# A Simple Parser Combinator Library in C++ (DRAFT)

Sven Eric Panitz

`www.panitz.name`

23rd December 2004

**Abstract**

Monadic Parser Combinators stem from functional programming. This paper exploits the ideas of parser combinators and applies them to the C++ programming language. The resulting library is extremely small, flexible and easy to use. The paper contains the complete source code of the resulting parser library. As an example a parser of N. Wirth's language PL/0 is given in terms of the parser library.

## 1 Introduction

Parser combinators are a technique developed for lazily evaluated functional programming languages. Some early impressive example can be found in [FL89]. A recent and very efficient implemention is the *parsec*-library for Haskell [LM01]. The idea of a parser combinator library is to construct more complex parsers by combining simpler parsers. The combination is done along the building rules of productions in a grammar. There are basically three kinds of parsers:

- atomic parsers, which accept exactly one token in the input stream.

- the sequence of two parsers.

- the alternative of two parsers.

In a production the alternative is expressed by a vertical bar |. For the sequence no special symbol is used.

Parser combinator libraries generally construct recursive descendant parsers with backtracking. They generally construct a list of results. Each result in this list constitutes a different pars for the grammar. An empty list denotes failure [Wad85].

Since a parser is a function, parser combinators are higher-order functions. In order to mimick the optical impression of a grammar, parser combinator libraries use overloaded operators for combinator functions.

C++ is a object-oriented language. However it is enormously flexible and allows almost every programming style. Most features known from functional languages, as operator overloading, generic types or higher order functions, are available in C++. Therefore defining a parser combinator library in C++ should be only some handcraft work. In the following section a simple such library is defined.

# 2   Parser Combinators in C++

In order to get a most flexible parser library we use generic types.[1] In C++ generic types are expressed via templates. The C++ compiler will instantiate type variables of a template in the source code. The instantiated source code then gets compiled. this is the so called *heterogenous* translation of generic code.

As a consequence we cannot compile general code for the parser library, but have to provide it in an header file, which needs to be included by applications.

─────── ParsLib.h ───────
```
1  #ifndef __PARSLIB_H
2  #define __PARSLIB_H
3
4  #include <vector>
5
6  using namespace std;
7  namespace name {namespace panitz{ namespace parser{
```

Within the library we define a hierarchical namespace in Java package style.

## 2.1   The Type of A parser

### 2.1.1   Constant applicative forms

Before we start to define the types of a parser, we define an auxiliary class. This will be necessary, when defining the parser of a recursive (or mutual recursive) production(s). The auxiliary class is used to express a memorizing constant function, a so called *constant applicative form*. In terms of the standard template library this is a so called *generator*. The class is generic over the result type of the function. It contains a field for the function pointer and a field for the result of the function. Two constructors are provided:

─────── ParsLib.h ───────
```
8   template <typename a>
9   class CAF {
10    private: a(*f)();
11    public:
12      a result;
13      a operator()(){
14        if (result==NULL) result=f();
15        return result;
16      }
17      CAF(a x):result(x){}
18      CAF(a(*f)()){this->f=f;result=NULL;}
19  };
```

Since the capability of type inference in C++ is limited and not applied for constructors, we define a function, which wraps the constructor call for `CAF`:

---

[1] In functional programming the term *polymorphic type* is used for *generic types*. However in object oriented programming *polymorphism* is used for something else.

```
──────── ParsLib.h ────────
20  template <typename a>
21  CAF<a> caf(a(*f)()){return CAF<a>(f);}
```

### 2.1.2  Results of a parser

A parser consumes some token of the input stream and produces some pars result. The result can be of any type, but generally it will be some abstract syntax tree. The token type can be of any type. Even the stream type of the input stream could be left most general, as long it has the typical*iterator* interface of retrieving the next token and a check for the end of the stream. However thoughout htis paper we will stick to the standard class `vector` from the *standard template library*.

A pars may fail. Then it usually consumes no token from the input stream. Unlike in functional implementations we will not express failure by an empty list, but have a flag for failure in the class `ParsResult`.

We keep the class for a pars results generic over two types: the token type and the type of the actual result. The class has three fields. One contains the further still to consume token, one the actual result and one the flag for failure:

```
──────── ParsLib.h ────────
22  template <typename a,typename b>
23  class ParsResult{
24    public:
25      a result;
26      vector<b> furtherToken;
27      bool failed;
28
29      ParsResult(vector<b>& furtherToken)
30        :furtherToken(furtherToken),failed(true){};
31      ParsResult(a result,vector<b>& furtherToken)
32        :result(result),furtherToken(furtherToken),failed(false){}
33  };
```

We provide a function which wraps the constructor for the result of a failed pars.

```
──────── ParsLib.h ────────
34  template <typename a,typename b>
35  ParsResult<a,b>* fail(vector<b>& furtherToken){
36    return new ParsResult<a,b>(furtherToken);
37  }
```

This enables us to rely on C++ type inference mechanism when constructing a failure object.

### 2.1.3  The general parser interface

The parser interface can be kept quite simple. It just contains one method. The method `parse` takes a vector of a generic token type and returns some `ParsResult`:

```
                        ──────── ParsLib.h ────────
38  template <typename a,typename b>
39  class Parser{
40    public: virtual ParsResult<a,b>* parse(vector<b>& xs)=0;
41  };
```

Different from libraries in functional programming just one result is returned. Not a list of results. If there is not a unique derivation in a grammar, then our implementation will just give one and not all of them. This is for effecency reasons. Unlike in lazily evaluated functional languages without further effort, always all derivations would be calculated. This is not desirable. However most grammars used in compiler construction have unique derivations.

In the following subsections the three basic kinds of parsers are defined.

## 2.2   Atomic Parsers

The most simple parser is a parser, which accepts exactly one token. This parser needs to be handcoded. This is done by way of the class `GetToken`. The class has two fields. An example of the token which is to be accepted and a equality function for this token type. The class is a subclass of `Parser` and as well a subclass of `CAF`. The result of the `CAF`-object is the object itself. This enables us to use `GetToken`-objects directly as parsers but as functions returning parsers as well.

```
                        ──────── ParsLib.h ────────
42  template <typename a>
43  class GetToken: public Parser<a,a>,public CAF<Parser<a,a>*>{
44    private:
45      a token;
46      bool(*eq)(a,a);
47    public:
48      GetToken(a token,bool(*eq)(a,a))
49        :CAF<Parser<a,a>*>(this),token(token),eq(eq){}
```

The implementation of the method `parse` is straightforward. Compare (if available) the next token with the token in questaion. Construct a successful or failure parser result object depending on the comparision.

```
                        ──────── ParsLib.h ────────
50      ParsResult<a,a>* parse(vector<a>& xs){
51        if (!xs.empty() && eq(token,xs[0])){
52          a tok = xs[0];
53          vector<a> restToken = vector<a>(xs.begin()+1,xs.end());
54          return new ParsResult<a,a>(tok,restToken);
55        }
56        return fail<a,a>(xs);
57      }
58  };
```

To ease use we again provide a function which wraps the constructor call:

```
                           ParsLib.h
59   template <typename a>
60   CAF<Parser<a,a>*> getToken(a token,bool(*eq)(a,a)){
61      return new GetToken<a>(token,eq);}
```

## 2.3   Sequence Operator

We can define classes for parser combination. The class `Seq` is use to combine to parser as a sequence. The sequence of to parsers denotes the following: first apply the first parser. In case of Success apply the second parser to the remaining token. in case of success combine the two partial results into a common result.

The most natural result for the sequence of two parser is a pair containing the two partial results. For this purpose we can use the class `pair` from the standrad template library. Thus we get the following class header:

```
                           ParsLib.h
62   template <typename a,typename b,typename c>
63   class Seq
64        :public Parser<pair<a,b>*,c >
65        ,public CAF<Parser<pair<a,b>*,c>*>{
```

As with class `GetToken`, we do not only extend the class parser, but also the class `CAF`, with the parser as its result type.

We will not directly combine two parser objects, but two `CAF` objects, which have a parser as result. We provide two internal fields for these objects. The fileds are initialized within the constructor:

```
                           ParsLib.h
66      private:
67          CAF<Parser<a,c>*> p1;
68          CAF<Parser<b,c>*> p2;
69
70      public:
71        Seq(CAF<Parser<a,c>*> p1
72           ,CAF<Parser<b,c>*> p2):
73           CAF<Parser<pair<a,b>*,c>*>(this)
74            ,p1(p1),p2(p2){}
```

Eventually the actual method `parse` needs to be defined. Again the implementation is straightforward. The only thing to take care of is storage management. Temporary results of the two combined parses need to be deleted.

```
                           ParsLib.h
75        virtual ParsResult<pair<a,b>*,c>* parse(vector<c>& xs){
76          ParsResult<a,c>* res1 = p1()->parse(xs);
77          if (!res1->failed) {
78             vector<c> further = res1->furtherToken;
79             a r1 = res1->result;
```

```
80
81          ParsResult<b,c>* res2 = p2()->parse(further);
82          if (!res2->failed){
83           b r2 = res2->result;
84           vector<c> further = res2->furtherToken;
85           delete res1;delete res2;
86           return
87           new ParsResult<pair<a,b>*,c>(new pair<a,b>(r1,r2),further);
88          }
89          delete res2;
90        }
91
92       delete res1;
93       return fail<pair<a,b>*,c>(xs);
94     }
95  };
```

As we are allready used to, a function which wraps the constructor is defined. This is done by way of overloading the comma operator.

```
                          ─── ParsLib.h ───
96  template <typename a,typename b,typename c>
97  CAF<Parser<pair<a,b>*,c>*> operator,
98              (CAF<Parser<a,c>*>p1,CAF<Parser<b,c>*> p2){
99    return ((Parser<pair<a,b>*,c>*)new Seq<a,b,c>(p1,p2));
100 }
```

## 2.4   Alternative Operator

In the same way as for the sequence combinator, we can define a class for the alternative combinator. However, we need to think of the result. In a combination of parsers, the two parser may have different result types. The result type of the combination is then either the result type of the first or of the other parser. In Haskell this can be easily expressed by the algebraic type `Either`. In C++ unfortunatly there is no corresponding standard class available. However it can be easily defined by way of C's *union* construct.

```
                          ─── ParsLib.h ───
101 template <typename A,typename B>
102 class Either{
103  public:
104     bool isLeft;
105     union LeftOrRight{A left;B right;} value;
106     Either(bool isLeft):isLeft(isLeft){};
107 };
```

We provide two simple functions to serve as constructors for this class. A function two construct an object for the left type:

```
                             ────── ParsLib.h ──────
108  template <typename A,typename B>
109  Either<A,B>* left(A v){
110    Either<A,B>* result=new Either<A,B>(true);
111    result->value.left=v;
112  }
```

And the corresponding function for the right type:

```
                             ────── ParsLib.h ──────
113  template <typename A,typename B>
114  Either<A,B>* right(B v){
115    Either<A,B>* result=new Either<A,B>(false);
116    result->value.right=v;
117  }
```

Now we are well prepared to define the alternative combinator. Its header can analogously
defined to the sequence combinator:

```
                             ────── ParsLib.h ──────
118  template <typename a,typename b,typename c>
119  class Alt:public Parser<Either<a,b>*,c>
120          ,public CAF<Parser<Either<a,b>*,c>*>{
121    private:
122      CAF<Parser<a,c>*> p1;
123      CAF<Parser<b,c>*> p2;
124    public:
125      Alt(CAF<Parser<a,c>*> p1, CAF<Parser<b,c>*> p2):
126          CAF<Parser<Either<a,b>*,c>*>(this),p1(p1),p2(p2){}
```

The implementation is again straighforward. First apply the first parser. In case of success
construct a left success result. Otherwise apply the second parser to the original input stream.
This is where backtracking is done. Once more we carefully need to delete intermediate results:

```
                             ────── ParsLib.h ──────
127    virtual ParsResult<Either<a,b>*,c>* parse(vector<c>& xs){
128      ParsResult<a,c>* res1 = p1()->parse(xs);
129      if (!res1->failed) {
130        vector<c> further = res1->furtherToken;
131        a r1 = res1->result;
132        delete res1;
133        return new ParsResult<Either<a,b>*,c>(left<a,b>(r1),further);
134      }
135      delete res1;
136      ParsResult<b,c>* res2 = p2()->parse(xs);
137      if (!res2->failed) {
138        b r2 = res2->result;
139        vector<c> further = res2->furtherToken;
140        delete res2;
```

```
141        return new ParsResult<Either<a,b>*,c>(right<a,b>(r2),further);
142      }
143      delete res2;
144      return fail<Either<a,b>*,c>(xs);
145    }
146  };
```

The wrapper for a constructor is again defined as an overloaded operator. The most natural choice for this operator is the vertical bar, which is used in the production rules of a grammer as well.

———— ParsLib.h ————
```
147  template <typename a,typename b,typename c>
148  CAF<Parser<Either<a,b>*,c>*>
149    operator|(CAF<Parser<a,c>*>p1,CAF<Parser<b,c>*> p2){
150    return ((Parser<Either<a,b>*,c>*)new Alt<a,b,c>(p1,p2));
151  }
```

## 2.5   Calculating the Result

Up to now we can express the rules of a grammar nicely. However, we can only construct parsers, which have pairs or `Either` objects as results. Generally we will want to construct some special result for certain production rules. In parser generators as yacc[Joh75] the grammar gets annotated by code, which will construct some result during the parses. In parser combinator libraries this code is attached to the production of a grammar by a further combinator. This combinator allow to express that for an successful pars a function will be applied to the result. This turns out to be a map.

We provide a further class in our library. It contains fields for a parser and a function:

———— ParsLib.h ————
```
152  template <typename a,typename b,typename c>
153  class Map:public Parser<b,c>,public CAF<Parser<b,c>*>{
154    private:
155      Parser<a,c>* p;
156      b(*f)(a);
157
158    public:
159      Map(b(*f)(a),Parser<a,c>* p)
160        :CAF<Parser<b,c>*>(this),f(f),p(p){}
```

The implementation of the actual method `parse` applies the inner parser to the input, and in case of success takes its actual result and applies the function to this, in order to constrcut the overall result. Again, intermediate results need to be deleted.

———— ParsLib.h ————
```
161      virtual ParsResult<b,c>* parse(vector<c>& xs){
162        ParsResult<a,c>* res1 = p->parse(xs);
163
164        if (!res1->failed){
```

```
165        a r1 = res1->result;
166        vector<c> further = res1->furtherToken;
167        delete res1;
168        return new ParsResult<b,c>((*f)(r1),further);
169      }
170      delete res1;
171      return fail<b,c>(xs);
172    }
173  };
```

An operator overloading is defined for construction of `Map` objects. We decided for the operator `<<`. It is used to attach some code to a rule of the grammar.

```
                          ── ParsLib.h ──
174  template <typename a,typename b,typename c>
175  CAF<Parser<b,c>*> operator<<(CAF<Parser<a,c>*> p,  b(*f)(a)){
176    return ((Parser<b,c>*)new Map<a,b,c>(f,p()));
177  }
```

With the new parser combinator `Map` we can express one further operator. For the alternativ combination of two parsers with the same result type, it is unnecessary to differentaite the two results through an `Either` object. We can provide a common function to extract the data stored in an `Either` object.

```
                          ── ParsLib.h ──
178  template <typename a>
179  a getLeftRight(Either<a,a>* either){return either->value.left;}
```

Now we provide an alternative combinator for two parsers with the same result type. We decided for the double vertical bar as this operator:

```
                          ── ParsLib.h ──
180  template <typename a,typename c>
181  CAF<Parser<a,c>*>
182    operator||(CAF<Parser<a,c>*> p1,CAF<Parser<a,c>*> p2){
183    return (p1|p2)<< getLeftRight<a>;
184  }
```

Finally we provide a class for the empty word production. It does not consume any token and always succeeds:

```
                          ── ParsLib.h ──
185  template <typename a,typename b>
186  class Result:public Parser<a,b>,public CAF<Parser<a,b>*>{
187    public:
188      a x;
189      Result(a x):CAF<Parser<a,b>*>(this),x(x){}
190      ParsResult<a,b>* parse(vector<b>& xs){
191        return new ParsResult<a,b>(x,xs);}
```

```
192  };
193  }}}//namespace
194  #endif
195
```

That's it. We defined a complete parser library for construction of recursive descendant parsers. As a consequence grammar transsscribed with our library to a parser may not contain left recursive productions.

## 3   Example

We give an example of how to use the parser library. As a language we will use Wirth's PL/0[Wir76]. The complete grammar is given in figure 1. It is not left recursive, such that we can directly transcribe it. However it is not left unique. There are different alternatives of a rule with common prefixes. This may lead to serious efficenciy problems.

In this example implementation we will not construct any reasonable result. A simple boolean value serves as result. Therefore we provide a generic unary function, which maps its argument to the value `true`. Furthermore an epsilon parser, which results `true` is given.

Our implementation does not differentiate between parser and lexer. The lexer is completely expressed within the parser. Productions for tokens are implemented in the parser. In order to deal with arbitrary whitespace, even a special prduction for consuming whitespace is provided. All this auxiliary code, together with some type synonous can be found in figure 2.

The tokenizer part of our parser is quite simple. Sequences of certain characters are consumed. The nested pair object is then simply mapped to the value `true`.[2] The call of the generic function `mkTrue` needs to be annotated with the concrete type instance, since the C++ type inference algorith is too weak to derive this type. the complete tokenizer part of the parser is given in figure 3.

Eventually we can define the parser for PL/0. We need to take care of recursive rules. For these a function returning the corresponding parser needs to be defined. The actual parser is then a genetor for this function. Having taken care of this, the grammar can directly be expressed in C++ code. The full implementation is given in figure 4.

We provide some `main` function for our parser.

```
                                 ParsePL0.cpp
143  int main(int argc,char** argv ){
144    FILE *fp;
145    fp = fopen(argv[1], "r");
146    vector<char> xs;
147
148    int c = getc(fp);
149    while (c != EOF) {
150      xs.push_back((char)c);
151      c = getc(fp);
```

---

[2]We neglected that the intermediate pair objects need to be deleted from the storage. This would have needed special versions of `mkTrue` and will blow up the code. There was not enough room for this within the paper.

$$
\begin{aligned}
program &\rightarrow\ block\ . &(1)\\
block &\rightarrow\ constDecl\ varDecl\ procDecls\ statement &(2)\\
constDecl &\rightarrow\ \texttt{CONST}\ constAssignmentList\ \texttt{;}\ |\ \epsilon &(3)\\
constAssignmentList &\rightarrow\ \texttt{ident = number} &(4)\\
&\quad |\ \texttt{ident = number ,}\ constAssignmentList\\
varDecl &\rightarrow\ \texttt{VAR}\ identList\ \texttt{;}\ |\ \epsilon &(5)\\
identList &\rightarrow\ \texttt{ident ,}\ identList\ |\ \texttt{ident} &(6)\\
procDecls &\rightarrow\ procDecl\ procDecls\ |\ \epsilon &(7)\\
procDecl &\rightarrow\ \texttt{PROCEDURE ident ;}\ block\ \texttt{;} &(8)\\
statement &\rightarrow\ blockSt\ |\ callSt\ |\ ifSt\ |\ whileSt\ |\ assignSt\ |\ \epsilon &(9)\\
assignSt &\rightarrow\ \texttt{ident := }\ expression &(10)\\
callSt &\rightarrow\ \texttt{CALL ident} &(11)\\
ifSt &\rightarrow\ \texttt{IF}\ condition\ \texttt{THEN}\ statement &(12)\\
whileSt &\rightarrow\ \texttt{WHILE}\ condition\ \texttt{DO}\ statement &(13)\\
blockSt &\rightarrow\ \texttt{BEGIN}\ statementList\ \texttt{END} &(14)\\
statementList &\rightarrow\ statement\ \texttt{;}\ statementList\ |\ statement &(15)\\
condition &\rightarrow\ \texttt{ODD}\ expression\ |\ expression\ compOp\ expression &(16)\\
compOp &\rightarrow\ \texttt{=}\ |\ \texttt{<>}\ |\ \texttt{<}\ |\ \texttt{>}\ |\ \texttt{<=}\ |\ \texttt{>=} &(17)\\
expression &\rightarrow\ term\ expression2 &(18)\\
&\quad |\ addOp\ term\ expression2\\
expression2 &\rightarrow\ addOp\ term\ expression2\ |\ \epsilon &(19)\\
addOp &\rightarrow\ \texttt{+}\ |\ \texttt{-} &(20)\\
term &\rightarrow\ factor\ term2 &(21)\\
term2 &\rightarrow\ multOp\ factor\ term2\ |\ \epsilon &(22)\\
multOp &\rightarrow\ \texttt{*}\ |\ \texttt{/} &(23)\\
factor &\rightarrow\ \texttt{ident}\ |\ \texttt{number}\ |\ \texttt{(}\ expression\ \texttt{)} &(24)
\end{aligned}
$$

Figure 1: Grammar of PL/0.

```
152     }
153
154    cout<<pl0::program()->parse(xs)->failed<<endl;
155    cout<<pl0::program()->parse(xs)->result<<endl;
156    cout<<pl0::program()->parse(xs)->furtherToken[0]<<endl;
157  }
158
```

The following simple program can be used as input for the parser:

```
                      Test1.pl0
1  CONST M = 7, N = 85;
2  VAR I,X,Y,Z,Q,R;
```

```
                                    ParsePL0.cpp
1   #include "../ParsLib.h"
2   #include <iostream>
3
4   using namespace name::panitz::parser ;
5
6   namespace pl0{
7
8   template <typename a>
9   bool mkTrue(a x){return true;}
10
11  typedef Parser<bool,char>* CP;
12  typedef CAF<CP> P;
13
14  typedef pair<bool,bool>* pb2;
15  typedef pair<pair<bool,bool>*,bool>* pb3;
16  typedef pair<pair<pair<bool,bool>*,bool>*,bool>* pb4;
17  typedef pair<pair<pair<pair<bool,bool>*,bool>*,bool>*,bool>* pb5;
18
19  bool charEq(char c1,char c2){
20  return c1==c2;}
21
22  P epsilon = new Result<bool,char>(true);
23
24  P gC(char c){ return getToken(c,charEq) << mkTrue<char>;}
25
26  P whiteChar = gC(' ')||gC('\n')||gC('\t');
27
28  CP getWhiteSpace();
29  P whiteSpace = caf(getWhiteSpace);
30  CP getWhiteSpace(){return ((whiteChar,whiteSpace)<<mkTrue<pb2> ||epsilon)() ;}
31
32  P gwC(char c){return (whiteSpace,gC(c)) << mkTrue<pb2>;}
```

Figure 2: Auxilliary definitions for PL/0 parser

```
3
4   PROCEDURE MULTIPLY; VAR A,B;
5   BEGIN   A := X; B := Y; Z := 0;
6     WHILE B > 0 DO BEGIN IF ODD B THEN Z := Z+A; A := 2*A; B := B/2; END
7   END;
8
9   PROCEDURE DIVIDE;
10  VAR W;
11  BEGIN   R := X; Q := 0; W := Y;
12    WHILE W = R DO W := 2*W;
13    WHILE W > Y DO
14     BEGIN Q := 2*Q; W := W/2; IF W = R THEN BEGIN R := R-W; Q := Q+1 END END
15  END;
16
17  PROCEDURE GCD; VAR F,G;
18  BEGIN   F := X; G := Y;
19    WHILE F <> G DO BEGIN IF F<G THEN G := G-F; IF G<F THEN F := F-G; END;
20    Z := F
21  END;
22
23  BEGIN   I:=2000;
24    WHILE I<>0 DO
25      BEGIN X := M; Y := N; CALL MULTIPLY;X := 25; Y := 3; CALL DIVIDE;
```

```
                              ─── ParsePL0.cpp ───
33  P ifT    = (gwC('I'),gC('F'))                        << mkTrue<pb2>;
34  P thenT  = (gwC('T'),gC('H'),gC('E'),gC('N'))        << mkTrue<pb4>;
35  P callT  = (gwC('C'),gC('A'),gC('L'),gC('L'))        << mkTrue<pb4>;
36  P whileT = (gwC('W'),gC('H'),gC('I'),gC('L'),gC('E'))<< mkTrue<pb5>;
37  P constT = (gwC('C'),gC('O'),gC('N'),gC('S'),gC('T'))<< mkTrue<pb5>;
38  P beginT = (gwC('B'),gC('E'),gC('G'),gC('I'),gC('N'))<< mkTrue<pb5>;
39  P oddT   = (gwC('O'),gC('D'),gC('D'))                << mkTrue<pb3>;
40  P endT   = (gwC('E'),gC('N'),gC('D'))                << mkTrue<pb3>;
41  P varT   = (gwC('V'),gC('A'),gC('R'))                << mkTrue<pb3>;
42  P doT    = (gwC('D'),gC('O'))                        << mkTrue<pb2>;
43
44  P procedureT = (gwC('P'),gC('R'),gC('O'),gC('C'),gC('E')
45                  ,gC('D'),gC('U'),gC('R'),gC('E'))
46     << mkTrue<pair<pair<pair<pair<pair<pair<pair<pair<bool,bool>*
47         ,bool>*,bool>*,bool>*,bool>*,bool>*,bool>*,bool>*>;
48
49  P addT   = gwC('+');
50  P subT   = gwC('-');
51  P mulT   = gwC('*');
52  P divT   = gwC('/');
53  P eqT    = gwC('=');
54  P gtT    = gwC('>');
55  P ltT    = gwC('<');
56  P neqT   = (gwC('<'),gC('>'))                        <<mkTrue<pb2>;
57  P geT    = (gwC('>'),gC('='))                        <<mkTrue<pb2>;
58  P leT    = (gwC('<'),gC('='))                        <<mkTrue<pb2>;
59  P dotT   = gwC('.');
60  P commaT = gwC(',');
61  P semicolonT   = gwC(';');
62  P lparT  = gwC('(');
63  P rparT  = gwC(')');
64  P assignT= (gwC(':'),gC('='))                        <<mkTrue<pb2>;
65
66  P alphaT = gC('A')||gC('B')||gC('C')||gC('D')||gC('E')||gC('F')||gC('G')
67            ||gC('H')||gC('I')||gC('J')||gC('K')||gC('L')||gC('M')||gC('N')
68            ||gC('O')||gC('P')||gC('Q')||gC('R')||gC('S')||gC('T')||gC('U')
69            ||gC('V')||gC('W')||gC('X')||gC('Y')||gC('Z');
70
71  P digitT = gC('0')||gC('1')||gC('2')||gC('3')||gC('4')||gC('5')||gC('6')
72            ||gC('7')||gC('8')||gC('9');
73
74
75  CP getNumber();
76  P numberT=caf(getNumber);
77  CP getNumber(){return ((digitT,numberT)<<mkTrue<pb2>||digitT)();}
78  P wnumberT=(whiteSpace,numberT)          << mkTrue<pb2>;
79
80  CP getIdent();
81  P identT=caf(getIdent);
82  CP getIdent(){return ((alphaT,identT)<<mkTrue<pb2>||alphaT)();}
83  P widentT=(whiteSpace,identT)          << mkTrue<pb2>;
```

Figure 3: Tokenizer part of PL/0 parser

```
26            X := 84; Y := 36; CALL GCD; I:=I-1;
27        END;
28  END.//end of everything
```

# 4 Conclusion

We have implemented a very simple parser combinator library in C++. The implementation could be made straighforward. This is not surprising, since C++ is flexible enough to allow many very different styles of programming. It has been once more shown that operator overloading and generic types are key features for implementation of flexible libraries. This has been pointed out several times e.g. very impressively in a talk by Guy L. Steele[Ste99].

The solution chosen for recursive productions is a bit unsatisfactory. A function needs to be defined which results the parser. The function then needs to be wrapped in an `CAF` object.

It is not very surprising that due to missing support of full type inference, and through explicit memory management the implementation is more complex than a corresponding Haskell implementation.

A more than prototypic implementation of a parser library will certainly provide more ways to express parsers, as e.g. repetitions or seperated lists.

Systematic performance tests have not been made for the library.

## 4.1 Related Work

A fully implementation of a parser combinator library in C++ does not seem to be available. The FC++ library[MS01] does not contain parser combinators. On the website of its descendants (*Boost.FC++*) it is noted as future work[3].

A much more ambitious work has begun by Claessen [Cla]. His implementation of a combinator parser library is done in plain C.

# References

[Cla]    Coen Claessen. Parser combinators in c. http://www.cs.chalmers.se/~koen/ParserComboC/parser-combo-c.html.

[FL89]   R. Frost and J. Launchbury. Constructing natural language interpreters in a lazy functional language. *The Computer Journal*, 32(2):108–121, 1989.

[Joh75]  Stephen C. Johnson. Yacc: Yet another compiler compiler. Technical Report 39, Bell Laboratories, 1975.

[LM01]   Daan Leijen and Erik Meijer. Parsec: Direct style monadic parser combinators for the real world. Technical Report UU-CS-2001-27, Department of Computer Science, Universiteit Utrecht, 2001.

[MS01]   Brian McNamara and Yannis Smaragdakis. Functional programming in c++ using the fc++ library. *SIGPLAN Notices*, 36(4):25–30, 2001.

[Ste99]  Guy L. Steele. Growing a language. *Higher-Order and Symbolic Computation*, 12(3):221–236, October 1999.

---

[3]As stated on the website http://www.cc.gatech.edu/ yannis/fc++/boostpaper/fcpp.sectlimitations.html in december 2004

[Wad85]  Phil Wadler. How to replace failure by a list of successes. In *Functional Programming Languages and Computer Architecture*, number 201 in Lecture Notes in Computer Science, pages 113–128. Springer, 1985.

[Wir76]  Niklaus Wirth. *Algorithms + Datastructures = Programs*. Prentice-Hall, 1976.

```
                                   ParsePL0.cpp
84   CP getExpression();                  P expression=caf(getExpression);
85   CP getTerm2();                       P term2 = caf(getTerm2);
86   CP getExpression2();                 P expression2 = caf(getExpression2);
87   CP getStatement();                   P statement = caf(getStatement);
88   CP getStatementList();               P statementList = caf(getStatementList);
89   CP getBlock();                       P block=caf(getBlock);
90   CP getProcDecls();                   P procDecls=caf(getProcDecls);
91   CP getIdentList();                   P identList=caf(getIdentList);
92   CP getConstAssignList();             P constAssignList=caf(getConstAssignList);
93
94   P addOp  = addT || subT;
95   P mulOp  = mulT || divT;
96   P compOp = eqT  || neqT || geT  || leT  || gtT || ltT  ;
97
98   P factor = widentT || wnumberT || (lparT,expression,rparT)<<mkTrue<pb3>;
99
100  CP getTerm2(){return ((mulOp,factor,term2)<<mkTrue<pb3> || epsilon)();}
101
102  P term = (factor,term2) << mkTrue<pb2>;
103
104  CP getExpression2(){return ((addOp,term,expression2)<<mkTrue<pb3> || epsilon)();}
105
106  CP getExpression(){
107   return( (term,expression2)                  <<mkTrue<pb2>
108         ||(addOp,term,expression2)            <<mkTrue<pb3>)();}
109
110  P condition = (oddT,expression)              <<mkTrue<pb2>
111              ||(expression,compOp,expression) <<mkTrue<pb3>;
112
113  CP getStatementList(){
114   return ( (statement,semicolonT,statementList)<< mkTrue<pb3>
115          ||statement                 )();}
116
117  P whileSt = (whileT,condition,doT,statement)  << mkTrue<pb4>;
118  P ifSt    = (ifT,condition,thenT,statement)   << mkTrue<pb4>;
119  P callSt   = (callT,widentT)                  << mkTrue<pb2>;
120  P assignSt= (widentT,assignT,expression)      << mkTrue<pb3>;
121  P blockSt = (beginT,statementList,endT)       << mkTrue<pb3>;
122
123  CP getStatement(){return (blockSt||callSt||ifSt||whileSt||assignSt||epsilon)();}
124
125  P procDecl = (procedureT,widentT,semicolonT,block,semicolonT) <<mkTrue<pb5>;
126
127  CP getProcDecls(){return ((procDecl,procDecls)<<mkTrue<pb2>||epsilon)();}
128
129  CP getIdentList(){return ((widentT,commaT,identList)<<mkTrue<pb3>||widentT)();}
130
131  P varDecl = (varT,identList,semicolonT)<<mkTrue<pb3>||epsilon;
132
133  CP getConstAssignList(){
134   return ((widentT,eqT,wnumberT,commaT,constAssignList)<<mkTrue<pb5>
135         ||(widentT,eqT,wnumberT)                       <<mkTrue<pb3>)();}
136
137  P constDecl = (constT,constAssignList,semicolonT)<<mkTrue<pb3> ||epsilon;
138
139  CP getBlock(){return ((constDecl,varDecl,procDecls,statement)<<mkTrue<pb4>)();}
140
141  P program = (block,dotT)<<mkTrue<pb2>;
142  }
```

Figure 4: PL/0 parser