token(k); }
    static public Token FromDouble(double d) { return new Token(d); }
}
```

Application of the construction method from Section 4.2 to the above grammar gives the parser skeleton shown in the following.

```
class ArithSkelParser {
    public void Parse(TokenStream ts) {
        void E(TokenStream ts) { T(ts); Eopt(ts); return; }
        void Eopt(TokenStream ts) {
            switch( ts.Current.kind ) {
                case Kind.PLUS: ts.MoveNext(); T(ts); Eopt(ts); return;
                case Kind.MINUS: ts.MoveNext(); T(ts); Eopt(ts); return;
                default: return;
            }
        }
        void T(TokenStream ts) { F(ts); Topt(ts); return; }
        void Topt(TokenStream ts) {
            switch( ts.Current.kind ) {
                case Kind.MULT: ts.MoveNext(); F(ts); Topt(ts); return;
                case Kind.DIV: ts.MoveNext(); F(ts); Topt(ts); return;
                default: return;
            }
        }
        void F(TokenStream ts) {
            switch( ts.Current.kind ) {
                case Kind.NUM: ts.MoveNext(); return;
                case Kind.LPAR:
                    ts.MoveNext();
                    E(ts);
                    if ( ts.Current.kind != Kind.RPAR )
                        throw new ApplicationException("Parse error: expected ')");
                    ts.MoveNext(); return;
                default:
                    throw new ApplicationException("Parse error: expected number or '('");
            }
        }
    }
}
```

7.3 A scanner for arithmetic expressions

An appropriate scanner is shown below. It ignores blanks and uses the scanner method `ScanReal` for real numbers defined in Section 5.4.

```
class Scanner : IScanner {
    static double ScanReal(TextReader reader) {
        double n = reader.Read() - '0';
        while ( Char.IsDigit((char) reader.Peek()) )
            n = 10 * n + reader.Read() - '0';
        if ( reader.Peek() == '.' ) {
            reader.Read();
            return ScanFrac(n, 0.1, reader);
        } else
            return n;
    }
    static double ScanFrac(double n, double wt, TextReader reader) {
        while( Char.IsDigit((char) reader.Peek()) ) {
            n += wt * (reader.Read() - '0');
            wt /= 10.0;
        }
        return n;
    }
}

public TokenStream Scan(TextReader reader) {
    while ( reader.Peek() != -1 ) {
        if ( Char.IsWhiteSpace((char) reader.Peek()) ) {
            reader.Read();
        } else if ( Char.IsDigit((char) reader.Peek()) ) {
            yield return Token.FromDouble(ScanReal(reader));
        } else {
            char c = (char) reader.Read();
            switch ( c ) {
                case '+': yield return Token.FromKind(Kind.PLUS); break;
                case '-': yield return Token.FromKind(Kind.MINUS); break;
                case '*': yield return Token.FromKind(Kind.MULT); break;
                case '/': yield return Token.FromKind(Kind.DIV); break;
                case '(': yield return Token.FromKind(Kind.LPAR); break;
                case ')': yield return Token.FromKind(Kind.RPAR); break;
                default:
                    throw new ApplicationException("Illegal character: '"+c+"'");
            }
        }
    }
    yield return Token.FromKind(Kind.EOF);
}
```

7.4 Evaluating arithmetic expressions

Now we extend the parser skeleton from Section 7.2 to evaluate the arithmetic expressions while parsing them. As observed previously, arithmetic expressions should be evaluated from left to right, so the resulting attributed parser below is similar to that in Example 11, except that many subexpressions have been simplified by hand.

```
class ArithEvalParser {
  public double Parse(TokenStream ts) {
    ts.MoveNext();
    double result = E(ts);
    switch( ts.Current.kind ) {
      case Kind.EOF: return result;
      default: throw new ApplicationException("Parse error: "+ts.Current);
    }
  }
  double E(TokenStream ts) { return Eopt(T(ts), ts); }
  double Eopt(double inval, TokenStream ts) {
    switch( ts.Current.kind ) {
      case Kind.PLUS:
        ts.MoveNext(); return Eopt(inval + T(ts), ts);
      case Kind.MINUS:
        ts.MoveNext(); return Eopt(inval - T(ts), ts);
      default:
        return inval;
    }
  }
  double T(TokenStream ts) { return Topt(F(ts), ts); }
  double Topt(double inval, TokenStream ts) {
    switch( ts.Current.kind ) {
      case Kind.MULT:
        ts.MoveNext(); return Topt(inval * F(ts), ts);
      case Kind.DIV:
        ts.MoveNext(); return Topt(inval / F(ts), ts);
      default:
        return inval;
    }
  }
  double F(TokenStream ts) {
    switch( ts.Current.kind ) {
      case Kind.NUM:
        double nval = ts.Current.nval;
        ts.MoveNext(); return nval;
      case Kind.LPAR:
        ts.MoveNext();
        double ev = E(ts);
        if ( ts.Current.kind != Kind.RPAR ) {
          throw new ApplicationException("Parse error: expected ')');");
        }
        ts.MoveNext(); return ev;
      default:
        throw new ApplicationException("Parse error: expected number or '('");");
    }
  }
}
```

8 Some background

8.1 History and notation

Formal grammars were developed within linguistics by Noam Chomsky around 1956, and were first used in computer science by John Backus and Peter Naur in 1960 to describe the Algol programming language. Their notation was subsequently called *Backus-Naur Form* or *BNF*. In the original BNF notation, our grammar from Example 4 would read:

```
<E>      ::= <T> <Eopt>
<Eopt>   ::=   | - <T> <Eopt>
<T>      ::= 0 | 1
```

This notation uses a different convention than ours: nonterminals are surrounded by angular brackets, and terminals are not quoted. Also, here the empty string Λ is denoted by nothing (empty space). In compiler books one may find still another notation:

```
E      → T Eopt
Eopt   → ε
Eopt   → - T Eopt
T      → 0
T      → 1
```

In this notation there is only one alternative per rule, so defining a nonterminal may require several rules. Also, ϵ is used instead of our Λ .

As can be seen, the actual notation used for grammars varies, and combinations of these notations exist also. However, the underlying idea of derivation is always the same.

8.2 Extended Backus-Naur Form

Our grammar notation is a simplification of the so-called *Extended Backus-Naur Form* or *EBNF*. The full EBNF notation contains more complicated forms of alternatives \mathbf{f} .

In EBNF, an *alternative* \mathbf{f} is a *sequence* $\mathbf{e}_1 \dots \mathbf{e}_m$ of elements, not just symbols. An *element* \mathbf{e} may be a symbol as before, or

- an *option* of form $[\mathbf{f}]$, which can derive zero or one occurrence of sequence \mathbf{f} , or
- a *repetition* of form $\{ \mathbf{f} \}$, which can derive zero, one, or more occurrences of \mathbf{f} , or
- a *grouping* of form (\mathbf{f}) , which can derive an occurrence of \mathbf{f} .

A grammar in EBNF notation using the new kinds of elements can be converted to a grammar in our notation. The conversion is done by introducing extra nonterminals and rules:

- an option $[\mathbf{f}]$ is replaced by a new nonterminal Optf with rule $\text{Optf} = \mathbf{f} \mid \Lambda$.
- a repetition $\{ \mathbf{f} \}$ is replaced by a new nonterminal Repf with rule $\text{Repf} = \mathbf{f} \text{Repf} \mid \Lambda$.
- a grouping (\mathbf{f}) is replaced by a new nonterminal Grpf with rule $\text{Grpf} = \mathbf{f}$.

This shows that our simple grammar notation can express everything that EBNF can, possibly at the expense of introducing more nonterminals.

8.3 Classes of languages

The parsing method described in Section 4 is called *recursive descent* parsing and is an example of a *top-down* parsing method. It works for a class of grammars called $LL(1)$: those that can be parsed by reading the input symbols from the *Left*, making derivations always from the *Leftmost* nonterminal, and using a lookahead of 1 input symbol. This class includes all grammars that satisfy the requirements in Figure 6.

Another well-known class of grammars, more powerful than $LL(1)$, is the $LR(1)$ class which can be parsed *bottom-up*, reading the input symbols from the *Left*, making derivations always from the *Rightmost* nonterminal, and using a lookahead of 1 input symbol. Construction of bottom-up parsers is complicated, and is seldom done by hand. A useful subclass of $LR(1)$ is the class $LALR(1)$ (for ‘lookahead LR ’), which can be parsed more efficiently, by smaller parsers. The Unix utility ‘*Yacc*’ is an automatic parser generator for $LALR(1)$ grammars. The $LR(1)$ grammars are sufficiently powerful for most computing problems, but as exemplified by Exercise 4 there are grammars for which there is no equivalent $LR(1)$ grammar (and consequently no $LALR(1)$ -grammar or $LL(1)$ -grammar).

The class of grammars defined in Figure 1 is properly called the *context-free grammars*. This is just one class in the hierarchy identified by Chomsky: (0) the *unrestricted*, (1) the *context-sensitive*, (2) the *context-free*, and (3) the *regular* grammars. The unrestricted grammars are more powerful than the context-sensitive ones, which are more powerful than the context-free ones, which are more powerful than the regular grammars.

The unrestricted grammars cannot be parsed in general; they are of theoretical interest but of little practical use in computing. All context-sensitive grammars can be parsed, but may take an excessive amount of time and space, and so are of little practical use. The context-free grammars are those defined in Figure 1; they are highly useful in computing, in particular the subclasses $LL(1)$, $LALR(1)$, and $LR(1)$ mentioned above. The regular grammars can be parsed very efficiently using a constant amount of memory, but they are rather weak; they cannot define parenthesized arithmetic expressions, for instance.

The following table summarizes the grammar classes:

Chomsky hierarchy	Example rules	Comments
0: Unrestricted	"a" B "b" \rightarrow "c"	Rewrite system
1: Context-sensitive	"a" B "b" \rightarrow "a" "c" "b"	Non-abbreviating rewrite system
2: Context-free	B \rightarrow "a" B "b"	As defined in Figure 1. Some interesting subclasses: <hr/> $LR(1)$ bottom-up parsing <hr/> $LALR(1)$ bottom-up, ‘ <i>Yacc</i> ’ <hr/> $LL(1)$ top-down, these notes
3: Regular	B \rightarrow "a" "a" B	parsing by finite automata

8.4 Further reading

A description (in Danish) of practical recursive descent parsing using Turbo Pascal is given by Kristensen [2], who provided Example 1 and other inspiration.

There is a rich literature on scanning and parsing in connection with compiler construction. The standard reference is Aho, Sethi, and Ullman [1]. More information on recursive descent parsing is found in Lewis, Rosenkrantz, and Stearns [3], and in Wirth [4, Chapter 5].

9 Exercises

Exercise 1 Write down a grammar for arrays of (unsigned) integers. For instance, the empty array of integers is $\{\}$. Other examples of arrays of integers are $\{117\}$, $\{2,3,5,7,11,13\}$. Show the derivations of $\{\}$ and $\{7, 9, 13\}$. \square

Exercise 2 Consider the grammar

$$\begin{aligned} E &= T "+" E \mid T "-" E \mid T . \\ T &= "0" \mid "1" . \end{aligned}$$

Left factorize it and find selection sets for the alternatives of the resulting grammar. \square

Exercise 3 Consider the grammar below, which is self left recursive:

$$S = S S \mid "0" \mid "1" .$$

Apply the technique for removing left recursion (Section 3.4). Find first-, follow-, and selection sets for the resulting grammar. Does it satisfy the grammar requirements?

What strings are derivable from this grammar? Find a grammar which generates the same strings and satisfies the requirements (this is quite easy). \square

Exercise 4 The grammar

$$P = "a" P "a" \mid "b" P "b" \mid \Lambda .$$

generates palindromes (strings which are equal to their reverse). Find first-, follow-, and selection sets for this grammar. Which requirement in Figure 6 is not satisfied? (In fact, there is no way to transform this grammar into one that satisfies the requirements). \square

Exercise 5 Consider the grammar in Exercise 2. Left factorize it. Construct a parser skeleton for the left factorized grammar, using the tokens '+', '-', '0', and '1'. \square

Exercise 6 The grammar

$$T = "0" \mid "1" \mid "(" T ")" .$$

describes simple expressions such as 1, (1), ((0)), etc. with well-balanced parentheses. Choose a suitable set of tokens to represent the terminal symbols, and construct a parser skeleton for the grammar. Test it on the expressions above, and on some ill-formed inputs. \square

Exercise 7 Write a grammar and construct a parser for parenthesized expressions such as 0, 0+(1), 1-(1+1), (0-1)-1, etc. \square

Exercise 8 Consider the grammar for polynomials from Example 3. (1) Remove the left recursion in the rule for **Poly**. (2) Left factorize the rule for **Term**. (3) Choose a suitable set of tokens to represent the terminal symbols. Note that **Natnum** in the grammar stands for a family of terminal symbols 0, 1, 2, (4) Construct a parser skeleton for the transformed grammar and test it. \square

Exercise 9 Show that the requirements in Figure 6 imply that for every grammar rule, and distinct alternatives f_i and f_j , it holds that $Select(f_i) \cap Select(f_j) = \{\}$. \square

Exercise 10 The input language for the scanner in Example 10 is described by the grammar:


```
input = "-" input | "0" input | "1" input | blank input | Λ .
blank = " " | "\t" | "\n" .
```

Make sure the grammar satisfies the requirements, then use the construction method of Section 4 to systematically make a scanner for it. Your scanner must check the form of the input, but need not return a list of terminals. □

Exercise 11 Extend the parser constructed in Exercise 5 to evaluate the parsed expression and return its value. You may decide yourself whether evaluation should be from left to right or right to left. □

Exercise 12 Extend the parser constructed in Exercise 5 to build an abstract syntax tree for the parsed expression, using the following classes:

```
abstract class Expr { }

class Zero : Expr {}

class One : Expr {}

class Minus : Expr {
    Expr E1, E2;

    public Minus(Expr e1, Expr e2)
    { E1 = e1; E2 = e2; }
}

class Plus : Expr {
    Expr E1, E2;

    public Plus(Expr e1, Expr e2)
    { E1 = e1; E2 = e2; }
}
```

What are the types of the attributed parsing methods? □

Exercise 13 Extend the scanner from Section 7.3 to recognize C# floating-point numbers with exponents such as '6.6256E34' or '3E8'. □

Exercise 14 What changes are necessary to make the parser in Example 13 build representations from right to left? □

Exercise 15 Check that the grammar at the end of Section 7.1 satisfies the grammar requirements. □

Exercise 16 Extend the grammar, scanner, and parser from Section 7 to handle arithmetic expressions with exponentiation, such that $3.0 * 4.0^2.0$ evaluates to 48, that is, 3 times the square of 4. Note that the exponentiation operator usually associates to the right and has higher precedence than multiplication and division, so $2.0^2.0^3.0$ is $2.0^{(2.0^3.0)}$ and evaluates to 256, not to 64.

What changes are necessary if the exponentiation operator were '**' instead of '^'? □

Exercise 17 The following classes may be used to represent the arithmetic expressions from Section 7:

```
abstract class Expr { }

class Binop : Expr {
  Expr E1, E2;
  char Op;

  public Binop(Expr e1, char op, Expr e2)
  { E1 = e1; Op = op; E2 = e2; }
}

class Real : Expr {
  double R;

  public Real(double r)
  { R = r; }
}
```

Write an attributed parser that builds abstract syntax trees of this form.

□

References

- [1] A.V. Aho, R. Sethi, and J.D. Ullman. *Compilers, Principles, Techniques, and Tools*. Addison-Wesley, 1986.
- [2] J.T. Kristensen. *Konstruktion af indlæseprogrammer*. Teknisk Forlag, 1990.
- [3] P.M. Lewis II, D.J. Rosenkrantz, and R.E. Stearns. *Compiler Design Theory*. The Systems Programming Series. Addison-Wesley, 1976.
- [4] Niklaus Wirth. *Algorithms + Data Structures = Programs*. Prentice-Hall, 1976.

Index

- \implies (derivation), 4
- Λ (empty sequence), 4
- abstract syntax tree, 30
- alternative, 4, 38
- arithmetic expressions, 34
- associate to the left, 27
- attribute, 27
- attributed parser, 27
- Backus-Naur Form, 38
- blank, 21
- BNF, 38
- bottom-up parsing, 39
- context-free grammar, 39
- context-sensitive grammar, 39
- derivation, 4
- derivation tree, 5
- EBNF, 38
- element, 38
- Extended Backus-Naur Form, 38
- factor, 34
- factorize, 9
- First*(**f**) (first-set), 12
- first-set, 11, 12
- Follow*(**A**) (follow-set), 12
- follow-set, 11, 12
- grammar, 3, 4
- grammar notation, 4
- grammar requirements, 12
- grammar rule, 4
- grouping, 38
- identifier, 22
- inherited attribute, 27
- keywords, 23
- LALR* grammar, 39
- language
 - generated by grammar, 5
- left associative, 27
- left factorization, 9
- left recursive, 10
- lexical analysis, 21
- LL* grammar, 39
- LR* grammar, 39
- name
 - scanning of, 22
- nonterminal symbol, 4
- numeral, 24
- option, 38
- Parse** (method), 17, 19, 27
- parse tree, 8
- parser, 3
- parser with attributes, 27
- parser construction, 17
- parser skeleton, 16
- parsing, 6
 - bottom-up, 39
 - recursive descent, 39
 - top-down, 6, 39
- parsing theory, 6
- precedence, 34
- recursive descent parsing, 39
- regular grammar, 39
- repetition, 38
- requirements on grammar, 12
- rule, 4
- Scan** (method), 21, 23, 25, 36
- ScanFrac** (method), 25
- ScanName** (method), 22
- Scanner** (method), 36
- scanner, 3, 21
 - IScanner** (interface), 21
- scanning, 21
- ScanReal** (method), 25
- Select*(**f**) (selection set), 11
- selection set, 11
- self left recursive, 10
- sequence, 4
- starting symbol, 4
- symbol, 4
- syntax, 27
- syntax analysis, 6

- synthesized attribute, 27
- term, 34
- terminal symbol, 4
- Token (struct), 16, 23, 35
- token, 16
- token stream, 16
 - TokenStream (type), 16, 21
- top-down parsing, 6, 39
- tree
 - derivation, 5
 - parse, 8
- unrestricted grammar, 39
- Yacc, 39