

Top Down Parsing

Copyright and AI Notice

These slides are intended for studying and teaching at Nazarbayev University in Astana, Kazakhstan.

If you want to redistribute or adapt these slides for different purposes, check the license first.

No generative AI was used in the preparation of these slides.

©Hans de Nivelles, 2025

Building a Parser

As with tokenizers, there are essentially two ways to build a parser:

- Write it by hand.
- Use a parser generator (CUP, Yacc, Bison, Maphoon).

If your grammar is big, likely to change, it is better to use a parser generator.

Tasks of the Parser

Determine the structure of the sequence of tokens that is produced by the tokenizer, and produce some representation that can be further processed.

Decide if the sequence of tokens is correct according to the grammar.

If yes, determine the value of the attribute. (It can be a tree representation, or a value.)

If not, give error messages that are as useful as possible.

Optional: Try to recover after the error.

Why do you do this?

Writing a Parser by Hand

For grammars that are not complicated, it is possible to write a parser by hand. In most cases, one needs to make changes in the grammar.

It turns out convenient to consider grammar rules where the right hand side is a regular expression.

The resulting parser is usually called **recursive descent** parser, because it recursively descends from the start symbol to the terminal symbols.

Example

We will construct a parser for the following context-free grammar:

$$\begin{aligned} E &\rightarrow E + E_1 & | & E - E_1 & | & E_1 \\ E_1 &\rightarrow E_1 \times E_2 & | & E_1 / E_2 & | & E_2 \\ E_2 &\rightarrow -E_2 & | & E_3 \\ E_3 &\rightarrow \text{int} & | & \text{double} & | & \text{ident} & | \\ & & & \text{ident}(A) & | & (E) \\ A &\rightarrow E & | & A, E \end{aligned}$$

$$\Sigma = \{E, E_1, E_2, E_3, A\} \cup$$

$$\{+, -, \times, /, \text{int}, \text{double}, \text{ident}, (,), ', '\}.$$

Start symbol is E .

Example

Possible words are:

$$f(a, b, c)$$

$$f(a + b, f(a + b), a + b \times c),$$

$$(1 + 2) \times (2.0 + 3)$$

What is the purpose of the different non-terminals E, E_1, E_2, E_3 ?

What is the purpose of A ?

We assume that we want to construct an AST, which will have type **syntaxtree**.

Assuming this, what will be the attributes of E, E_1, E_2, E_3, A ?

Implementing a Top Down Parser

In order to obtain a top down parser, one has to write a parsing function for every non-terminal symbol:

parseE, **parseE₁**, **parseE₂**, **parseE₃**, and **parseA**.

The functions must have access to the tokenizer. We come back to this later.

They must return the attribute of their non-terminal symbol.

Interface to the Tokenizer, General Structure

I recommend that you write a class `parser` with a field `std::optional< symbol > lookahead` and the following two methods:

```
// Returning a non-const reference allows to move out
// the attribute:

const symbol& topdown::getlookahead( ) {
    if( !lookahead. has_value( ))
        lookahead = source. get( ); // From the tokenizer
    return lookahead. value( );
}

void topdown::resetlookahead( ) {
    lookahead. reset( );
}
```

parseE₂

Method parseE₂ is easy to implement:

```
if getlookahead( ) is '-' then  
    resetlookahead( )  
     $a = \text{parseE}_2( )$   
    return neg}(a)  
else  
    return parseE}_3( )
```

There are two possible rules for E_2 . If the lookahead is '-', we know it is the first rule. Otherwise, it must be the second.

$\text{neg}(a)$ is a tree whose top node is **neg**.

Problems with parseE

When writing parseE, we run into problems: There are three rules to choose from. Two start with an E , and one starts with E_1 .

At the beginning, we don't know if we must call parseE or parseE₁.

We could try to read ahead and see if there is a '+' or '-' further in the input. This will not work. For example, in $a * -b$, the minus sign is unary, and in $a * (b-c)$, it's in parentheses. We still have to apply rule $E \rightarrow E_1$.

We could try to look at the first token with which a word generated from E or E_1 can start. Unfortunately both can start with $-$, int, double, ident, ' (', so that won't help.

Method parseE_3

Method parseE_3 has a similar problem, but not so bad. There are two rules that start with `ident`: $E_3 \rightarrow \text{ident}$ and $E_3 \rightarrow \text{ident}(A)$.

In this case, one can postpone the decision until after the `ident` :

```
if getlookahead( ) is ident then
     $s$  = the attribute of getlookahead( )
    resetlookahead( )
    if getlookahead( ) is '(' then
        resetlookahead( )
         $a$  =  $\text{parseA}( )$ 
        if getlookahead( ) is not )' then syntax error: expected )'
        resetlookahead( )
        return  $s(a)$ 
    else
        return  $s$ 
```

Rewrite the Grammar

Symbol E rewrites to the following words

$$E_1, E_1 + E_1, E_1 - E_2, E_1 + E_1 + E_1, \dots, E_1 - E_1 - E_1.$$

One can write a regular expression for all words obtained by the rules $E \rightarrow E + E_1$, $E \rightarrow E - E_1$, $E \rightarrow E_1$.

$$E_1 ((' + ' | ' - ') E_1)^*.$$

All problematic groups of rules have to be replaced by regular expressions.

This is usually possible for realistic programming languages.

There is no general algorithm, so you have to use your insight and creativity.

Rewrite the Grammar (2)

The result is:

$$E \rightarrow E_1 (('+' | '-') E_1)^*$$

$$E_1 \rightarrow E_2 (('\times' | '/') E_2)^*$$

$$E_2 \rightarrow ('-')^* E_3$$

$$E_3 \rightarrow \text{int} \mid \text{double} \mid '(' E ') \mid \\ \text{ident} (\epsilon \mid '(' A '))$$

$$A \rightarrow E (',' E)^*$$

$$\Sigma = \{ E, E_1, E_2, E_3 \} \cup$$

$$\{ +, -, \times, /, \text{int}, \text{double}, \text{ident}, (,), ', ' \}.$$

Start symbol is still E .

Choice '|' can be implemented with **if**. Repetition '*' can be implemented by using **while**.

Method parseE can now be implemented as follows:

```
a = parseE1( )
while getlookahead( ) is '-' or '+' do
    t = getlookahead( ).
    resetlookahead( )
    a2 = parseE1( )
    if t = '+' then
        a = add(a, a2)
    else
        a = sub(a, a2)
return a
```

The other functions can be written in similar way.

Interface to the Tokenizer in Java

Because a previous version of this lecture used Java, I show how to write the code in Java.

In Java, every class Object is implicitly optional, because it can be null.

```
class TopDown
{
    Lexer lex;
    java_cup.runtime.Symbol lookahead;

    public TopDown( Lexer lex )
        { this. lex = lex; }
```

```
public Symbol getLookahead( ) {
    if( lookahead == null )
        lookahead = lex. NextToken( );
    return lookahead;
}

public void resetLookahead( ) {
    lookahead = null;
}
};
```

Errors

Errors must be taken seriously, and must be included in the design from the beginning.

Creating good error messages is difficult.

Giving up after an error, is acceptable only for very simple programs. Realistic programs need to continue as good as possible, in order to collect as much information as possible in one run. This is called **error recovery**.

Recovery

The ability to recover from errors is an important feature for the usefulness of a parser.

In most cases, one ignores the input, until a synchronization symbol is found. (For example ; or a closing))

One has to create a special attribute for the garbage read, and this attribute must not result in further errors.

Summary

For simple grammars, it is possible to implement a top down parser by hand. Such parsers are called recursive descent parsers.

In most cases, you have to modify the grammar rules in order to make the parse functions deterministic. Often, this can be achieved by using regular expressions.

You have to write a parse function for every non-terminal, which returns the attribute of this non-terminal.

Errors must be taken seriously. Recovering from errors is difficult. It is easier with a parser generator.